**Fairwinds**

KUBERNETES BENCHMARK REPORT 2024

# Cost, Reliability and Security Workload Results

See how your Kubernetes workloads
compare to other organizations

# INTRODUCTION

Kubernetes security, cost efficiency and reliability are of top concerns to cloud native users. Unfortunately, many DevOps team leaders lack visibility into what's happening within clusters. The result is unnecessary risk, cloud cost overruns and lost customers due to poor application performance.

The Kubernetes Benchmark Report 2024 evaluates over 330,000 workloads and hundreds of organizations on their cluster health. The report serves as a tool for Kubernetes users to benchmark their clusters against efficiency, reliability and security.

Overall we've seen improvements in efficiency and reliability trends over the last three years. While improvements, there are still changes needed. 30% of organizations need container rightsizing to improve efficiency for example. A large majority, greater than 65% of organizations, are missing liveness and/or readiness probes and many organizations still rely on cached images.

Security remains bimodal. Either organizations are putting security configurations in place or they are missing the mark. There are still too many organizations running insecure capabilities (28%) on 90% or more workloads. Likewise 71% or more of their workloads are allowed to run with root access. Security is a crucial factor in Kubernetes so must be addressed with policies to ensure proper configuration.

New this year, the Kubernetes Benchmark Report looks at NSA Hardening Guideline checks providing a snapshot into how organizations are doing. While some benchmark checks see low numbers of workloads impacted, we do see a lot of workloads missing network policies.

Use this report to understand your cluster deficiencies, where to make investments and how to configure Kubernetes to have a positive business impact. Data from the benchmark results is anonymous and sourced from users of Fairwinds Insights.

Fairwinds is the trusted partner for Kubernetes cost optimization and policy enforcement.

**With Fairwinds, customers ship cloud native applications faster, more cost-effectively and with less risk.**

## 30%
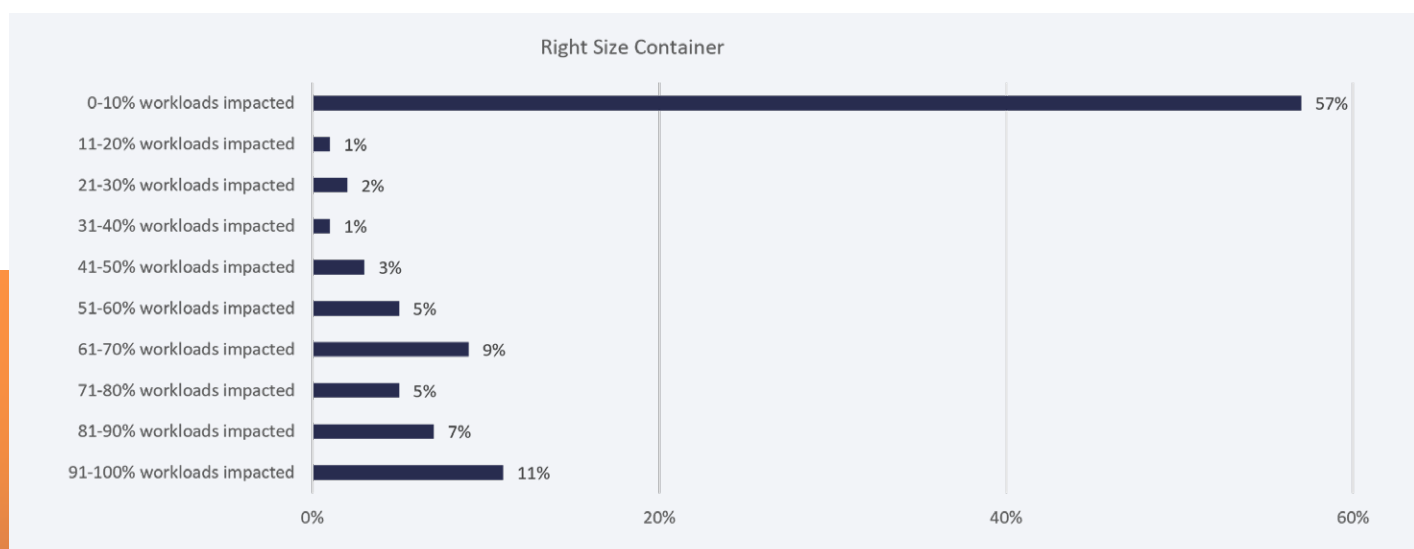of organizations need container rightsizing to improve efficiency

## <65%
of organizations are missing liveness and/or readiness probes and many organizations still rely on cached images

# COST EFFICIENCY

## Rightsize Container

### GOLDILOCKS

Historically, we broke out the data by CPU and Memory, which has been interesting to track since right-sizing can go in one of two directions: Increasing resources to improve reliability, or lowering resources to improve utilization and efficiency. However, as companies grow their container usage, they recognize that right-sizing needs to become part of their workflow and only need to answer a simple question: "Does the container need to be right-sized or not?" As a result, we've simplified our reporting to align with this question and have revealed a somewhat bi-modal distribution: 57% of organizations only have 10% or fewer workloads requiring right-sizing, but 37% of organizations have 50% or more requiring some level of investigation.

### Right Size Container

| Workloads impacted | Percentage |
| --- | --- |
| 0-10% workloads impacted | 57% |
| 11-20% workloads impacted | 1% |
| 21-30% workloads impacted | 2% |
| 31-40% workloads impacted | 1% |
| 41-50% workloads impacted | 3% |
| 51-60% workloads impacted | 5% |
| 61-70% workloads impacted | 9% |
| 71-80% workloads impacted | 5% |
| 81-90% workloads impacted | 7% |
| 91-100% workloads impacted | 11% |

**57%** of organizations only have 10% or fewer workloads requiring right-sizing

**37%** of organizations have 50% or more requiring some level of investigation.
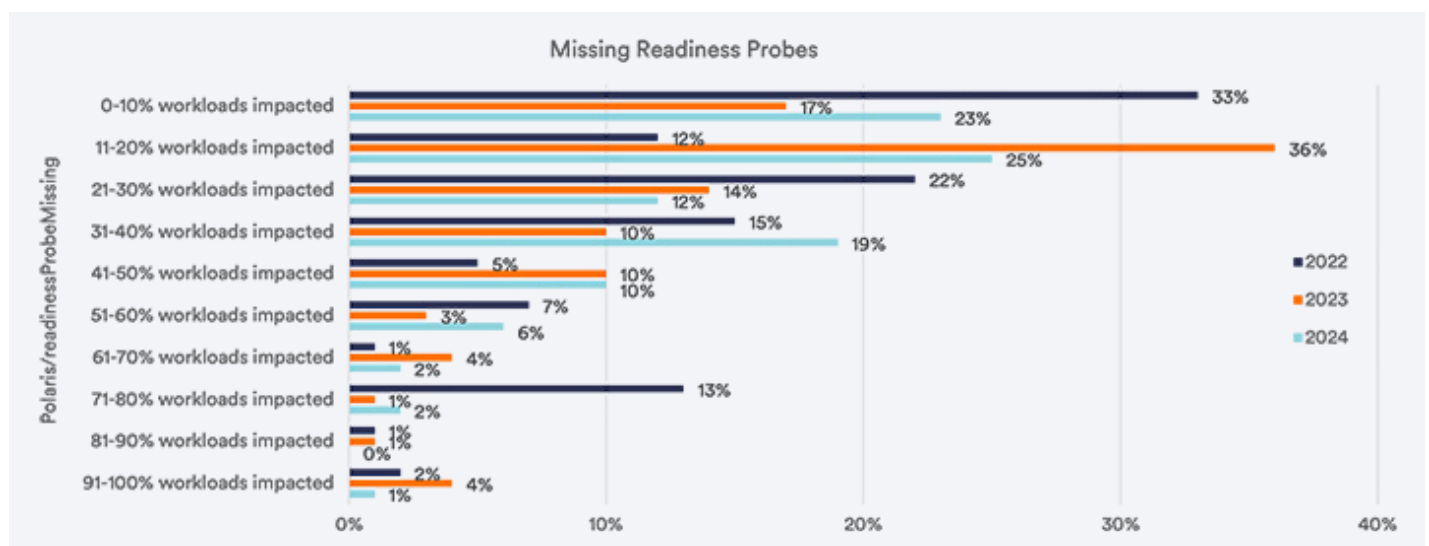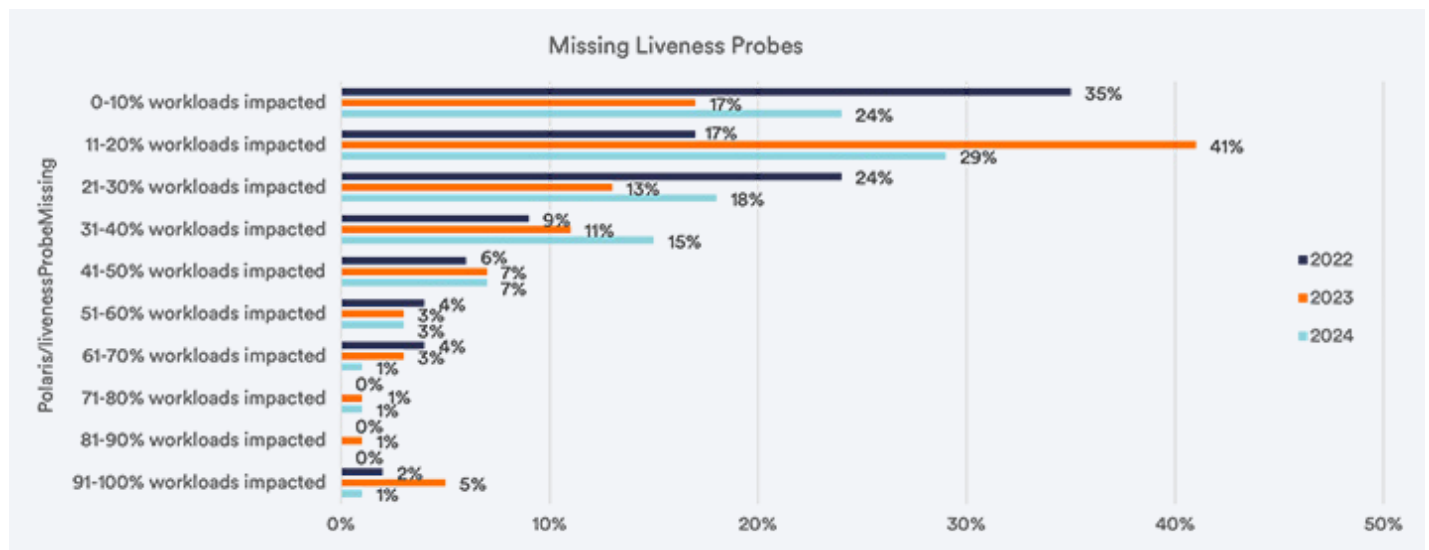
# RELIABILITY

## Missing Liveness and Readiness Probes

**POLARIS**

A liveness probe indicates whether or not the container is running, a fundamental indicator of how a Kubernetes workload should function. If this liveness probe moves into a failing state, then Kubernetes automatically sends a signal to restart the container in an attempt to restore your service to an operational state. If each container in the pod does not have a liveness probe, then a faulty or non-functioning pod will continue to run indefinitely, using up valuable resources and potentially causing application errors.

Liveness and readiness probes make a big difference in the reliability of an application. While we don't see large amounts of workloads impacted i.e. 50% or greater, we do see that 69% of organizations have between 11-50% of workloads missing liveness probes. Similarly we see 66% of workloads with 11-50% of workloads missing readiness probes.

Kubernetes readiness and liveness probes are an important part of making it work properly in the event of a container failure. Be sure to check your workloads for this configuration setting.
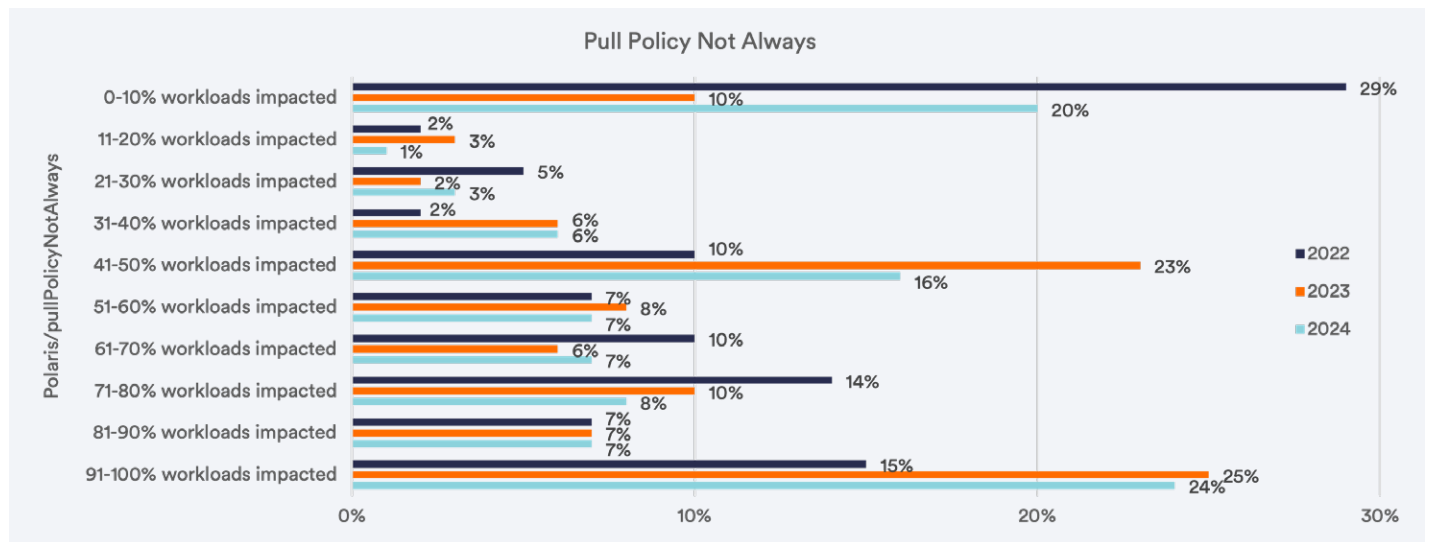
### Missing Liveness Probes

| Polaris/livenessProbeMissing | 2022 | 2023 | 2024 |
|---|---|---|---|
| 0-10% workloads impacted | 35% | 17% | 24% |
| 11-20% workloads impacted | 17% | 41% | 29% |
| 21-30% workloads impacted | 24% | 13% | 18% |
| 31-40% workloads impacted | 9% | 11% | 15% |
| 41-50% workloads impacted | 6% | 7% | 7% |
| 51-60% workloads impacted | 4% | 3% | 3% |
| 61-70% workloads impacted | 4% | 3% | 1% |
| 71-80% workloads impacted | 0% | 1% | 1% |
| 81-90% workloads impacted | 0% | 1% | 1% |
| 91-100% workloads impacted | 2% | 5% | 1% |

### Missing Readiness Probes

| Polaris/readinessProbeMissing | 2022 | 2023 | 2024 |
|---|---|---|---|
| 0-10% workloads impacted | 33% | 17% | 23% |
| 11-20% workloads impacted | 12% | 36% | 25% |
| 21-30% workloads impacted | 22% | 14% | 12% |
| 31-40% workloads impacted | 15% | 10% | 19% |
| 41-50% workloads impacted | 5% | 10% | 10% |
| 51-60% workloads impacted | 7% | 3% | 6% |
| 61-70% workloads impacted | 1% | 4% | 2% |
| 71-80% workloads impacted | 13% | 1% | 2% |
| 81-90% workloads impacted | 1% | 0% | 1% |
| 91-100% workloads impacted | 2% | 4% | 1% |

# Pull Policy Not Always

## POLARIS

Relying on cached versions of a Docker image can become a reliability issue. By default, an image will be pulled if it isn't already cached on the node attempting to run it. This issue can cause variations in images that are running per node, or potentially provide a way to gain access to an image without having direct access to the ImagePullSecret.

Again, we see an increase this year in all workloads impacted. While the distribution of percentage varies, we see that 24% of organizations are relying on cached images for more than 90% of workloads, impacting the reliability of applications..
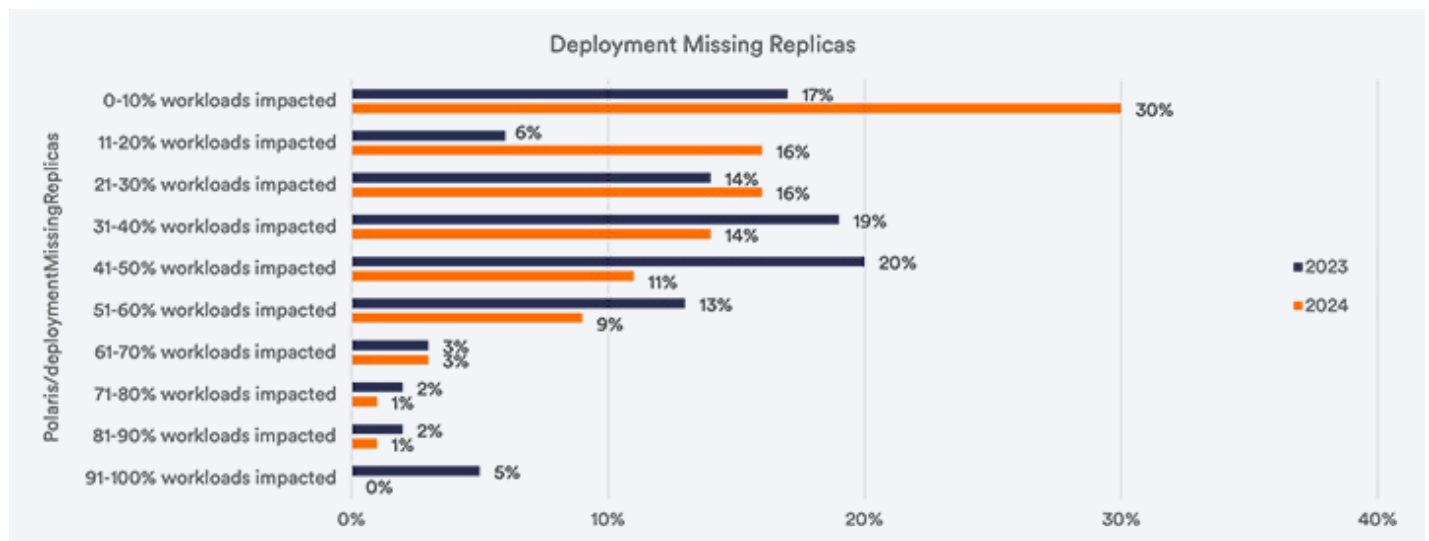
### Pull Policy Not Always

| | 2022 | 2023 | 2024 |
|---|---|---|---|
| 0-10% workloads impacted | 29% | 10% | 20% |
| 11-20% workloads impacted | 2% | 3% | 1% |
| 21-30% workloads impacted | 5% | 2% | 3% |
| 31-40% workloads impacted | 2% | 6% | 6% |
| 41-50% workloads impacted | 10% | 23% | 16% |
| 51-60% workloads impacted | 7% | 8% | 7% |
| 61-70% workloads impacted | 10% | 6% | 7% |
| 71-80% workloads impacted | 14% | 10% | 8% |
| 81-90% workloads impacted | 7% | 7% | 7% |
| 91-100% workloads impacted | 15% | 25% | 24% |

# Deployment Missing Replicas

## POLARIS

For the second year we look at when a deployment is missing replicas. In many cases workloads have not been configured with a replica configuration. 55% of organizations have between more than 21% of workloads missing replicas. On the plus side, there has been improvements for some organizations with 30% having less than 10% of workloads impacted.

Deployments help maintain the stability and high availability of containers. Without these in place, if a node crashes, a deployment or ReplicaSet will not replace failed pods. This could be dangerous for the reliability of an application.
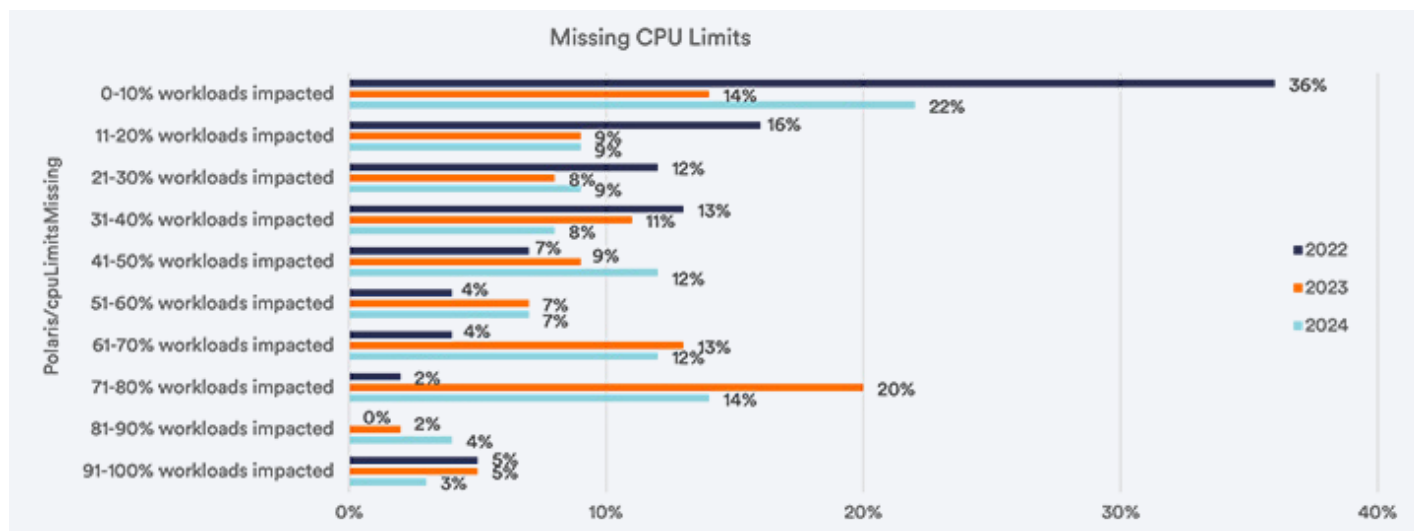
### Deployment Missing Replicas

| | 2023 | 2024 |
|---|---|---|
| 0-10% workloads impacted | 17% | 30% |
| 11-20% workloads impacted | 6% | 16% |
| 21-30% workloads impacted | 14% | 16% |
| 31-40% workloads impacted | 19% | 14% |
| 41-50% workloads impacted | 20% | 11% |
| 51-60% workloads impacted | 13% | 9% |
| 61-70% workloads impacted | 3% | 3% |
| 71-80% workloads impacted | 2% | 1% |
| 81-90% workloads impacted | 2% | 1% |
| 91-100% workloads impacted | 5% | 0% |

# CPU Limits Missing

## POLARIS

More organizations are ensuring CPU limits are set. 22% of organizations have less than 10% of workloads missing CPU limits. We've seen drops across the board showing that organizations are adding them and making improvements to the cost-efficiency and reliability of applications.

If you do not specify CPU limit then the container will not have any upper bound. This can impact reliability as the CPU intensive container slows down and could exhaust all CPU available on the node.



# CPU Requests Missing

## POLARIS

Similarly to CPU limits, we've seen an increase of organizations with missing CPU requests. Previously we saw that only 50% and 78% of organizations were missing requests on at least 10% of their workloads. We've seen a reduction to 67% of organizations with greater than 11% of workloads impacted. Organizations are realizing the value of ensuring CPU requests are set.

If a single pod is allowed to consume all of the node CPU and memory, then other pods will be starved for resources. Setting resource requests increases reliability by guaranteeing the pod will have access to those resources—and preventing other pods from consuming all of the available resources on a node (this is referred to as the "noisy neighbor problem").

# SECURITY

Kubernetes is not secure out-of-the-box. By default, everything you deploy to your cluster is able to communicate with everything else. Pod security is important because pods can run as root, and even if you disable privileged mode, there are still plenty of pod settings that an attacker inside your cluster can enable to get root access to the underlying node.
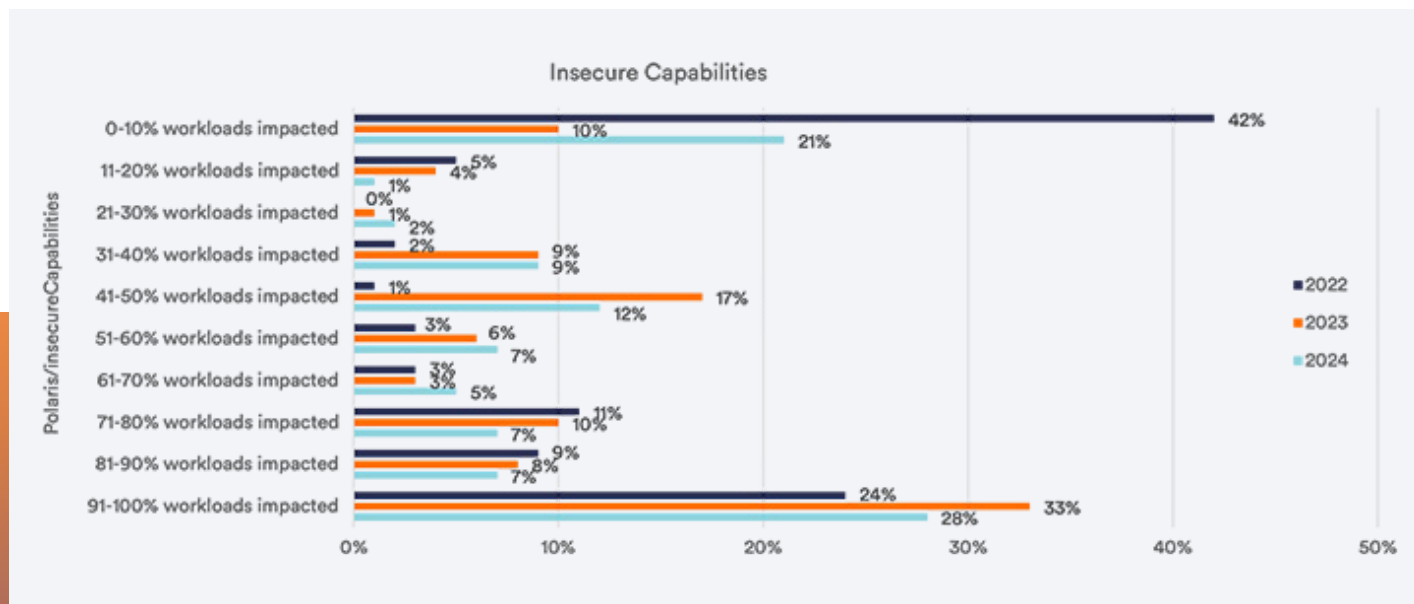
The checks below all help to ensure organizations are configuring Kubernetes with security in mind. New this year, we've also added NSA hardening check benchmarks.

## Insecure Capabilities

### POLARIS

Certain Linux capabilities are enabled by default for Kubernetes workloads, though most workloads don't really need these capabilities. The effort organizations made to pare back these capabilities has dropped. Whereas in 2022 42% of organizations were turning off these capabilities, that number dropped to 10% in 2023, but is rebounding at 21% this year.

Now we see 28% of organizations with more than 90% of workloads running with insecure capabilities. We need to see the majority of organizations turning off these capabilities vs. only slight improvements.
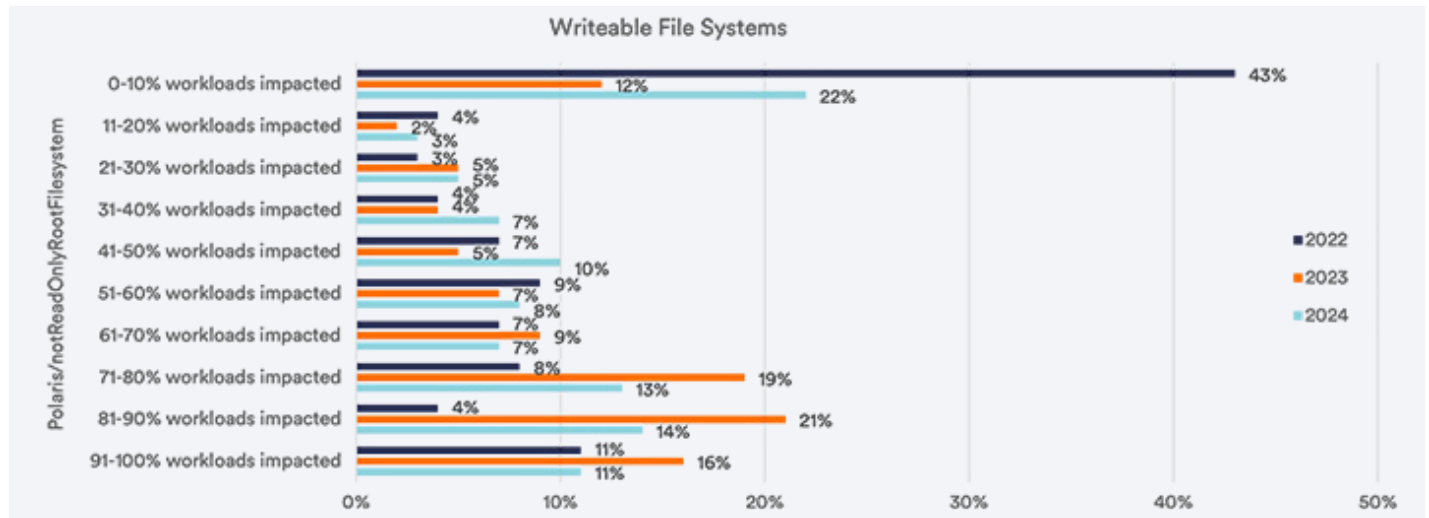


Insecure Capabilities

| Polaris/insecureCapabilities | 2022 | 2023 | 2024 |
|---|---|---|---|
| 0-10% workloads impacted | 42% | 10% | 21% |
| 11-20% workloads impacted | 5% | 4% | 1% |
| 21-30% workloads impacted | 0% | 1% | 2% |
| 31-40% workloads impacted | 2% | 9% | 9% |
| 41-50% workloads impacted | 1% | 17% | 12% |
| 51-60% workloads impacted | 3% | 6% | 7% |
| 61-70% workloads impacted | 3% | 3% | 5% |
| 71-80% workloads impacted | 11% | 10% | 7% |
| 81-90% workloads impacted | 9% | 8% | 7% |
| 91-100% workloads impacted | 24% | 33% | 28% |

## 28% of organizations with more than 90% of workloads running with insecure capabilities

# Writeable File Systems

## POLARIS

readOnlyRootFilesystem is a security setting that controls whether a container is able to write into its filesystem. It is a feature most organizations want enabled in the event of a hack. If an attacker gets in, they will not be able to tamper with the application or write foreign executables to disk. Unfortunately, Kubernetes workloads do not set this to true by default, which means teams need to explicitly ensure it happens to get the most secure configuration possible.

Whereas we saw in 2022 a binary distribution of organizations locking down filesystems inside their containers, and in 2023 a huge increase of workloads impacted, this year we see some improvements. There was a drop in the number of organizations working to override the insecure defaults for 71-100% of their workloads (56% of orgs in 2023, down to 38% in 2024).
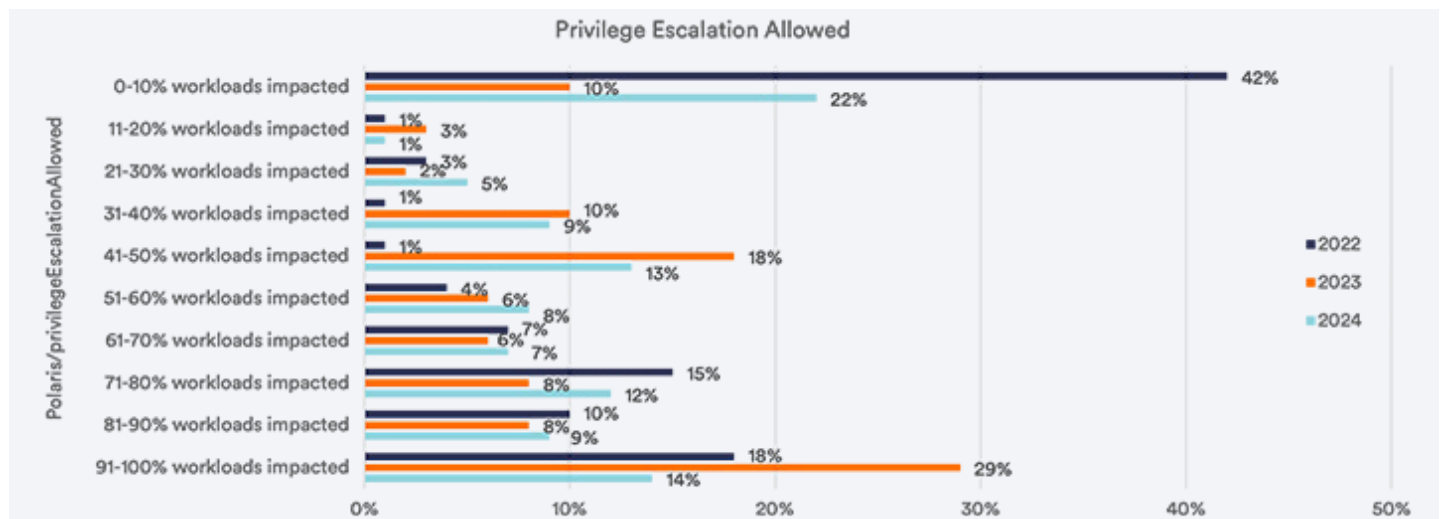


# Privilege Escalation Allowed

## POLARIS

Under particular configurations, a container may be able to escalate its privileges. Setting allowPrivilegeEscalation to false will set the no_new_privs flag on the container process, preventing setuid binaries from changing the effective user ID. Setting this flag is particularly important when using runAsNonRoot, which can otherwise be circumvented. Because this, too, is not set by default, security-conscious teams need to explicitly set it.

There is good news in this area as organizations improve. We've seen a 15% reduction in organizations running 90% of their workloads impacted with privilege escalation (29% of orgs in 2023 to 14% in 2024). We also see that 22% of organizations (up by 12%) have less than 10% of workloads impacted.
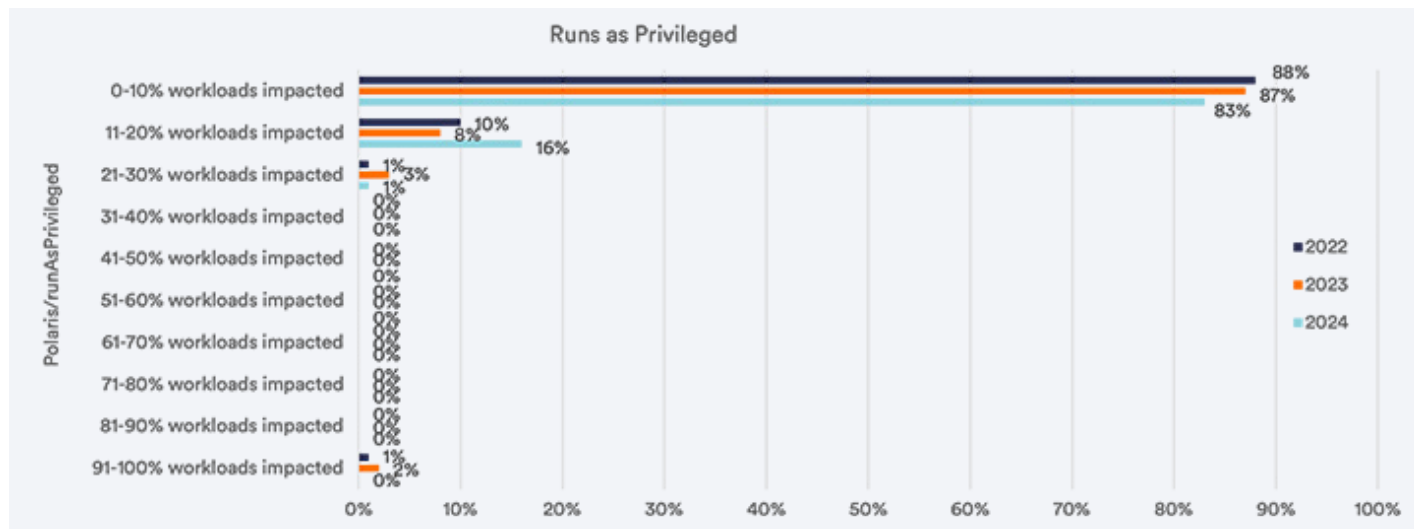
# Runs as Privileged

## POLARIS

The "privileged" key determines if any container in a pod can enable privileged mode. By default, a container is not allowed to access any devices on the host, but a privileged container is given access to all devices on the host. This feature allows the container nearly all the same access as processes running on the host, which is useful for containers looking to use Linux capabilities, like manipulating the network stack and accessing devices.

There has been a slight dip in organizations that have this covered - 87% in 2023 to 83% in 2024. We've also seen doubling of 11-20% of workloads impacted.



# Run as Root Allowed

## POLARIS

As a general rule, containers should not be configured to run as a root user in a Kubernetes cluster. For example, a malicious user or container could take advantage of the root privileges to compromise the system or access sensitive data.

While the results are binary, we are seeing organizations get a handle on workloads running as root allowed. Down from 44% in 2023 to 30% in 2024, we see 71% or more of their workloads as allowing root access. We do see a 10% improvement on organizations with less than 10% of workloads impacted.
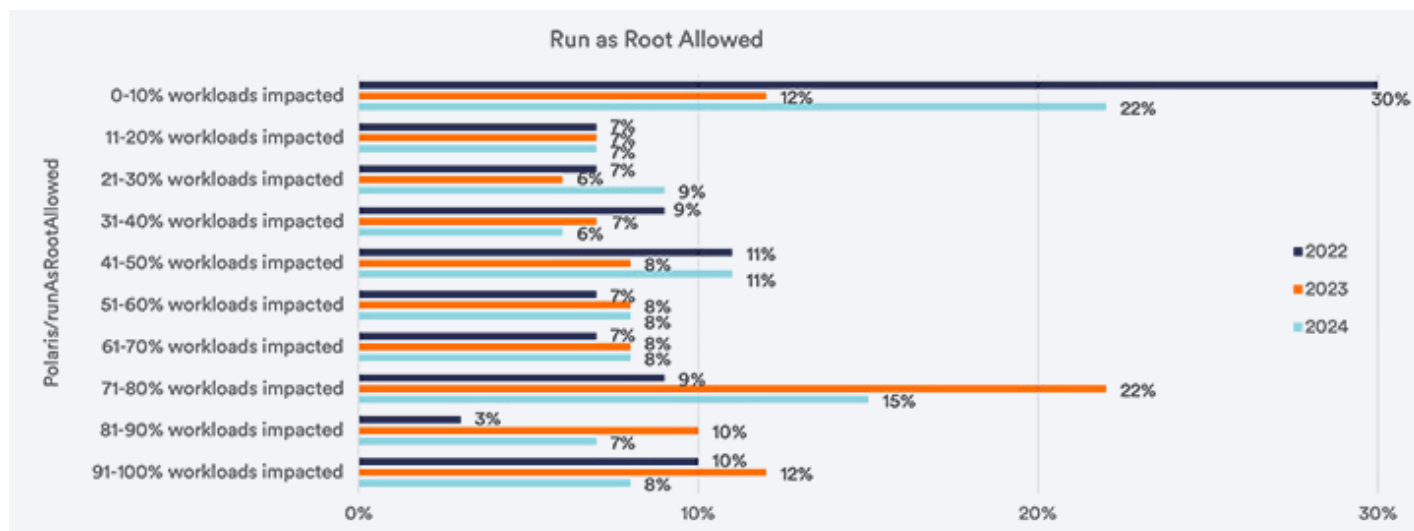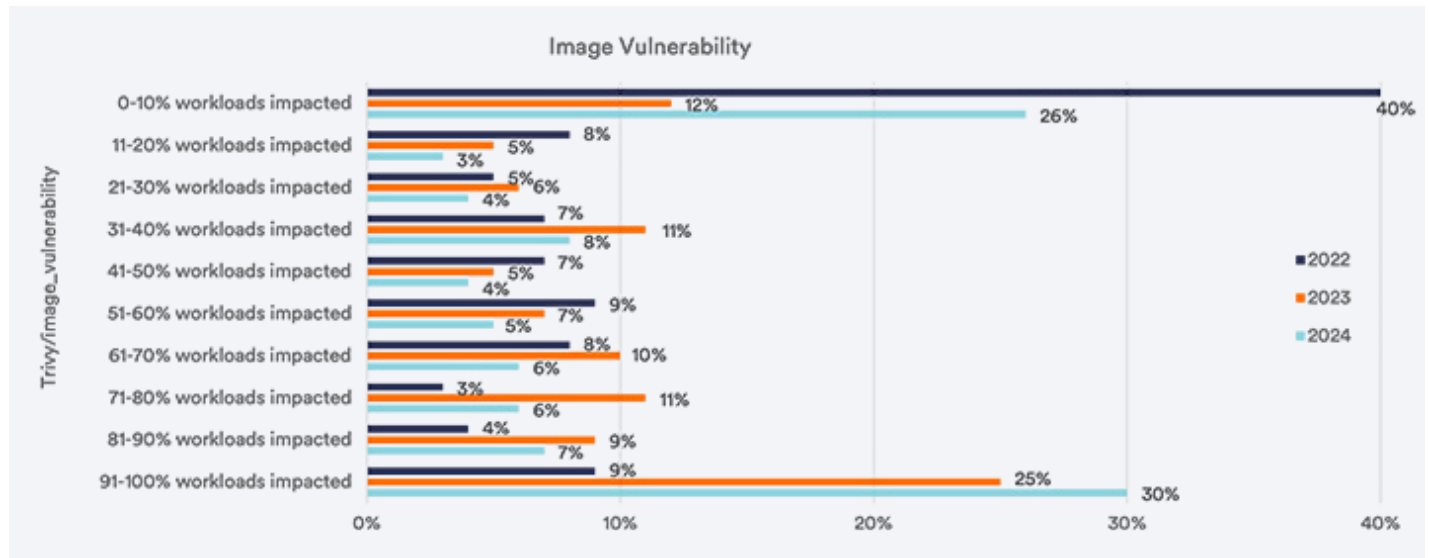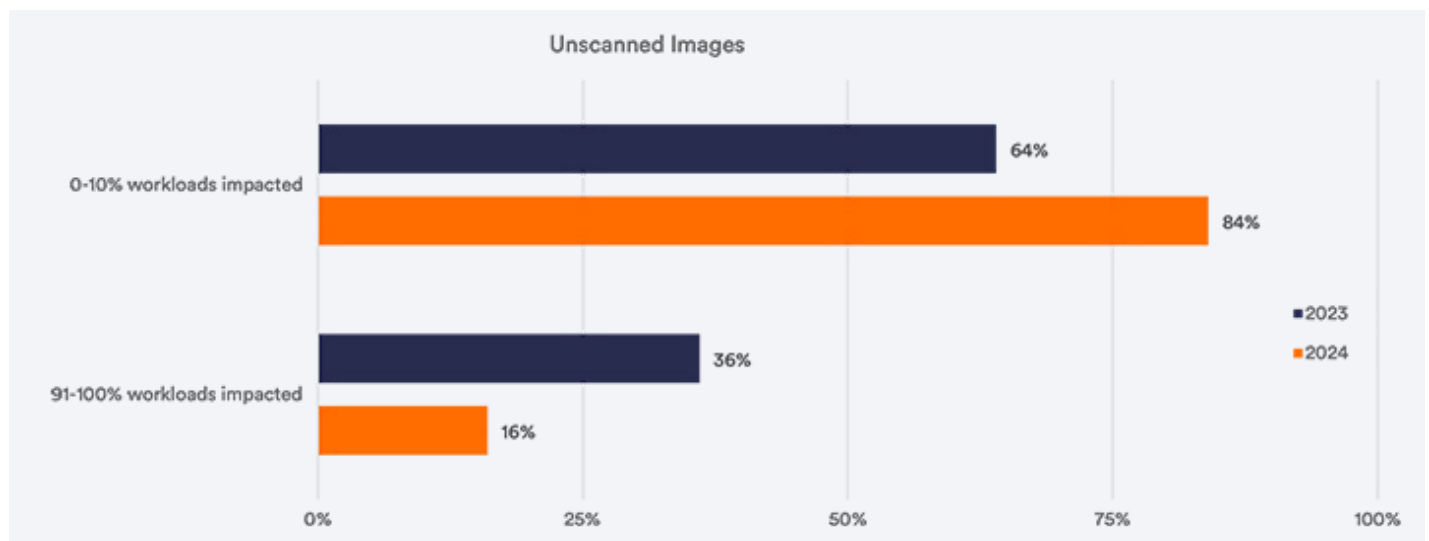
# Image Vulnerability

## TRIVY

Since tracking image vulnerabilities over the last three years, we've seen a huge increase of more than 90% of workloads impacted (9% in 2022 to 30% in 2024). Known vulnerabilities can be exploited by malicious actors and need to be patched/remediated. While workloads impacted vary across percentages impacted, we have seen an improvement in the last year (12 to 26%) of workloads with less than 10% having known image vulnerabilities. This data shows that organizations that adopt regular image patching processes can keep the percentage of workloads affected by vulnerabilities low.



# Unscanned Images

## TRIVY

Organizations have made improvements to scanning images. 20% more organizations have less than 10% of workloads impacted by unscanned images (64% in 2023 to 84% in 2024). On the flip side, now only 16% of organizations are not scanning images (more than 91% or greater impacted). With easily accessible open source scanning tools, every organization should have images scanned for vulnerabilities.

# Outdated Helm Charts

## NOVA

Outdated Helm charts are a pervasive issue across most organizations. This year the percentage of organizations running workloads with outdated Helm charts has increased. While the amount of workloads impacted varies, organizations appear to be getting a handle on tracking outdated Helm charts. More than 70% of organizations have 11% or more of workloads impacted with outdated Helm charts. 36% have 50% or more of workloads impacted. The add-ons running your cluster are most likely installed by Helm. Each add-on has its own release cadence, and some updates may come with critical security patches. Helm charts must be updated but they are increasingly hard to monitor and predict. That's why so many teams are running outdated Helm charts. Nova cross-checks Helm charts running in the cluster with the latest version available.

### Outdated Helm Chart

| Workloads impacted | 2022 | 2023 | 2024 |
|---|---|---|---|
| 0-10% workloads impacted | 37% | 16% | 27% |
| 11-20% workloads impacted | 9% | 4% | 5% |
| 21-30% workloads impacted | 10% | 15% | 12% |
| 31-40% workloads impacted | 6% | 11% | 12% |
| 41-50% workloads impacted | 5% | 8% | 9% |
| 51-60% workloads impacted | 7% | 15% | 12% |
| 61-70% workloads impacted | 8% | 7% | 8% |
| 71-80% workloads impacted | 5% | 13% | 7% |
| 81-90% workloads impacted | 2% | 4% | 3% |
| 91-100% workloads impacted | 11% | 7% | 6% |

**70%** of organizations have 11% or more of workloads impacted with outdated Helm charts
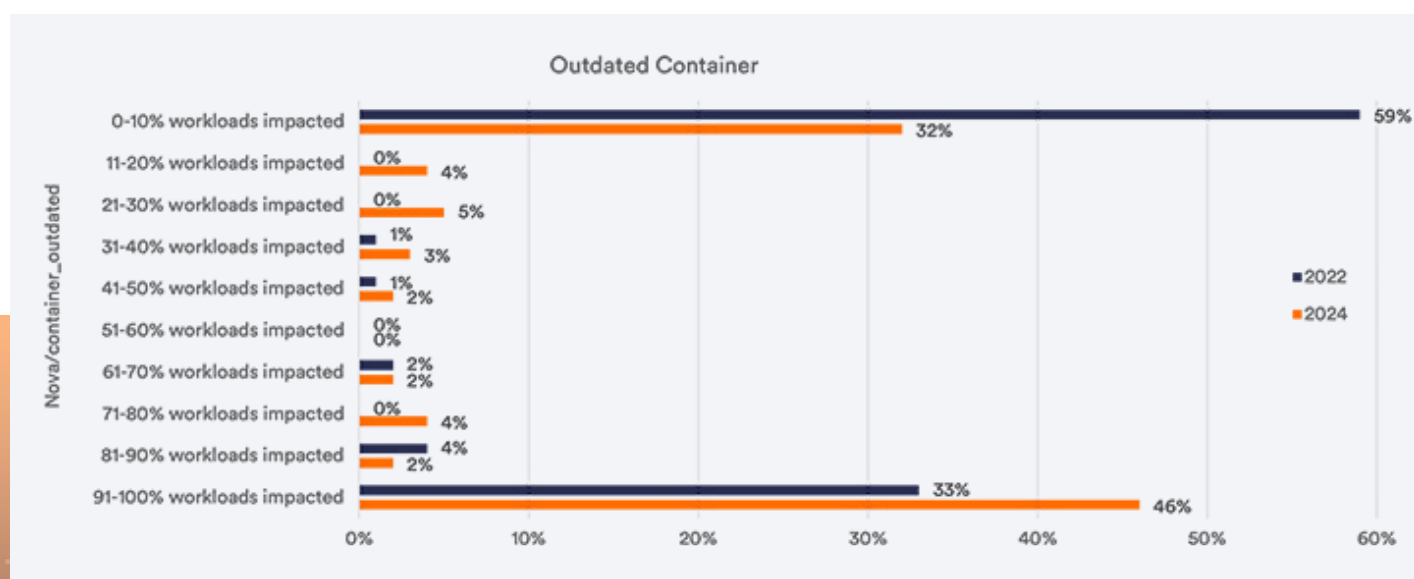
# Outdated Container Images

## NOVA

For the second year, we benchmark how many organizations are running outdated container images. We continue to see extremes: either less than 10% or greater than 90% of workloads are impacted. This year we've seen a spike in 90% or more of workloads impacted - up by 13% (33% in 2023, 46% in 2024).

Nova users can run a flag called "containers", which examines all container images in a Kubernetes cluster and notifies users if a new version is available. Nova provides three different alternatives for updating images: the absolute latest version, the latest minor version and the most recent patch version.

When benchmarking yourself against these findings, the important takeaway is to action outdated helm charts or container images. While you might be on the lower spectrum of impacted workloads, it still could introduce avoidable risk.

**Outdated Container**

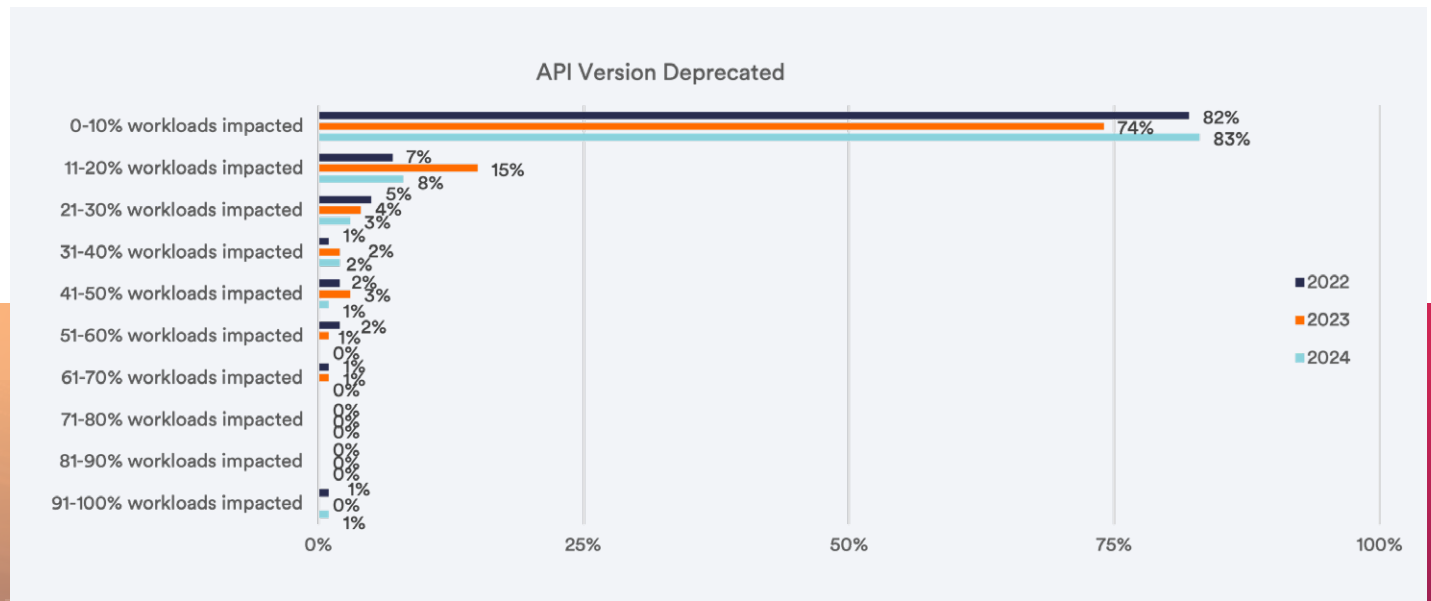| Nova/container_outdated | 2022 | 2024 |
|---|---|---|
| 0-10% workloads impacted | 59% | 32% |
| 11-20% workloads impacted | 0% | 4% |
| 21-30% workloads impacted | 0% | 5% |
| 31-40% workloads impacted | 1% | 3% |
| 41-50% workloads impacted | 1% | 2% |
| 51-60% workloads impacted | 0% | 0% |
| 61-70% workloads impacted | 2% | 2% |
| 71-80% workloads impacted | 0% | 4% |
| 81-90% workloads impacted | 4% | 2% |
| 91-100% workloads impacted | 33% | 46% |

**90%** This year we've seen a spike in 90% or more of workloads impacted - up by 13% (33% in 2023, 46% in 2024

# API Version Deprecated

Most organizations appear to have only a few workloads with deprecated API versions. However, deprecated APIs issues need to be resolved to ensure workloads continue to function after Kubernetes upgrades. Therefore, even monitoring for deprecated APIs - even if it's just 10% of your workloads - remains an important step in de-risking Kubernetes upgrades.



API Version Deprecated

| Workloads impacted | 2022 | 2023 | 2024 |
|---|---|---|---|
| 0-10% workloads impacted | 82% | 74% | 83% |
| 11-20% workloads impacted | 7% | 15% | 8% |
| 21-30% workloads impacted | 5% | 4% | 3% |
| 31-40% workloads impacted | 1% | 2% | 2% |
| 41-50% workloads impacted | 2% | 3% | 1% |
| 51-60% workloads impacted | 0% | 1% | 2% |
| 61-70% workloads impacted | 1% | 1% | 0% |
| 71-80% workloads impacted | 0% | 0% | 0% |
| 81-90% workloads impacted | 0% | 0% | 0% |
| 91-100% workloads impacted | 1% | 0% | 1% |

**10%** even monitoring for deprecated APIs - even if it's just 10% of your workloads - remains an important step in de-risking Kubernetes upgrades

## NSA Hardening Checks

The National Security Agency (NSA) and the Cybersecurity and Infrastructure Security Agency (CISA) released Kubernetes hardening guide in August 2021 and continues to update it. The checks make recommendations for hardening Kubernetes clusters and outlines a strong defense-in-depth approach to ensure that when an attacker compromises your cluster, the blast radius will be as small as possible.
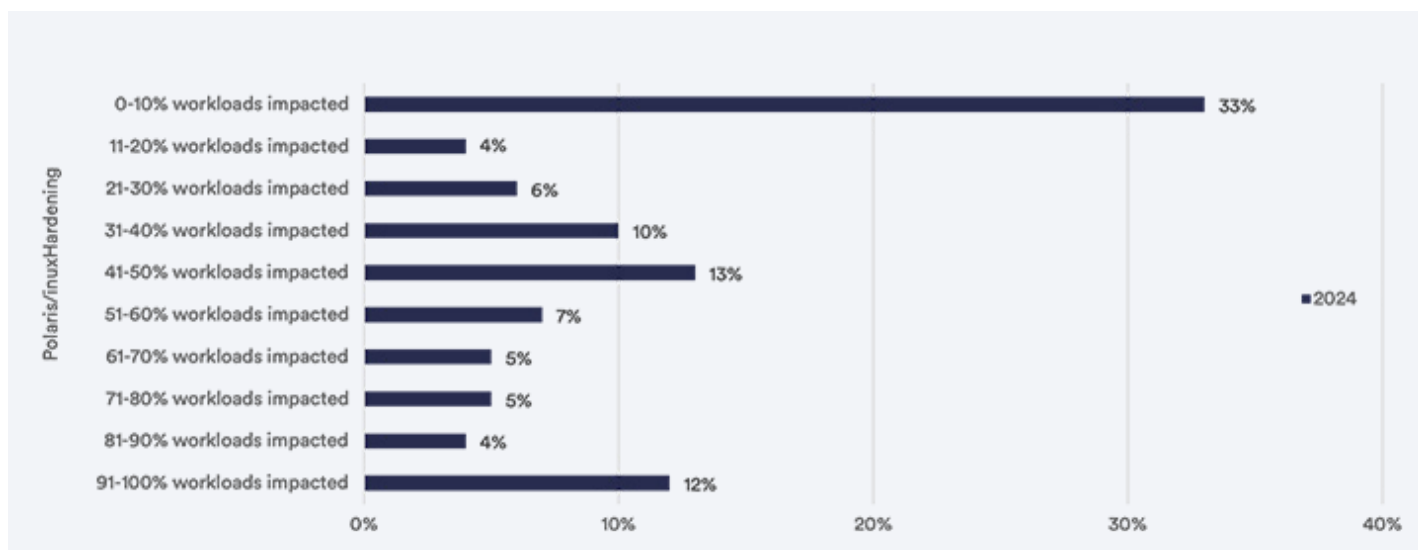
Fairwinds provides support for a number of these recommendations, but added these additional checks in response to the NSA guidelines.

**LINUX HARDENING**

**POLARIS**

By default, a workload likely has more ability than is needed to run its application. Tuning these privileges minimizes the radius, speed, and impact of a container compromise. This check ensures a workload uses either AppArmor, dropping Linux capabilities, SELinux, or a seccomp profile to grant workloads the minimal privileges needed to function. This is important because minimizing a workload's privileges also minimizes an attacker's ability to gain access to other workloads or to your cluster.

33% of organizations have more than 50% of workloads that has too much privilege. These organizations must look to identify these workloads to implement Linux hardening.

# Missing Network Policy

**POLARIS**

A NetworkPolicy controls Pod network communication with other Pods, IP addresses, or namespaces.

The Pod does not have an accompanying NetworkPolicy that limits its egress and ingress traffic. This may allow undesired access to external resources, or access from other Pods.

37% of organizations have protected workloads with network policy. However, 58% of organizations have workloads missing network policy. That is surprisingly high for organizations as setting network policy is an essential step to securing containers.
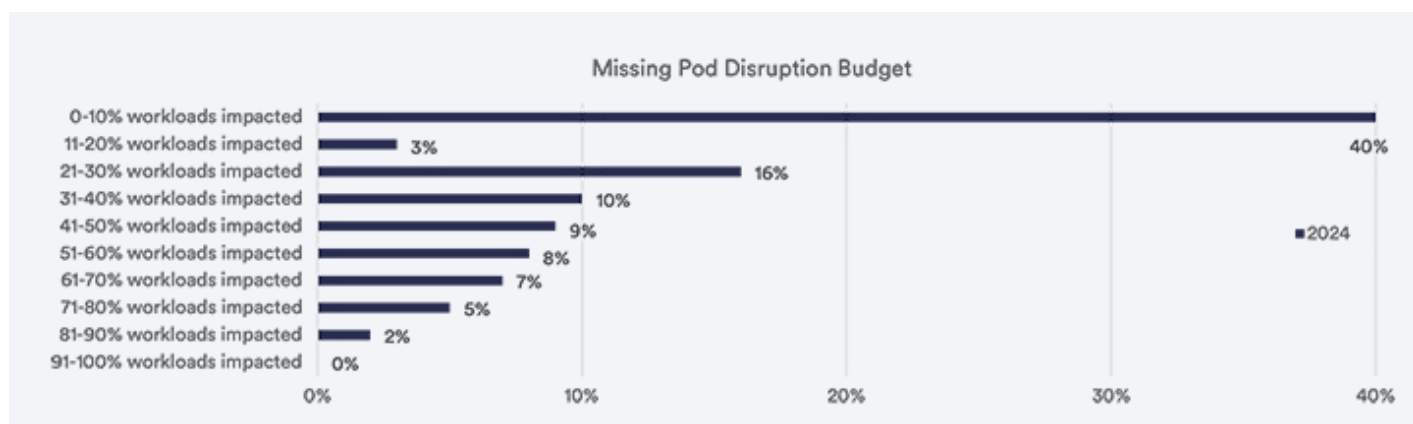


Missing Network Policy

| Workloads impacted | 2024 |
|---|---|
| 0-10% workloads impacted | 37% |
| 11-20% workloads impacted | 1% |
| 21-30% workloads impacted | 0% |
| 31-40% workloads impacted | 1% |
| 41-50% workloads impacted | 3% |
| 51-60% workloads impacted | 15% |
| 61-70% workloads impacted | 11% |
| 71-80% workloads impacted | 14% |
| 81-90% workloads impacted | 10% |
| 91-100% workloads impacted | 8% |

# Missing Pod Disruption Budget

**POLARIS**

A PodDisruptionBudget is very helpful for ensuring zero-downtime for your applications, especially when going through cluster upgrades. It ensures that a certain number of Pods remain available when Nodes need to be cycled out.

Things improve with pod distribution budget set. 40% of organizations have this covered for 90% of their workloads.
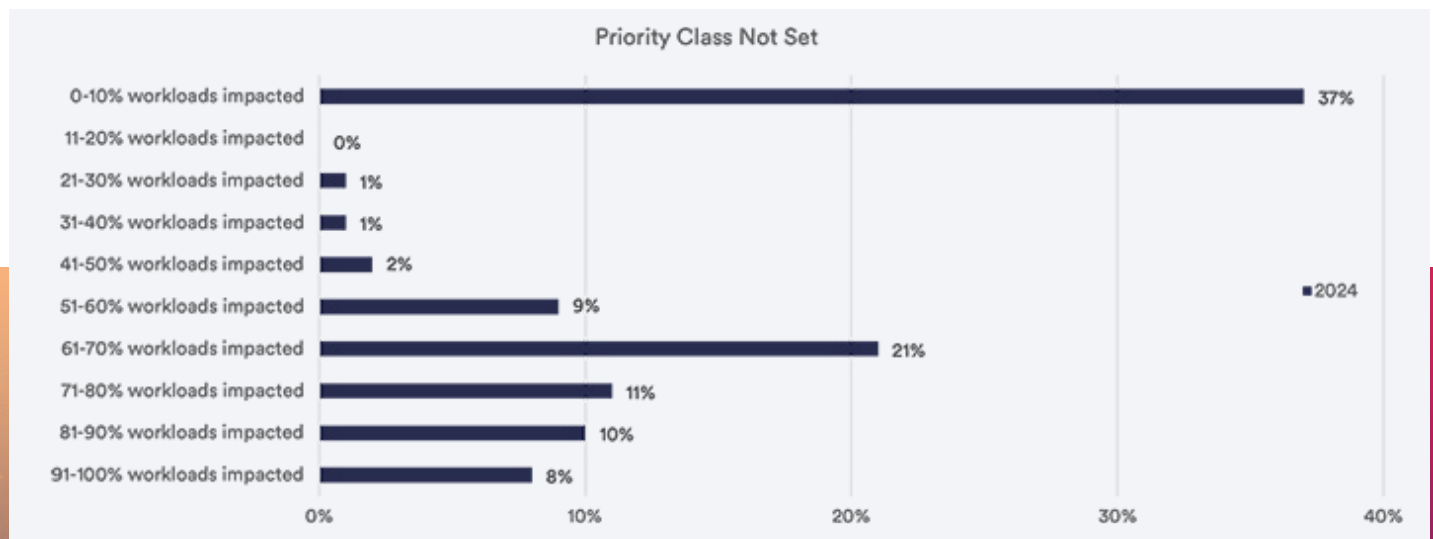


Missing Pod Disruption Budget

| Workloads impacted | 2024 |
|---|---|
| 0-10% workloads impacted | 40% |
| 11-20% workloads impacted | 3% |
| 21-30% workloads impacted | 16% |
| 31-40% workloads impacted | 10% |
| 41-50% workloads impacted | 9% |
| 51-60% workloads impacted | 8% |
| 61-70% workloads impacted | 7% |
| 71-80% workloads impacted | 5% |
| 81-90% workloads impacted | 2% |
| 91-100% workloads impacted | 0% |

# Priority Class Not Set

## POLARIS

Kubernetes allows users to specify the priority of different pods. This helps to make sure mission-critical pods get scheduled ahead of lower-priority pods. Taking advantage of this feature can help increase the reliability of your applications.

21% of organizations have 61-70% of workloads without priority class set compared to 37% that have it on less than 10% of workloads. Overall 58% of organization have 50% or more workloads impacted.

### Priority Class Not Set

| Workloads impacted | 2024 |
|---|---|
| 0-10% workloads impacted | 37% |
| 11-20% workloads impacted | 0% |
| 21-30% workloads impacted | 1% |
| 31-40% workloads impacted | 1% |
| 41-50% workloads impacted | 2% |
| 51-60% workloads impacted | 9% |
| 61-70% workloads impacted | 21% |
| 71-80% workloads impacted | 11% |
| 81-90% workloads impacted | 10% |
| 91-100% workloads impacted | 8% |

**21%** of organizations have 61-70% of workloads without priority class set

**58%** of organization have 50% or more workloads impacted

# CONCLUSION

Fairwinds will continue to update the results of this benchmark data to help the cloud native community understand how they stack up against peers. The important takeaway is this: as organizations expand Kubernetes across multiple teams and expand clusters, it is hard to enforce standardization. Without it, you risk wasting cloud resources and introducing risk. You don't want your platform and DevOps teams spending all their time as a helpdesk.

Fairwinds Insights is for platform teams running Kubernetes to standardize and enable development best practices.

✓ Decrease friction

✓ Improve development experience

✓ Increase development velocity

✓ Accelerate time to market and revenue generation

# FAIRWINDS INSIGHTS

Platform teams use Fairwinds Insights to implement standards, enforce development best practices for developers to self-service and do the right thing. With Insights, platform teams enable development to increase velocity and accelerate revenue generating activities.

Use Insights to make platform teams' lives easier, give developers relief, keep security and compliance teams happy and demonstrate responsible spending.

DevOps teams using Insights can stop serving as a Kubernetes help desk.

### Gain Visibility into Kubernetes

Your platform team owns Kubernetes. Fairwinds Insights integrates seamlessly across the entire development lifecycle. Insights continuously scans workloads against your mission critical policies, identifies issues and automates remediation.

**Centralized**
One view into clusters and IaC

**Consistent**
Standardize and enforce policy

**Prioritized**
Team focus on what matters most

**Fairwinds**

## WHY FAIRWINDS

Fairwinds is your trusted partner for Kubernetes governance and guardrails. With Fairwinds, customers ship cloud-native applications faster, more cost effectively, and with less risk. We provide a unified view between dev, sec, and ops, removing friction between those teams with software that simplifies complexity. Fairwinds Insights is built on Kubernetes expertise and integrates our leading open source tools to help you save time, reduce risk, and deploy with confidence.

**WWW.FAIRWINDS.COM**