

CHAPTER 1

Revolution in the Cloud

There was never a time when the world began, because it goes round and round like a circle, and there is no place on a circle where it begins.

—Alan Watts

There's a revolution going on. Actually, three revolutions.

The first revolution is the creation of *the cloud*, and we'll explain what that is and why it's important. The second is the dawn of *DevOps*, and you'll find out what that involves and how it's changing operations. The third revolution is the wide adoption of *containers*. Together, these three waves of change are creating a new software world: the *cloud native* world. The operating system for this world is called *Kubernetes*.

In this chapter, we'll briefly recount the history and significance of these revolutions, and explore how the changes are affecting the way we all deploy and operate software. We'll outline what *cloud native* means, and what changes you can expect to see in this new world if you work in software development, operations, deployment, engineering, networking, or security.

Thanks to the effects of these interlinked revolutions, we think the future of computing lies in cloud-based, containerized, distributed systems, dynamically managed by automation, on the *Kubernetes* platform (or something very like it). The art of developing and running these applications—*cloud native DevOps*—is what we'll explore in the rest of this book.

If you're already familiar with all of this background material, and you just want to start having fun with *Kubernetes*, feel free to skip ahead to [Chapter 2](#). If not, settle down comfortably, with a cup of your favorite beverage, and we'll begin.

The Creation of the Cloud

In the beginning (well, the 1960s, anyway), computers filled rack after rack in vast, remote, air-conditioned datacenters, and users would never see them or interact with them directly. Instead, developers submitted their jobs to the machine remotely and waited for the results. Many hundreds or thousands of users would all share the same computing infrastructure, and each would simply receive a bill for the amount of processor time or resources they used.

It wasn't cost-effective for each company or organization to buy and maintain its own computing hardware, so a business model emerged where users would share the computing power of remote machines, owned and run by a third party.

If that sounds like right now, instead of last century, that's no coincidence. The word *revolution* means "circular movement," and computing has, in a way, come back to where it began. While computers have gotten a lot more powerful over the years—today's Apple Watch is the equivalent of about three of the mainframe computers shown in [Figure 1-1](#)—shared, pay-per-use access to computing resources is a very old idea. Now we call it the cloud, and the revolution that began with timesharing mainframes has come full circle.



Figure 1-1. Early cloud computer: the IBM System/360 Model 91, at NASA's Goddard Space Flight Center

Buying Time

The central idea of the cloud is this: instead of buying a *computer*, you buy *compute*. That is, instead of sinking large amounts of capital into physical machinery, which is hard to scale, breaks down mechanically, and rapidly becomes obsolete, you simply buy time on someone else's computer and let them take care of the scaling, maintenance, and upgrading. In the days of bare-metal machines—the “Iron Age,” if you like—computing power was a capital expense. Now it's an operating expense, and that has made all the difference.

The cloud is not just about remote, rented computing power. It is also about distributed systems. You may buy raw compute resources (such as a Google Compute *instance*, or an AWS Lambda *function*) and use it to run your own software, but increasingly you also rent *cloud services*: essentially, the use of someone else's software. For example, if you use PagerDuty to monitor your systems and alert you when something is down, you're using a cloud service (sometimes called *software as a service*, or SaaS). The success of these SaaS services is partly due to this recent revolution of the cloud. Now almost anyone can create a new app or website, host it on a public cloud provider, and scale it up to a global audience if they find some success.

Infrastructure as a Service

When you use cloud infrastructure to run your own services, what you're buying is *infrastructure as a service* (IaaS). You don't have to expend capital to purchase it, you don't have to build it, and you don't have to upgrade it. It's just a commodity, like electricity or water. Cloud computing is a revolution in the relationship between businesses and their IT infrastructure.

Outsourcing the hardware is only part of the story; the cloud also allows you to outsource the *software* that you don't write: operating systems, databases, clustering, replication, networking, monitoring, high availability, queue and stream processing, and all the myriad layers of software and configuration that span the gap between your code and the CPU. Managed services can take care of almost all of this *undifferentiated heavy lifting* for you (you'll find out more about the benefits of managed services in [Chapter 3](#)).

The revolution in the cloud has also triggered another revolution in the people who use it: the DevOps movement.

The Dawn of DevOps

Before DevOps, developing and operating software were essentially two separate jobs, performed by two different groups of people. *Developers* wrote software, and they passed it on to *operations* staff, who ran and maintained the software *in production*.

(that is to say, serving real users, instead of merely running internally for testing or feature development purposes). Like the massive mainframe computers that needed their own floor of the building, this separation had its roots in the middle of the last century. *Software development* was a very specialist job, and so was *computer operation*, and there was very little overlap between these two roles.

The two departments had quite different goals and incentives, which often conflicted with each other. Developers tended to focus on shipping new features quickly, while operations teams cared mostly about making services stable and reliable over the long term. In some cases there would be security policies in place that prevented software developers from even having access to the logs or metrics for their own applications running in production. They would need to ask permission from the operations team to debug the application and deploy any fixes. And it was often the operations team who were blamed anytime there was an issue with an application, regardless of the cause.

As cloud computing became more popular, the industry changed. Distributed systems are complex, and the internet is very big. The technicalities of these distributed systems—when it comes to recovering from failures, handling timeouts, smoothly upgrading versions—are not so easy to separate from the design, architecture, and implementation of the system.

Further, “the system” is no longer just your software: it comprises in-house software, cloud services, network resources, load balancers, monitoring, content distribution networks, firewalls, DNS, and so on. All these things are intimately interconnected and interdependent. The people who write the software have to understand how it relates to the rest of the system, and the people who operate the system have to understand how the software works and fails.

Improving Feedback Loops

The origins of the DevOps movement lie in attempts to bring these two groups together: to collaborate, to share understanding, to share responsibility for systems reliability and software correctness, and to improve the scalability of both the software systems and the teams of people who build them.

DevOps is about improving the feedback loops and handoff points that exist between various teams when writing code, building apps, running tests, and deploying changes to ensure that things are running smoothly and efficiently.

What Does DevOps Mean?

DevOps has occasionally been a controversial term to define, both with people who insist it’s nothing more than a modern label for existing good practice in software

development, and with those who reject the need for greater collaboration between development and operations at all.

There is also widespread misunderstanding about what DevOps actually is: A job title? A team? A methodology? A skill set? The influential DevOps writer John Willis has identified four key pillars of DevOps, which he calls *culture, automation, measurement, and sharing* (CAMS). Organizations practicing DevOps have a culture that embraces collaboration, rejects siloing knowledge between teams, and comes up with ways of measuring how they can be constantly improving. Another way to break it down is what Brian Dawson has called the DevOps trinity: *people and culture, process and practice, and tools and technology*.

Some people think that cloud and containers mean that we no longer need DevOps—a point of view sometimes called *NoOps*. The idea is that since all IT operations are outsourced to a cloud provider, or another third-party service, businesses don't need full-time operations staff.

The NoOps fallacy is based on a misapprehension of what DevOps work actually involves:

With DevOps, much of the traditional IT operations work happens before code reaches production. Every release includes monitoring, logging, and A/B testing. CI/CD pipelines automatically run unit tests, security scanners, and policy checks on every commit. Deployments are automatic. Controls, tasks, and non-functional requirements are now implemented before release instead of during the frenzy and aftermath of a critical outage.

—Jordan Bach ([AppDynamics](#))

For now, you will find lots of job postings for the title of DevOps Engineer and a huge range of what is expected of that role, depending on the organization. Sometimes it will look more like a traditional “sysadmin” role and have little interaction with software engineers. Sometimes the role will be embedded alongside developers building and deploying their own applications. It is important to consider what DevOps means to you and what you want it to look like at an organization.

The most important thing to understand about DevOps is that it is primarily an organizational, human issue, not a technical one. This accords with Jerry Weinberg’s Second Law of Consulting:

No matter how it looks at first, it’s always a people problem.

—Gerald M. Weinberg, *The Secrets of Consulting*

And DevOps does really work. Studies regularly suggest that companies that adopt DevOps principles release better software faster, react better and faster to failures and problems, are more agile in the marketplace, and dramatically improve the quality of their products:

DevOps is not a fad; rather it is the way successful organizations are industrializing the delivery of quality software today and will be the new baseline tomorrow and for years to come.

—Brian Dawson, *CloudBees*

Infrastructure as Code

Once upon a time, developers dealt with software, while operations teams dealt with hardware and the operating systems that run on that hardware.

Now that hardware is in the cloud, everything, in a sense, is software. The DevOps movement brings software development skills to operations: tools and workflows for rapid, agile, collaborative building of complex systems. This is often referred to as *infrastructure as code* (IaC).

Instead of physically racking and cabling computers and switches, cloud infrastructure can be automatically provisioned by software. Instead of manually deploying and upgrading hardware, operations engineers have become the people who write the software that automates the cloud.

The traffic isn't just one-way. Developers are learning from operations teams how to anticipate the failures and problems inherent in distributed, cloud-based systems, how to mitigate their consequences, and how to design software that degrades gracefully and fails safely.

Learning Together

Both development teams and operations teams are learning how to work together. They're learning how to design and build systems, how to monitor and get feedback on systems in production, and how to use that information to improve the systems. Even more importantly, they're learning to improve the experience for their users, and to deliver better value for the business that funds them.

The massive scale of the cloud and the collaborative, code-centric nature of the DevOps movement have turned operations into a software problem. At the same time, they have also turned software into an operations problem. All of which raises these questions:

- How do you deploy and upgrade software across large, diverse networks of different server architectures and operating systems?
- How do you deploy to distributed environments, in a reliable and reproducible way, using largely standardized components?

Enter the third revolution: the *container*.

The Coming of Containers

To deploy a piece of software, you need not only the software itself, but its *dependencies*. That means libraries, interpreters, subpackages, compilers, extensions, and so on.

You also need its *configuration*: settings, site-specific details, license keys, database passwords—everything that turns raw software into a usable service.

The State of the Art

Earlier attempts to solve this problem include using *configuration management* systems, such as Puppet or Ansible, which consist of code to install, run, configure, and update the software.

Another solution is the *omnibus package*, which, as the name suggests, attempts to cram everything the application needs inside a single file. An omnibus package contains the software, its configuration, its dependent software components, *their* configuration, *their* dependencies, and so on. (For example, a Java omnibus package would contain the Java runtime as well as all the Java Archive [JAR] files for the application.)

Some vendors have even gone a step further and included the entire computer system required to run it, as a *virtual machine image* (VM image), but these are large and unwieldy, time-consuming to build and maintain, fragile to operate, slow to download and deploy, and vastly inefficient in performance and resource footprint.

From an operations point of view, not only do you need to manage these various kinds of packages, but you also need to manage a fleet of servers to run them on.

Servers need to be provisioned, networked, deployed, configured, kept up-to-date with security patches, monitored, managed, and so on.

This all takes a significant amount of time, skill, and effort just to provide a platform to run software on. Isn't there a better way?

Thinking Inside the Box

To solve these problems, the tech industry borrowed an idea from the shipping industry: the *container*. In the 1950s, a truck driver named **Malcolm McLean** proposed that, instead of laboriously unloading goods individually from the truck trailers that brought them to the ports and loading them onto ships, the trucks themselves—or rather, the truck bodies—could be loaded onto the ship.

A truck trailer is essentially a big metal box on wheels. If you can separate the box—the container—from the wheels and chassis used to transport it, you have something that is very easy to lift, load, stack, and unload, and can go right onto a ship or another truck at the other end of the voyage. Containers also use standard

dimensions, which allows the entire shipping industry including boats, trains, and trucks, to know what to expect when it comes to moving them from place to place. (Figure 1-2).

McLean's container shipping firm, Sea-Land, became very successful by using this system to ship goods far more cheaply, and **containers quickly caught on**. Today, hundreds of millions of containers are shipped every year, carrying trillions of dollars' worth of goods.



Figure 1-2. Standardized containers dramatically cut the cost of shipping bulk goods (photo by [Lucarelli](#), licensed under Creative Commons)

Putting Software in Containers

The software container is exactly the same idea: a standard packaging and distribution format that is generic and widespread, enabling greatly increased carrying capacity, lower costs, economies of scale, and ease of handling. The container format contains everything the application needs to run, baked into an *image file* that can be executed by a *container runtime*.

How is this different from a virtual machine image? That, too, contains everything the application needs to run—but a lot more besides. A typical VM image is around 1 GiB.¹ A well-designed container image, on the other hand, might be a hundred times smaller.

Because the virtual machine contains lots of unrelated programs, libraries, and things that the application will never use, most of its space is wasted. Transferring VM images across the network is far slower than optimized containers.

Even worse, virtual machines are *virtual*: the underlying physical CPU effectively implements an *emulated* CPU, which the virtual machine runs on. The virtualization layer has a dramatic, negative effect on **performance**: in tests, virtualized workloads run about 30% slower than the equivalent containers.

In comparison, containers run directly on the real CPU, with no virtualization overhead, just as ordinary binary executables do.

And because containers only hold the files they need, they’re much smaller than VM images. They also use a clever technique of addressable filesystem *layers*, which can be shared and reused between containers.

For example, if you have two containers, each derived from the same Debian Linux base image, the base image only needs to be downloaded once, and each container can simply reference it.

The container runtime will assemble all the necessary layers and only download a layer if it’s not already cached locally. This makes very efficient use of disk space and network bandwidth.

Plug and Play Applications

Not only is the container the unit of deployment and the unit of packaging; it is also the unit of *reuse* (the same container image can be used as a component of many different services), the unit of *scaling*, and the unit of *resource allocation* (a container can run anywhere sufficient resources are available for its own specific needs).

Developers no longer have to worry about maintaining different versions of the software to run on different Linux distributions, against different library and language versions, and so on. The only thing the container depends on is the operating system kernel (Linux, for example).

Simply supply your application in a container image, and it will run on any platform that supports the standard container format and has a compatible kernel.

¹ The *gibibyte* (GiB) is the International Electrotechnical Commission (IEC) unit of data, defined as 1,024 *mebibytes* (MiB), and *kibibyte* (KiB), is defined as 1,024 bytes. We’ll use IEC units (GiB, MiB, KiB) throughout this book to avoid any ambiguity.

Kubernetes developers Brendan Burns and David Oppenheimer put it this way in their paper “[Design Patterns for Container-Based Distributed Systems](#)”:

By being hermetically sealed, carrying their dependencies with them, and providing an atomic deployment signal (“succeeded”/“failed”), [containers] dramatically improve on the previous state of the art in deploying software in the datacenter or cloud. But containers have the potential to be much more than just a better deployment vehicle—we believe they are destined to become analogous to objects in object-oriented software systems, and as such will enable the development of distributed system design patterns.

Conducting the Container Orchestra

Operations teams, too, find their workload greatly simplified by containers. Instead of having to maintain a sprawling estate of machines of various kinds, architectures, and operating systems, all they have to do is run a *container orchestrator*: a piece of software designed to join together many different machines into a *cluster*. A container orchestrator is a kind of unified compute substrate, which appears to the user as a single very powerful computer on which containers can run.

The terms *orchestration* and *scheduling* are often used loosely as synonyms. Strictly speaking, though, orchestration in this context means coordinating and sequencing different activities in service of a common goal (like the musicians in an orchestra). Scheduling means managing the resources available and assigning workloads where they can most efficiently be run. (Not to be confused with scheduling in the sense of *scheduled jobs*, which execute at preset times.)

A third important activity is *cluster management*: joining multiple physical or virtual servers into a unified, reliable, fault-tolerant, apparently seamless group.

The term *container orchestrator* usually refers to a single service that takes care of scheduling, orchestration, and cluster management.

Containerization (using containers as your standard method of deploying and running software) offered obvious advantages, and a de facto standard container format has made possible all kinds of economies of scale. But one problem still stood in the way of the widespread adoption of containers: the lack of a standard container orchestration system.

As long as several different tools for scheduling and orchestrating containers competed in the marketplace, businesses were reluctant to place expensive bets on which technology to use. But all that was about to change.

Kubernetes

Google was running containers at scale for production workloads long before anyone else. Nearly all of Google’s services run in containers: Gmail, Google Search, Google

Maps, Google App Engine, and so on. Because no suitable container orchestration system existed at the time, Google was compelled to invent one.

From Borg to Kubernetes

To solve the problem of running a large number of services at global scale on millions of servers, Google developed a private, internal container orchestration system it called **Borg**.

Borg is essentially a centralized management system that allocates and schedules containers to run on a pool of servers. While very powerful, Borg is tightly coupled to Google's own internal and proprietary technologies, difficult to extend, and impossible to release to the public.

In 2014, Google founded an open source project named Kubernetes (from the Greek word κυβερνήτης, meaning “helmsman, pilot”) that would develop a container orchestrator that everyone could use, based on the lessons learned from Borg and its successor, **Omega**.

The rise of Kubernetes was meteoric. While other container orchestration systems existed before Kubernetes, none have had quite the same widespread adoption that Kubernetes has found. With the advent of a truly free and open source container orchestrator, adoption of both containers and Kubernetes grew at a phenomenal rate.

Kubernetes continues to grow in popularity and is becoming the norm for running containerized applications. According to a report published by **Datadog**:

Kubernetes has become the de facto standard for container orchestration. Today, half of organizations running containers use Kubernetes, whether in self-managed clusters, or through a cloud provider service... Kubernetes adoption has more than doubled since 2017, and continues to grow steadily, without any signs of slowing down.

Much like containers standardized the way software is packaged and deployed, Kubernetes is standardizing the platform on which to run those containers.

Why Kubernetes?

Kelsey Hightower, a staff developer advocate at Google, coauthor of **Kubernetes Up & Running** (O'Reilly), and all-around legend in the Kubernetes community, has put it this way:

Kubernetes does the things that the very best system administrator would do: automation, failover, centralized logging, monitoring. It takes what we've learned in the DevOps community and makes it the default, out of the box.

—Kelsey Hightower

Many of the traditional sysadmin tasks like upgrading servers, installing security patches, configuring networks, and running backups are less of a concern in the cloud native world. Kubernetes can automate these things for you so that your team can concentrate on doing its core work.

Some of these features, like *load balancing* and *autoscaling*, are built into the Kubernetes core; others are provided by add-ons, extensions, and third-party tools that use the Kubernetes API. The Kubernetes ecosystem is large, and growing all the time.

Kubernetes makes deployment easy

Ops staff love Kubernetes for these reasons, but there are also some significant advantages for developers. Kubernetes greatly reduces the time and effort it takes to deploy. Zero-downtime deployments are common, because Kubernetes does rolling updates by default (starting containers with the new version, waiting until they become healthy, and then shutting down the old ones).

Kubernetes also provides facilities to help you implement continuous deployment practices such as *canary deployments*: gradually rolling out updates one server at a time to catch problems early (see “[Canary Deployments](#)” on page 247). Another common practice is *blue-green* deployments: spinning up a new version of the system in parallel, and switching traffic over to it once it’s fully up and running (see “[Blue/Green Deployments](#)” on page 246).

Demand spikes will no longer take down your service, because Kubernetes supports autoscaling. For example, if CPU utilization by a container reaches a certain level, Kubernetes can keep adding new replicas of the container until the utilization falls below the threshold. When demand falls, Kubernetes will scale down the replicas again, freeing up cluster capacity to run other workloads.

Because Kubernetes has redundancy and failover built in, your application will be more reliable and resilient. Some managed services can even scale the Kubernetes cluster itself up and down in response to demand so that you’re never paying for a larger cluster than you need at any given moment (see “[Autoscaling](#)” on page 104). That does mean that your applications need to be designed in a way to run in a dynamic environment, but Kubernetes gives you standard ways to leverage that sort of infrastructure.

The business will love Kubernetes too, because it cuts infrastructure costs and makes much better use of a given set of resources. Traditional servers, even cloud servers, are mostly idle most of the time. The excess capacity that you need to handle demand spikes is essentially wasted under normal conditions.

Kubernetes takes that wasted capacity and uses it to run workloads, so you can achieve much higher utilization of your machines—and you get scaling, load balancing, and failover for free too.

While some of these features, such as autoscaling, were available before Kubernetes, they were always tied to a particular cloud provider or service. Kubernetes is *provider-agnostic*: once you've defined the resources you use, you can run them on any Kubernetes cluster, regardless of the underlying cloud provider.

That doesn't mean that Kubernetes limits you to the lowest common denominator. Kubernetes maps your resources to the appropriate vendor-specific features: for example, a load-balanced Kubernetes service on Google Cloud will create a Google Cloud load balancer; on Amazon, it will create an Amazon Web Services (AWS) load balancer. Kubernetes abstracts away the cloud-specific details, letting you focus on defining the behavior of your application.

Just as containers are a portable way of defining software, Kubernetes resources provide a portable definition of how that software should run.

Will Kubernetes Disappear?

Oddly enough, despite the current excitement around Kubernetes, we may not be talking much about it in years to come. Many things that once were new and revolutionary are now so much part of the fabric of computing that we don't really think about them: microprocessors, the mouse, the internet.

Kubernetes, too, is likely to fade into the background and become part of the plumbing. It's boring, in a good way! Once you learn what you need to know to deploy your application to Kubernetes, you can spend your time focusing on adding features to your application.

Managed service offerings for Kubernetes will likely do more and more of the heavy lifting behind running Kubernetes itself. In 2021, Google Cloud Platform (GCP) released a new offering to their existing Kubernetes service called Autopilot that handles cluster upgrades, networking, and scaling the VMs up and down depending on the demand. Other cloud providers are also moving in that direction and offering Kubernetes-based platforms where developers only need to worry about running their application and not focus on the underlying infrastructure.

Kubernetes Is Not a Panacea

Will all software infrastructure of the future be entirely Kubernetes-based? Probably not. Is it incredibly easy and straightforward to run any and all types of workloads? Not quite.

For example, running databases on distributed systems requires careful consideration as to what happens around restarts and how to ensure that data remains consistent.

Orchestrating software in containers involves spinning up new interchangeable instances without requiring coordination between them. But database replicas are not interchangeable; they each have a unique state, and deploying a database replica requires coordination with other nodes to ensure things like schema changes happen everywhere at the same time.

—Sean Loiselle, Cockroach Labs

While it's perfectly possible to run stateful workloads like databases in Kubernetes with enterprise-grade reliability, it requires a large investment of time and engineering that it may not make sense for your company to make (see “[Run Less Software](#)” [on page 45](#)). It's usually more cost-effective to use a managed database service instead.

Secondly, some things may not actually need Kubernetes, and can run on what are sometimes called *serverless* platforms, better named *functions as a service* (FaaS) platforms.

Cloud functions

AWS Lambda, for example, is a FaaS platform that allows you to run code written in Go, Python, Java, Node.js, C#, and other languages without you having to compile or deploy your application at all. Amazon does all that for you. Google Cloud has similar offerings with Cloud Run and Functions, and Microsoft also offers Azure Functions.

Because you're billed for the execution time in increments of milliseconds, the FaaS model is perfect for computations that only run when you need them to, instead of paying for a cloud server, which runs all the time whether you're using it or not.

These cloud functions are more convenient than containers in some ways (though some FaaS platforms can run containers as well). But they are best suited to short, standalone jobs (AWS Lambda limits functions to 15 minutes of run time, for example), especially those that integrate with existing cloud computation services, such as Azure Cognitive Services or the Google Cloud Vision API.

These types of event-driven platforms are often called “serverless” models. Technically, there is still a server involved: it's just somebody else's server. The point is that you don't have to provision and maintain that server; the cloud provider takes care of it for you.

Not every workload is suitable for running on FaaS platforms, by any means, but it is still likely to be a key technology for cloud native applications in the future.

Nor are cloud functions restricted to public FaaS platforms such as Lambda Functions or Azure Functions: if you already have a Kubernetes cluster and want to run

FaaS applications on it, open source projects like [OpenFaaS](#) and [Knative](#) make this possible.

Some of these Kubernetes serverless platforms encompass both long-running containers and event-driven short-lived functions, which may mean that in the future the distinction between these types of compute may blur or disappear altogether.

Cloud Native

The term *cloud native* has become an increasingly popular shorthand way of talking about modern applications and services that take advantage of the cloud, containers, and orchestration, often based on open source software.

Indeed, the [Cloud Native Computing Foundation \(CNCF\)](#) was founded in 2015 to, in their words, “foster a community around a constellation of high-quality projects that orchestrate containers as part of a microservices architecture.”

Part of the Linux Foundation, the CNCF exists to bring together developers, end users, and vendors, including the major public cloud providers. The best-known project under the CNCF umbrella is Kubernetes itself, but the foundation also incubates and promotes other key components of the cloud native ecosystem: Prometheus, Envoy, Helm, Fluentd, gRPC, and many more.

So what exactly do we mean by *cloud native*? Like most such things, it means different things to different people, but perhaps there is some common ground.

First, *cloud* does not necessarily mean a public cloud provider, like AWS or Azure. Many organizations run their own internal “cloud” platforms, often while also simultaneously using one or multiple public providers for different workloads. The term *cloud* loosely means the platform of servers used to run software infrastructure, and that can take many forms.

So what makes an application cloud native? Just taking an existing application and running it on a cloud compute instance does not make it cloud native. Neither is it just about running it in a container, or using cloud services such as Azure’s Cosmos DB or Google’s Pub/Sub, although those may well be important aspects of a cloud native application.

So let's look at a few of the characteristics of cloud native systems that most people can agree on:

Automatable

If applications are to be deployed and managed by machines, instead of humans, they need to abide by common standards, formats, and interfaces. Kubernetes provides these standard interfaces in a way that means application developers don't even need to worry about them.

Ubiquitous and flexible

Because they are decoupled from physical resources such as disks, or any specific knowledge about the compute node they happen to be running on, containerized microservices can easily be moved from one node to another, or even one cluster to another.

Resilient and scalable

Traditional applications tend to have single points of failure: the application stops working if its main process crashes, or if the underlying machine has a hardware failure, or if a network resource becomes congested. Cloud native applications, because they are inherently distributed, can be made highly available through redundancy and graceful degradation.

Dynamic

A container orchestrator such as Kubernetes can schedule containers to take maximum advantage of available resources. It can run many copies of containers to achieve high availability, and perform rolling updates to smoothly upgrade services without ever dropping traffic.

Observable

Cloud native apps, by their nature, are harder to inspect and debug. So a key requirement of distributed systems is *observability*: monitoring, logging, tracing, and metrics all help engineers understand what their systems are doing (and what they're doing wrong).

Distributed

Cloud native is an approach to building and running applications that takes advantage of the distributed and decentralized nature of the cloud. It's about how your application works, not where it runs. Instead of deploying your code as a single entity (known as a *monolith*), cloud native applications tend to be composed of multiple, cooperating, distributed *microservices*. A microservice is simply a self-contained service that does one thing. If you put enough microservices together, you get an application.

It's not just about microservices

However, microservices are also not a panacea. Monoliths are easier to understand, because everything is in one place, and you can trace the interactions of different parts. But it's hard to scale a monolith, both in terms of the code itself and the teams of developers who maintain it. As the code grows, the interactions between its various parts grow exponentially, and the system as a whole grows beyond the capacity of a single brain to understand it all.

A well-designed cloud native application is composed of microservices, but deciding what those microservices should be, where the boundaries are, and how the different services should interact is no easy problem. Good cloud native service design consists of making wise choices about how to separate the different parts of your architecture. However, even a well-designed cloud native application is still a distributed system, which makes it inherently complex, difficult to observe and reason about, and prone to failure in surprising ways.

While cloud native systems tend to be distributed, it's still possible to run monolithic applications in the cloud, using containers, and gain considerable business value from doing so. This may be a step on the road to gradually migrating parts of the monolith outward to modern microservices, or a stopgap measure pending the redesign of the system to be fully cloud native.

The Future of Operations

Operations, infrastructure engineering, and system administration are highly skilled jobs. Are they at risk in a cloud native future? We think not.

Instead, these skills will only become more important. Designing and reasoning about distributed systems is hard. Networks and container orchestrators are complicated. Every team developing cloud native applications will need operations skills and knowledge. Automation frees up staff from boring, repetitive, manual work to deal with more complex, interesting, and fun problems that computers can't yet solve for themselves.

That doesn't mean all current operations jobs are guaranteed. Sysadmins used to be able to get by without coding skills, except maybe cooking up the odd simple shell script. In the cloud native world, that won't be enough to succeed.

In a software-defined world, the ability to write, understand, and maintain software becomes critical. If you don't want to learn new skills, the industry will leave you behind—and it has always been that way.

Distributed DevOps

Rather than being concentrated in a single operations team that services other teams, ops expertise will become distributed among many teams.

Each development team will need at least one ops specialist, responsible for the health of the systems or services the team provides. They will be a developer as well, but they will also be the domain expert on networking, Kubernetes, performance, resilience, and the tools and systems that enable the other developers to deliver their code to the cloud.

Thanks to the DevOps revolution, there will no longer be room in most organizations for devs who can't ops, or ops who don't dev. The distinction between those two disciplines is obsolete and is rapidly being erased altogether. Developing and operating software are merely two aspects of the same thing.

Some Things Will Remain Centralized

Are there limits to DevOps? Or will the traditional central IT and operations team disappear altogether, dissolving into a group of roving internal consultants, coaching, teaching, and troubleshooting ops issues?

We think not, or at least not entirely. Some things still benefit from being centralized. It doesn't make sense for each application or service team to have its own way of detecting and communicating about production incidents, for example, or its own ticketing system, or deployment tools. There's no point in everybody reinventing their own wheel.

Developer Productivity Engineering

The point is that self-service has its limits, and the aim of DevOps is to speed up development teams, not slow them down with unnecessary and redundant work.

Yes, a large part of traditional operations can and should be devolved to other teams, primarily those that deploy code and respond to code-related incidents. But to enable that to happen, there needs to be a strong central team building and supporting the DevOps ecosystem in which all the other teams operate.

Instead of calling this *team operations*, we like the name *developer productivity engineering*. Some organizations call this role *platform engineer* or maybe even *DevOps engineer*. The point is that these teams do whatever is necessary to help other software engineering teams do their work better and faster: operating infrastructure, building tools, busting problems.

And while developer productivity engineering remains a specialist skill set, the engineers themselves may move outward into the organization to bring that expertise where it's needed.

Lyft engineer Matt Klein has suggested that, while a pure DevOps model makes sense for startups and small firms, as an organization grows, there is a natural tendency for infrastructure and reliability experts to gravitate toward a central team. But he says that team can't be scaled indefinitely:

By the time an engineering organization reaches ~75 people, there is almost certainly a central infrastructure team in place starting to build common substrate features required by product teams building microservices. But there comes a point at which the central infrastructure team can no longer both continue to build and operate the infrastructure critical to business success, while also maintaining the support burden of helping product teams with operational tasks.

—Matt Klein

At this point, not every developer can be an infrastructure expert, just as a single team of infrastructure experts can't service an ever-growing number of developers. For larger organizations, while a central infrastructure team is still needed, there's also a case for embedding *site reliability engineers* (SREs) into each development or product team. They bring their expertise to each team as consultants and also form a bridge between product development and infrastructure operations.

You Are the Future

If you're reading this book, it means you are a part of this new cloud native future. In the remaining chapters, we'll cover all the knowledge and skills you'll need as a developer or operations engineer working with cloud infrastructure, containers, and Kubernetes.

Some of these things will be familiar, and some will be new, but we hope that when you've finished the book you'll feel more confident in your own ability to acquire and master cloud native skills. Yes, there's a lot to learn, but it's nothing you can't handle. You've got this!

Now read on.

Summary

We've necessarily given you a rather quick tour of the landscape, including the history of DevOps, cloud computing, and the emerging standard of using containers and Kubernetes for running cloud native applications. We hope it's enough to bring you up to speed with some of the challenges in this field and how they're likely to change the IT industry.

A quick recap of the main points before we move on to meet Kubernetes in person in the next chapter:

- Cloud computing frees you from the expense and overhead of managing your own hardware, making it possible for you to build resilient, flexible, scalable distributed systems.
- DevOps is a recognition that modern software development doesn't stop at shipping code: it's about closing the feedback loop between those who write the code and those who use it.
- DevOps also brings a code-centric approach and good software engineering practices to the world of infrastructure and operations.
- Containers allow you to deploy and run software in small, standardized, self-contained units. This makes it easier and cheaper to build large, diverse, distributed systems, by connecting together containerized microservices.
- Orchestration systems take care of deploying your containers, scheduling, scaling, networking, and all the things that a good system administrator would do, but in an automated, programmable way.
- Kubernetes is the de facto standard container orchestration system, and it's ready for you to use in production right now, today. It is still a fast-moving project, and all of the major cloud providers are offering more managed services to handle the underlying core Kubernetes components automatically.
- “Serverless” event-driven computing is also becoming popular for cloud native applications, often using containers as the runtime. Tools are available to run these types of functions on Kubernetes clusters.
- Cloud native is a useful shorthand for talking about cloud-based, containerized, distributed systems, made up of cooperating microservices, dynamically managed by automated infrastructure as code.
- Operations and infrastructure skills, far from being made obsolete by the cloud native revolution, are and will become more important than ever.
- What will go away is the sharp distinction between software engineers and operations engineers. It's all just software now, and we're all engineers.