# CS6216 Advanced Topics in Machine Learning (Systems)

# Cloud systems for AI

Yao LU

06 Nov 2024

National University of Singapore
School of Computing

# From LLMs to the cloud



Chef
(LLM)

Restaurant
(serving systems)

Disney world
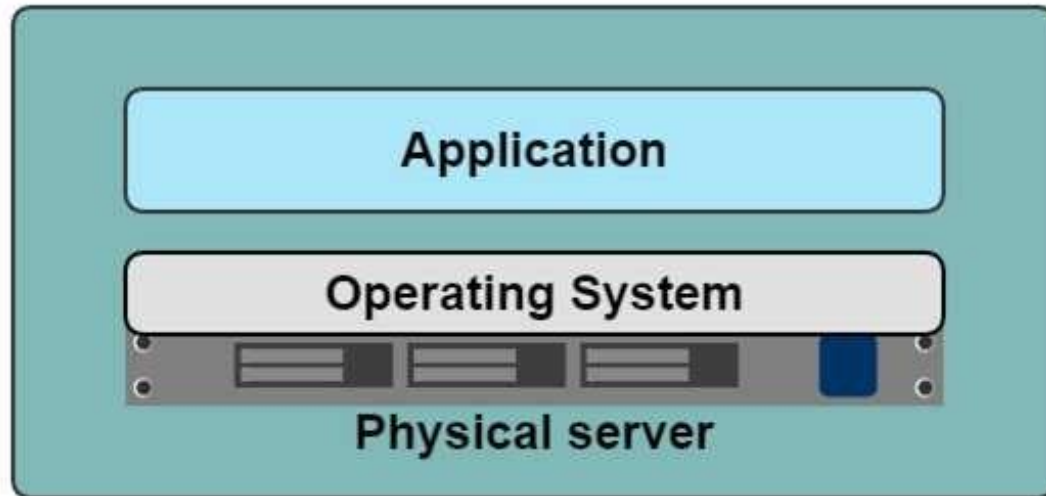(cloud systems)

From serving to cloud systems:

- **Multi-tenancy**: from scaling-up to scaling-out (models, users, applications, tasks etc.)

- **Operations of** large-scale, heterogeneous infrastructures

# Outline

- Brief history of cloud computing

- Cloud native technologies

- Current practice and opportunities of AI on cloud
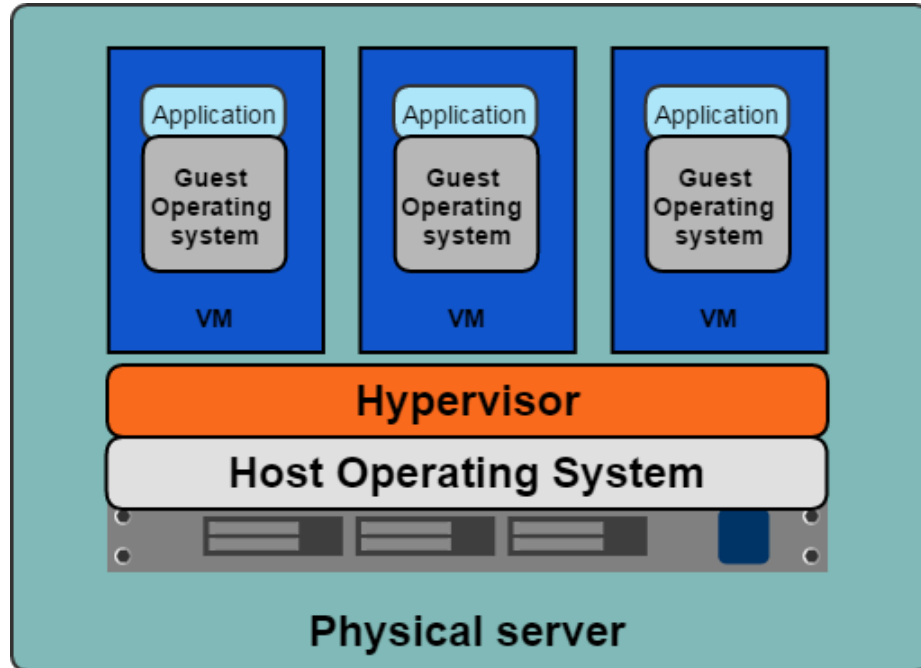
# A history lesson

## In the Dark Ages



- Slow deployment times
- Huge costs & wasted resources
- Difficult to scale & migrate
- Vendor lock in

One application on one physical server

# A history lesson

## Hypervisor-based Virtualization



- One physical server can contain multiple applications
- Each application runs in a virtual machine (VM)

- Better resource pooling
  - One physical machine divided into multiple virtual machines
- Easier to scale
- VMs in the cloud
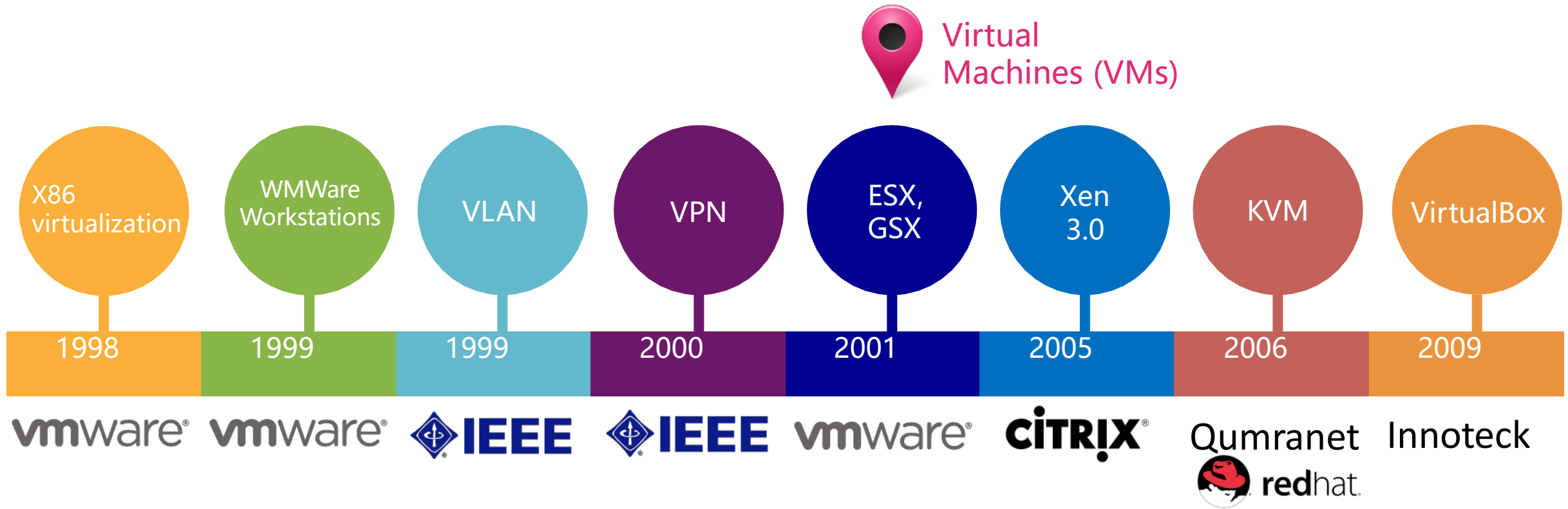  - Rapid elasticity
  - Pay as you go model

# Brief history of cloud computing

| X86 virtualization | WMWare Workstations | VLAN | VPN | ESX, GSX | Xen 3.0 | KVM | VirtualBox |
|---|---|---|---|---|---|---|---|

Virtual Machines (VMs)

| 1998 | 1999 | 1999 | 2000 | 2001 | 2005 | 2006 | 2009 |
|---|---|---|---|---|---|---|---|
| vmware® | vmware® | IEEE | IEEE | vmware® | CiTRIX® | Qumranet redhat | Innoteck |

**Mature of virtualization:** no.1 important technology

# Cloud computing offerings

| IaaS | | PaaS | | Open-source IaaS | Open-source PaaS | | FaaS |
|------|------|------|------|------|------|------|------|
| S3, EC2 | Google App Engine | First PaaS | Microsoft Azure cloud service | Open-source IaaS | Cloud Foundry | Google Cloud Engine | Lambda |
| 2006 | 2008 | 2009 | 2010 | 2011 | 2011 | 2013 | 2014 |
| amazon web services | Google | HEROKU | Microsoft | openstack | Pivotal | Google | amazon web services |

$5,000,000

**2005**  Team of engineers

Months of development

$5,000

**2017**  An engineer

Weeks development

# Cloud computing offerings



| Traditional On-Premises IT | Colocation | Hosting | IaaS | PaaS | SaaS |
|---|---|---|---|---|---|
| Data | Data | Data | Data | Data | Data |
| Application | Application | Application | Application | Application | Application |
| Databases | Databases | Databases | Databases | Databases | Databases |
| Operating System | Operating System | Operating System | Operating System | Operating System | Operating System |
| Virtualization | Virtualization | Virtualization | Virtualization | Virtualization | Virtualization |
| Physical Servers | Physical Servers | Physical Servers | Physical Servers | Physical Servers | Physical Servers |
| Network & Storage | Network & Storage | Network & Storage | Network & Storage | Network & Storage | Network & Storage |
| Data Center | Data Center | Data Center | Data Center | Data Center | Data Center |

■ Provider-Supplied   ■ Self-Managed

# However,

- Each VM stills requires
  - CPU allocation
  - Storage
  - RAM
  - An entire guest operating system

- The more VMs you run, the more resources you need

- Guest OS means wasted resources

- Application portability not guaranteed

# Looking for all kinds of solutions…

| | Development VM | QA Server | Single Prod Server | Onsite Cluster | Public Cloud | Contributor's laptop | Customer Servers |
|---|---|---|---|---|---|---|---|
| **Static website** | ? | ? | ? | ? | ? | ? | ? |
| **Web frontend** | ? | ? | ? | ? | ? | ? | ? |
| **Background workers** | ? | ? | ? | ? | ? | ? | ? |
| **User DB** | ? | ? | ? | ? | ? | ? | ? |
| **Analytics DB** | ? | ? | ? | ? | ? | ? | ? |
| **Queue** | ? | ? | ? | ? | ? | ? | ? |

Too many to consider

# An analogy: cargo transportation

# What are the possibilities

# The challenge continued

# Shipping containers



Multiplicity of Goods

A standard container that is loaded with virtually any goods, and stays sealed until it reaches final delivery.

Do I worry about how goods interact (e.g. coffee beans next to spices)

...in between, can be loaded and unloaded, stacked, transported efficiently over long distances, and transferred from one mode of transport to another

Multiplicity of methods for transporting/storing

Can I transport quickly and smoothly (e.g. from boat to train to truck)
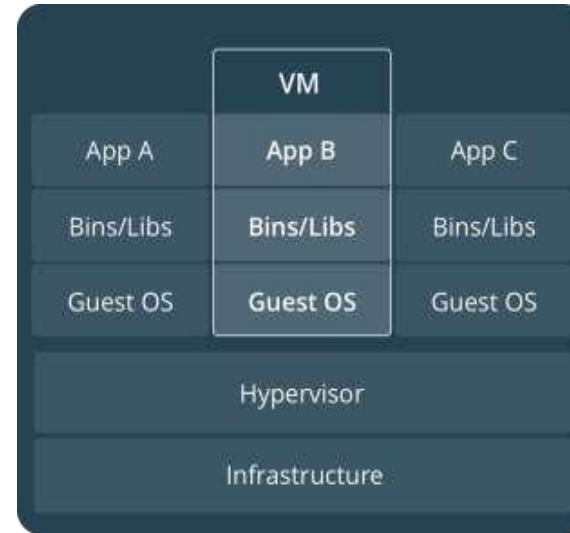
# Container for code?



- **Speed:** share the same OS kernel. No OS to boot = applications online in seconds

- **Portability:** Standardized software packaging. Less dependencies between process layers = ability to move between infrastructure & OS

- **Efficiency:** Less OS overhead & improved VM density
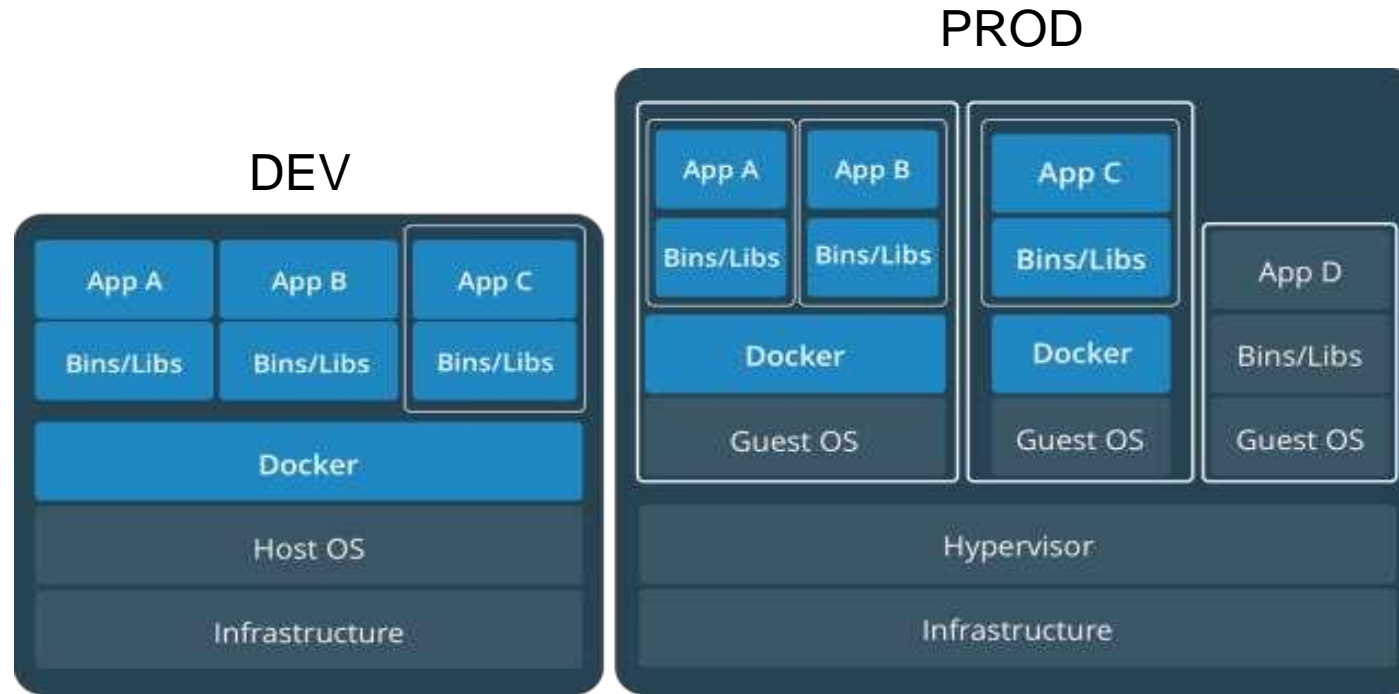
# Comparing containers and VMs



Containers are an app level construct

VMs are an infrastructure level construct to turn one machine into many servers
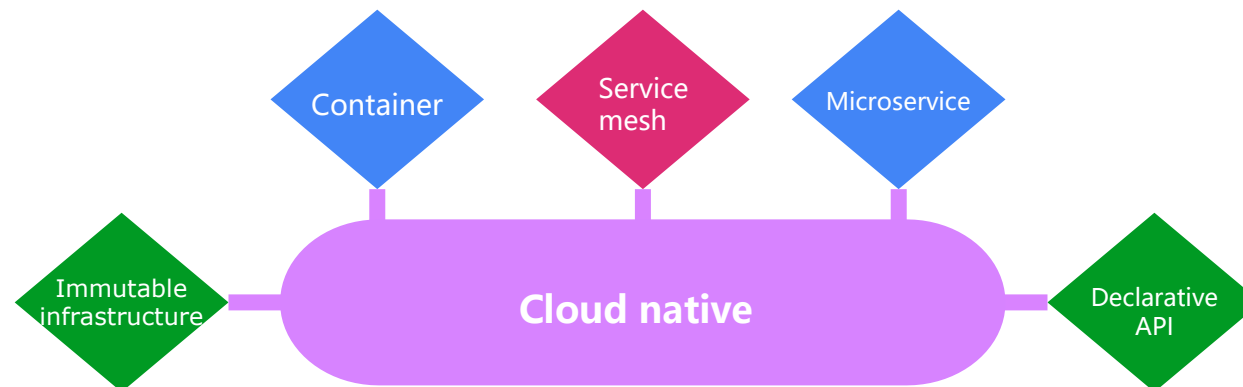
# Containers and VMs together



Containers and VMs together provide a tremendous amount of flexibility for IT to optimally deploy and manage apps.
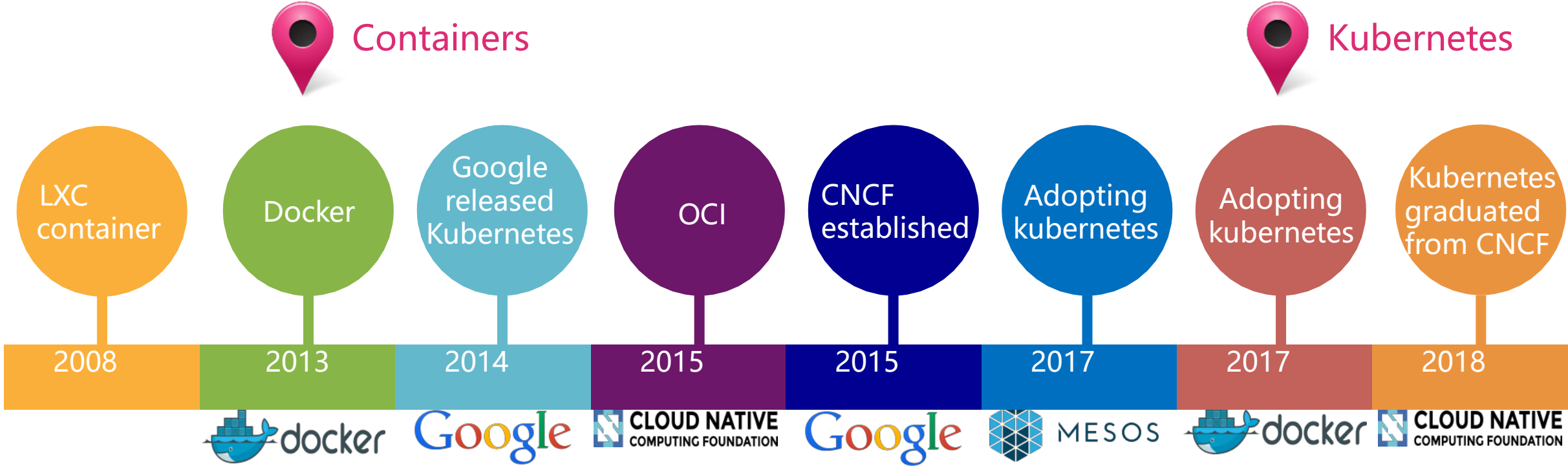
# Cloud native technologies

**Definitions by Cloud Native Computing Foundation (CNCF) :**

**CLOUD NATIVE**
**COMPUTING FOUNDATION**

- Cloud native practices empower organizations to develop, build, and deploy workloads in computing environments (public, private, hybrid cloud) to meet their organizational needs at scale in a programmatic and repeatable manner. It **is characterized by loosely coupled systems that interoperate in a manner that is secure, resilient, manageable, sustainable, and observable.**

- Cloud native technologies and architectures typically consist of some combination of containers, service meshes, multi-tenancy, microservices, immutable infrastructure, serverless, declarative APIs etc.

- Combined with robust automation, cloud native practices allow organizations to make high-impact changes frequently, predictably, with minimal toil and clear separation of concerns.
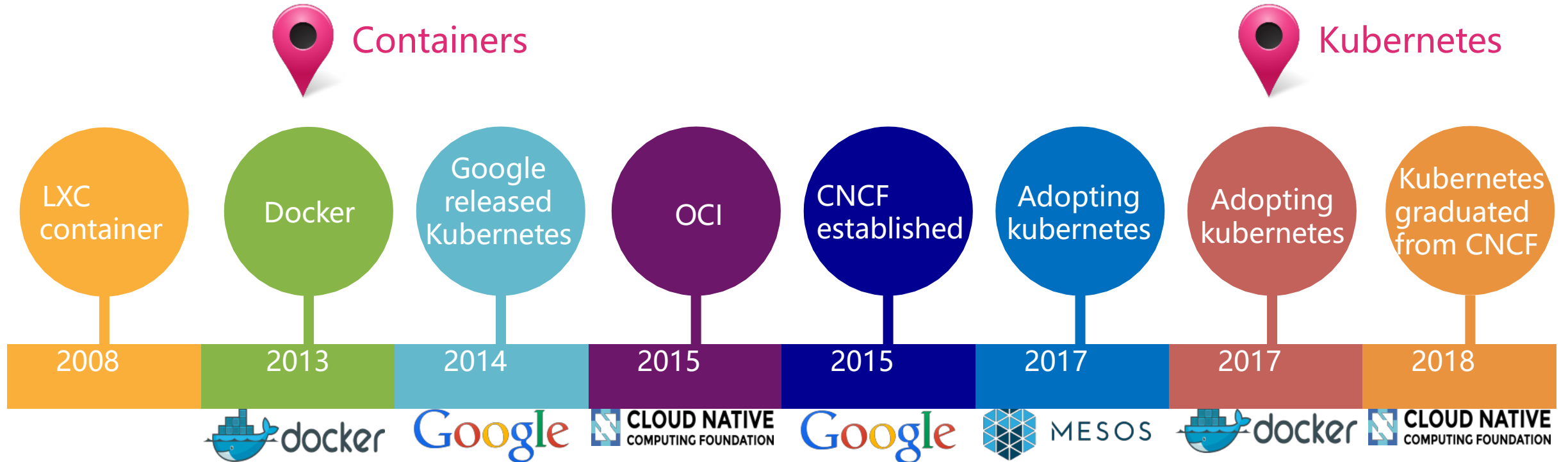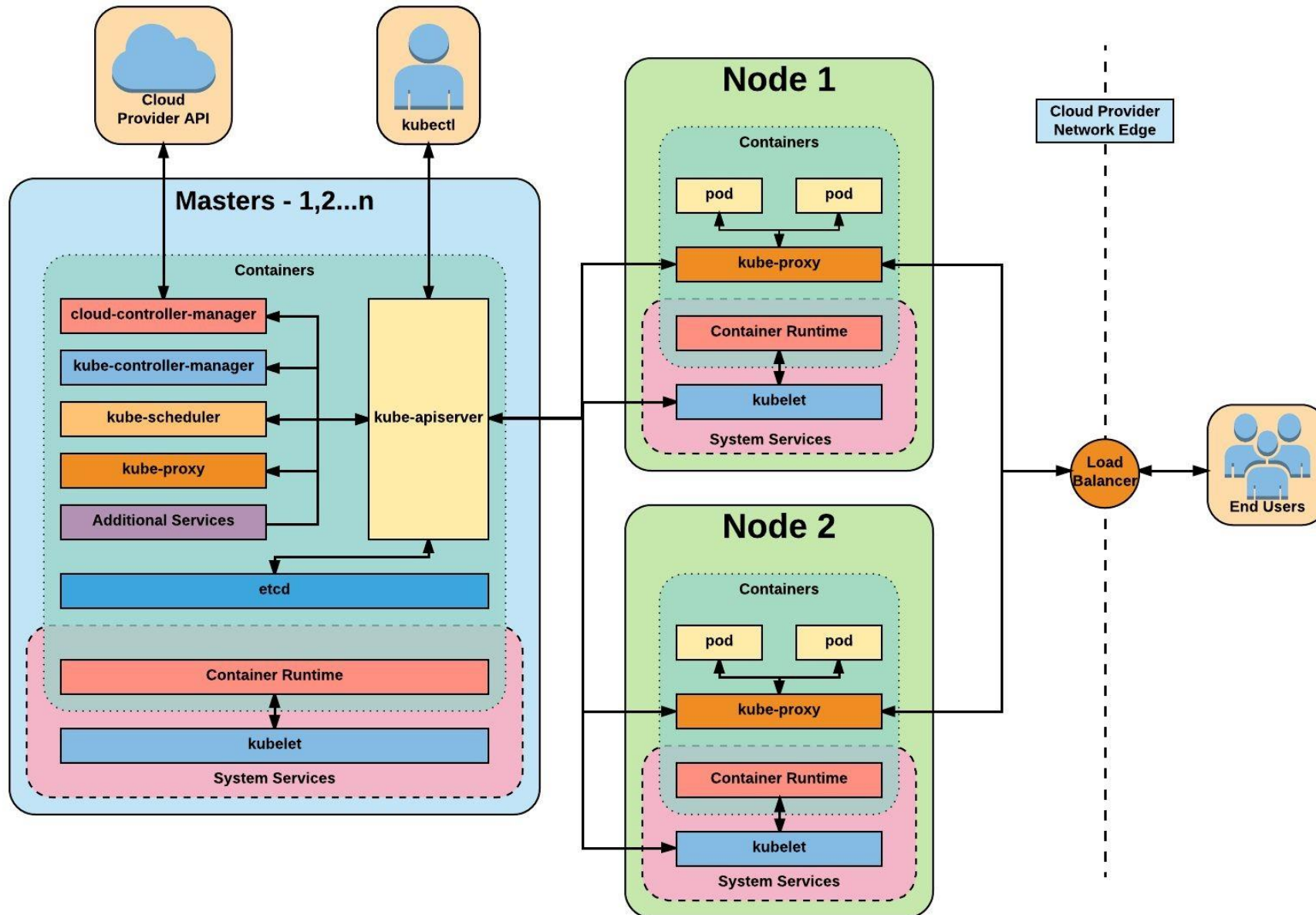
# Rise of containers and Kubernetes (K8s)

Containers

Kubernetes

| LXC container | Docker | Google released Kubernetes | OCI | CNCF established | Adopting kubernetes | Adopting kubernetes | Kubernetes graduated from CNCF |
|---|---|---|---|---|---|---|---|
| 2008 | 2013 | 2014 | 2015 | 2015 | 2017 | 2017 | 2018 |
| | docker | Google | CLOUD NATIVE COMPUTING FOUNDATION | Google | MESOS | docker | CLOUD NATIVE COMPUTING FOUNDATION |

**Kubernetes** or **K8s** is a project spun out of Google as a open source next-gen container scheduler

# Rise of containers and Kubernetes (K8s)

Containers

Kubernetes

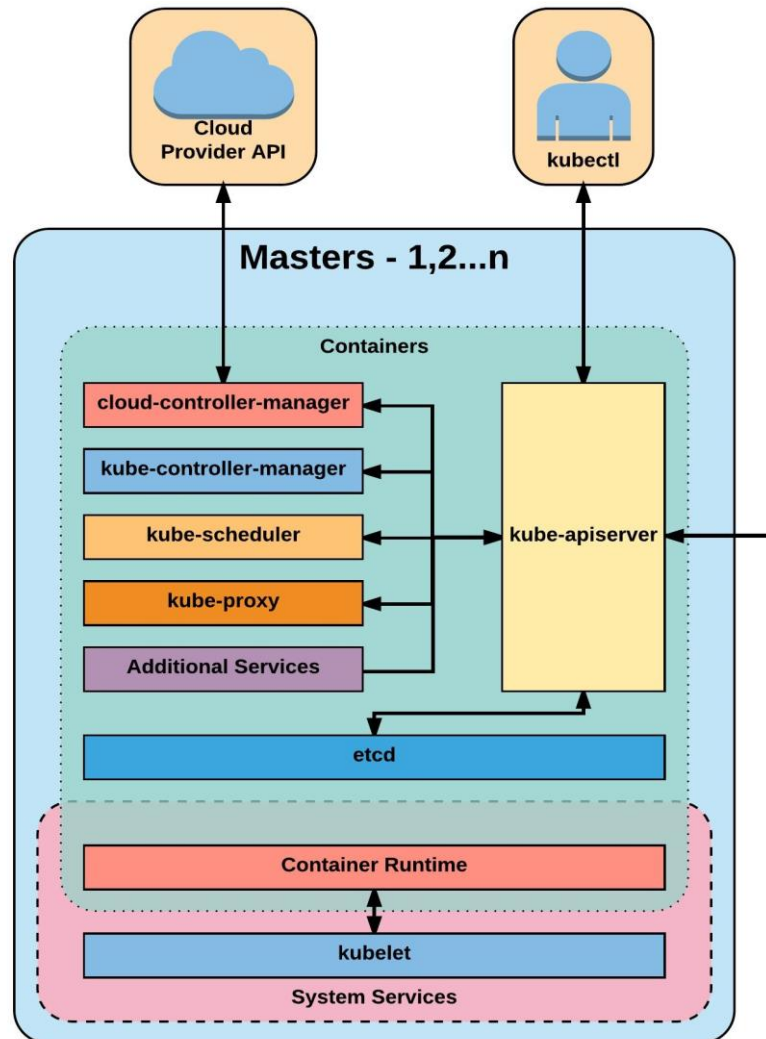| LXC container | Docker | Google released Kubernetes | OCI | CNCF established | Adopting kubernetes | Adopting kubernetes | Kubernetes graduated from CNCF |
|---|---|---|---|---|---|---|---|
| 2008 | 2013 | 2014 | 2015 | 2015 | 2017 | 2017 | 2018 |
| | docker | Google | CLOUD NATIVE COMPUTING FOUNDATION | Google | MESOS | docker | CLOUD NATIVE COMPUTING FOUNDATION |

- K8s is an orchestration tool for managing distributed services or containerized applications across a distributed cluster of nodes.

- K8s follows a **client-server architecture with a master and worker nodes**. Core concepts in Kubernetes include pods, services (logical pods with a stable IP address) and deployments (a definition of the desired state for a pod or replica set).

- K8s **users define rules** for how container management should occur, and then K8s handles the rest
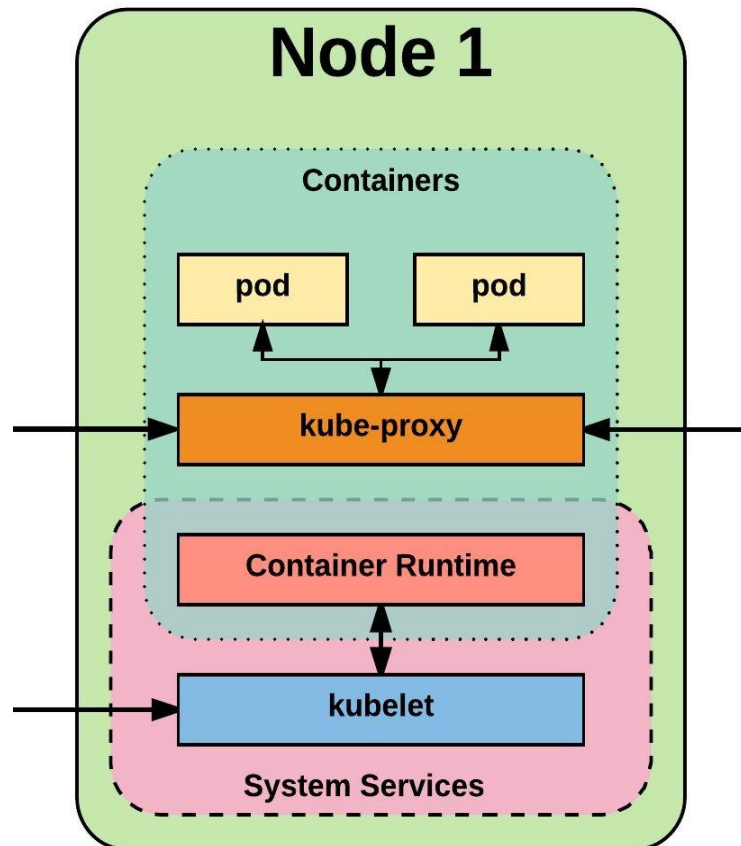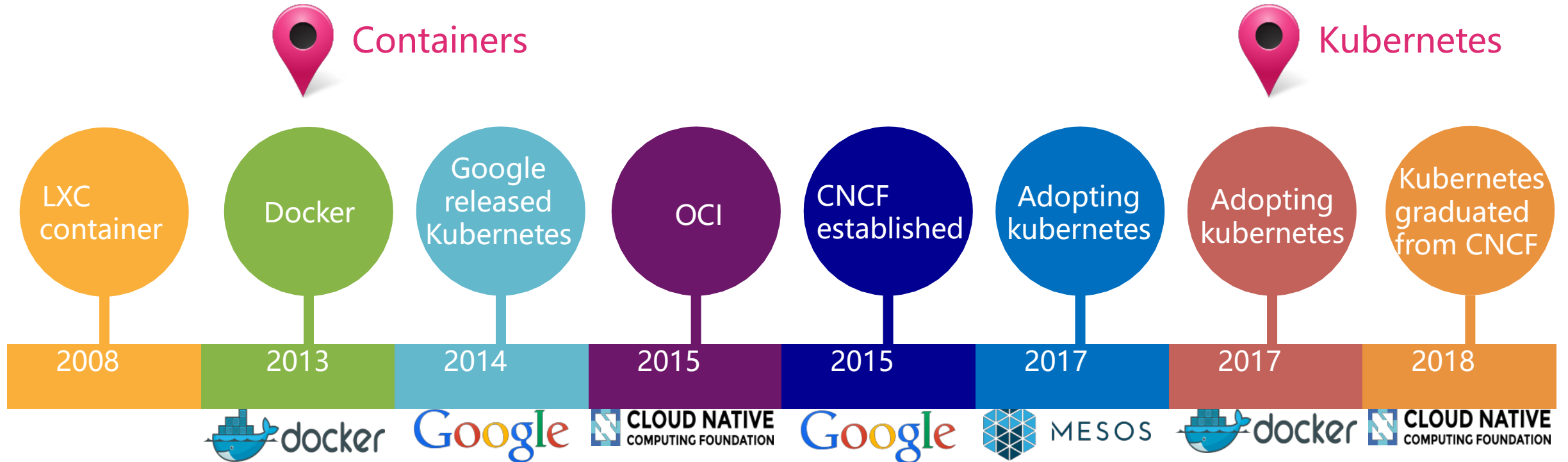
# Architecture overview

# Master components



- **Kube-apiserver:** provides REST interface into the K8s control plane and datastore.

- **Etcd:** the cluster datastore; providing a strong, consistent and highly available key-value store used for persisting cluster state

- **Kube-controller-manager:** manages all core component control loops; monitors and steers the cluster towards the desired state.

- **Cloud-controller-manager:** provides cloud-provider specific knowledge and integration capability.

- **Kube-scheduler:** evaluates workload resource requirements and place it on a matching resource.

# Node components



- **kubelet:** node agent for managing pod lifecycle on its host.
- **kube-proxy:** managing the network rules on each node and performs connection forwarding or load balancing.
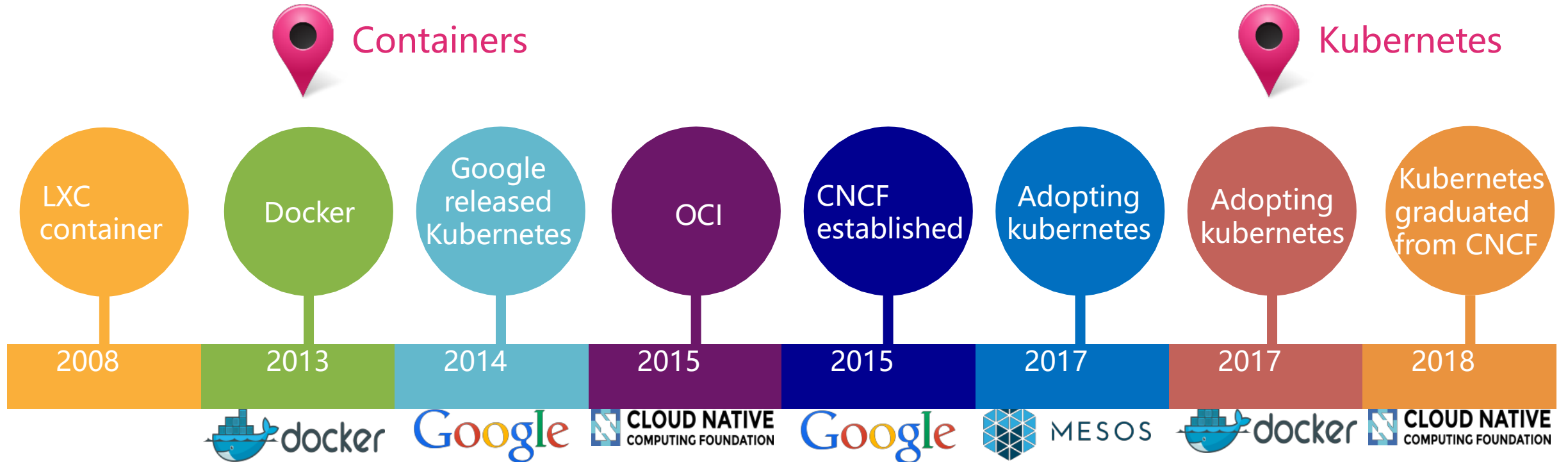- **container runtime:** executes and manages containers.

# Rise of containers and Kubernetes (K8s)

Containers

Kubernetes

| 2008 | 2013 | 2014 | 2015 | 2015 | 2017 | 2017 | 2018 |
|------|------|------|------|------|------|------|------|
| LXC container | Docker | Google released Kubernetes | OCI | CNCF established | Adopting kubernetes | Adopting kubernetes | Kubernetes graduated from CNCF |
| | docker | Google | CLOUD NATIVE COMPUTING FOUNDATION | Google | MESOS | docker | CLOUD NATIVE COMPUTING FOUNDATION |

Advantages of using K8s in reliable & efficient software deployment

- **Velocity:** fast to deploy while maintaining availability by immutable infrastructures & declarative configurations

- **Scaling:** fast and auto scaling of software and develop team

- **Infrastructure abstraction:** applications-infrastructure separation & portability

- **Efficiency:** lower costs of running a server, develop/deploy/test software

# Rise of containers and Kubernetes (K8s)

Containers

Kubernetes

| 2008 | 2013 | 2014 | 2015 | 2015 | 2017 | 2017 | 2018 |
|------|------|------|------|------|------|------|------|
| LXC container | Docker | Google released Kubernetes | OCI | CNCF established | Adopting kubernetes | Adopting kubernetes | Kubernetes graduated from CNCF |
|  | docker | Google | CLOUD NATIVE COMPUTING FOUNDATION | Google | MESOS | docker | CLOUD NATIVE COMPUTING FOUNDATION |

**Cloud evolvement in the last two decades**

- From physical machines to virtual machines to containers

- Different offerings: IaaS, PaaS, SaaS, Caas, FaaS on Public /private / hybrid cloud

- Kubernetes becoming standard

- High-available service on low-available hardware

# Outline

- Brief history of cloud computing
- Cloud native technologies
- Current practice and opportunities of AI on cloud

# Placement and load balancing (PLB)

**Question:** put 18KB into [A: 10KB | B: 20KB | C: 19KB | D: 25KB | E: 30KB]

- **Placement**

  The overall goal is to reduce violation to users' Service Level Agreements (SLAs), given that resource usages are **dynamic**.

  - **First Fit:** the first one that fits → B: 20KB

  - **Best Fit:** the one that just fits → C: 19KB

  - **Worst first:** the one that has the most resource → E: 30KB

  - **More advanced:**

    - **Multi-resource:** memory, disk, CPU, etc.
    - **More complex policies:** constrains, leave-one-out, leave-two-out and so on

# Placement and load balancing (PLB)

- **The real problem:** usages can change → no theoretical guarantee for optimal placement, since everything is data-driven

# Placement and load balancing (PLB)

- **Load balancing:** migrate to "make" some room



**But migration is not free, often very expensive: stateful VMs, DBs, etc.**

- Better migration mechanisms: cache compression, disaggregated memory etc.

- Resource usage prediction: placement & balancing based on predicted usages

# Failovers

- Failovers ensure a robust & highly available service
  - Duplica on independent resources to avoid simultaneous failure



- Failovers are useful for hot software patching / updates

- Failovers can co-exist with PLB which makes it a lot more complex

# Failovers

- **A fast failover** involves efficient context switch & recovery

| | | |
|---|---|---|
| App 3 | App 4 | |
| App 2 | App 3 | App 2 |
| App 1 | App 1 | App 4 |

- Route requests to duplica
- Recover main from duplica
  - Reinstate using logs
  - Live migration

Failover due to:
- Unexpected: software/hardware failure
- Scheduled: Software update

# Serverless computing

- Some "rewrap" of ideas, but many cloud-native techniques are the same underneath

**No servers to provision or manage.** User describes application; system finds out best provision.

**Scales with usage.** System expanse and shrinks automatically with actual usage.

**Build-in availability and fault tolerance.** System also provides safety belts at no cost to the users.

**No pay for idle.** Billing model – user pays only for actual usage; financial risks at the operator.

Containers & orchestrators

(New) Serverless functions & overscriptions

# Serverless functions, or Function-as-a-service (FaaS)

- FaaS is an example of serverless computing to simply deployment of event-driven function calls

- Examples
  - **Netflix:** media encoding, thumbnailing, content recommendation
  - **Airbnb:** user authentication, process booking requests, payment
  - **Coca-cola:** supply-chain operations, personalized promotions

# Resource oversubscription



100 seats, sell 105 tickets

Overbooking of cloud resources

CPU RAM

- (Almost) direct revenue boost , given the base at **$100B!**
- But still, new technologies needed

# Resource oversubscription

- **Technology prerequisites:**
    - Virtualization to cut CPU/disk/memory into fine granularity
    - Quick allocation / migration
    - Multi-tenancy over shared resources

- **Key problem for oversubscription:**
    - Increase oversubscription rate, while reduce/prevent violations w/ user SLAs
      SLA can be latency of query, service availability, etc.

- **Mechanisms for an oversubscribed system:**
    - Similar PLB but at high resource usages

# Outline

- Brief history of cloud computing

- Cloud native technologies

- Current practice and opportunities of AI on cloud

# Practice and opportunities of AI on GPU cloud

- Cloud native technologies are open-box solutions for may AI use cases
- **Key idea:** containerizing your models

# Practice and opportunities of AI on GPU cloud

- But, some cloud-native ideas couldn't be applied

  - GPU virtualization and oversubscription

  - Fine-grained scheduling and operations

    - Each container is a big black-box

# Resource oversubscription for GPUs?

User / App 1 ⬌  GPU(s) 1

User / App 2 ⬌  GPU(s) 2

User / App 3 ⬌  GPU(s) 3

*Tight coupling* between
resource specification and allocation

This means it's hard to switch context / allocation, even when the resource is in idle.

# Breaking the tight coupling between apps & allocation

User / App 1

User / App 2

User / App 3

Resource pool

vGPU  vGPU

vGPU  vGPU

vGPU  vGPU

Compilers, Runtime & Resource manager

GPU 1

GPU 2

GPU 3

"Virtulizating" GPUs into thread and memory blocks,
But, big engineering challenges

# Looking forward

- Multi-tenancy in AI workloads
  - Users, apps, tasks, models, adapters, prompts, …

- Breaking the black-boxes
  - Co-design of AI and cloud systems
  - Cloud-native → AI-native



```
lequnchen@rhiannon: ~/download/punica-checkout
(punica-demo) lequnchen@rhiannon:~/download/punica-checkout$
```

Punica: Multi-Tenant LoRA Serving

# Looking forward

- New opportunities in cloud AI
  - **New cloud** with heterogeneous, ephemeral infra
    spot instances, intermittent/green power

  - **New services:** multi-agent AI, RAGs etc.

  - **New hardware & architecture:** RISC-V AI chips, AISC chips

  - **New applications:** AI-for-X

ChatDev

# What we have covered & not covered

- MLsys foundations

- Automatic differentiation

- Hardware acceleration

- Parallelism and training techniques

- Transformers, attention and optimizations

- Serving LLMs

- Fine-tuning and alignment techniques

- AI for systems

- Application systems

- ML compilers

- Cloud systems for AI

<div align="center">Covered</div>

- Compute graph optimization

- Heterogeneous runtime

- Serving multi-modal models

- Serving mixture-of-experts models

- ML ops


- Many more..

<div align="center">Not covered</div>

Wish you all the best in your PhD / Masters journey!