

# CS4262/5462 Machine Learning Systems

## Cloud systems & AI

Yao LU

09 Apr 2026

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:

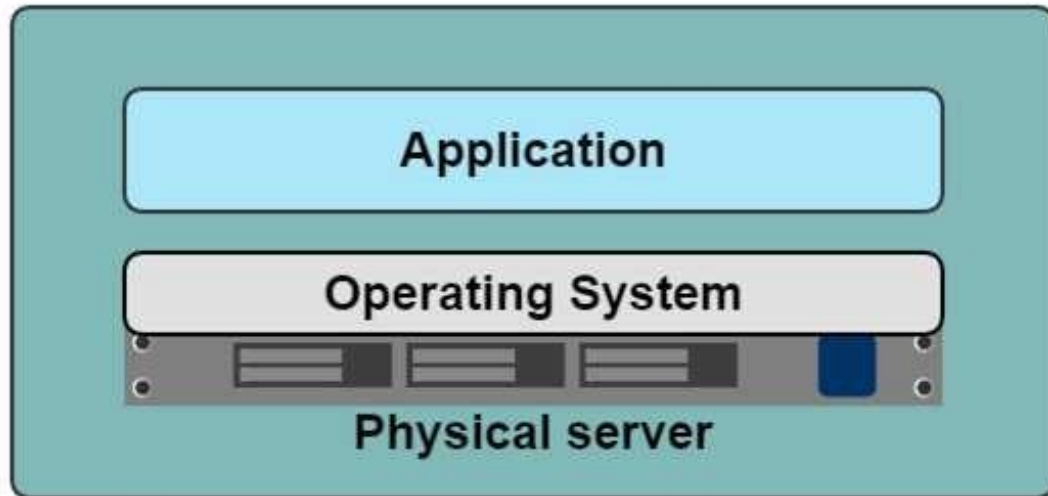
- **Multi-tenancy**: from scaling-up to scaling-out (models, users, applications, tasks etc.)
- **Operations of** large-scale, heterogeneous infrastructures

# Outline

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

# A history lesson

In the Dark Ages (1960s)

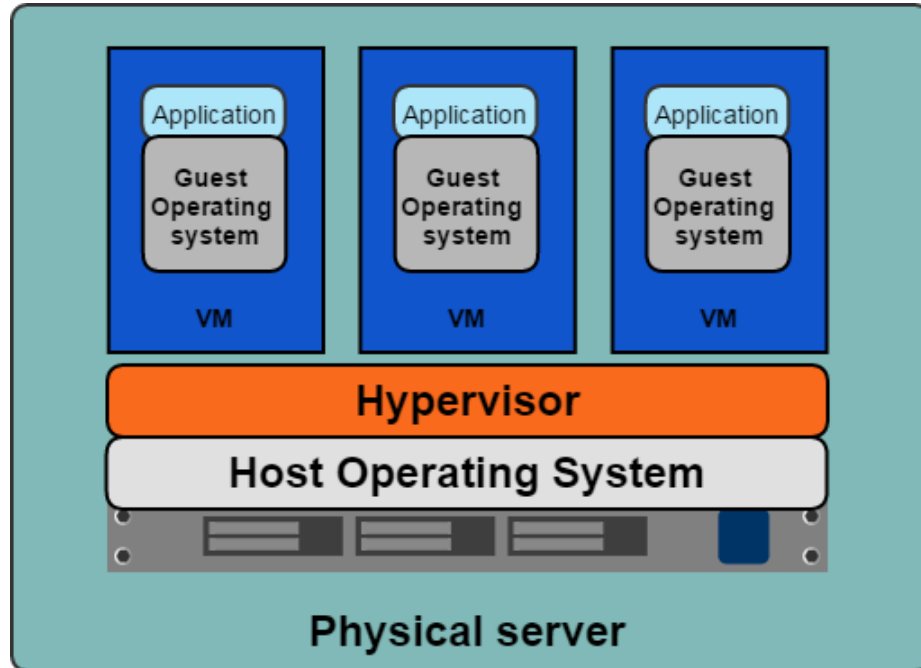


- Slow deployment times
- Huge costs & wasted resources
- Difficult to scale & migrate
- Vendor lock in

One application on one physical server

# A history lesson

## Hypervisor-based Virtualization (1970s)



- 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

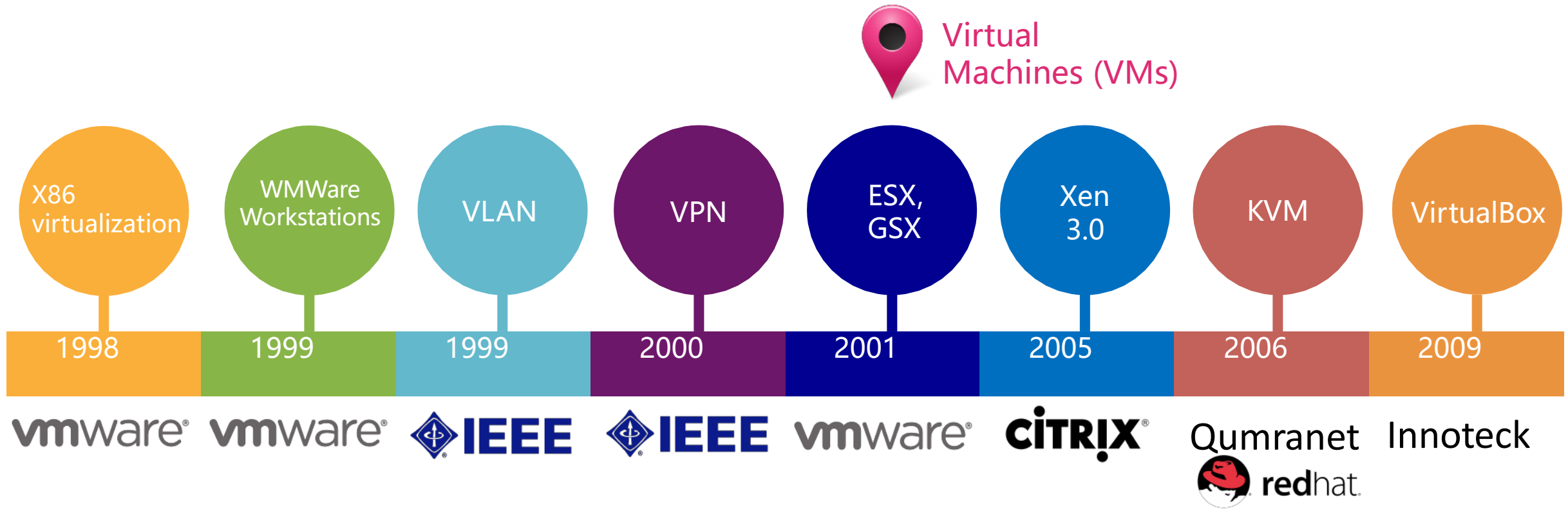
- One physical server can contain multiple applications
- Each application runs in a virtual machine (VM)

 Microsoft Azure

 amazon  
web services™

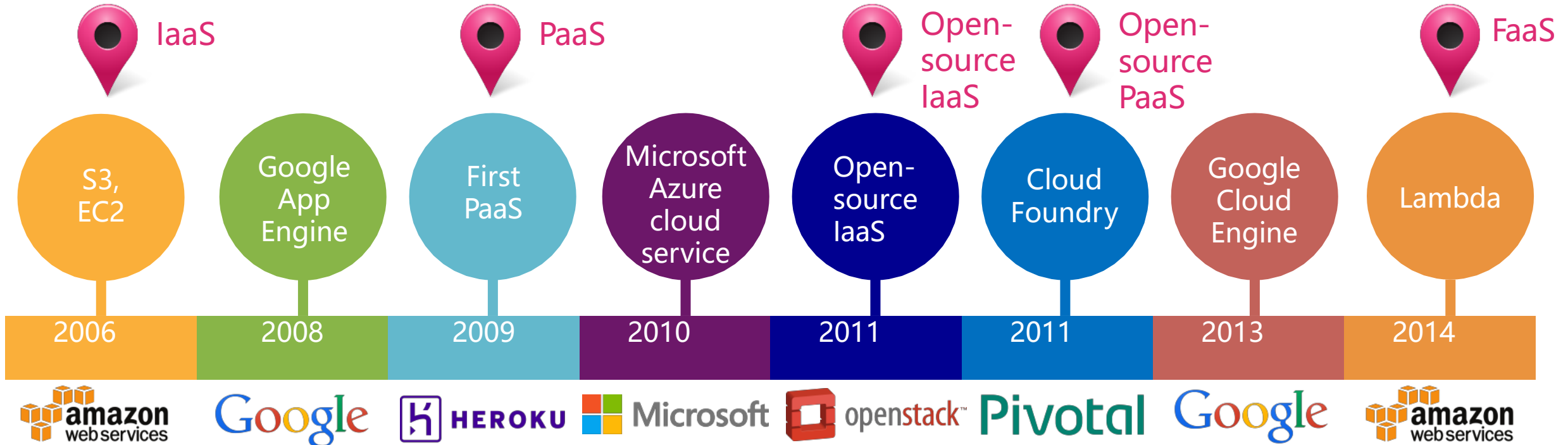
 vmware®

# Brief history of cloud computing



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

# Cloud computing offerings



\$5,000,000

\$5,000

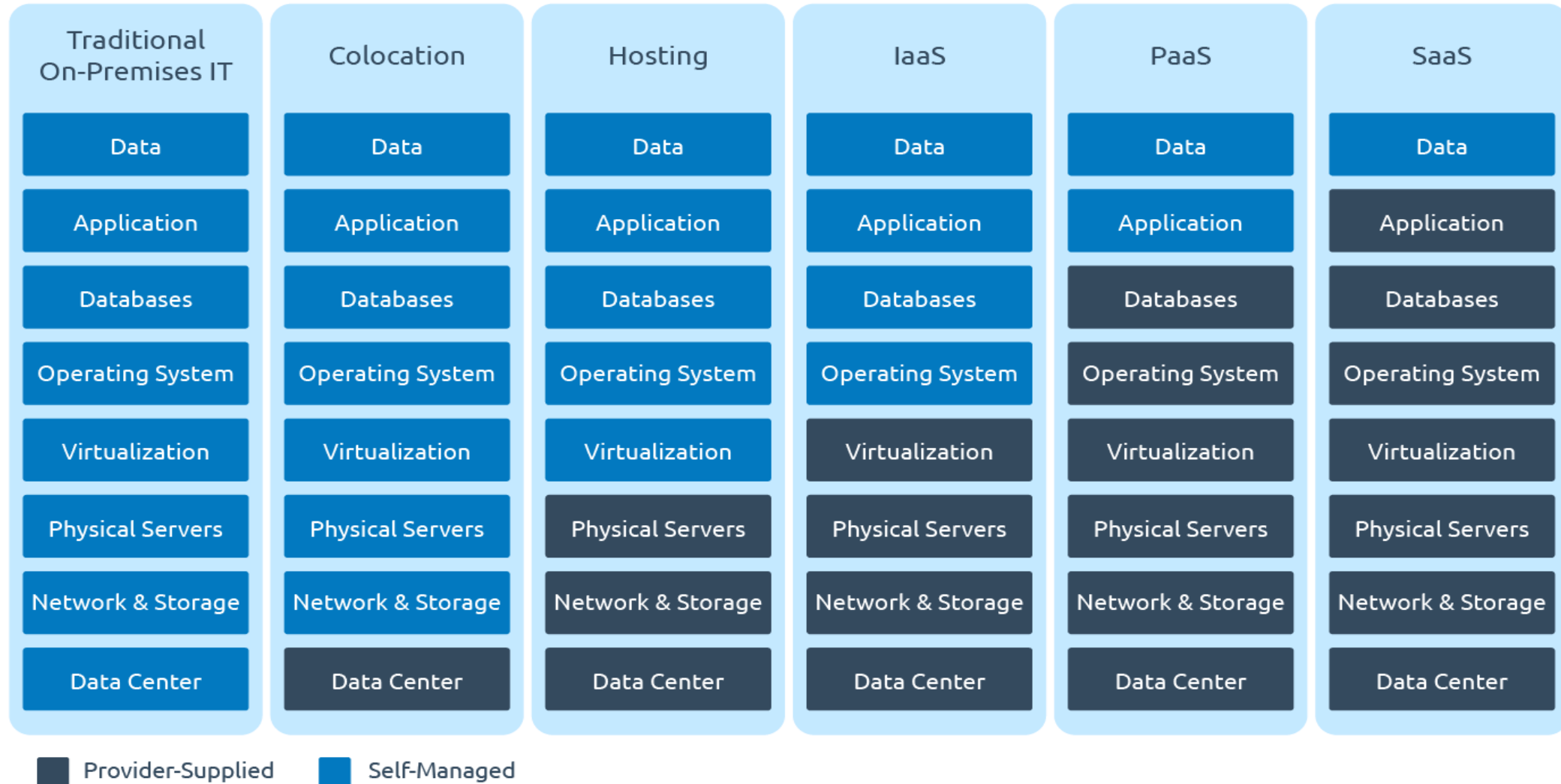
**2005** Team of engineers

**2017** An engineer

Months of development

Weeks development

# Cloud computing offerings



# 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



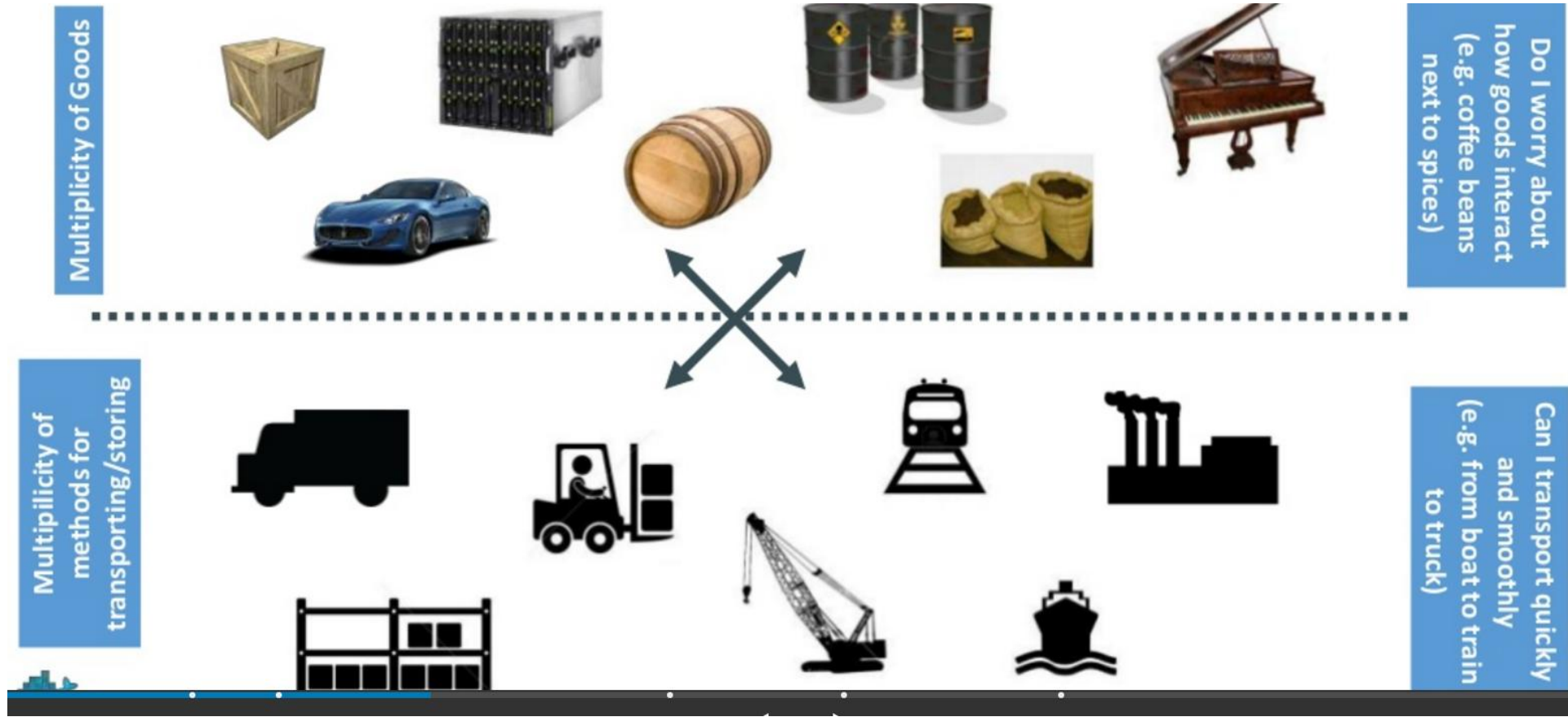
# Looking for all kinds of solutions...

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
















Too many to consider

# An analogy: cargo transportation



# What are the possibilities

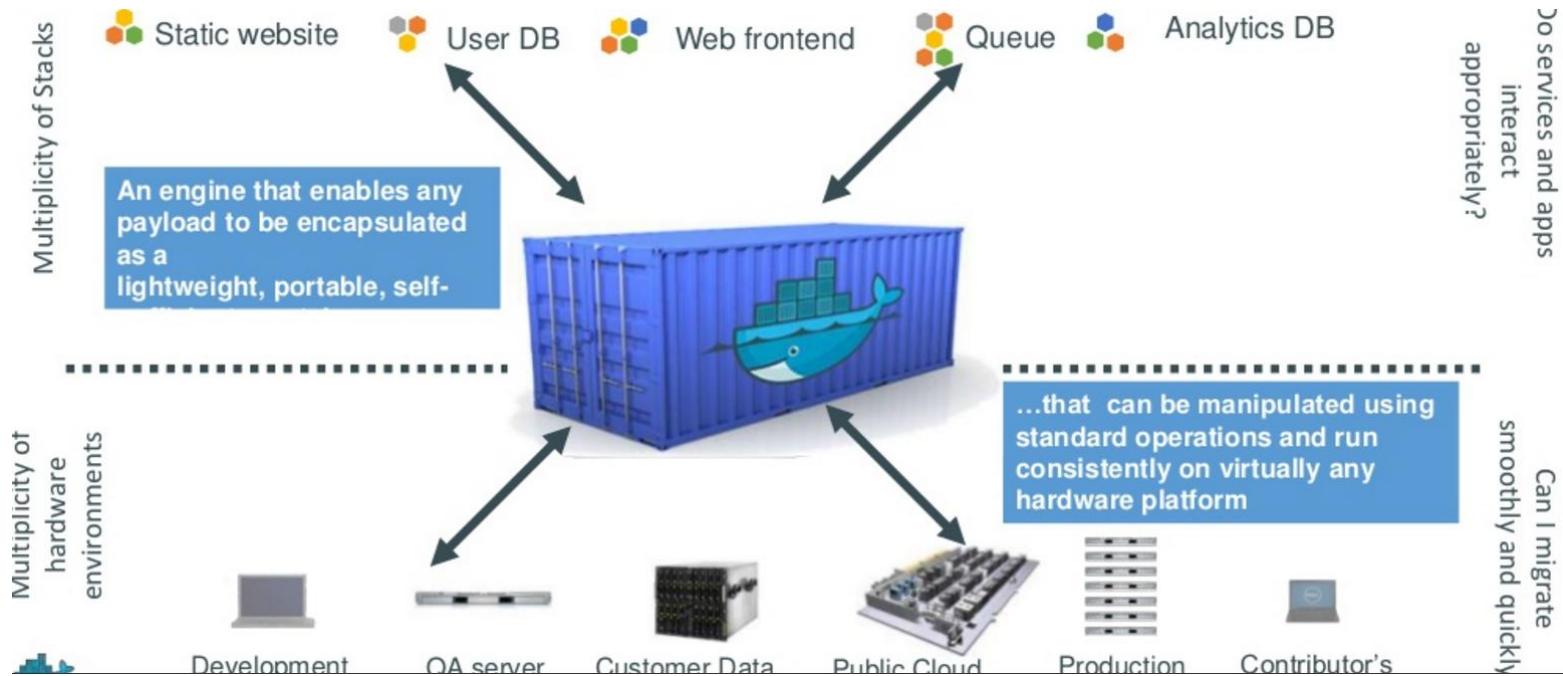
	?	?	?	?	?	?	?
	?	?	?	?	?	?	?
	?	?	?	?	?	?	?
	?	?	?	?	?	?	?
	?	?	?	?	?	?	?
	?	?	?	?	?	?	?
							

Any ideas?

# Shipping containers



# Container for code?

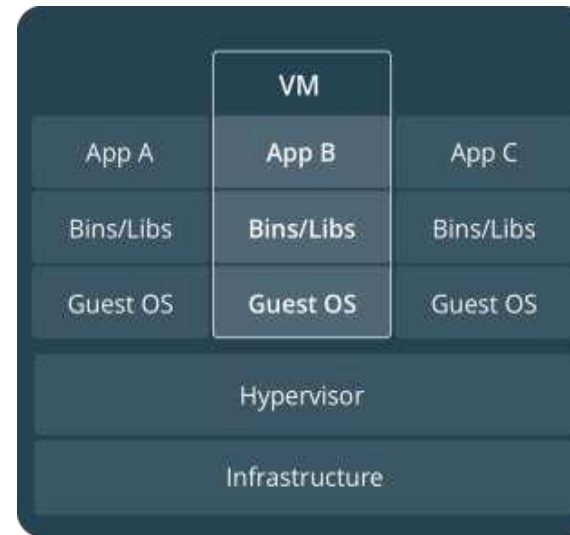


- **Speed:** share the same OS kernel. No OS to boot = applications online in seconds
- **Portability:** Standardized software packaging. Less dependencies between process layers = ability to move between infrastructure & OS
- **Efficiency:** Less OS overhead & improved VM density
- **Isolation:** Achieved by Linux namespace and Cgroups

# Comparing containers and VMs

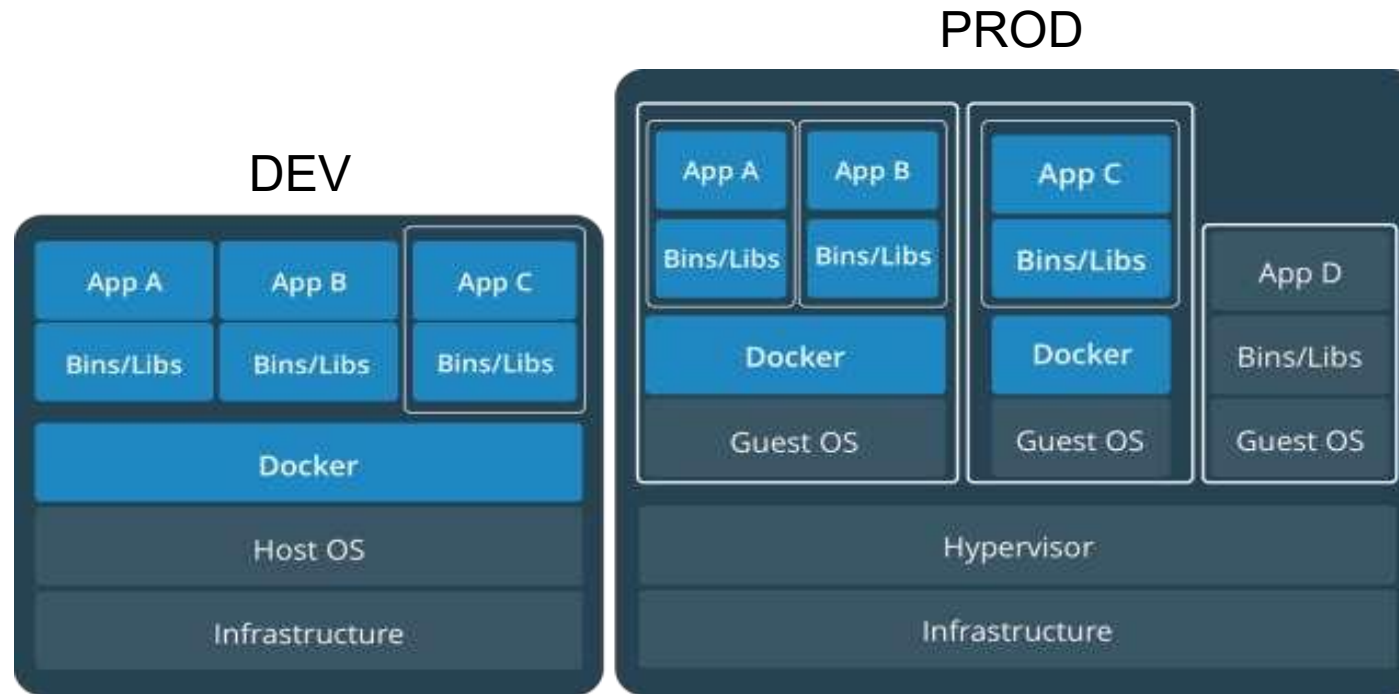


Containers are an app level construct



VMs are an infrastructure level construct to turn one machine into many servers

# Containers and VMs together



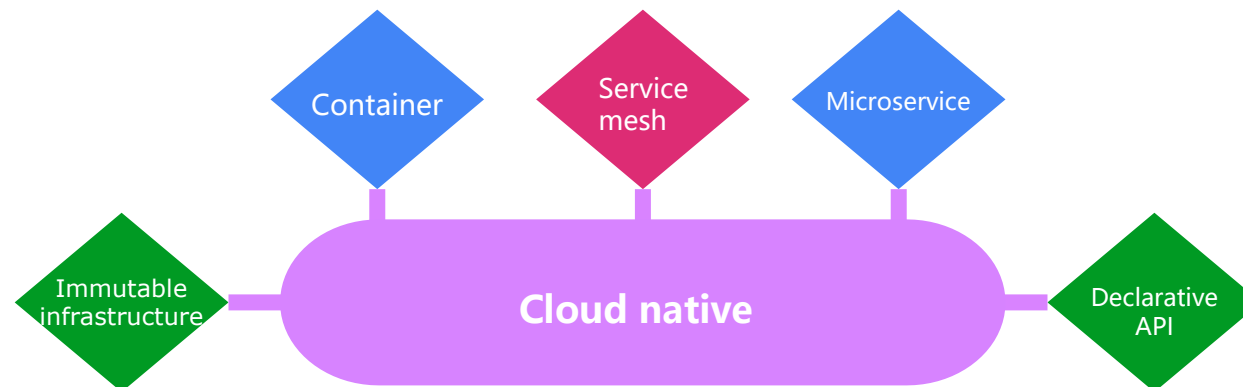
Containers and VMs together provide a tremendous amount of flexibility for IT to optimally deploy and manage apps.

# Cloud native technologies

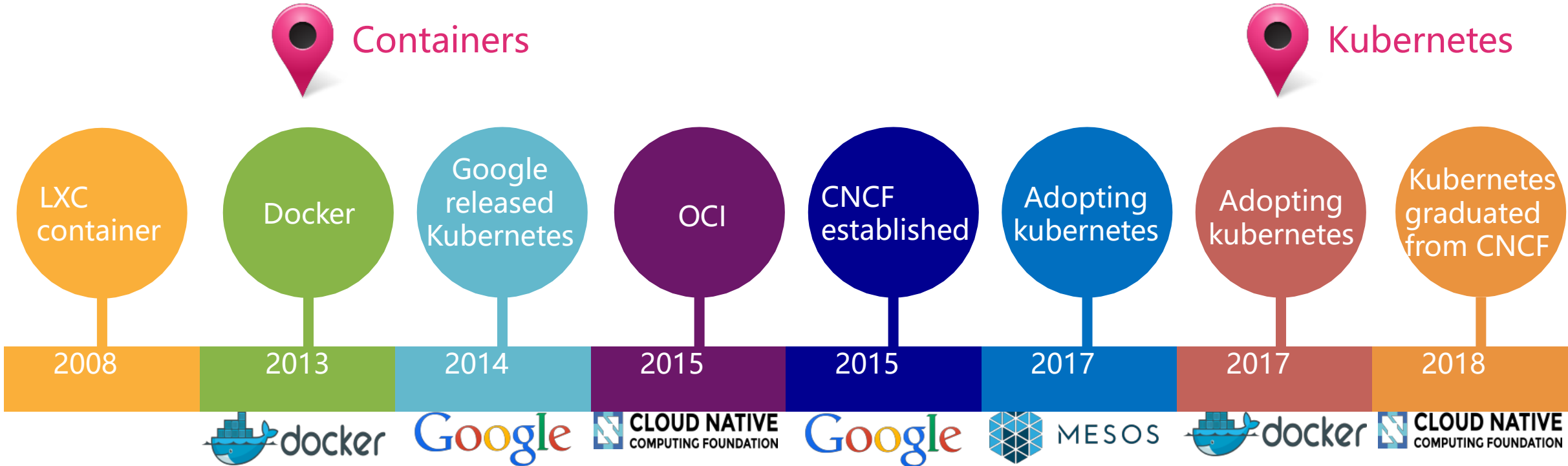
## Definitions by Cloud Native Computing Foundation (CNCF) :



- Cloud native 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.

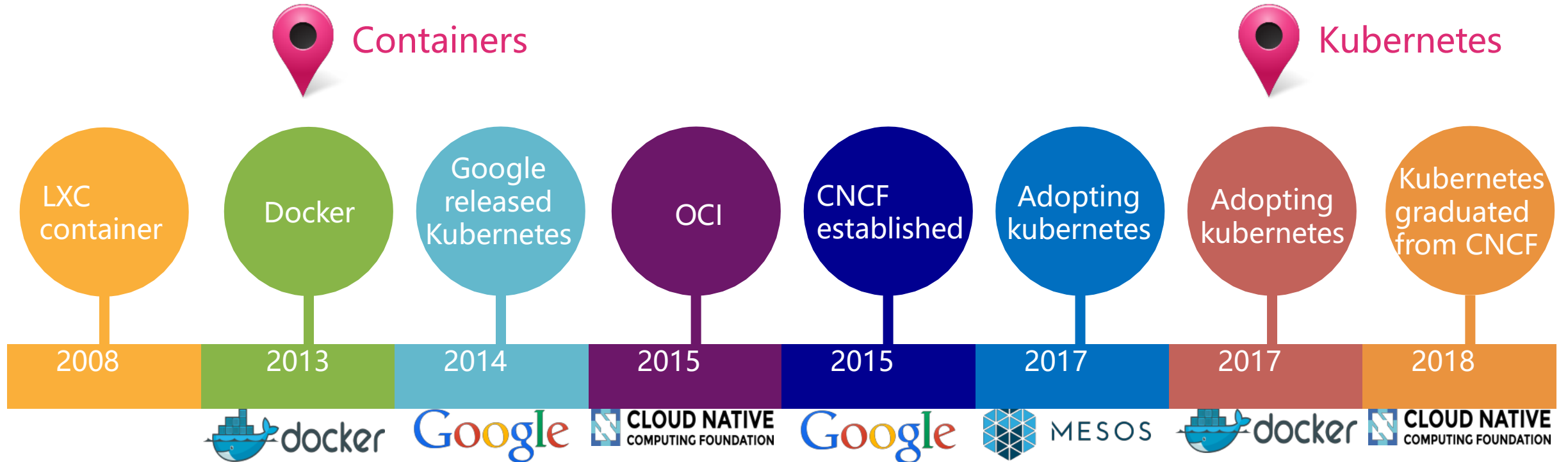


# Rise of containers and Kubernetes (K8s)



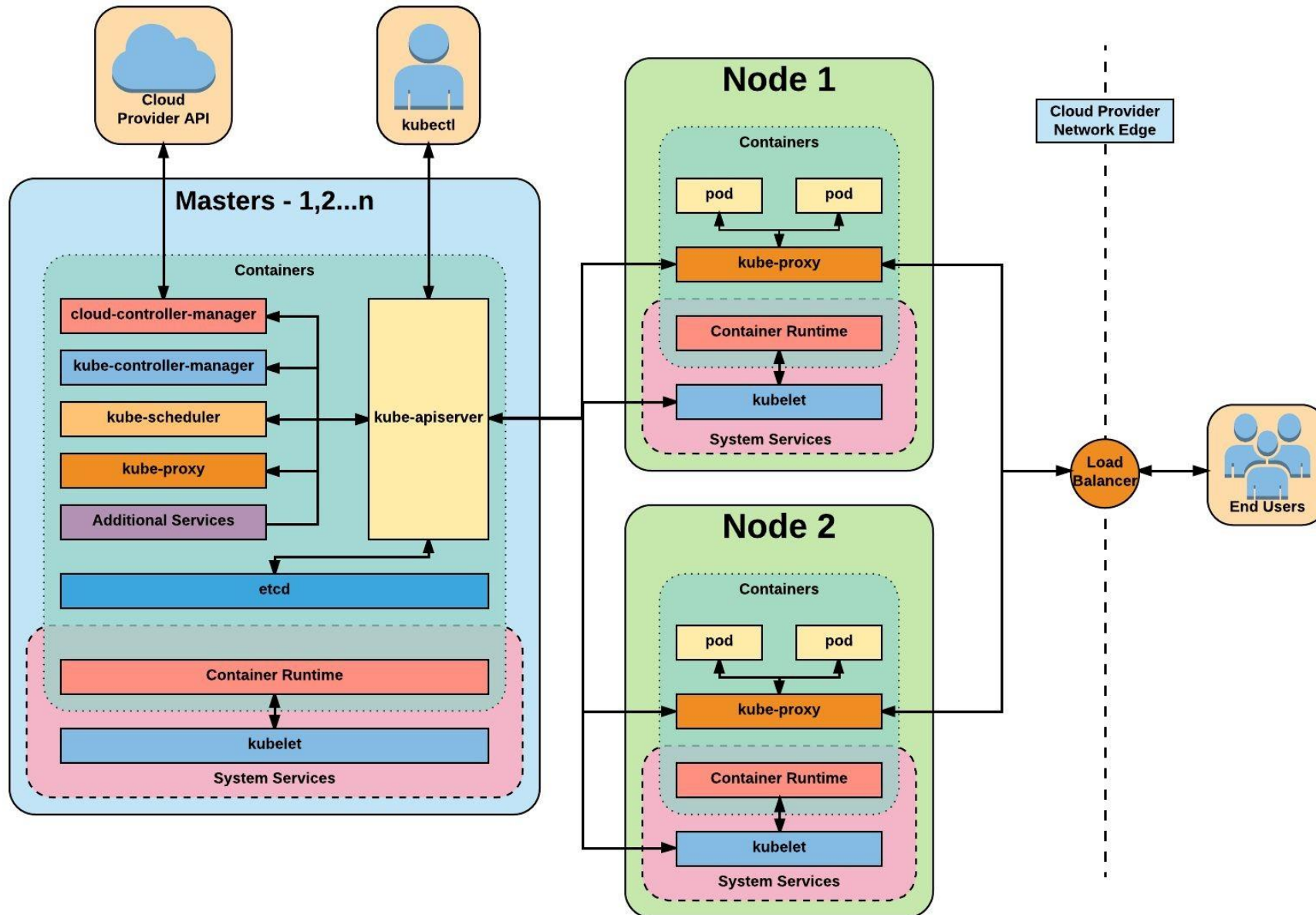
**Kubernetes** or **K8s** is a project spun out of Google as an open source next-gen container scheduler

# Rise of containers and Kubernetes (K8s)

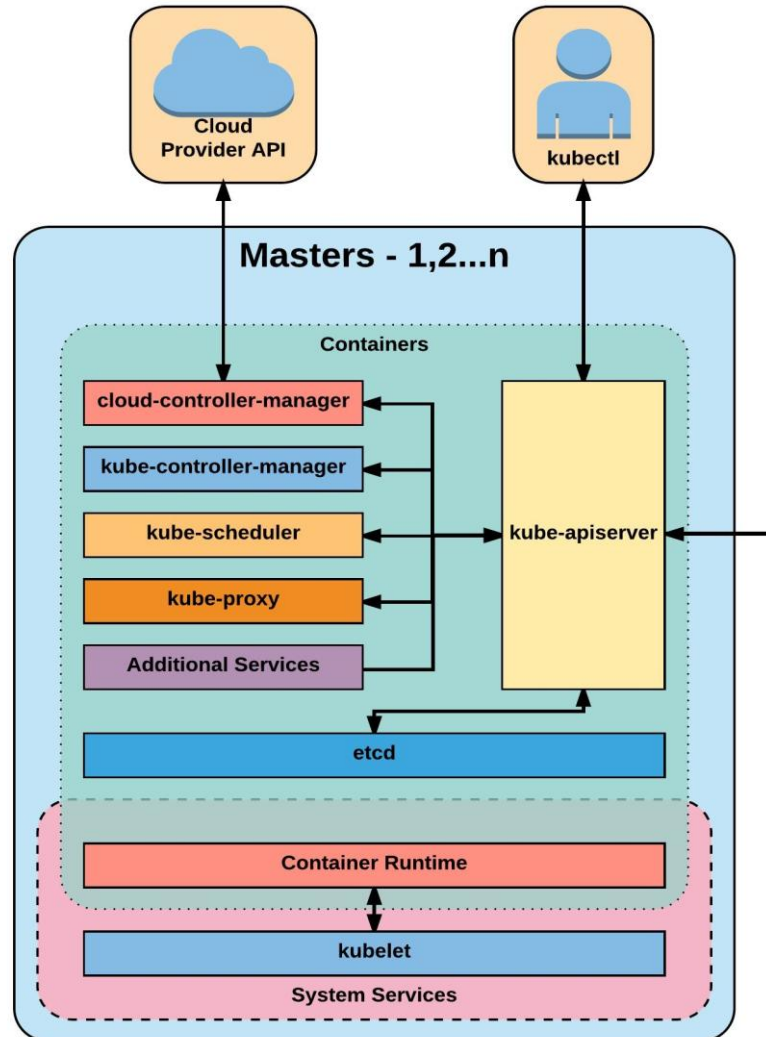


- K8s is an orchestration tool for managing distributed services or containerized applications across a distributed cluster of nodes.
- K8s follows a **client-server architecture with a master and worker nodes**. Core concepts in Kubernetes include pods, services (logical pods with a stable IP address) 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

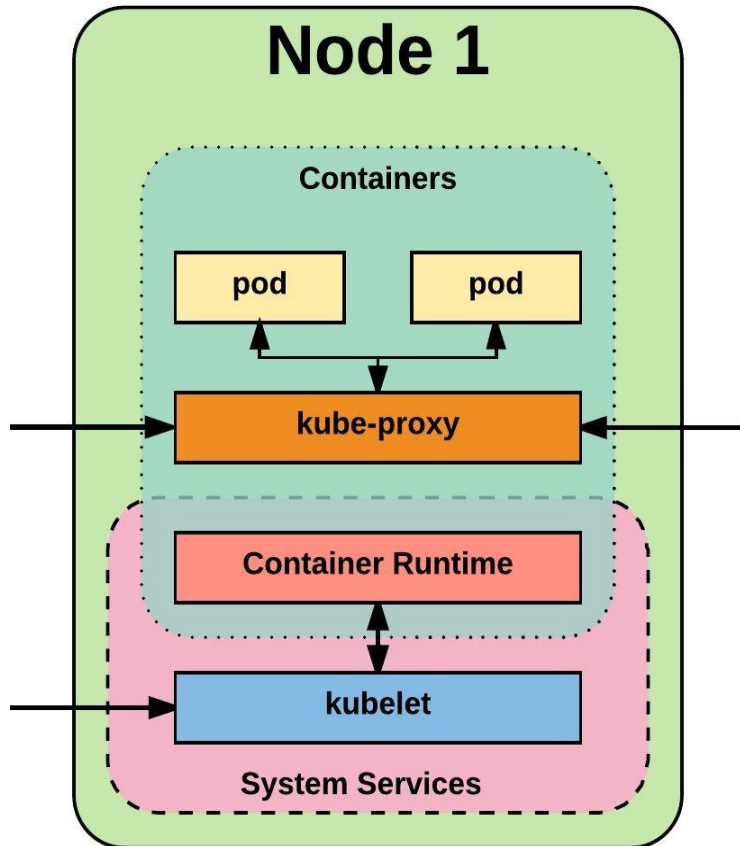


# Master components



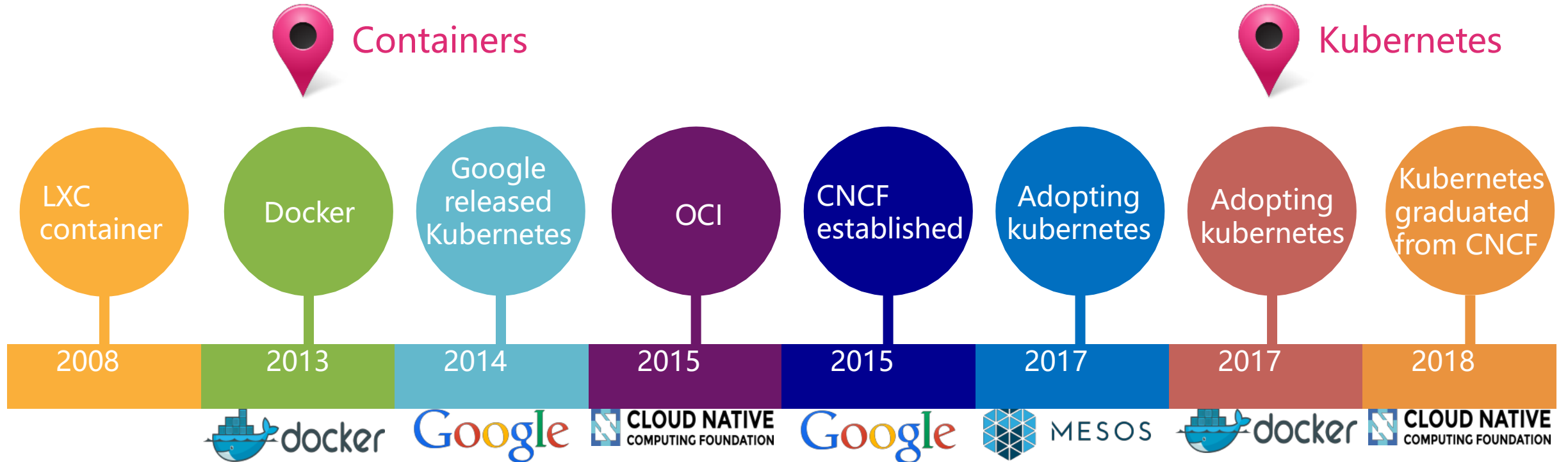
- **Kube-apiserver:** provides REST interface into the K8s control plane and datastore.
- **Etcd:** the cluster datastore; providing a strong, consistent and highly available key-value store used for persisting cluster state
- **Kube-controller-manager:** manages all core component control loops; monitors 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.

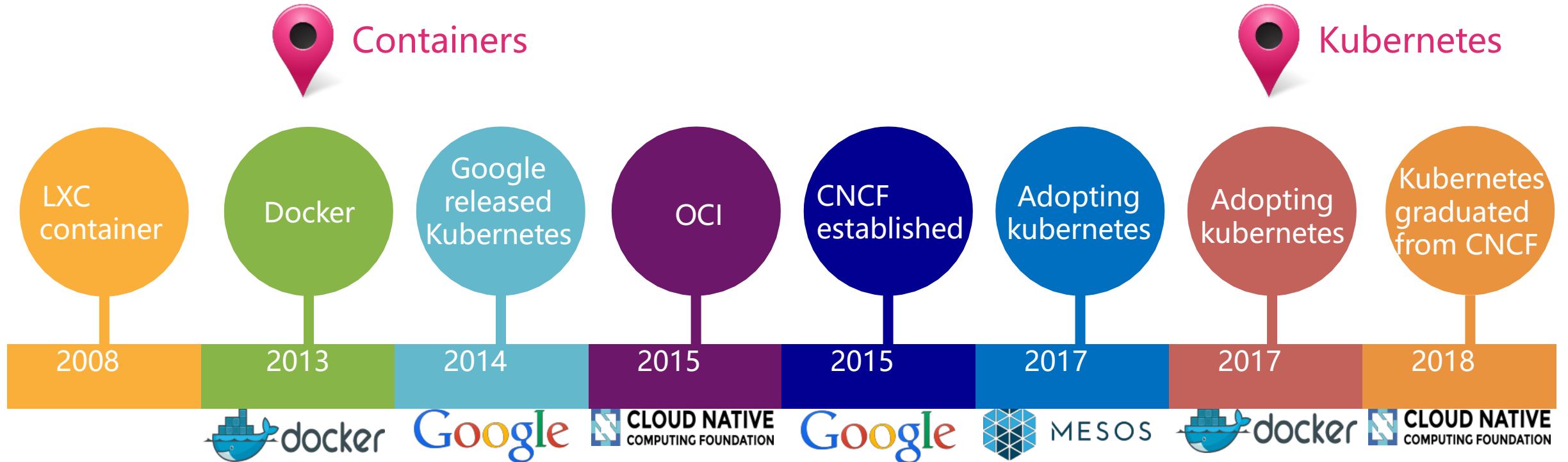
# Rise of containers and Kubernetes (K8s)



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)



## Cloud evolution 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 (what can be the problems?)**
- Current practice and opportunities of AI on cloud



Chef  
(LLM)



Restaurant  
(serving systems)



Disney world  
(cloud systems)

# Placement and load balancing (PLB)

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

- **Placement**

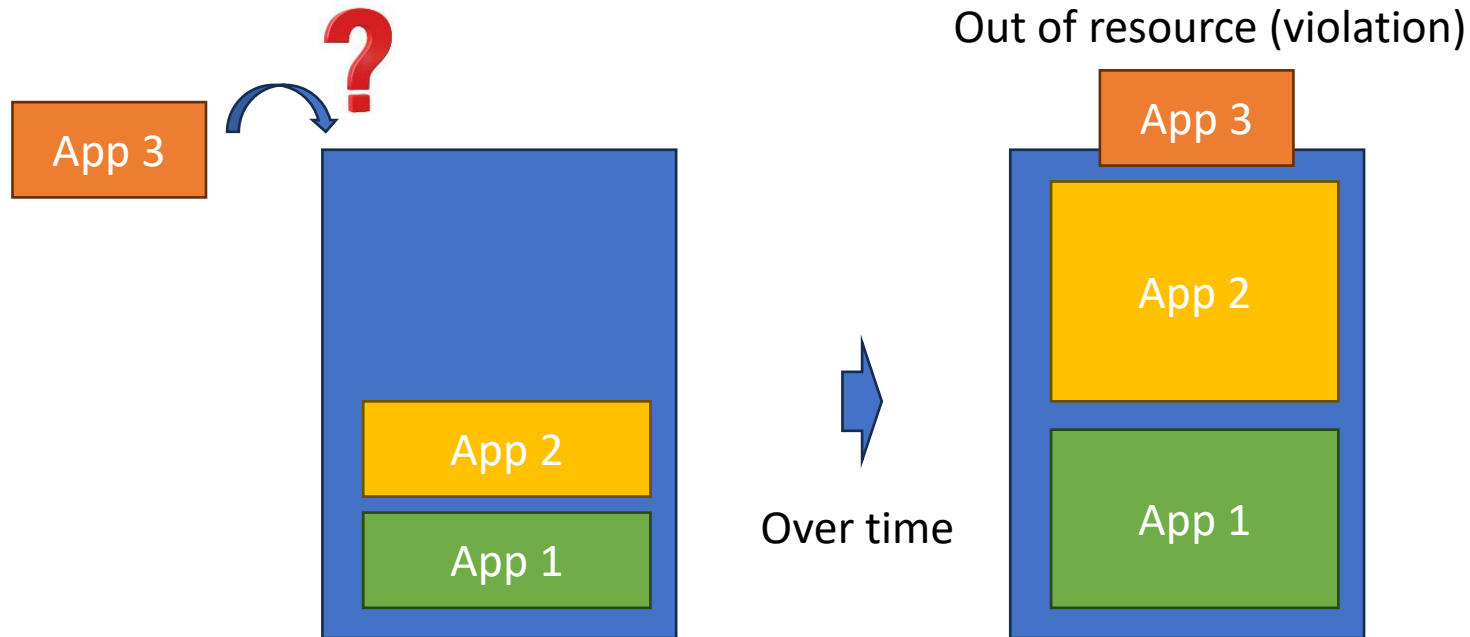
The overall goal is to reduce violation to users' Service Level Agreements (SLAs), given that resource usages are dynamic.

What can be the placement strategies?

- **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. [TETRIS, SIGCOMM 2014]
  - **More complex policies:** constrains, leave-one-out, leave-two-out and so on

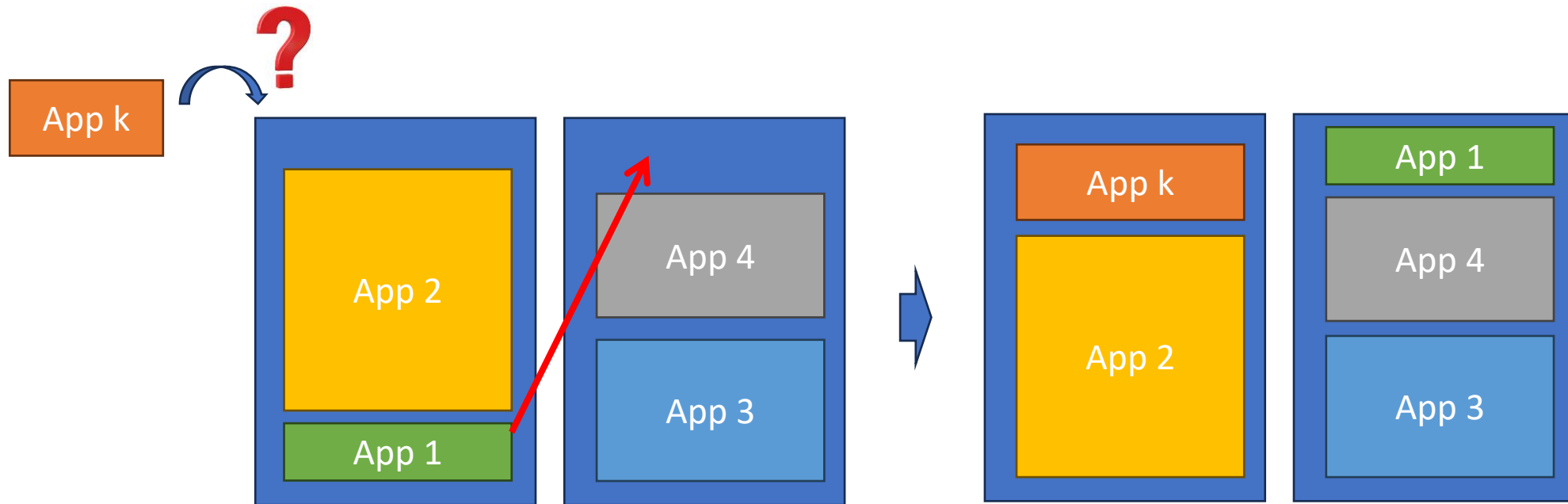
# Placement and load balancing (PLB)

- **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

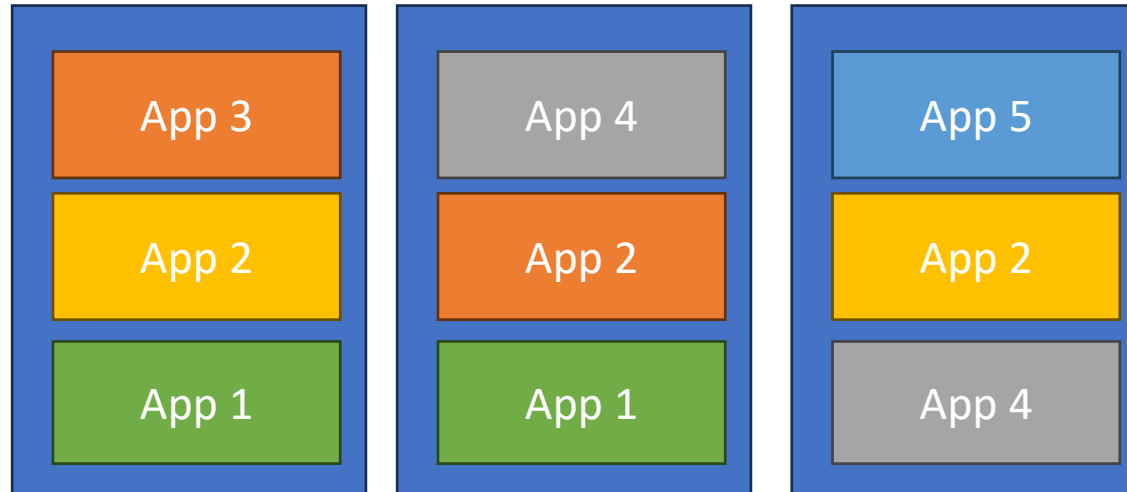


**But migration is not free, often very expensive: stateful VMs, DBs, etc.**

- Better migration mechanisms: cache compression, disaggregated memory etc.
- Resource usage prediction: placement & balancing based on predicted usages

# Failovers

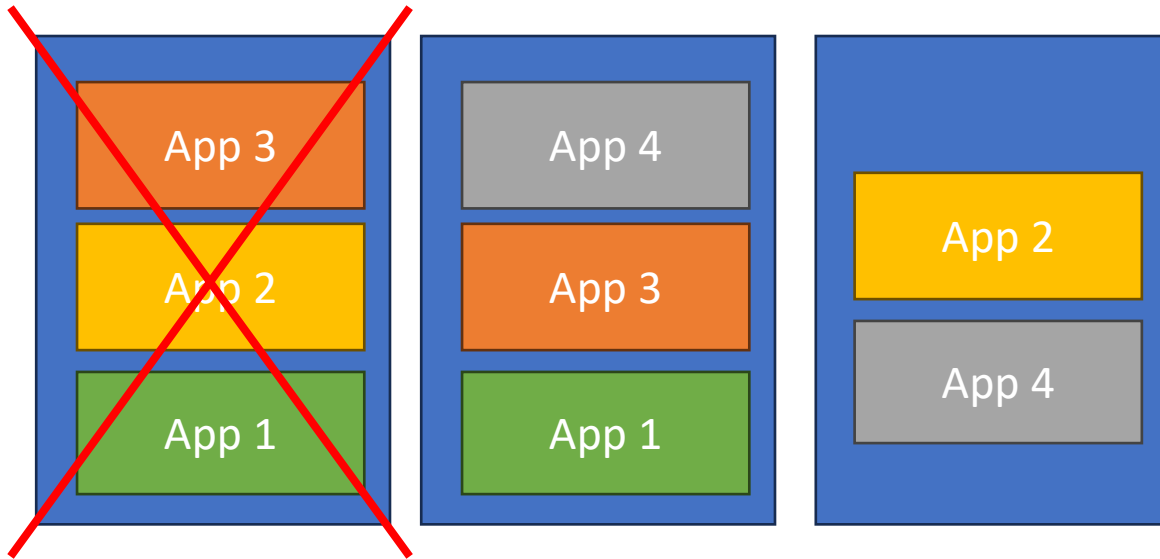
- Failovers ensure a robust & highly available service
  - Duplication of independent resources to avoid simultaneous failure



- 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



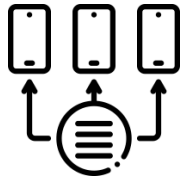
- Route requests to duplica
- Recover main from duplica
  - Reinstall using logs
  - Live migration

Failover due to:

- Unexpected: software/hardware failure
- Scheduled: Software update

# Serverless computing

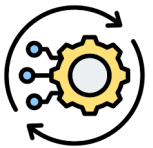
- Some “rewrap” of ideas, but many cloud-native techniques are the same underneath



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



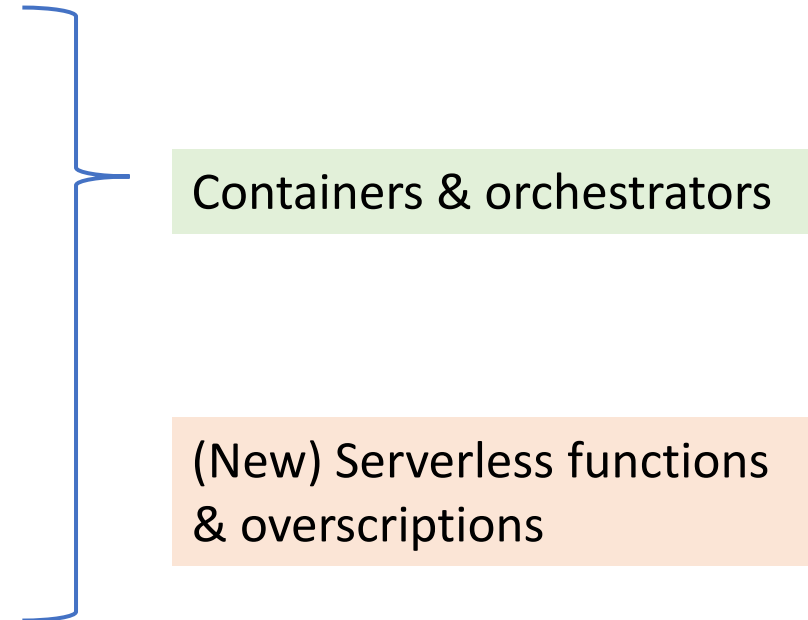
**Scales with usage.** System expands and shrinks automatically with actual usage.



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



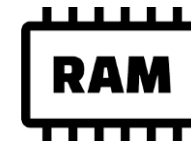
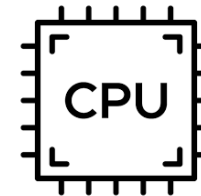
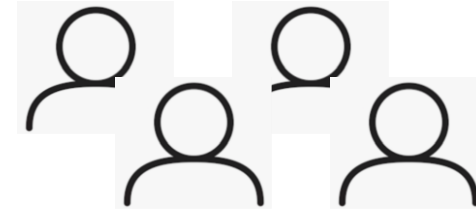
**No pay for idle.** Billing model – user pays only for actual usage; financial risks at the operator.



# Resource oversubscription



100 seats, sell 105 tickets



- (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

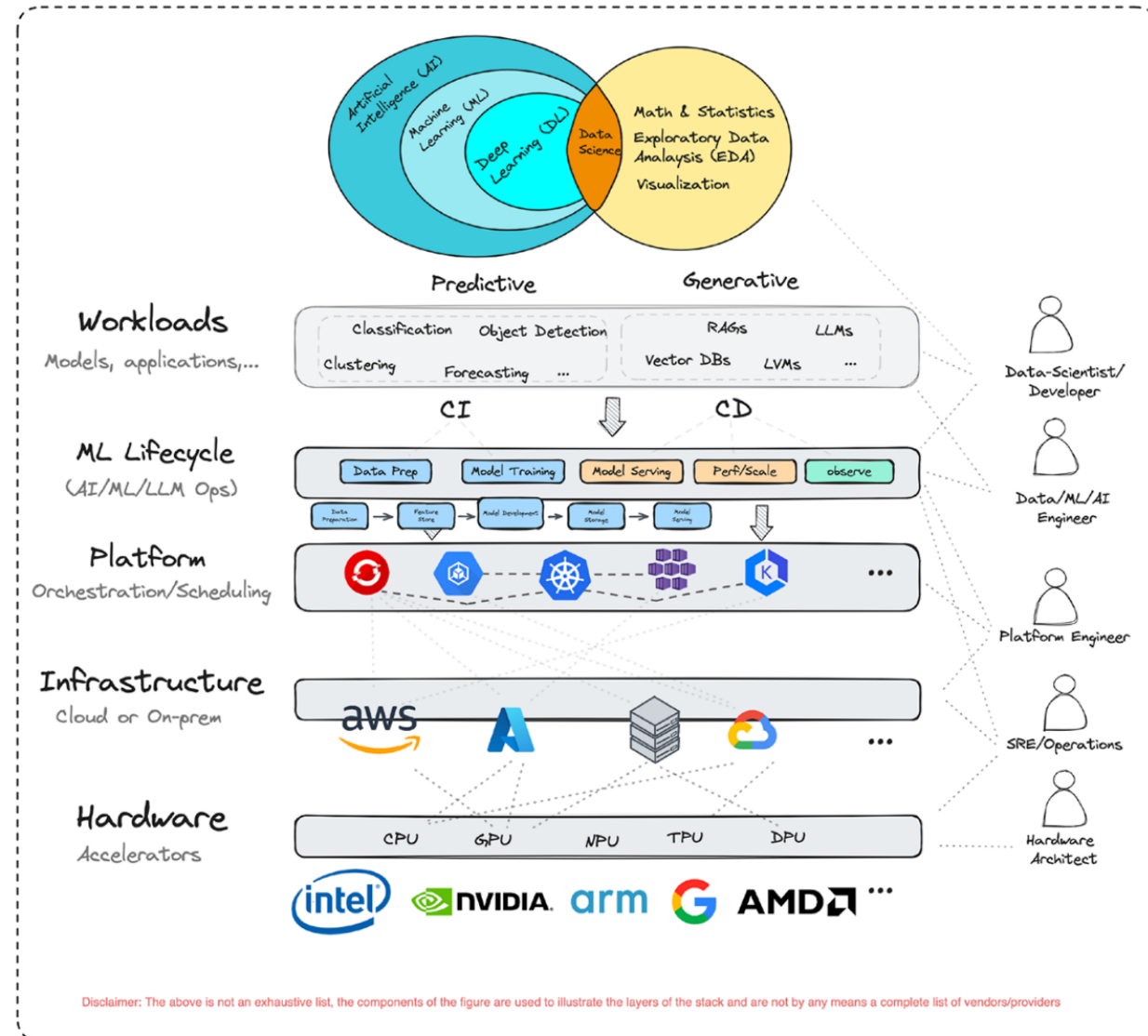
# Outline

- 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 AI

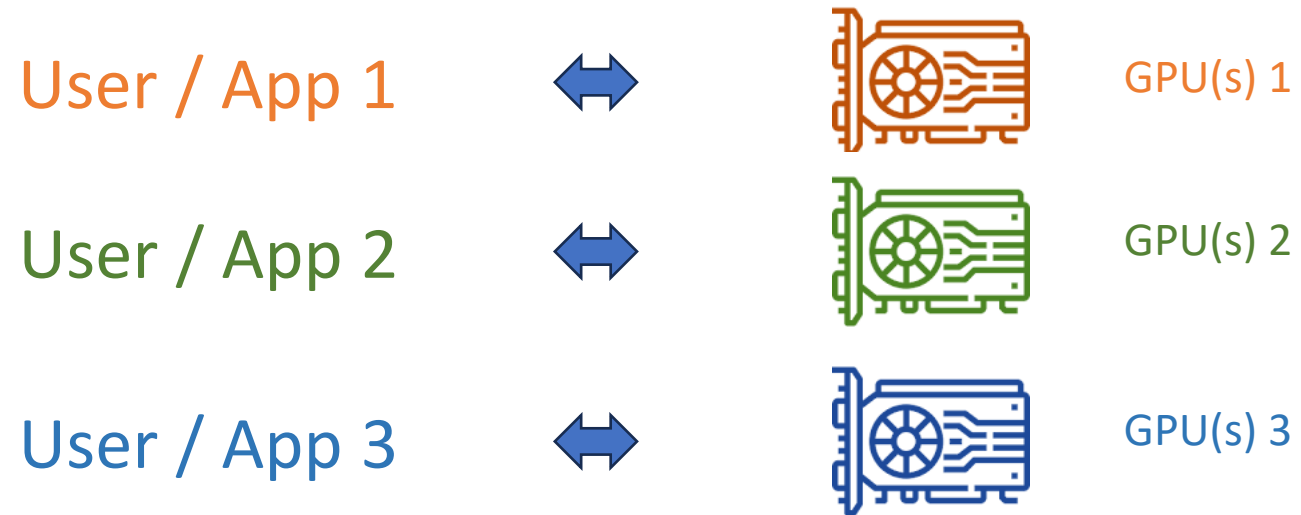
- Cloud native technologies are open-box solutions for many AI use cases
- **Key idea:** containerizing your models



# 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

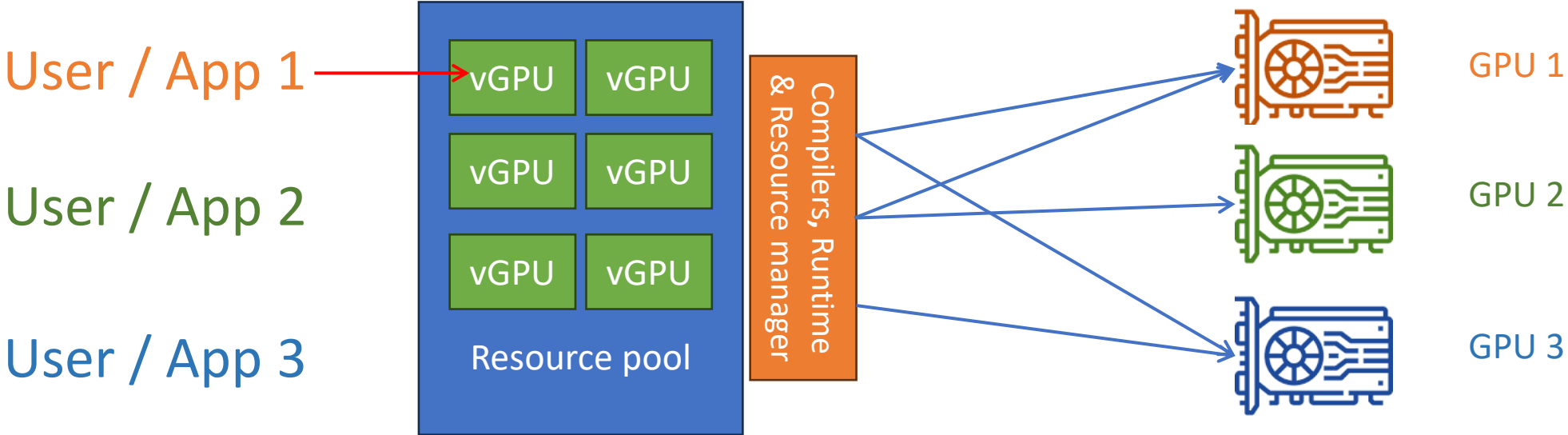
# Resource oversubscription for GPUs?



*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



“Virtulizing” GPUs into thread and memory blocks,  
But, big engineering challenges

# Looking forward

- New opportunities in cloud AI
  - **New cloud** with heterogeneous, ephemeral infra spot instances, intermittent/green power
  - **New services:** RL, agentic workflows, etc.
  - **New hardware & architecture:** RISC-V AI chips, AISC chips
  - **New applications:** AI-for-X

#GPUs:	ANY	0X	1X	2X	4X	8X	8X+	On-Demand	Any GPU	Planet Earth	Auto Sort
m:12140	host:73118	Spain, ES	ROME2D32GM	↑2758 Mbps ↓5998 Mbps	100 ports	verified Max Duration 9 days	\$0.468/hr				
	1x RTX 4090	PCIe 4.0,8x	12.7 GB/s	AMD EPYC 7642 ...	nvme	73.7 DLPerf	Reliability 99.68%	RENT			
<small>vast.ai Max CUDA: 12.2 24 GB 3479.4 GB/s</small>											
m:12959	host:57805	North Carolina, US	WRX80 Creator R2	↑5451 Mbps ↓3953 Mbps	75 ports	verified Max Duration 1 mon, 29d	\$0.503/hr				
	1x RTX 4090	PCIe 4.0,16x	24.4 GB/s	AMD Ryzen Thre...	nvme	75.7 DLPerf	Reliability 99.59%	RENT			
<small>vast.ai Max CUDA: 12.2 24 GB 3568.8 GB/s</small>											
m:11541	host:73118	Spain, ES	ROME2D32GM	↑2758 Mbps ↓5998 Mbps	100 ports	verified Max Duration 9 days	\$0.478/hr				
	1x RTX 4090	PCIe 4.0,8x	12.7 GB/s	AMD EPYC 7642 ...	nvme	73.7 DLPerf	Reliability 99.68%	RENT			
<small>vast.ai Max CUDA: 12.0 24 GB 3223.7 GB/s</small>											
m:9599	host:25384	Japan, JP	WRX80 Creator R2	↑5451 Mbps ↓3953 Mbps	75 ports	verified Max Duration 2 mon.	\$0.556/hr				
	1x RTX 4090	PCIe 4.0,16x	24.4 GB/s	AMD Ryzen Thre...	nvme	75.7 DLPerf	Reliability 99.90%	RENT			
<small>vast.ai Max CUDA: 12.2 24 GB 3576.5 GB/s</small>											





# 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
- Application systems & AI agents
- LLM safety
- ML Operations
- Data & Cloud systems for AI
- Compute graph optimization
- Serving multi-modal models
- Serving mixture-of-experts models
- Many more..

Covered

Not covered

# Logistics for next week

- LLM benchmark – Leaderboard: submit your result before next Wednesday
- In-person presentation
  - 1 Slide briefing + Demo (project 2,3), present to other students
  - TA & I will grade also
- Peer review & voting
  - Review 3 other groups' projects
  - Submit a short description for how another project can help in yours
- Reflection
  - In which way your solution can be improved by others' ideas.
- Final report
  - Research style: background, motivation, solution, evaluation, analysis.
  - Review + reflect + demo (project 2,3)