

Microservices and Serverless Computing

Corso di Sistemi Distribuiti e Cloud Computing A.A. 2023/24

Valeria Cardellini

Laurea Magistrale in Ingegneria Informatica

Microservices

- A “new” emerging architectural style for distributed applications that structures an application as a **collection of loosely coupled services**
- Not so new: deriving from **SOA** and **Web services**
 - But with some significant differences
- Address how to build, manage, and evolve architectures out of **small, self-contained** units
 - *Modularization*: decompose app into a set of **independently deployable services**, that are **loosely coupled** and **cooperating** and can be **rapidly deployed and scaled**
 - Services equipped with **memory persistence tools** (e.g., relational databases and NoSQL data stores)

The ancestors: Service Oriented Architecture

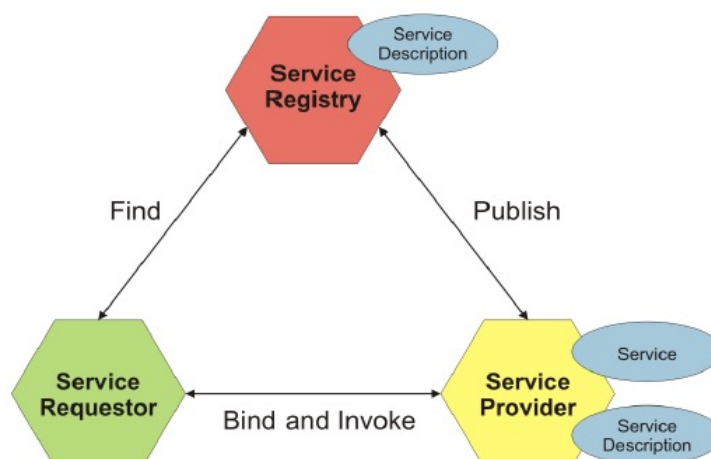
- **Service Oriented Architecture (SOA)**: architectural paradigm for designing loosely coupled distributed sw systems
- Definition (by OASIS docs.oasis-open.org/soa-rm/v1.0/soa-rm.pdf)
SOA is a paradigm for organizing and utilizing distributed capabilities that may be under the control of different ownership domains. It provides a uniform means to **offer**, **discover**, **interact with** and **use** capabilities to produce desired effects consistent with **measurable** preconditions and expectations
- Properties of SOA (by W3C www.w3.org/TR/ws-arch)
 - Logical view
 - Message orientation and description orientation
 - Service granularity, network orientation
 - Platform neutral

Valeria Cardellini – SDCC 2023/24

2

Service Oriented Architecture

- 3 interacting entities
 1. **Service requestor** or **consumer**: requests service execution
 2. **Service provider**: provides service and makes it available
 3. **Service registry**: offers publication and search tools to discover the services offered by providers



Valeria Cardellini – SDCC 2023/24

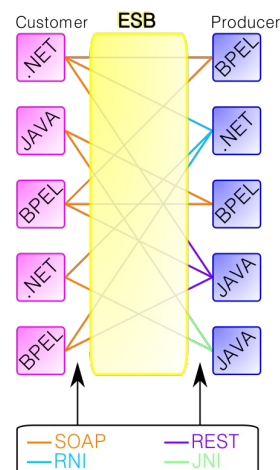
3

Web services

- Web services: implementation of SOA
- Definition (by W3C www.w3.org/TR/ws-arch)
 - Web service: software system designed to support interoperable machine-to-machine (M2M) interaction over a network
 - Web service interface described in a machine-processable format
 - Other systems interact with web service in a manner prescribed by its description using SOA messages, typically conveyed using HTTP

Web services

- More than 60 standards and specifications, most used:
 - To describe: **WSDL** (Web Service Description Language)
 - To communicate: **SOAP** (Simple Object Access Protocol)
 - To register: **UDDI** (Universal Description, Discovery and Integration)
 - To define business processes: **BPEL** (Business Process Execution Language), **BPMN** (Business Process Model and Notation)
 - To define SLA: **WSLA**
- A variety of technologies
 - Including **ESB** (Enterprise Service Bus), integration platform that provides fundamental interaction and communication services for complex applications



SOA vs. microservices

- Heavyweight vs. lightweight technologies
 - SOA tends to rely strongly on heavyweight middleware (e.g., ESB), while **microservices** rely on **lightweight** technologies
- Protocol families
 - SOA is often associated with web services protocols
 - SOAP, WSDL, and WS-* family of standards
 - **Microservices** typically rely on **REST** and **HTTP**
- Views
 - SOA mostly viewed as integration solution
 - **Microservices** are typically applied to build **individual software applications**

Microservices and containers

- Microservices as ideal complementation of **container-based virtualization**
 - “**Microservice instance per container**”: package each microservice as container image and deploy each microservice instance as container
 - Manage each container at runtime scaling and/or migrating it
- Pros and cons:
 - ✓ Scaling out/in microservice instance by changing number of container replicas
 - ✓ Scaling up/down microservice instance assigning more/less resources
 - ✓ Isolate microservice instance
 - ✓ Apply resource limits on microservice instance
 - ✓ Build and start rapidly
 - ✗ Require **container orchestration** to manage multi-container app

Microservices: benefits

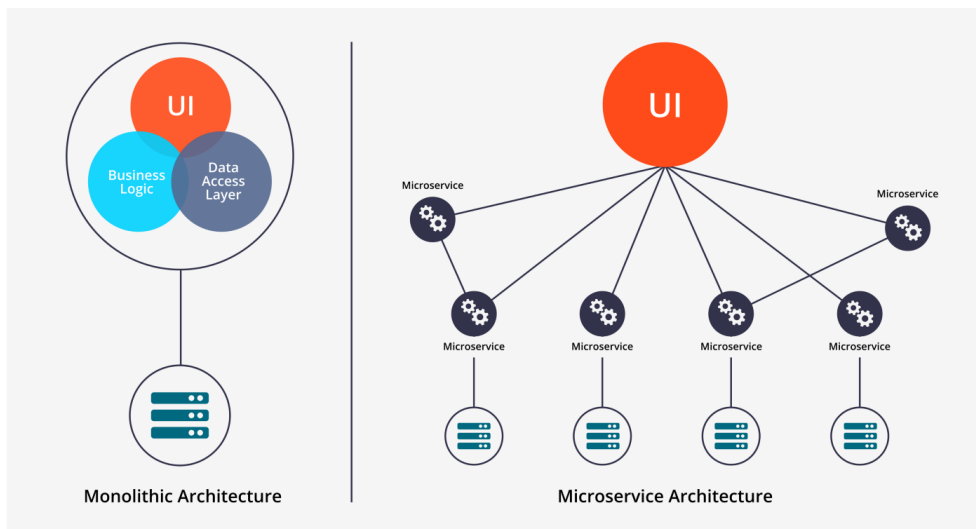
- Increased software agility
 - Microservice: **independent unit** of development, deployment, operations, versioning, and scaling
 - All interactions with a microservice happen via its API, which encapsulates its implementation details
 - Exploit **container-based virtualization**
- Improved scalability and fault isolation
- Increased reusability across different areas of business
- Improved data security
- Faster development and delivery
- Greater autonomy
 - Teams can be be more autonomous

Microservices: concerns

- Increased network traffic
 - Inter-service calls over a network cost more in terms of network latency
- Higher complexity
 - Increased operational complexity (e.g., deployment)
 - Global testing and debugging is more complicated

How to decompose the application

- How to decompose a complex app (implemented as *monolithic* application) into microservices?
- Mostly an art, no winner strategy but rather a number of strategies microservices.io/patterns



How to decompose the application

- Main decomposition patterns
 - Let's consider an e-commerce application that takes orders from customers, verifies inventory and available credit, and ships them
 1. Decompose by **business capability** and define services corresponding to business capabilities
 - Business capability: something that a business does in order to generate value
 - E.g., *Order Management* is responsible for orders
 2. Decompose by **domain-driven design (DDD) subdomain**
 - A domain consists of multiple subdomains; each subdomain corresponds to a different part of the business
 - E.g., *Order Management*, *Inventory*, *Product Catalogue*, *Delivery*

How to decompose the application

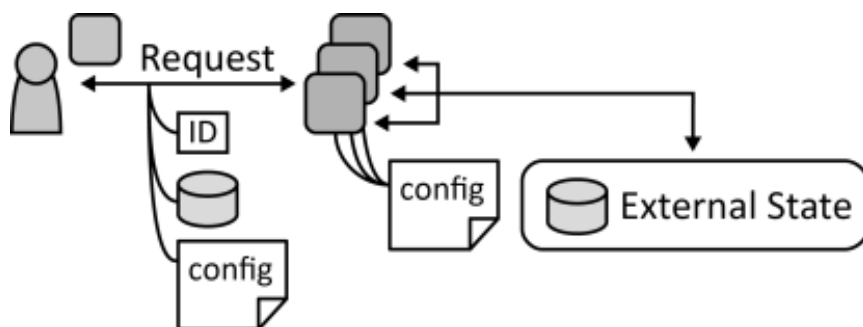
- Other decomposition patterns
 3. Decompose by **use case** and define services that are responsible for particular actions
 - E.g., *Shipping Service* is responsible for shipping complete orders
 4. Decompose by **nouns or resources** and define a service that is responsible for all operations on entities/resources of a given type
 - E.g., *Account Service* is responsible for managing user accounts

Microservices and scalability

- How to achieve microservice scalability?
 - Use **multiple instances of same microservice** and **load balance** requests across multiple instances
- How to improve microservice scalability?
 - State is complex to manage and scale
 - Prefer **stateless services**: scale better and faster than stateful services

Stateless service

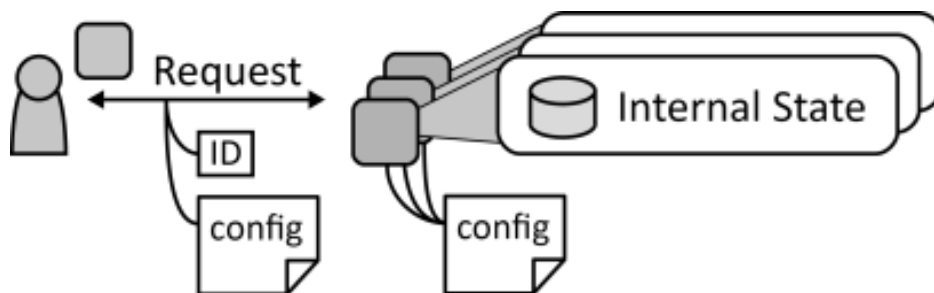
- **Stateless service**: state is handled external of service to ease its scaling out and to make application more tolerant to service failures



Cloud Computing Patterns,
www.cloudcomputingpatterns.org/stateless_component/

Stateful service

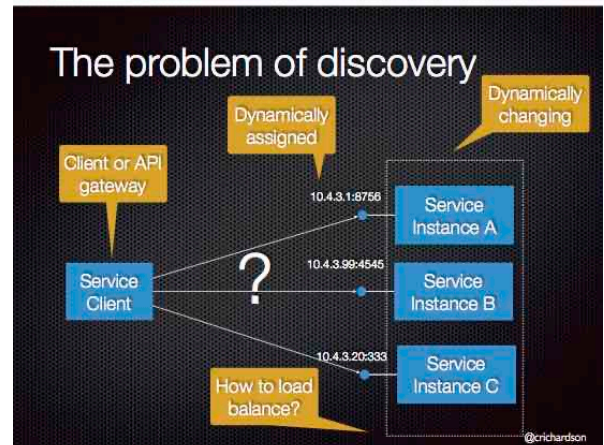
- **Stateful service**: multiple instances of scaled-out service need to synchronize their internal state to provide a unified behavior
- Issue: how can a scaled-out stateful service maintain a synchronized internal state?



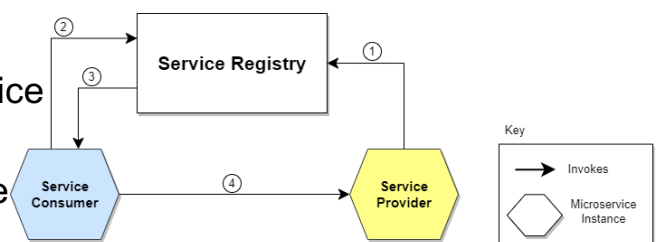
Cloud Computing Patterns,
www.cloudcomputingpatterns.org/stateful_component/

Service discovery

- We also need **service discovery**
 - Microservice instances have dynamically assigned network locations (IP address and port) and their set changes dynamically because of auto-scaling, failures, and upgrades



- Service discovery provides
 - a mechanism for a microservice instance to register
 - a way to find the service once it has registered

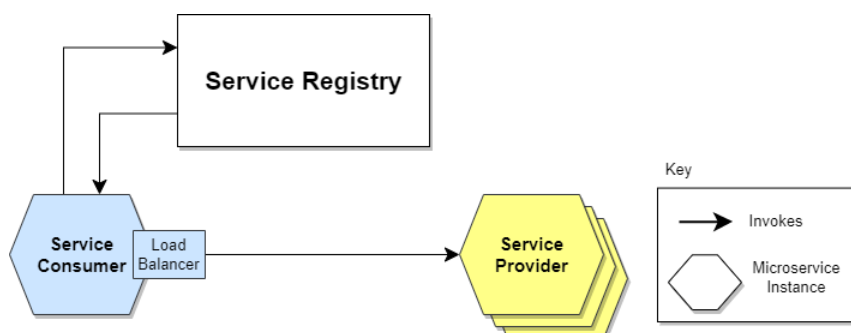


Patterns for service discovery

1. Client-side service discovery

- Client of service is responsible for determining network locations of available service instances and load balancing requests among them
- Client queries Service Register, then it uses a load-balancing algorithm to choose one of the available service instances and performs a request

microservices.io/patterns/server-side-discovery.html

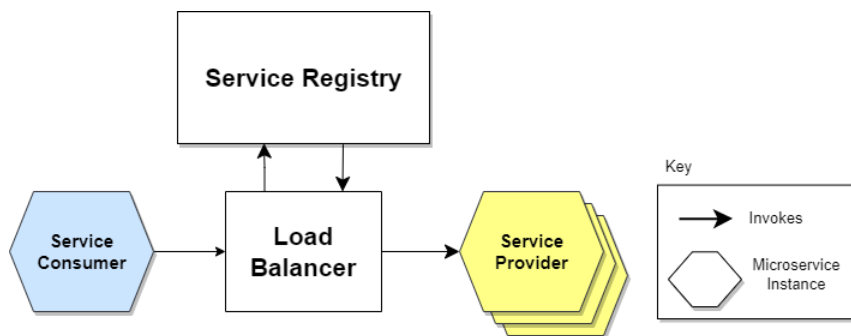


Patterns for service discovery

2. Server-side service discovery

- Client uses an intermediary that acts as *Load Balancer* and runs at a well known location
- Client makes a request to a service via a load balancer. The load balancer queries the Service Registry and routes each request to an available service instance

microservices.io/patterns/server-side-discovery.html



Integration of microservices

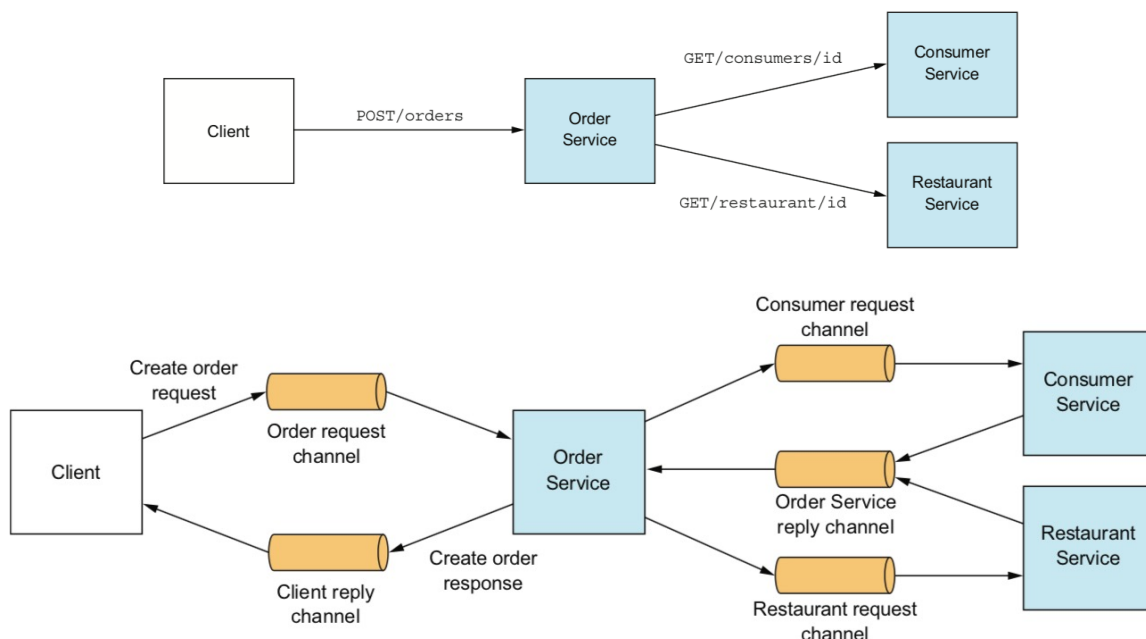
- Let's consider two issues related to integration of microservices
 - Synchronous vs. asynchronous communication
 - Orchestration vs. choreography

Synchronous vs. asynchronous

- Should communication be synchronous or asynchronous?
 - Synchronous: request/response style of communication
 - Asynchronous: event-driven style of communication
- Synchronous communication
 - Synchronous request/response-based communication mechanisms, such as HTTP-based REST or gRPC
- Asynchronous communication
 - Asynchronous, message-based communication mechanisms such as pub-sub systems, message queues and related protocols
 - Interaction style can be one-to-one or one-to-many
- Synchronous communication may reduce availability

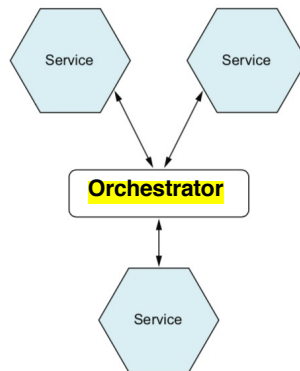
Synchronous vs. asynchronous

- Example of synchronous communication vs. asynchronous communication



Orchestration and choreography

- Microservices can interact among them according to 2 patterns:
 - Orchestration
 - Choreography
- **Orchestration**: **centralized** approach
 - A single centralized process (*orchestrator*, *conductor* or *message broker*) coordinates interaction
 - Orchestrator is responsible for invoking and combining services, which can be unaware of composition

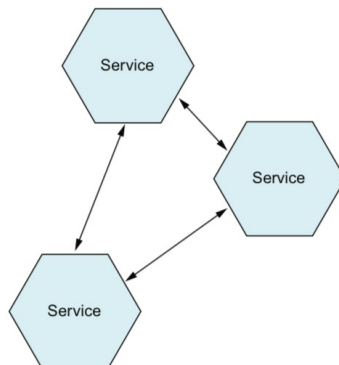


Valeria Cardellini – SDCC 2023/24

22

Orchestration and choreography

- **Choreography**: **decentralized** approach
 - A global description of participating services, which is defined by exchange of messages, rules of interaction and agreements between two or more endpoints
 - Services can exchange messages directly

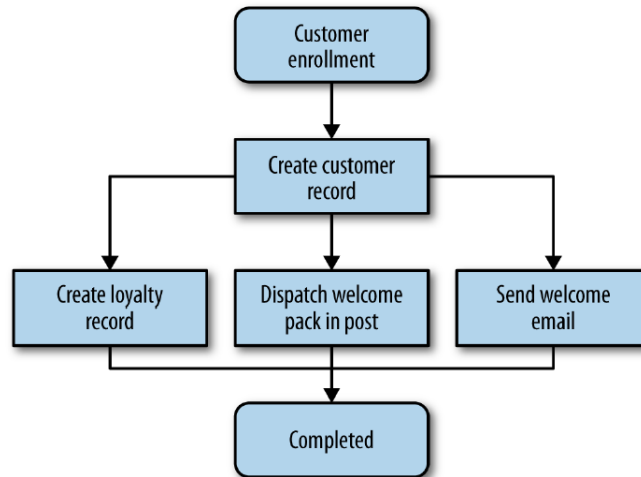


Valeria Cardellini – SDCC 2023/24

23

Example: orchestration and choreography

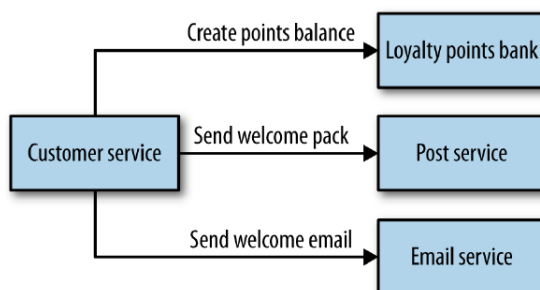
- Example: workflow for customer creation, i.e., process for creating a new customer



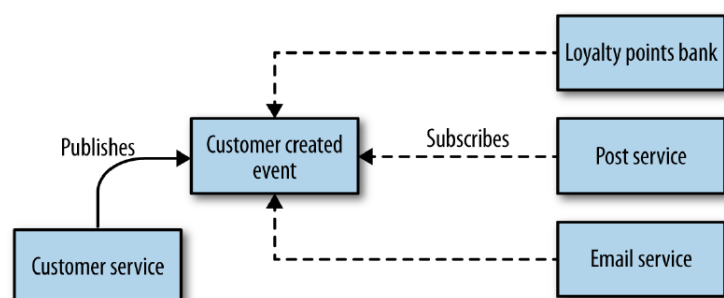
From: S. Newman, "Building Microservices", O'Really, 2015.

Example: orchestration and choreography

Orchestration



Choreography



Orchestration vs choreography

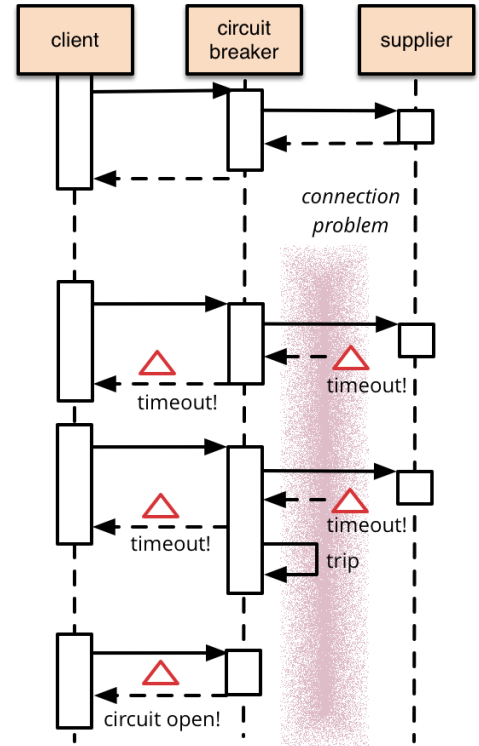
- Orchestration:
 - ✓ Simpler and more popular
 - ✗ SPoF and performance bottleneck
 - ✗ Tight coupling
 - ✗ Higher network traffic and latency
- Choreography
 - ✓ Lower coupling, less operational complexity, and increased flexibility and ease of changing
 - ✗ Services need to know about each other's locations
 - ✗ Extra work to monitor and track services
 - ✗ Implementing mechanisms such as guaranteed delivery is more challenging

Design patterns for microservice-based applications

- Let's examine some design patterns
 1. Circuit breaker
 2. Database per service
 3. Saga (and event sourcing)
 4. CQRS
 5. Log aggregation
 6. Distributed request tracing

Patterns: Circuit breaker

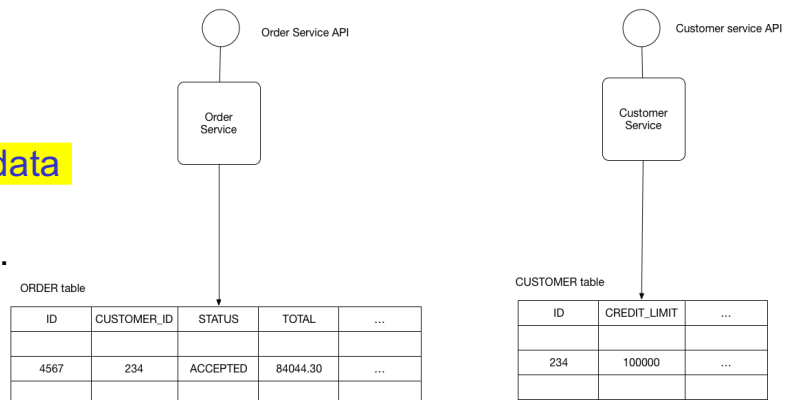
- **Problem:** How to prevent a network or service failure from cascading to other services?
- **Solution:** A service client invokes a remote service via a proxy that functions in a similar fashion to an **electrical circuit breaker**
 - When the *number of consecutive failures* crosses a threshold, the circuit breaker trips, and for the duration of a *timeout* period all attempts to invoke the remote service will fail immediately
 - After the timeout expires, the circuit breaker allows a limited number of test requests to pass through. If those requests succeed, the circuit breaker resumes normal operation. Otherwise, in case of failure the timeout period begins again



microservices.io/patterns/reliability/circuit-breaker.html

Patterns: Database per service

- **Problem:** which database architecture?
- **Solution:** **keep each microservice's persistent data private** to that service and accessible only via its API. Service transactions only involve its DB

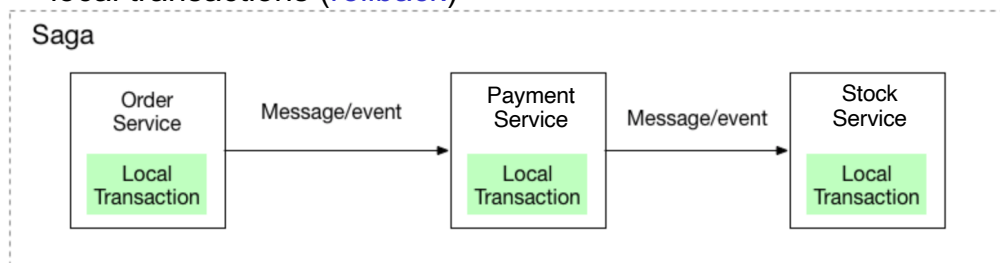


- Pros and cons
 - ✓ Helps ensure that services are loosely coupled
 - ✓ Each service can use the most convenient DB type (e.g., NoSQL)
 - ✗ More complex to implement transactions that span multiple services
 - ✗ Complexity of managing multiple DBs
 - But do not need to provision a DB server for each service
 - Options: private-tables-per-service, schema-per-service, database-server-per-service

microservices.io/patterns/data/database-per-service.html

Patterns: Saga

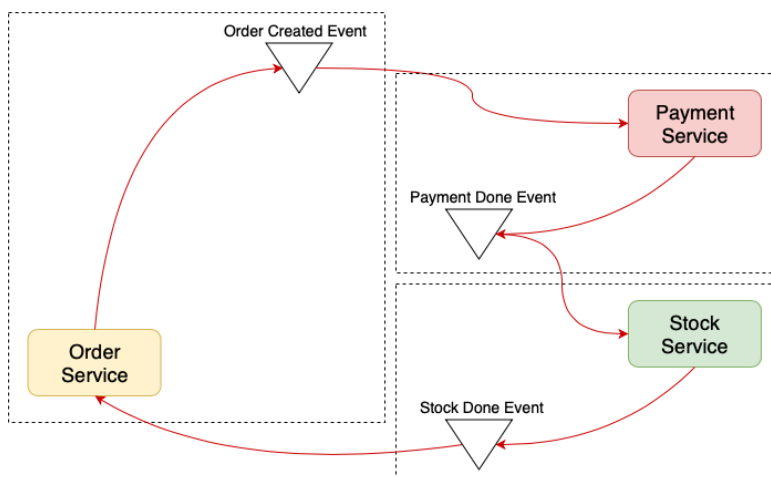
- **Problem:** each service has its own DB, however some transactions span multiple services: how to maintain data consistency across services without using distributed transactions (2PC protocol)?
- **Solution:** implement each transaction that spans multiple services as a saga
- **Saga:** **sequence of local transactions**
 - Each local transaction updates its DB and publishes a message or event to trigger the next local transaction in the saga
 - If local transaction fails, then saga executes a series of compensating transactions that undo changes made by preceding local transactions (**rollback**)



Patterns: Saga

- 2 ways to coordinate saga:
 - **Choreography:** each local transaction publishes events that trigger local transactions in other services
 - **Orchestration:** an orchestrator tells participants what local transactions to execute

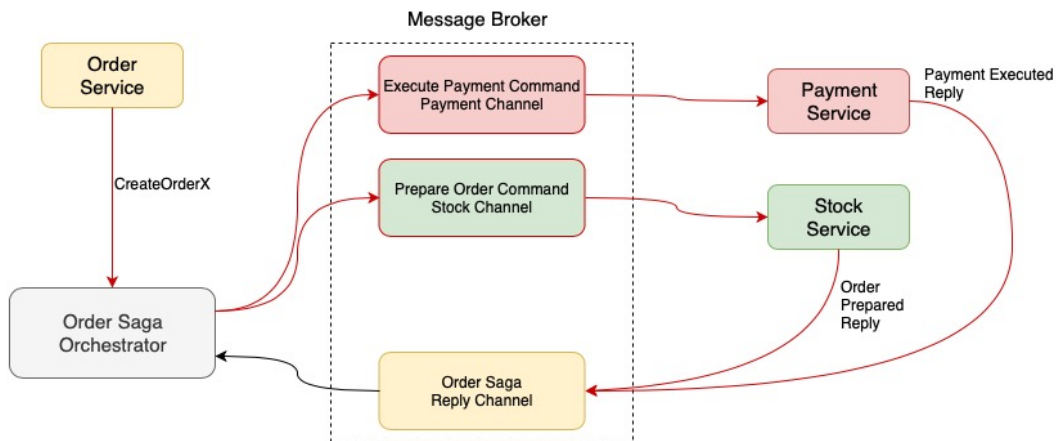
Choreography



Patterns: Saga

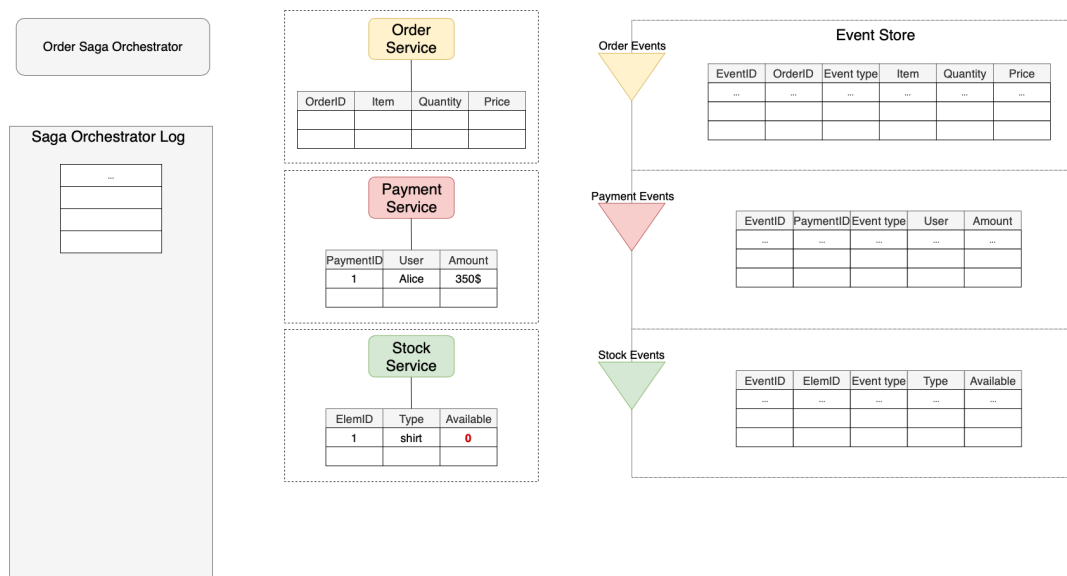
- 2 ways to coordinate saga:
 - **Choreography**: each local transaction publishes events that trigger local transactions in other services
 - **Orchestration**: orchestrator tells participants what local transactions to execute

Orchestration



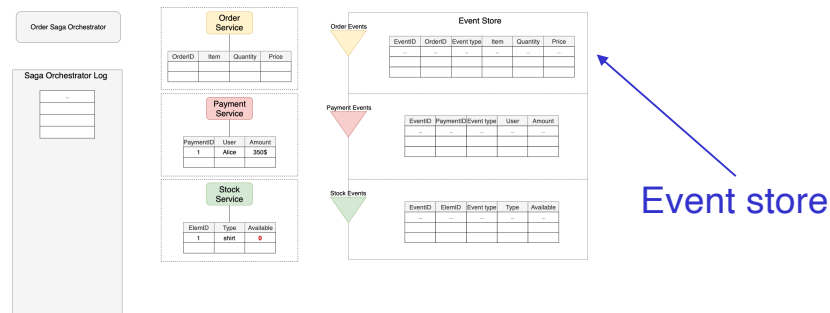
Patterns: Saga

- Example: orchestration-based saga
 - Source: MSc thesis by Andrea Cifola, see [Microservice SAGAexample.pdf](#)



Patterns: Saga

- Example: orchestration-based saga



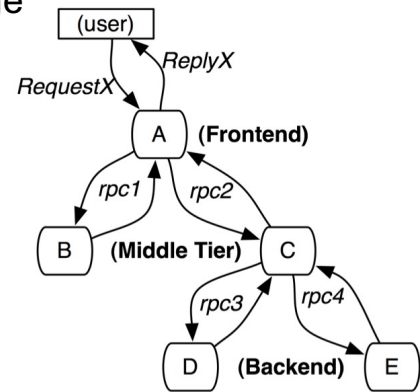
- We also use another pattern: [event sourcing](https://microservices.io/patterns/data/event-sourcing.html)
microservices.io/patterns/data/event-sourcing.html
 - **Problem:** a service that participates in a saga needs to atomically update the DB and sends messages/events in order to avoid data inconsistencies
 - **Solution:** persist a sequence of domain events that represent state changes; each event in the sequence is stored in an *append-only event store* (a DB of events)

Patterns: CQRS

- **Problem:** How to implement a query that retrieves data from multiple services in a microservice architecture? How to separate read and write load allowing you to scale each independently?
- **Solution:** define a view DB, which is a read-only replica that is designed to support that query
 - Application keeps replica updated by subscribing to [Domain events](#) published by the service that owns data
- Called Command Query Responsibility Segregation (CQRS), i.e., separate read and update operations for a data store
microservices.io/patterns/data/cqrs.html

Monitoring microservices

- Service distribution, even at large scale: difficult to monitor microservice apps and capture causal and temporal relationships among microservices
 - Aka **microservices observability** challenge
- Why do we need monitoring?
 - To debug the application
 - To analyze performance and latency, including **tail latency**
 - To analyze **service dependencies**
 - To identify **root cause** of anomalies, which requires to:
 - Construct a service dependency graph that outlines the sequence and structure of microservices that are invoked
 - Localize the root cause microservices using the graph, traces, logs, and KPIs



A request is passed through multiple microservices with different functionalities

Monitoring microservices

- Let's examine 2 patterns to monitor microservices
 1. Log aggregation
 2. Distributed request tracing

Patterns: Log aggregation

- **Problem:** How to understand application behavior and troubleshoot problems?
- **Solution:** Use a **centralized logging service** that aggregates logs from each microservice instance
 - DevOps team can search and analyze logs and configure alerts that are triggered when certain messages appear in logs
 - E.g., AWS CloudWatch
- ✗ Centralized (if physical, not only logical)
- ✗ Handling large volume of logs requires substantial infrastructure

microservices.io/patterns/observability/application-logging.html

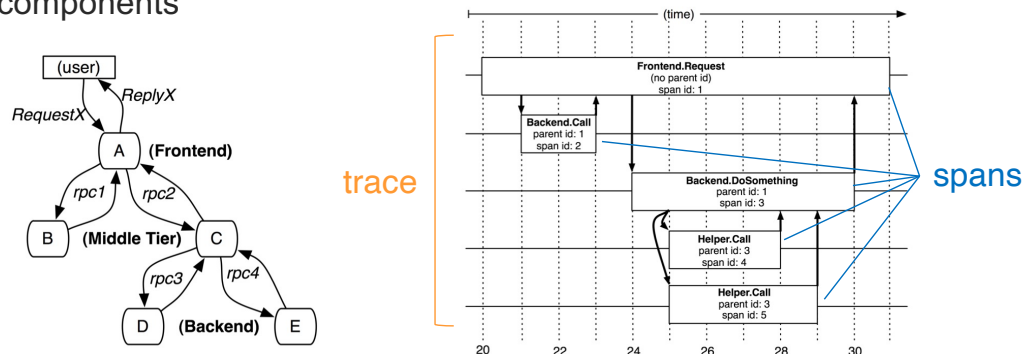
Patterns: Distributed tracing

- **Problem:** How to understand **complex** app behavior and troubleshoot problems?
- **Solution:** **Instrument microservices with code** that
 - Assigns to each user request a unique request id (aka **trace id**), that allows to track that request through the microservices it traverses
 - Passes trace id to each microservice involved in handling the user request
 - Includes trace id in log messages
 - Records **trace context** (e.g., start time, operation, duration) in a (distributed) data store
- ✗ Storing and aggregating traces can require significant infrastructure

microservices.io/patterns/observability/distributed-tracing.html

Monitoring microservices: tools

- Dapper
 - Google's production distributed systems tracing infrastructure
 - Introduces the concept of *traces* and *spans*
 - **Span**: individual unit of work (e.g., HTTP request, call to DB) in application; must have an operation name, start time, and duration
 - **Trace**: collection/list of spans connected in a parent/child relationship (can also be thought of as DAG of spans); traces specify how requests are propagated through services and other components



Barroso et al., [Dapper, a large-scale distributed systems tracing infrastructure](#), 2010

Valeria Cardellini – SDCC 2023/24

40

Monitoring microservices: tools

- Dapper
 - Traces are sampled using an **adaptive sampling rate**, why?
 - Storing everything would require too much storage and network traffic, as well as introducing too much application overhead
 - Span data is written to local log files, then pulled from there by Dapper daemons, sent over a collection infrastructure, and finally traces are stored into BigTable, with one row in a trace table dedicated to each trace id

Valeria Cardellini – SDCC 2023/24

41

Monitoring microservices: tools

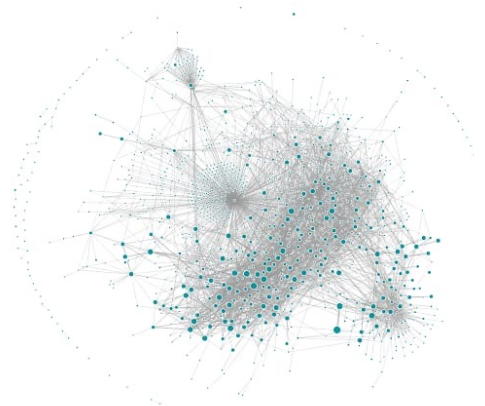
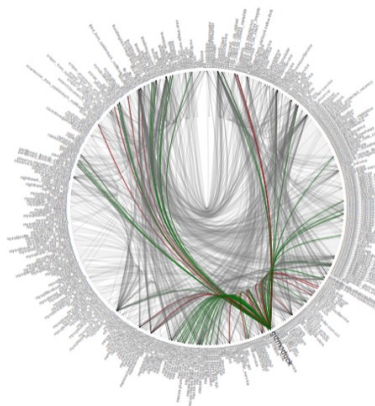
- Open-source tools for distributed tracing
 - [Jaeger](#)
 - Uses Spark/Flink for aggregate trace analysis
 - [Zipkin](#)
 - [OpenTelemetry](#)
 - Broad language support
 - Integrated with popular frameworks and libraries
- Need for standards to support interoperability between different tracing tools
 - W3C defines [Trace Context](#) specification, a standardized format for unifying tracing data

Some large-scale examples

- Netflix, Twitter, Uber: 500+ microservices



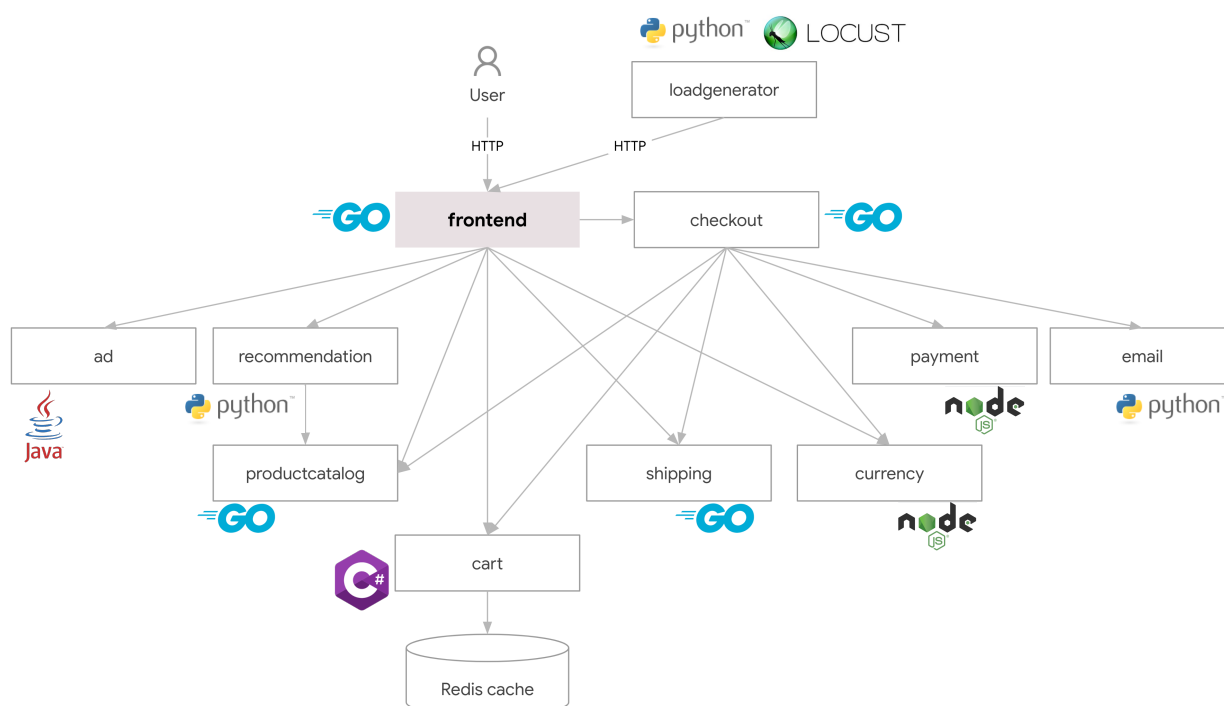
NETFLIX



Example of microservices app

- Let's examine a microservices app: Google's Online Boutique
github.com/GoogleCloudPlatform/microservices-demo
- Online store composed of 11 microservices written in different languages
 - Programming languages are silos, but in the last 15 years renaissance in programming language diversity: need for **polyglot programmers**
- How to realize a polyglot application?
 1. **REST and JSON** as message interchange format
 2. **gRPC and protocol buffers** as IDL and message interchange format: see Online Boutique

Online Boutique: architecture

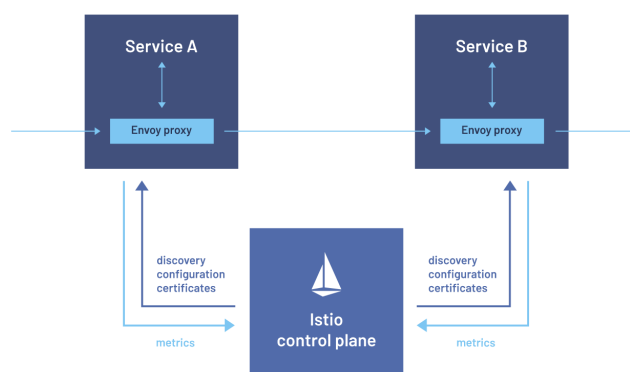


Online Boutique: features

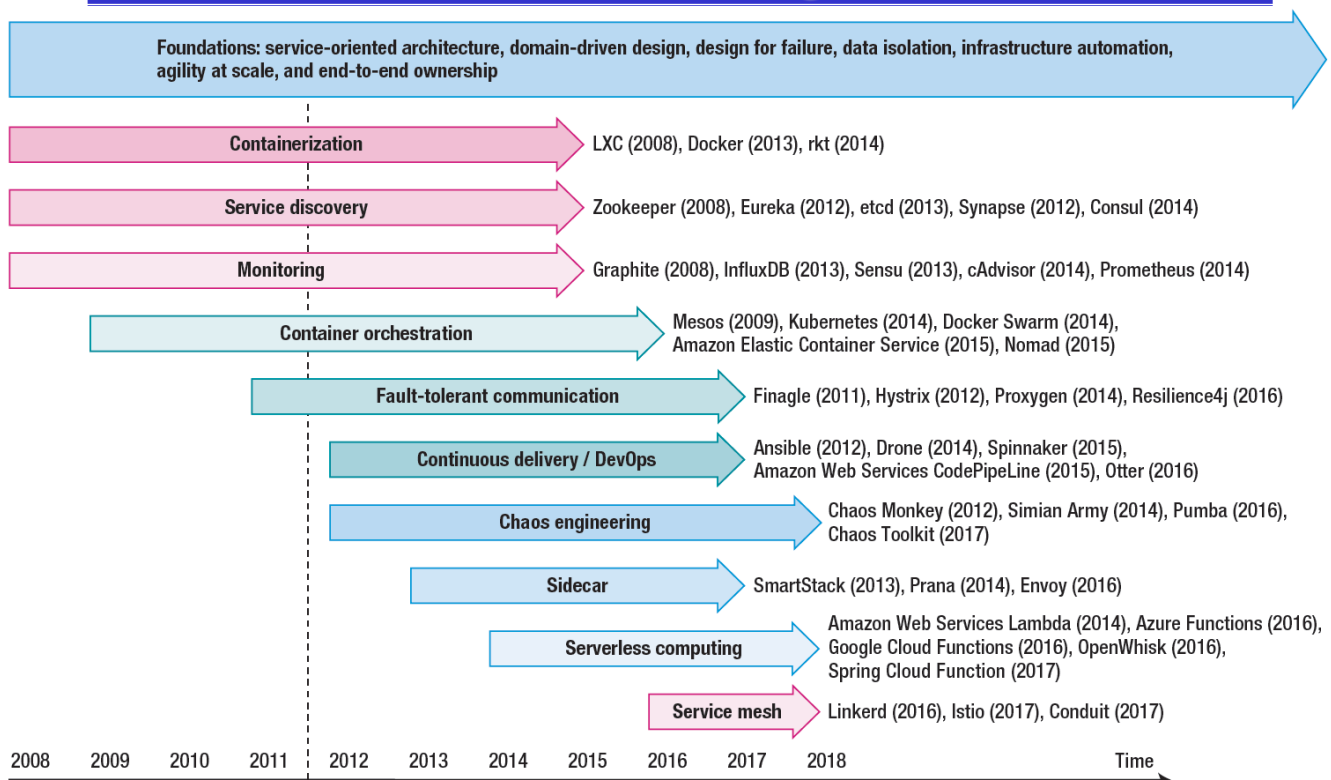
- Composed of 11 microservices written in different languages that communicate using gRPC
- Used by Google to demonstrate use of many technologies:
 - [Kubernetes](#) and Google Kubernetes Engine ([GKE](#))
 - gRPC: we know it 😊
 - [Istio](#): service mesh
 - [Cloud Operations](#): integrated monitoring, logging, and trace managed services for apps and systems running on Google Cloud
 - [Skaffold](#): command line tool that facilitates continuous development for Kubernetes applications
 - [Locust](#): load generator

Service mesh

- Dedicated infrastructure layer added to microservice app for facilitating service-to-service communications between microservices using a proxy
- Goals: transparently add capabilities like **observability**, **networking** and traffic management (including load balancing), **security**, without adding them to application code
- Products: [Istio](#), [Linkerd](#), [Envoy](#)



Microservice technologies timeline



The first use of "microservices" as a common architectural approach

From "Microservices: The Journey So Far and Challenges Ahead".

Valeria Cardellini – SDCC 2023/24

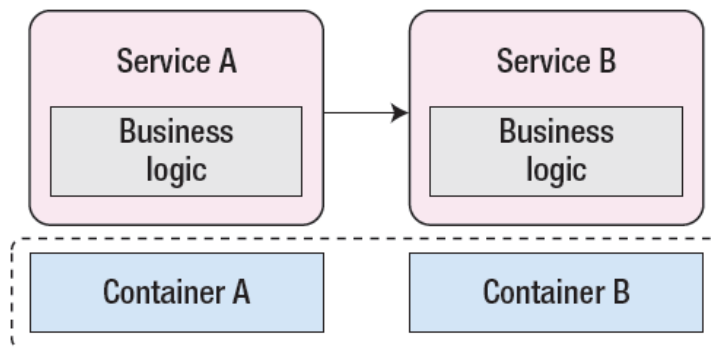
48

Generations: at the beginning

- 4 generations of microservice architectures
- **1st generation** based on:
 - **Container-based virtualization**
 - **Service discovery** (e.g., [Eureka](#), [etcd](#), [ZooKeeper](#))
 - [etcd](#): distributed reliable key-value store (e.g., used by Kubernetes as primary data store)
 - [Eureka](#): REST based service developed by Netflix; used in AWS cloud for locating services for load balancing and failover of middle-tier servers
 - **Monitoring** (e.g., [Graphite](#), [InfluxDB](#) and [Prometheus](#))
 - Enable runtime monitoring and analysis of microservice resources behavior at different levels of detail

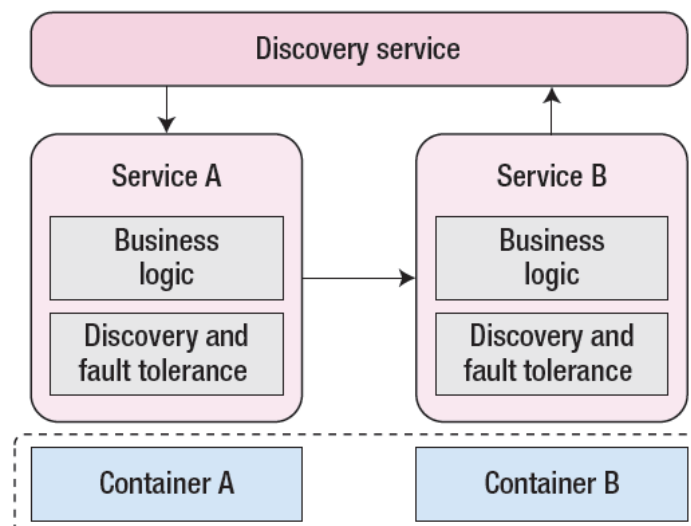
Generations: container orchestration

- Then, **container orchestration**
 - E.g., [Kubernetes](#), [Docker Swarm](#)
 - Automate container allocation and management tasks, abstracting away underlying physical or virtual infrastructure from service developers
 - But application-level failure-handling mechanisms are still implemented in services code



Generations: service discovery and fault tolerance

- **2nd generation** based on **discovery services** and **fault-tolerant communication libraries**
 - Let services communicate more efficiently and reliably

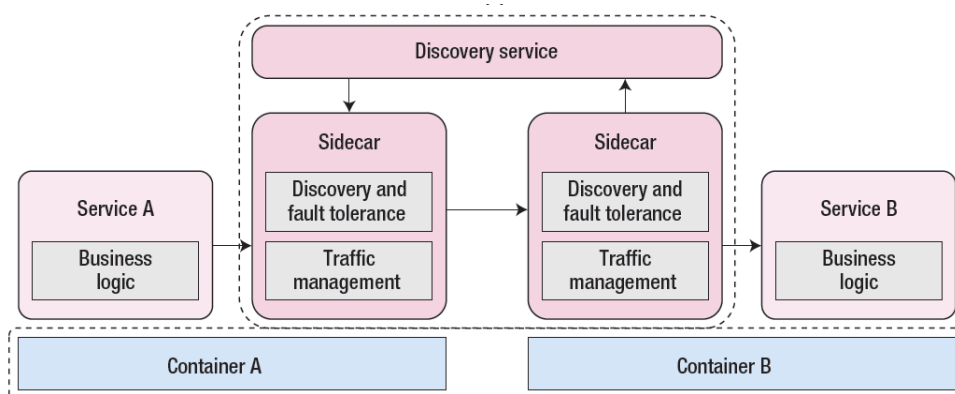


Generations: service discovery and fault tolerance

- Examples of discovery services and fault-tolerant communication libraries
 - [Consul](#): service discovery (now a more powerful service mesh, also to securely connect applications on Kubernetes)
 - [Finagle](#): fault tolerant, protocol-agnostic RPC system for the JVM, designed and used in production at Twitter (and other organizations)
 - [Hystrix](#): latency and fault tolerance library (currently in maintenance mode) designed by Netflix to:
 - isolate points of access to remote systems, services and 3rd party libraries
 - stop cascading failures (circuit breaker pattern)
 - enable resilience in complex DS

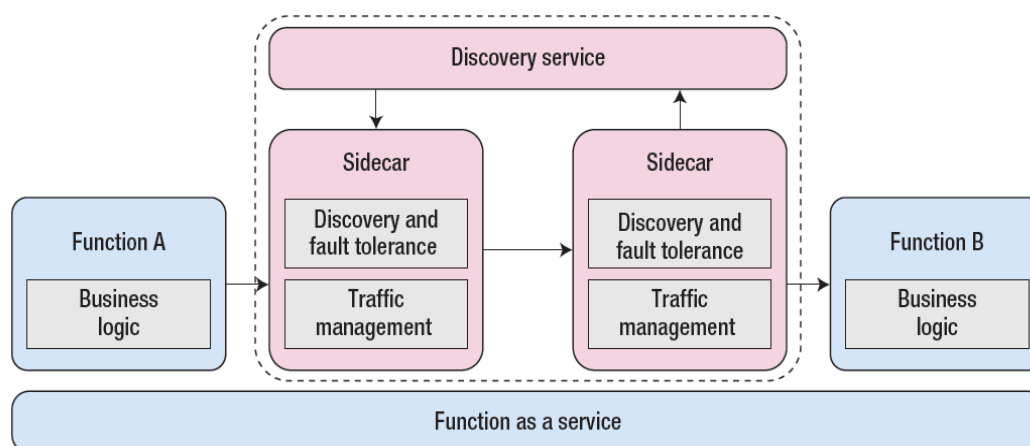
Generations: sidecar and service mesh

- **3rd generation** based on **sidecar** (or service proxy) and **service mesh** technologies (e.g., [Envoy](#) and [Istio](#))
 - Encapsulate communication-related features such as service discovery and use of protocol-specific and fault-tolerant communication libraries
 - Goal: abstract them from service developers, improve sw reusability and provide homogeneous interface



Generations: serverless

- **4th generation** based on **Function as a Service (FaaS)** and **serverless computing** to further simplify microservice development and delivery



Serverless computing

- Cloud computing model which aims to abstract server management and low-level infrastructure decisions away from users by means of **full automation**
- Users can develop, run and manage application code (i.e., **functions**), without no worry about provisioning, managing and scaling computing resources
- Runtime environment is **fully managed** by Cloud (or private infrastructure and platform) provider
- Serverless: functions still run on “servers” somewhere but we don’t care
- Function as a Service (**FaaS**) often as synonym
 - Still some discussion

Serverless through an analogy

- Services for moving homes

		Packaging	Delivery	Operations	Legal	Financial	Personnel
serverless	Modern movers	All objects	Any route	All decisions	All covered	Fine-grained Utilization-based	Small team
IaaS/PaaS cloud	Traditional movers	Limited support	Major roads	Basic	Basic	Coarse-grained	Large team
self-hosting	Moving it yourself (with family and friends)	Yourself	Yourself	Yourself	Yourself	Yourself	Yourself

Serverless: many definitions

van Eyk et al., Serverless is More: From PaaS to Present Cloud Computing (135 cit.), IEEE IC, May 2018 [47]:

“Serverless Computing is a form of cloud computing which allows users to run event-driven and granularly billed applications, without having to address the operational logic. Function-as-a-Service (FaaS) is a form of serverless computing where the cloud provider manages the resources, lifecycle, and event-driven execution of user-provided functions.”

Hellerstein et al., Serverless Computing: One Step Forward, Two Steps Back (360 cit.). CIDR, Jan 2019 [23]:

“Serverless computing offers the attractive notion of a platform in the cloud where developers simply upload their code, and the platform executes it on their behalf as needed at any scale. Developers need not concern themselves with provisioning or operating servers, and they pay only for the compute resources used when their code is invoked... Serverless is not only FaaS. It is FaaS supported by a “standard library”: the various multi-tenanted, autoscaling services provided by the vendor. In the case of AWS, this includes S3 (large object storage), DynamoDB (key-value storage), SQS (queuing services), SNS (notification services), and more.”

Castro et al., The Rise of Serverless Computing (196 cit.), CACM, Dec 2019 [10]:

“Serverless computing is a platform that hides server usage from developers and runs code on-demand automatically scaled and billed only for the time the code is running. This definition captures the two key features of serverless computing: (a) Cost—billed only for what is running (pay-as-you-go)...; serverless essentially supports “scaling to zero” and avoids the need to pay for idle servers. (b) Elasticity—scaling from zero to “infinity.” ... The main differentiators of serverless platforms is transparent autoscaling and fine-grained resource charging only when code is running. Function-as-a-Service is a serverless computing platform where the unit of computation is a function that is executed in response to triggers such as events or HTTP requests. Mobile Backend as-a-Service (MBaaS) or more generalized Backend as-a-Service (BaaS) bears a close resemblance to serverless computing.”

Jonas, ..., Patterson et al., Cloud Programming Simplified: A Berkeley View on Serverless Computing (485 cit.), arXiv, Feb 2019 [26], refined in a follow up publication by the same authors What Serverless Computing Is and Should Become: The Next Phase of Cloud Computing (75 cit.), CACM, May 2021 [41]

“In serverless computing, programmers create applications using high level abstractions offered by the cloud provider... They may also use serverless object storage, message queues, key-value store databases, mobile client data sync, and so on, a group of services offerings known collectively as Backend-as-a-Service (BaaS). Managed cloud function services are also called Function-as-a-Service (FaaS) and collectively Serverless Cloud Computing today = FaaS + BaaS. Three essential qualities of serverless computing are: 1. Providing an abstraction that hides the servers and the complexity of programming and operating them. 2. Offering a pay-as-you-go cost model instead of a reservation-based model, so there is no charge for idle resources. 3. Automatic, rapid, and unlimited scaling resources up and down to match demand closely, from zero to practically infinite.”

Serverless: many definitions

Kounev et al., Serverless Computing: What It Is, and What It Is Not?, Comm. ACM, 2023

Serverless computing is a cloud computing paradigm encompassing a class of cloud computing platforms that allow one to develop, deploy, and run applications (or components thereof) in the cloud **without allocating and managing virtualized servers and resources or being concerned about other operational aspects**.

The responsibility for operational aspects, such as fault tolerance or the **elastic scaling of computing, storage, and communication resources** to match varying application demands, is offloaded to the cloud provider.

Providers apply **utilization-based billing**: they charge cloud users with fine granularity, in proportion to the resources that applications actually consume from the cloud infrastructure, such as computing time, memory, and storage space.

Serverless: features

- Ephemeral compute resources
 - May only last for one function invocation
 - ✗ **Cold start**: when a request arrives and no container/microVM is ready to serve it, function execution must be delayed until a new container/microVM is launched
- Automated (i.e., zero configuration) elasticity
 - Compute resources auto-scale transparently from zero to peak load and back in response to workload shifts
- True pay-per-use: fine-grained and utilization-based
 - Pay only for consumed time (ms), rather than on pre-purchased units of capacity, e.g., AWS Lambda pricing

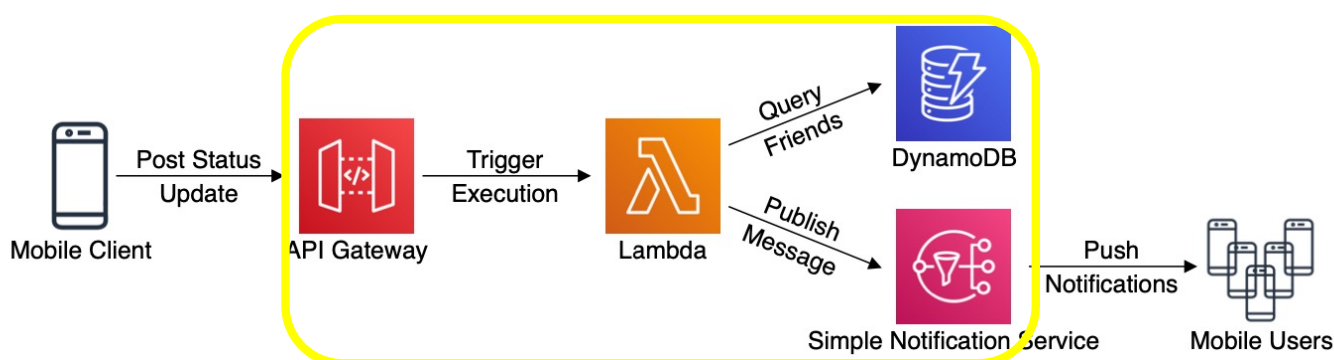
Architecture	Duration	Requests
x86 Price		
First 6 Billion GB-seconds / month	\$0.0000166667 for every GB-second	\$0.20 per 1M requests
Next 9 Billion GB-seconds / month	\$0.000015 for every GB-second	\$0.20 per 1M requests
Over 15 Billion GB-seconds / month	\$0.0000133334 for every GB-second	\$0.20 per 1M requests

Serverless: features

- Event-driven
 - When an event is triggered (e.g., file uploaded to storage, message ready in queue), a piece of infrastructure is allocated dynamically to execute the function code
- NoOps: simplifies the process of deploying code into production
 - Scaling, capacity planning and maintenance operations are hidden from developers or operators
- Supports diverse applications: from enterprise automation to scientific computing

Serverless application: example

- Mobile backend for a social media app
 1. Users compose status update and send it using mobile clients
 2. Platform orchestrates ops needed to propagate update inside the social media platform and to user's friends using serverless (AWS Lambda) and other cloud services
 3. Each friend receives updates on their social media app



Serverless Cloud services

- Several Cloud providers offer serverless computing on their public clouds as fully managed service
 - [AWS Lambda](#)
 - See [hands-on course](#)
 - Includes functions at the edge ([Lambda@Edge](#))
 - [Azure Functions](#)
 - [Google Cloud Functions](#)
 - [IBM Cloud Functions](#)
- Limited knobs to control performance of functions
 - Developers can configure only amount of memory allocated to function: amount of virtual CPU is proportional to amount of memory
- Cloud platforms also offer other supporting services (e.g., event notification, storage, message queue, DB) that are necessary for operating a serverless ecosystem

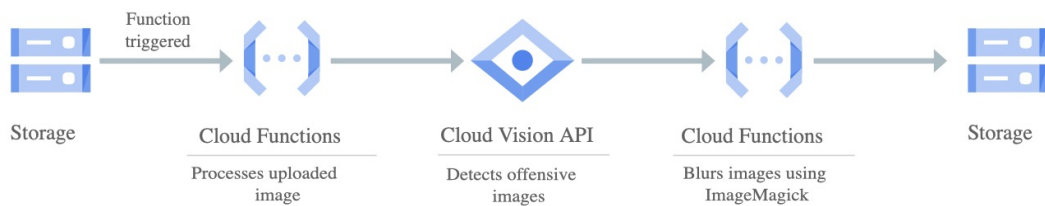
Example: Google Cloud Functions

- “Hello World” FaaS example from Google using Go
 - HTTP response that displays “Hello, World!”

```
// helloGet is an HTTP Cloud Function.  
func helloGet(w http.ResponseWriter, r *http.Request) {  
    fmt.Fprint(w, "Hello, World!")  
}
```

- Plus some initialization code, see full example cloud.google.com/functions/docs/tutorials

Example: Google Cloud Functions



- A more complex example
 - Function execution is triggered from storage when an image is uploaded to a Cloud Storage bucket
 - Function uses Cloud Vision API to detect violent or adult content
 - When violent or adult content is detected in an uploaded image, a second function is called to download the offensive image: it uses ImageMagick to blur the image, and then uploads the blurred image to the output bucket

Code: cloud.google.com/functions/docs/tutorials/imagemagick

Serverless: state

- Stateless functions are easy to manage (horizontal scalability, fast recovery, ...)
 - But stateless functions are not enough for a broad range of applications and algorithms
- How to handle stateful processing?
 - State can be external (e.g., handed over to external DB)
 - Issues to address:
 - Efficient access to shared state, so to keep auto-scaling benefits
 - Programming support, e.g., [Azure Durable Functions](#)
- How to handle transactions?

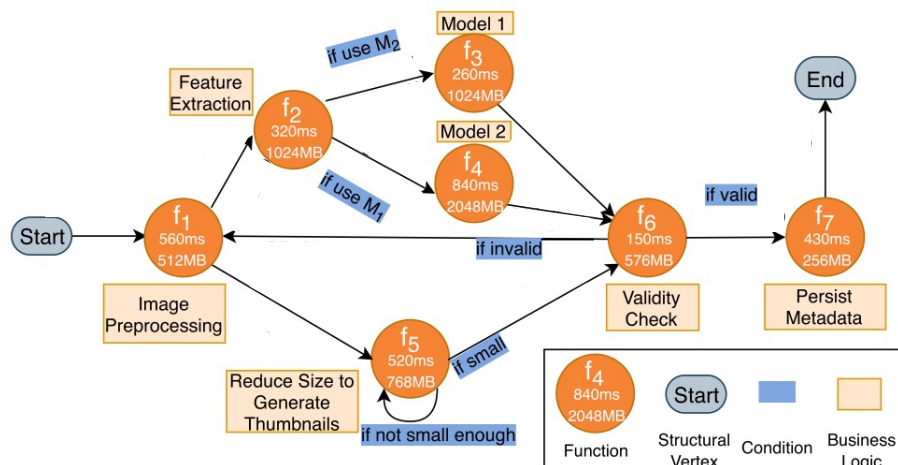
Serverless: challenges and limitations

- Performance
 - Cold starts
 - “The initial deployment of a function may take several minutes, while the underlying infrastructure is provisioned. Redeploying an existing function is faster, and incoming traffic is automatically migrated to the newest version.” (Google Cloud Functions)
 - Autoscaling: requires proactive one
- Runtime and language support
 - E.g., on Google: Node.js, Python, Go, Java, C#, Ruby, PHP
 - Language runtime impacts on performance and cost of serverless functions
- Resource limits
 - E.g., on AWS memory between 128 and 10240 MB per function
- Security
- Lack of standards and risk of vendor lock-in

Lower flexibility

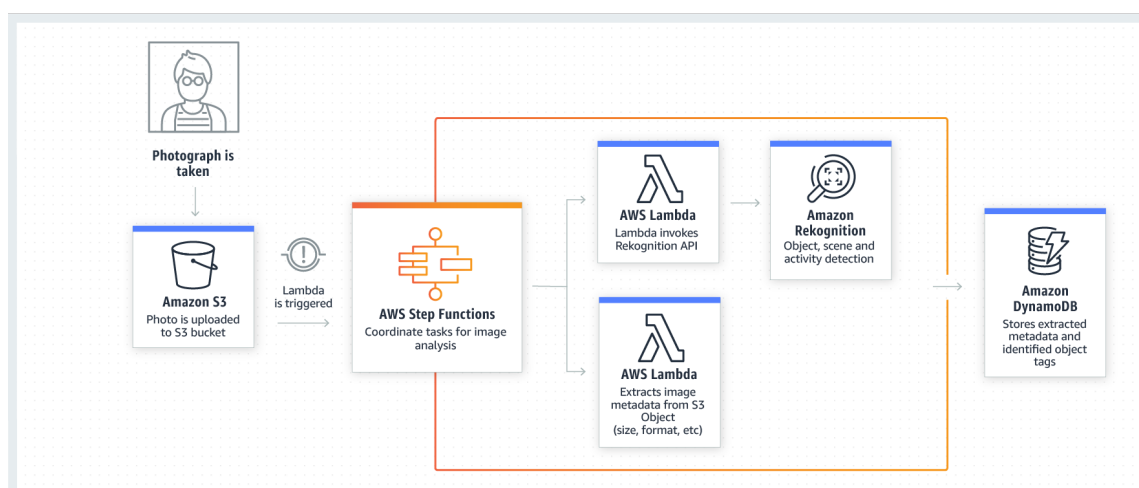
Composition of serverless functions

- Write small, simple, stateless functions
 - Complex functions are hard to understand, debug, and maintain
 - Separate code from data structures
- Then **compose** them in a **workflow**



Example: AWS Step Functions

- AWS Step Functions: serverless orchestration service that allows developers to coordinate multiple Lambda functions into workflows
- Example: process photo after its upload in S3

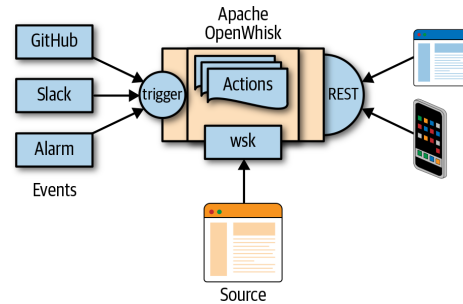


Open-source FaaS platforms

- Can run on commodity hardware
- Most platforms rely on Kubernetes for orchestration and management of serverless functions
 - Configuration management of containers
 - Container scheduling and service discovery
 - Elasticity management
- Prominent platforms
 - [Apache OpenWhisk](#)
 - [OpenFaaS](#)
 - [Fission](#)
 - [Knative](#)
 - [Nuclio](#)

OpenWhisk

- Open-source, distribute serverless platform that executes functions in response to events at any scale
openwhisk.apache.org



- Based on Docker containers
- Multiple container orchestration frameworks

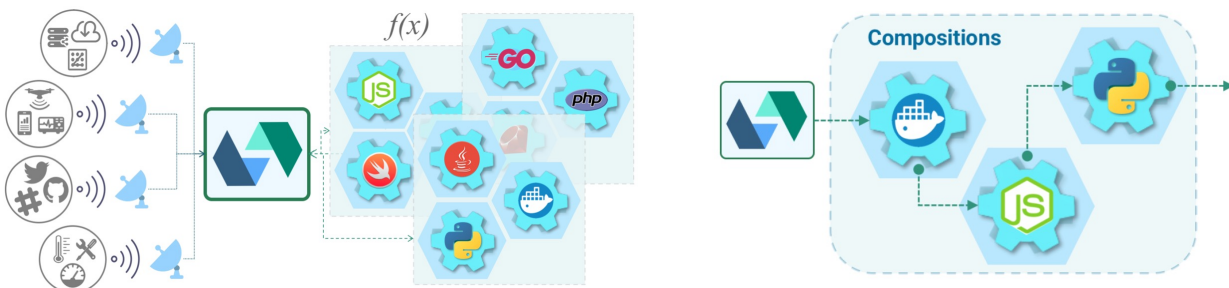


Valeria Cardellini – SDCC 2023/24

70

OpenWhisk

- Developers write functional logic (called **actions**)
 - In any supported programming language
 - Dynamically scheduled and run in response to associated events (via **triggers**) from external sources (**feeds**) or from HTTP requests
- Functions can be combined into **compositions**

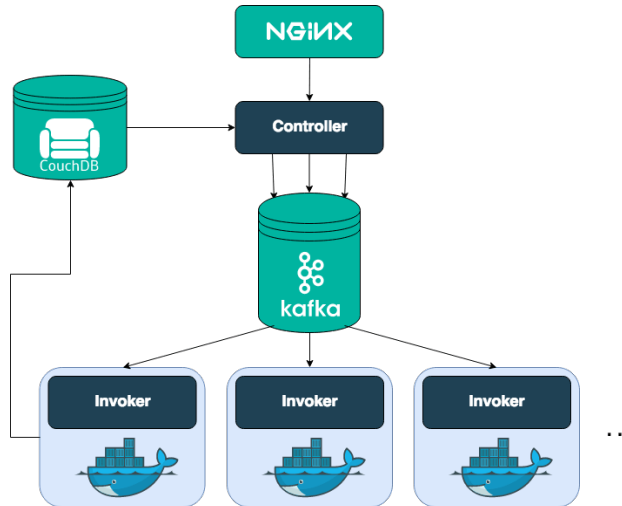


Valeria Cardellini – SDCC 2023/24

71

OpenWhisk architecture

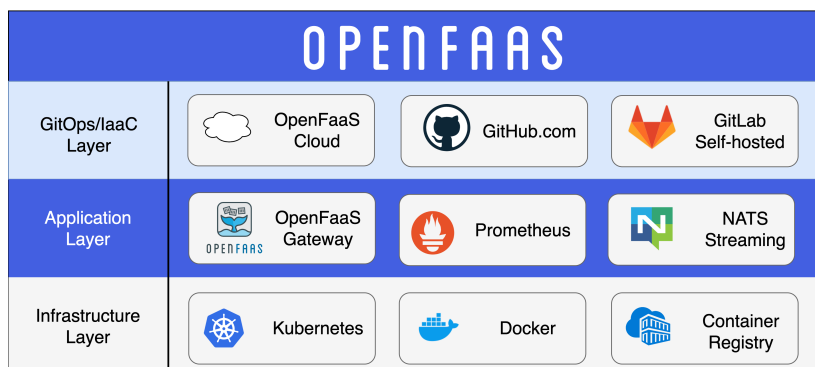
- Internal architecture powered by multiple distributed frameworks: Apache Kafka, [CouchDB](#), Docker and [NGINX](#)



OpenFaaS

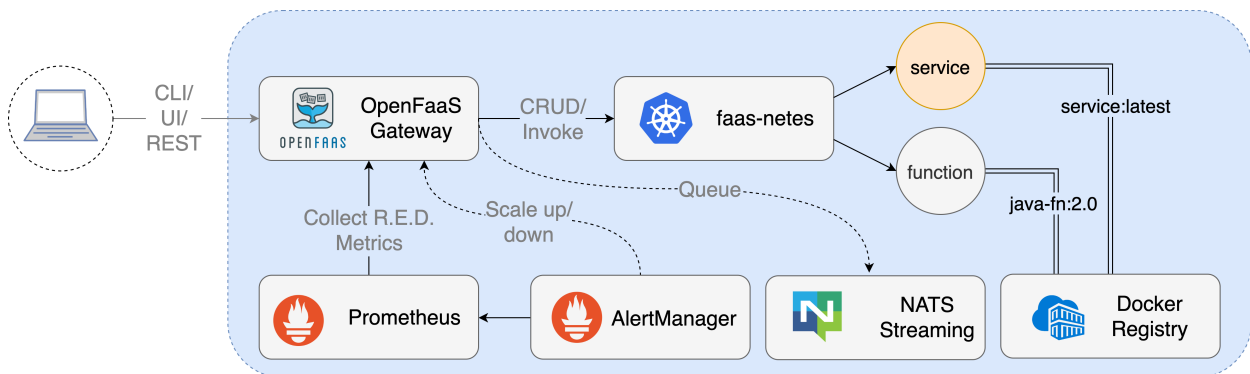


- Open-source FaaS framework for building functions on top of Docker and Kubernetes www.openfaas.com
- OpenFaaS stack
 - Gateway provides an external route into functions, collects metrics and scale functions
 - [Prometheus](#) provides metrics and enables auto-scaling
 - [NATS](#) provides asynchronous execution and queuing



OpenFaaS

- Conceptual workflow
 - OpenFaaS Gateway can be accessed through its REST API, via CLI or through UI
 - Prometheus collects metrics made available via Gateway's API and used for auto-scaling
 - NATS Streaming enables long-running tasks or function invocations to run in background

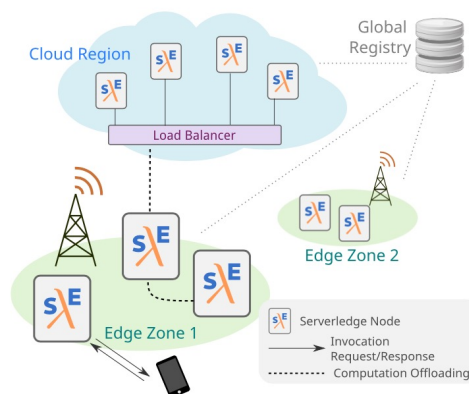


Valeria Cardellini – SDCC 2023/24

74

Serverless for the Edge-Cloud Continuum

- Open-source FaaS frameworks are unsuitable for Edge-Cloud continuum
- We are working on a decentralized FaaS framework called **Serverledge**: **thesis opportunities!**



Russo Russo et al., Decentralized Function-as-a-Service for the Edge-Cloud Continuum, Percom 2023 <https://github.com/grussorusso/serverledge>

Valeria Cardellini – SDCC 2023/24

75

Microservices: References and resources

- Lewis and Fowler, [Microservices](#)
- Lewis and Fowler, [Microservice Guides](#)
- Richardson, [Microservice Architecture](#)
- Jamshidi et al., [Microservices: The Journey So Far and Challenges Ahead](#), *IEEE Software*, 2018

Serverless: References and resources

- Roberts, [Serverless Architectures](#)
- Schleier-Smith et al., [What serverless computing is and should become: the next phase of cloud computing](#), *Comm. ACM*, 2021
- Kounev et al., [Serverless Computing: What It Is, and What It Is Not?](#), *Comm. ACM*, 2023