# Google Cloud

## Microservice Design and Architecture

# Learning objectives

- Decompose monolithic applications into microservices.

- Recognize appropriate microservice boundaries.

- Architect stateful and stateless services to optimize scalability and reliability.

- Implement services using 12-factor best practices.

- Build loosely coupled services by implementing a well-designed REST architecture.

- Design consistent, standard RESTful service APIs.

Google Cloud

# Agenda
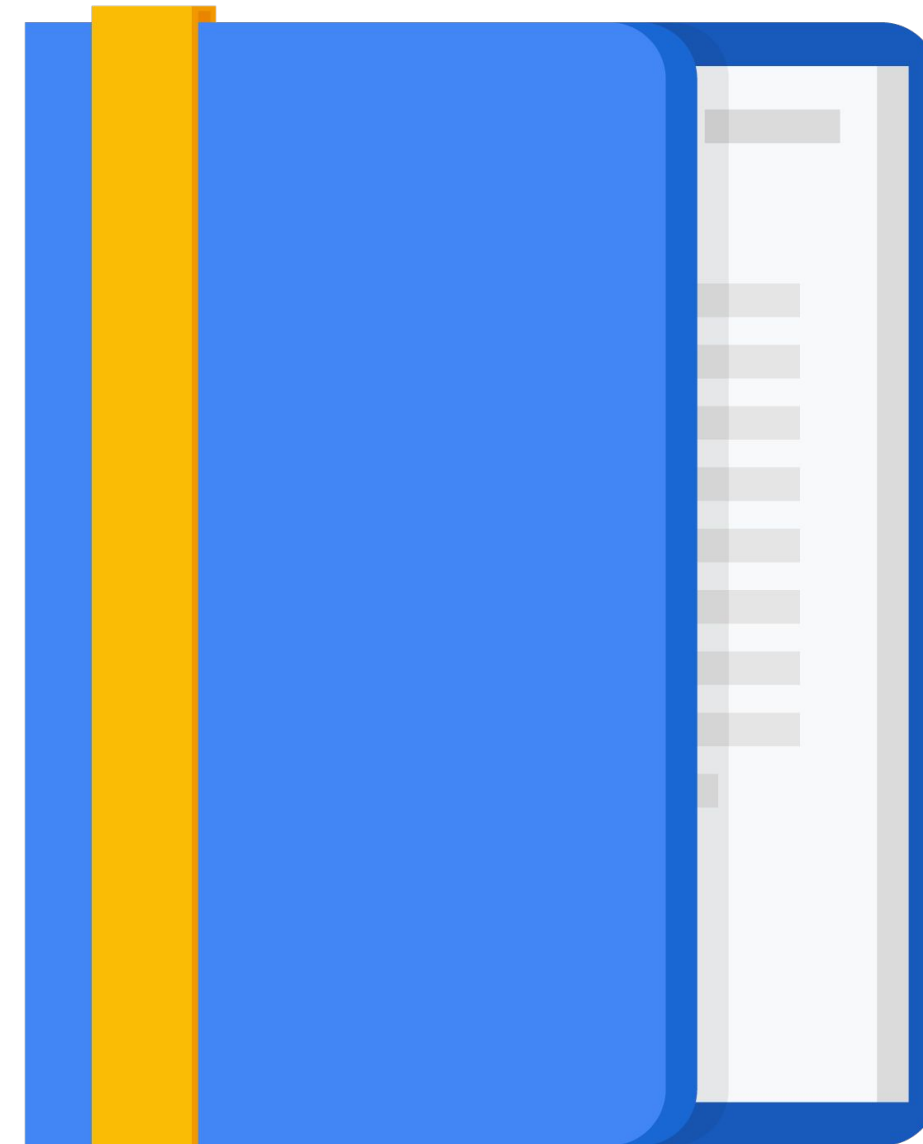
**Microservices**

Microservice Best Practices

Design Activity #4

REST

APIs

Design Activity #5

Google Cloud
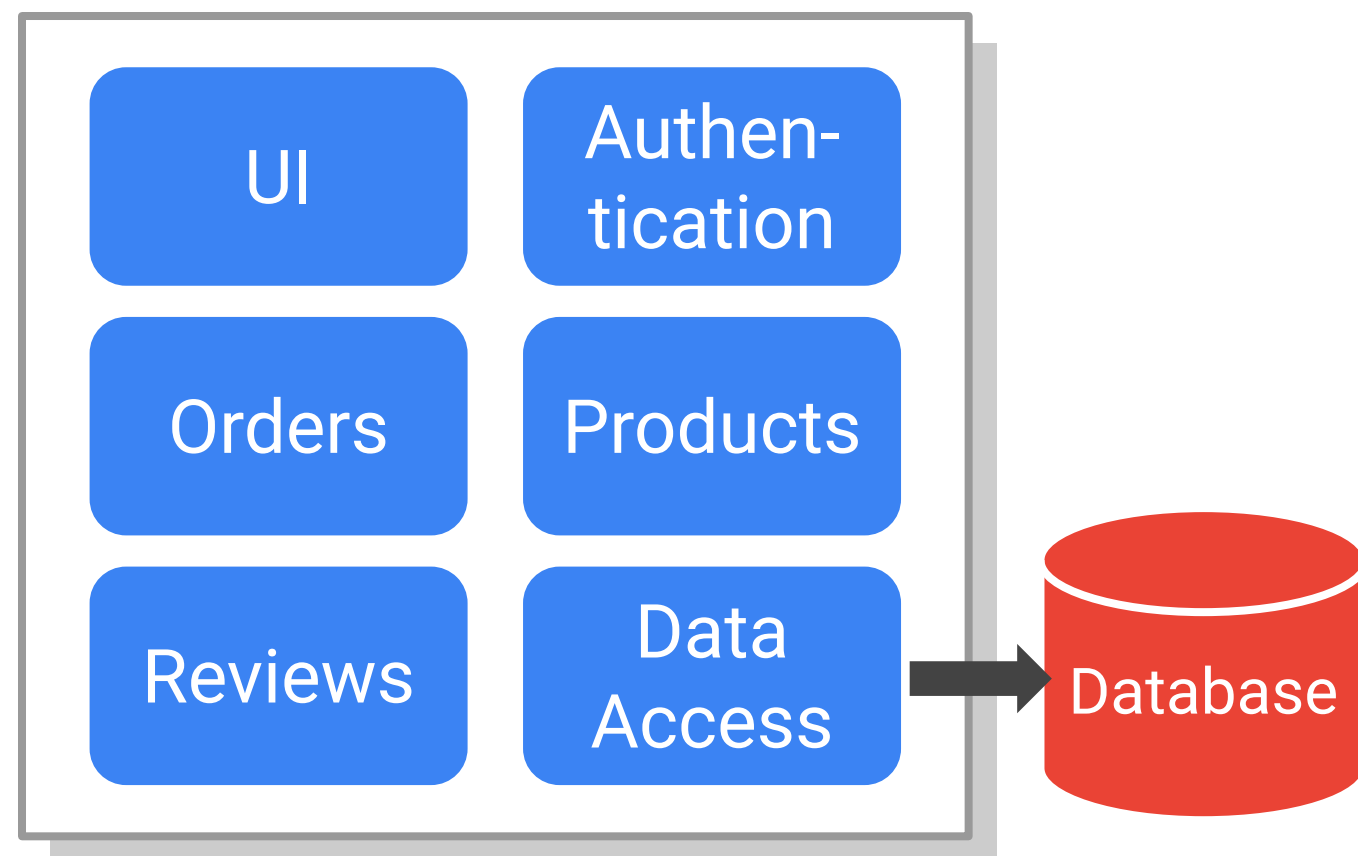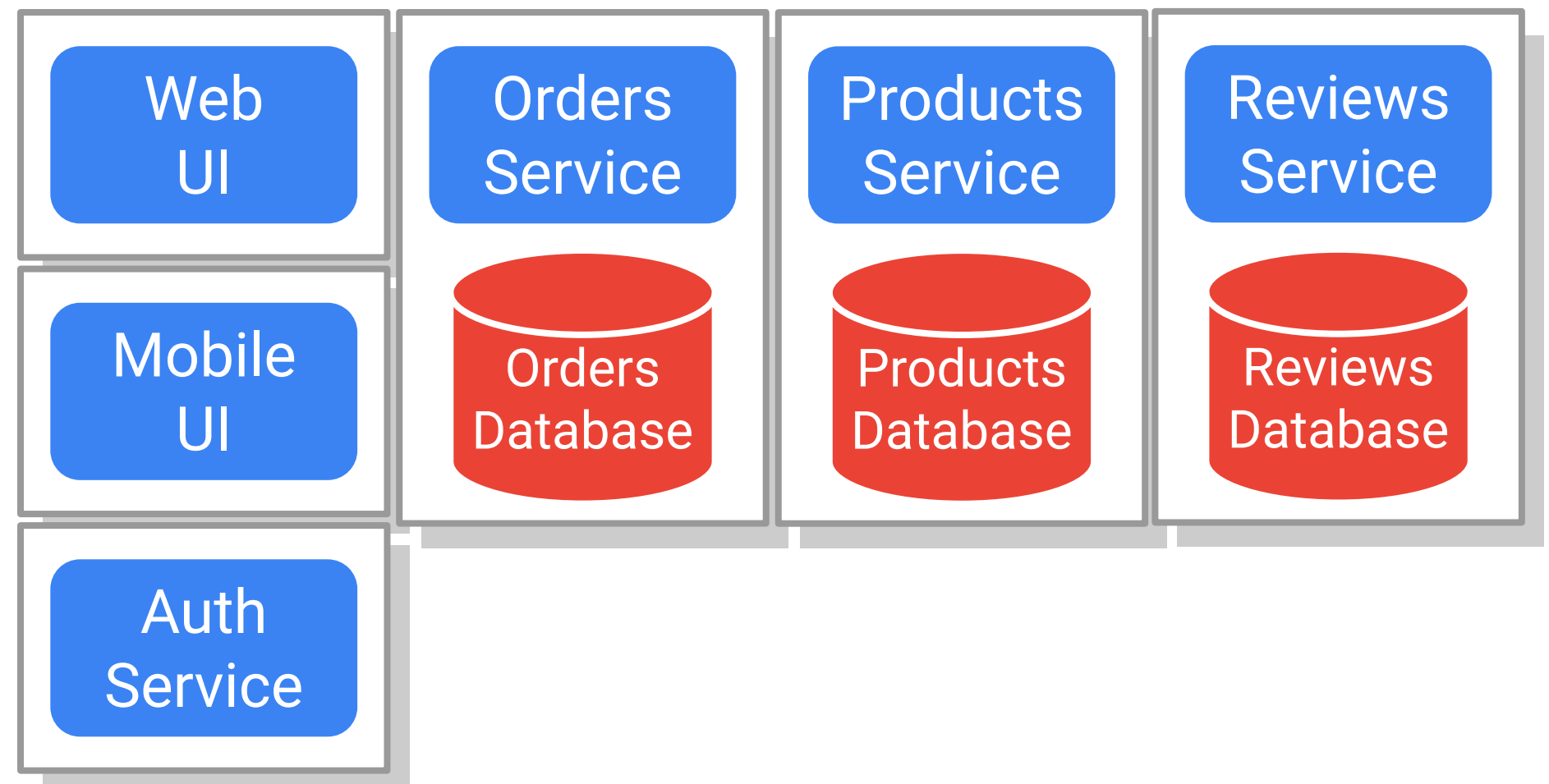
# Microservices divide a large program into multiple smaller, independent services

Monolithic applications implement all features in a single code base with a database for all data.

Microservice have multiple code bases, and each service manages its own data.

| UI | Authen-tication |
|----|----|
| Orders | Products |
| Reviews | Data Access |

→ Database

| Web UI | Orders Service | Products Service | Reviews Service |
|----|----|----|----|
| Mobile UI | Orders Database | Products Database | Reviews Database |
| Auth Service | | | |

Google Cloud

# Pros and cons of microservice architectures...

- Easier to develop and maintain
- Reduced risk when deploying new versions
- Services scale independently to optimize use of infrastructure
- Faster to innovate and add new features
- Can use different languages and frameworks for different services
- Choose the runtime appropriate to each service

- Increased complexity when communicating between services
- Increased latency across service boundaries
- Concerns about securing inter-service traffic
- Multiple deployments
- Need to ensure that you don't break clients as versions change
- Must maintain backward compatibility with clients as the microservice evolves

Google Cloud

# The key to architecting microservice applications is recognizing service boundaries

| Decompose applications by feature to minimize dependencies | Organize services by architectural layer | Isolate services that provide shared functionality |
|---|---|---|
| • Reviews service<br>• Orders service<br>• Products service<br>• Etc. | • Web, Android, and iOS user interfaces<br>• Data access services | • Authentication service<br>• Reporting service<br>• Etc. |

Google Cloud

# Stateful services have different challenges than stateless ones

Stateful services manage stored data over time

- Harder to scale
- Harder to upgrade
- Need to back up
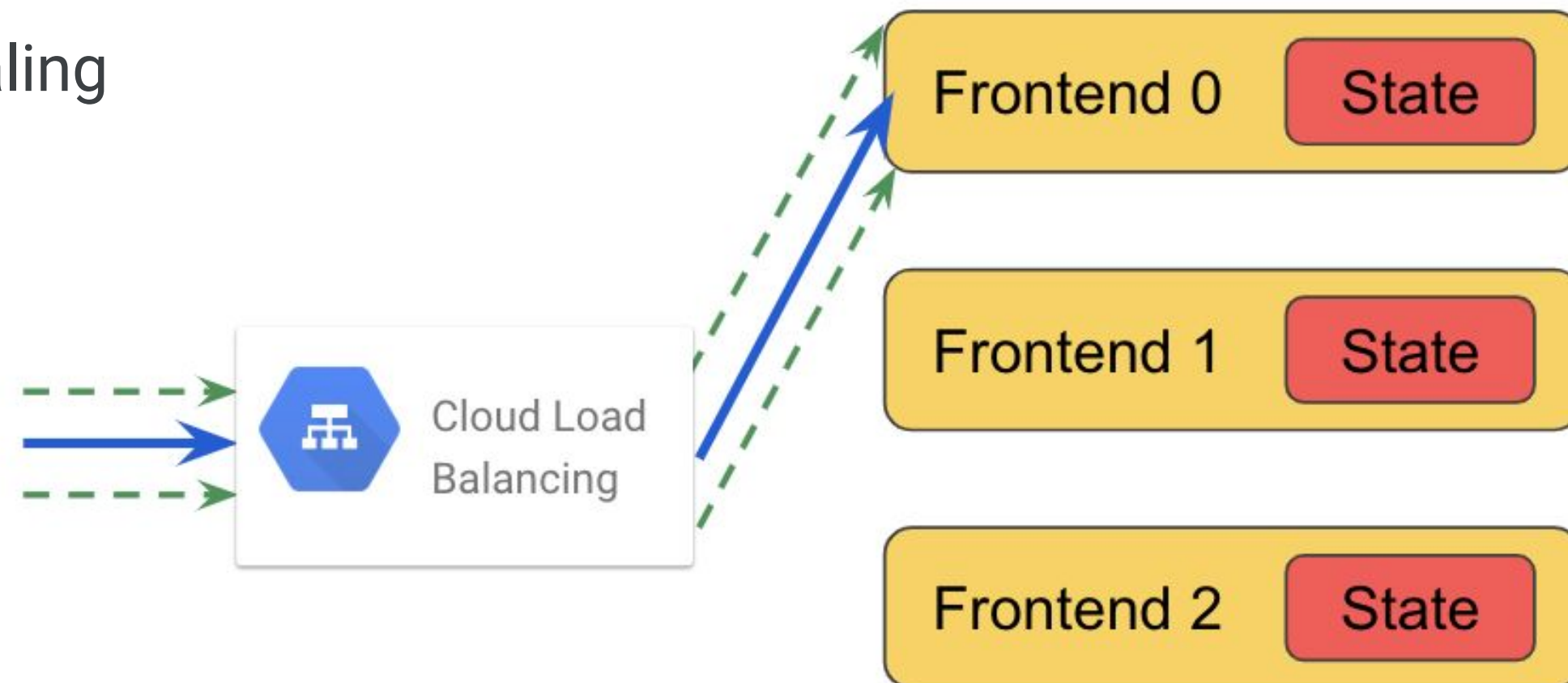
**Reviews Service**

**Reviews Database**

Stateless services get their data from the environment or other stateful services

- Easy to scale by adding instances
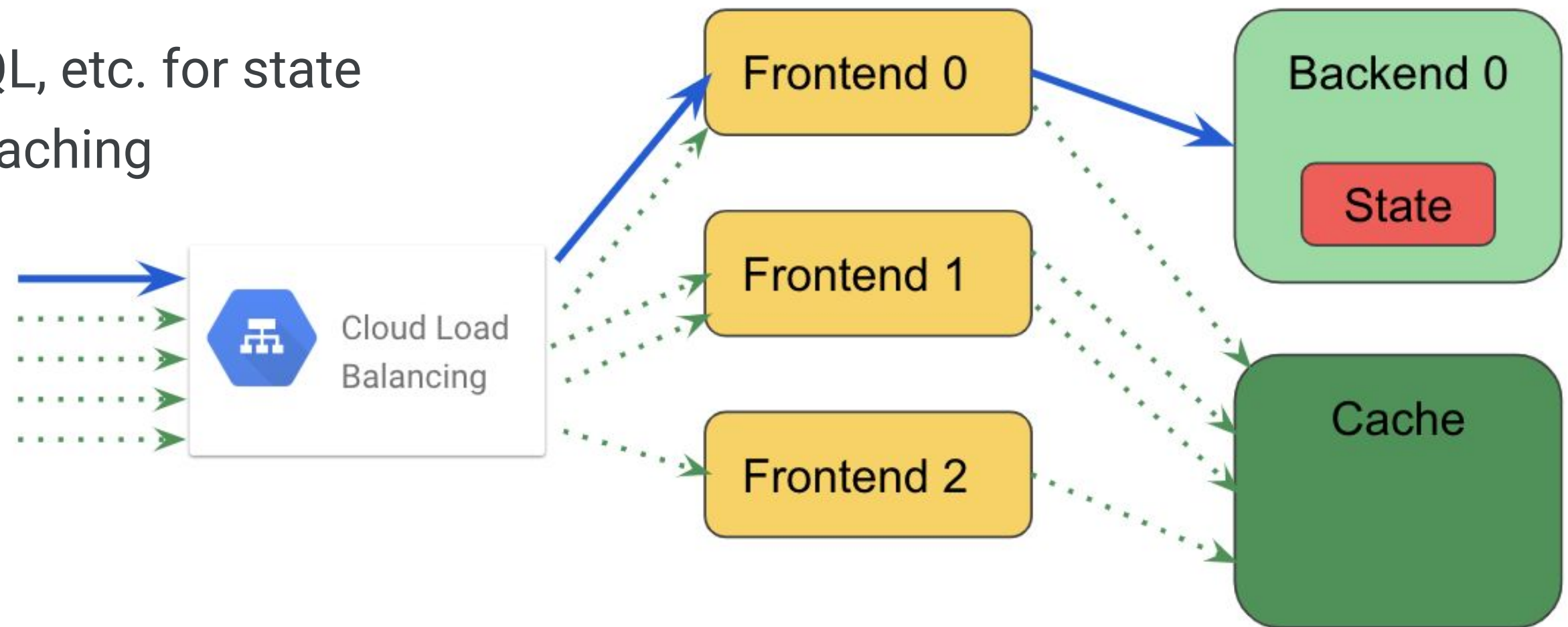- Easy to migrate to new versions
- Easy to administer

**Web UI**

Google Cloud

# Avoid storing shared state in-memory on your servers

- Requires sticky sessions to be set up in the load balancer
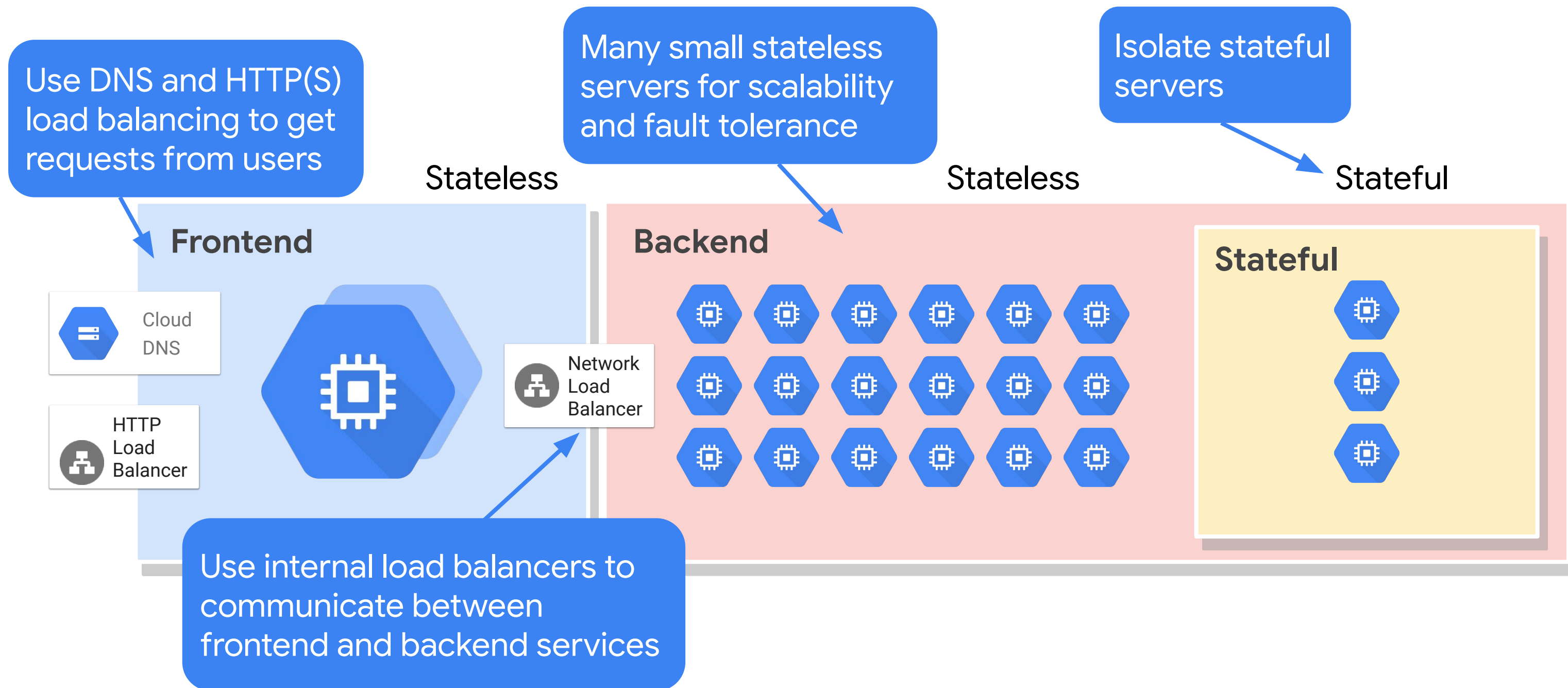- Hinders elastic autoscaling

# Store state using backend storage services shared by the frontend server

- Cache state data for faster access
- Take advantage of Google Cloud–managed data services
  - Firestore, Cloud SQL, etc. for state
  - Memorystore for caching



Google Cloud

# A general solution for large-scale cloud-based systems

Use DNS and HTTP(S) load balancing to get requests from users

Many small stateless servers for scalability and fault tolerance

Isolate stateful servers

Stateless

Stateless

Stateful

**Frontend**

**Backend**

**Stateful**

Cloud DNS

HTTP Load Balancer

Network Load Balancer

Use internal load balancers to communicate between frontend and backend services
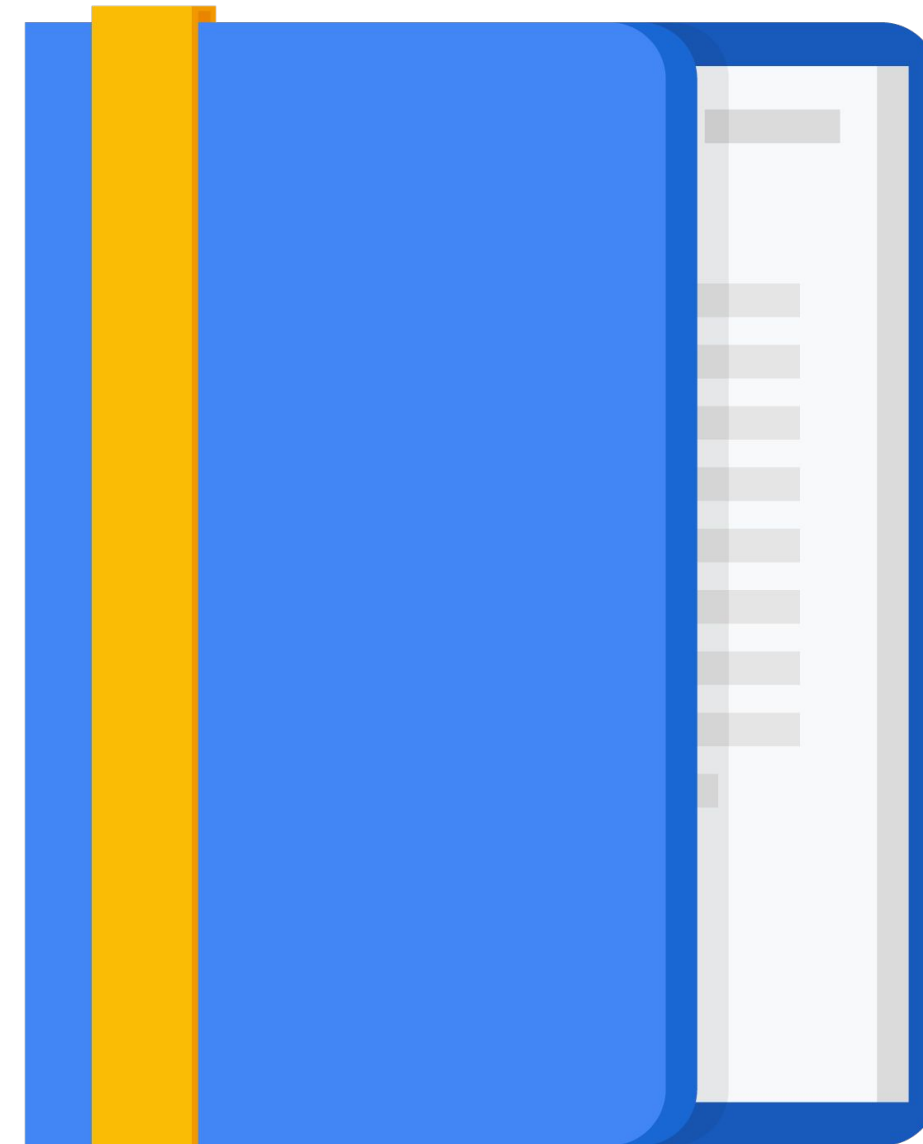
Google Cloud

# Agenda

Microservices

Design Activity #4

REST

APIs

Design Activity #5

Google Cloud

# The 12-factor app is a set of best practices for building web or software-as-a-service applications

- Maximize portability
- Deploy to the cloud
- Enable continuous deployment
- Scale easily


THE TWELVE-FACTOR APP

Google Cloud

# The 12 factors

| | |
|---|---|
| **1. Codebase**<br>One codebase tracked in revision control, many deploys | • Use a version control system like Git.<br>• Each app has one code repo and vice versa. |
| **2. Dependencies**<br>Explicitly declare and isolate dependencies | • Use a package manager like Maven, Pip, NPM to install dependencies.<br>• Declare dependencies in your code base. |
| **3. Config**<br>Store config in the environment | • Don't put secrets, connection strings, endpoints, etc., in source code.<br>• Store those as environment variables. |
| **4. Backing services**<br>Treat backing services as attached resources | • Databases, caches, queues, and other services are accessed via URLs.<br>• Should be easy to swap one implementation for another. |

Google Cloud

# The 12 factors (continued)

| | |
|---|---|
| 5. **Build, release, run**<br>Strictly separate build and run stages | • Build creates a deployment package from the source code.<br>• Release combines the deployment with configuration in the runtime environment.<br>• Run executes the application. |
| 6. **Processes**<br>Execute the app as one or more stateless processes | • Apps run in one or more processes.<br>• Each instance of the app gets its data from a separate database service. |
| 7. **Port binding**<br>Export services via port binding | • Apps are self-contained and expose a port and protocol internally.<br>• Apps are not injected into a separate server like Apache. |
| 8. **Concurrency**<br>Scale out via the process model | • Because apps are self-contained and run in separate process, they scale easily by adding instances. |

# The 12 factors (continued)

| | |
|---|---|
| 9. **Disposability**<br>Maximize robustness with fast startup and graceful shutdown | • App instances should scale quickly when needed.<br>• If an instance is not needed, you should be able to turn it off with no side effects. |
| 10. **Dev/prod parity**<br>Keep development, staging, and production as similar as possible | • Container systems like Docker makes this easier.<br>• Leverage infrastructure as code to make environments easy to create. |
| 11. **Logs**<br>Treat logs as event streams | • Write log messages to standard output and aggregate all logs to a single source. |
| 12. **Admin processes**<br>Run admin/management tasks as one-off processes | • Admin tasks should be repeatable processes, not one-off manual tasks.<br>• Admin tasks shouldn't be a part of the application. |

Google Cloud

# Activity 4: Designing microservices for your application

Refer to your Design and Process Workbook.

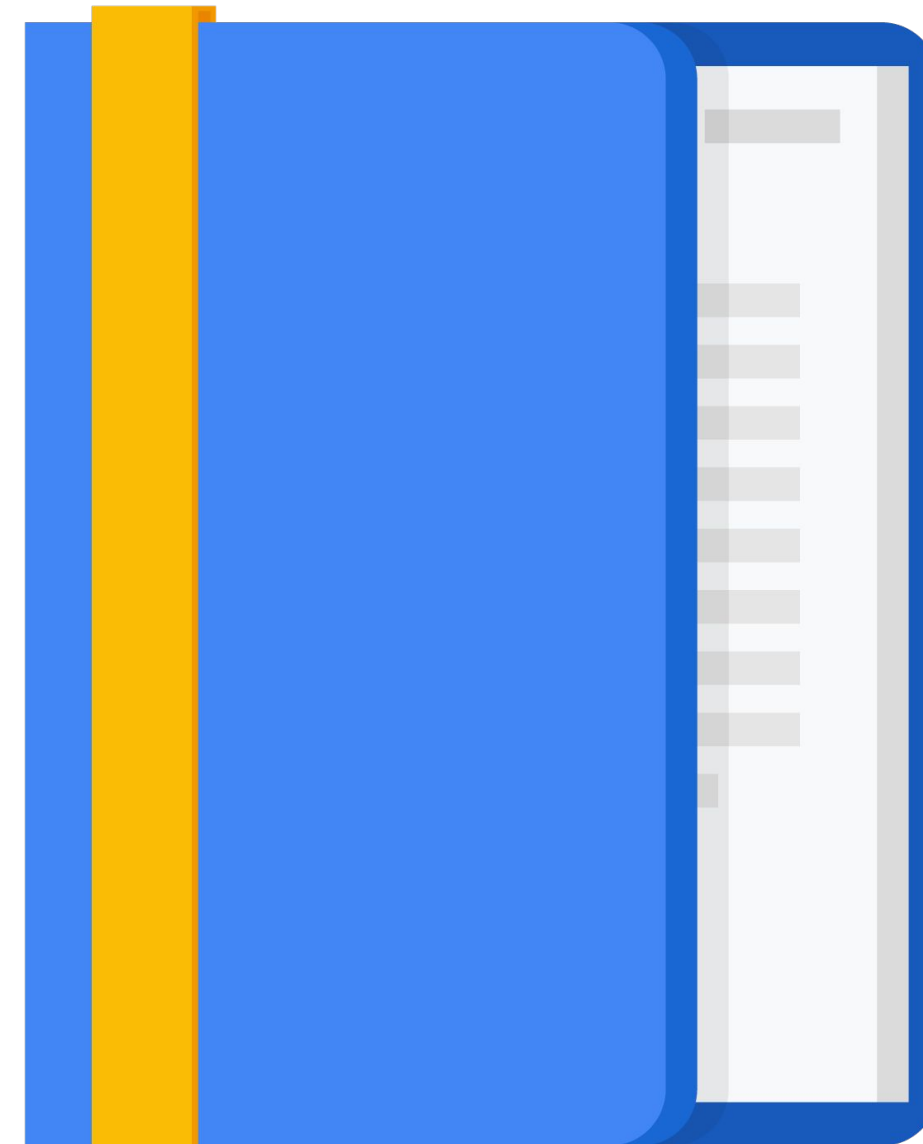- Diagram the microservices required by your case-study application.



Google Cloud

# Agenda

Microservices

Microservice Best Practices

Design Activity #4

REST

APIs

Design Activity #5

Google Cloud

# A good microservice design is loosely coupled

- Clients should not need to know too many details of services they use
- Services communicate via HTTPS using text-based payloads
  - Client makes GET, POST, PUT, or DELETE request
  - Body of the request is formatted as JSON or XML
  - Results returned as JSON, XML, or HTML
- Services should add functionality without breaking existing clients
  - Add, but don't remove, items from responses

*If microservices aren't loosely coupled, you'll end up with a really complicated monolith.*

Google Cloud

# REST architecture supports loose coupling

- REST stands for *Representational State Transfer*

- Protocol independent
  - HTTP is most common
  - Others possible like gRPC

- Service endpoints supporting REST are called *RESTful*

- Client and Server communicate with Request – Response processing

Google Cloud

# RESTful services communicate over the web using HTTP(S)

- URIs (or endpoints) identify resources
  - Responses return an immutable representation of the resource information

- REST applications provide consistent, uniform interfaces
  - Representation can have links to additional resources

- Caching of immutable representations is appropriate

Google Cloud

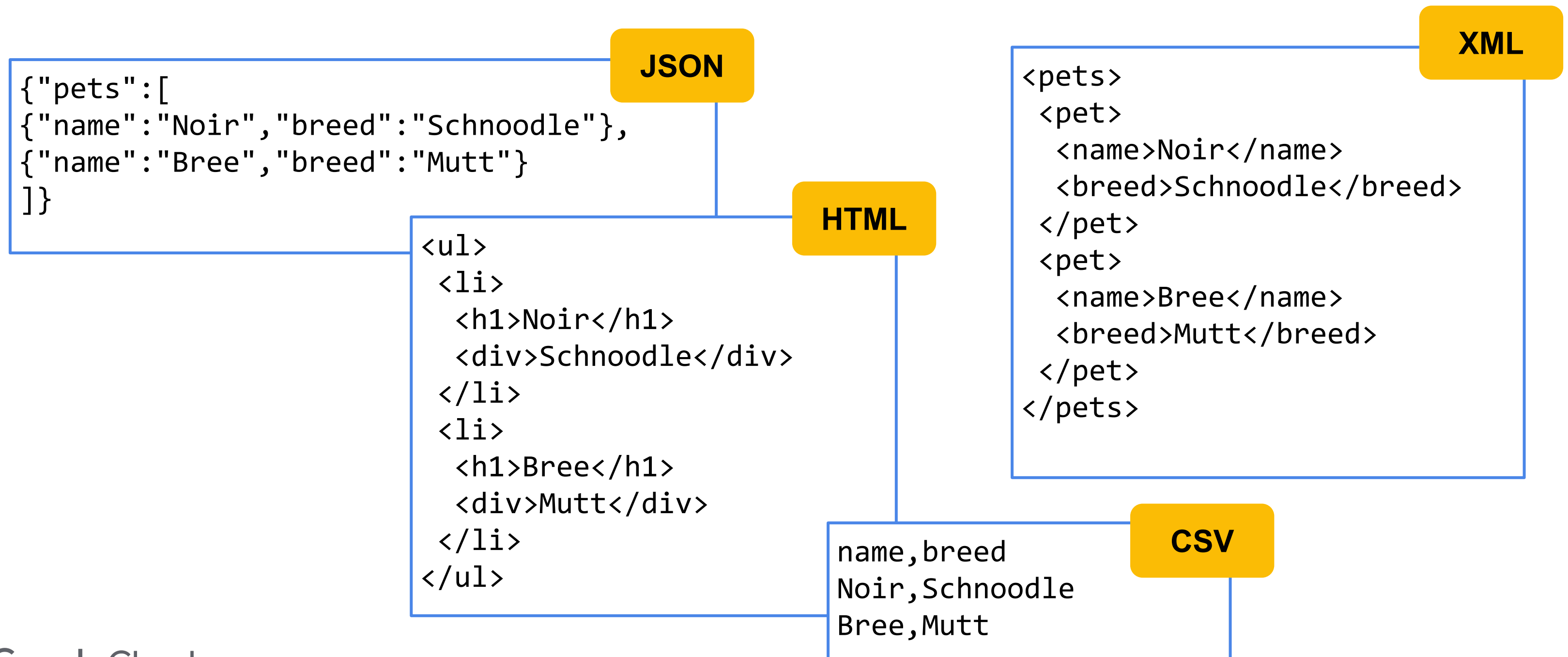# Resources and representations

- Resource is an abstract notion of information

- Representation is a copy of the resource information
  - Representations can be single items or a collection of items



Google Cloud

# Passing representations between services is done using standard text-based formats

**JSON**

```
{"pets":[
{"name":"Noir","breed":"Schnoodle"},
{"name":"Bree","breed":"Mutt"}
]}
```

**XML**

```
<pets>
 <pet>
  <name>Noir</name>
  <breed>Schnoodle</breed>
 </pet>
 <pet>
  <name>Bree</name>
  <breed>Mutt</breed>
 </pet>
</pets>
```

**HTML**

```
<ul>
 <li>
  <h1>Noir</h1>
  <div>Schnoodle</div>
 </li>
 <li>
  <h1>Bree</h1>
  <div>Mutt</div>
 </li>
</ul>
```

**CSV**

```
name,breed
Noir,Schnoodle
Bree,Mutt
```

Google Cloud

# Clients access services using HTTP requests

| <VERB> | <URI> | <HTTP Version> |
|---|---|---|
| <Request Header> | | |
| <Request Body> | | |

- VERB: GET, PUT, POST, DELETE

- URI: Uniform Resource Identifier (endpoint)

- Request Header: metadata about the message
  - Preferred representation formats (e.g., JSON, XML)

- Request Body: (Optional) Request state
  - Representation (JSON, XML) of resource
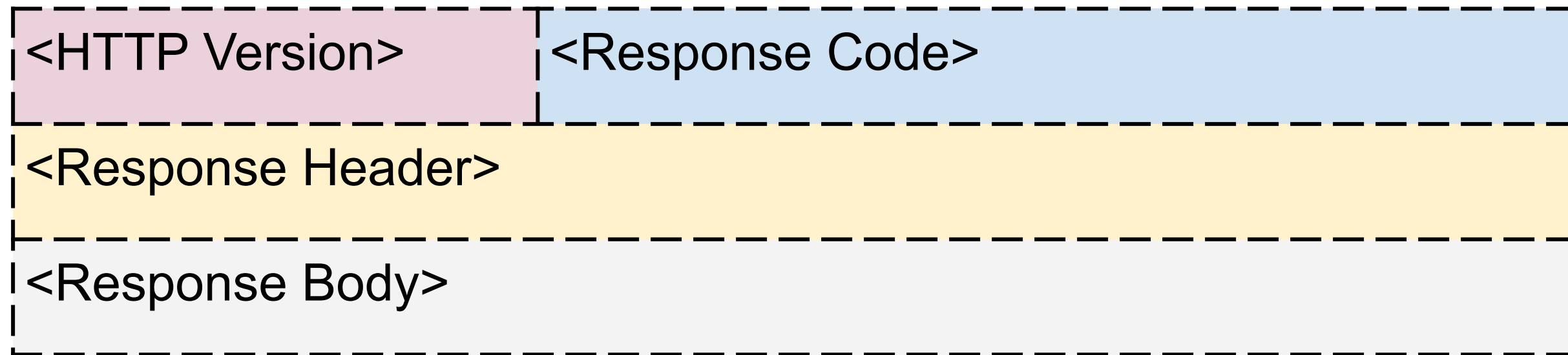
Google Cloud

# HTTP requests are simple and text-based

```
GET / HTTP/1.1
Host: pets.drehnstrom.com
```

```
POST /add HTTP/1.1
Host: pets.drehnstrom.com
Content-Type: json
Content-Length: 35

{"name":"Noir","breed":"Schnoodle"}
```

# The HTTP verb tells the server what to do

- **GET** is used to retrieve data

- **POST** is used to create data

  - Generates entity ID and returns it to the client

- **PUT** is used to create data or alter existing data

  - Entity ID must be known
  - *PUT should be* idempotent*, which means that whether the request is made once or multiple times, the effects on the data are exactly the same*

- **DELETE** is used to remove data

Google Cloud

# Services return HTTP responses

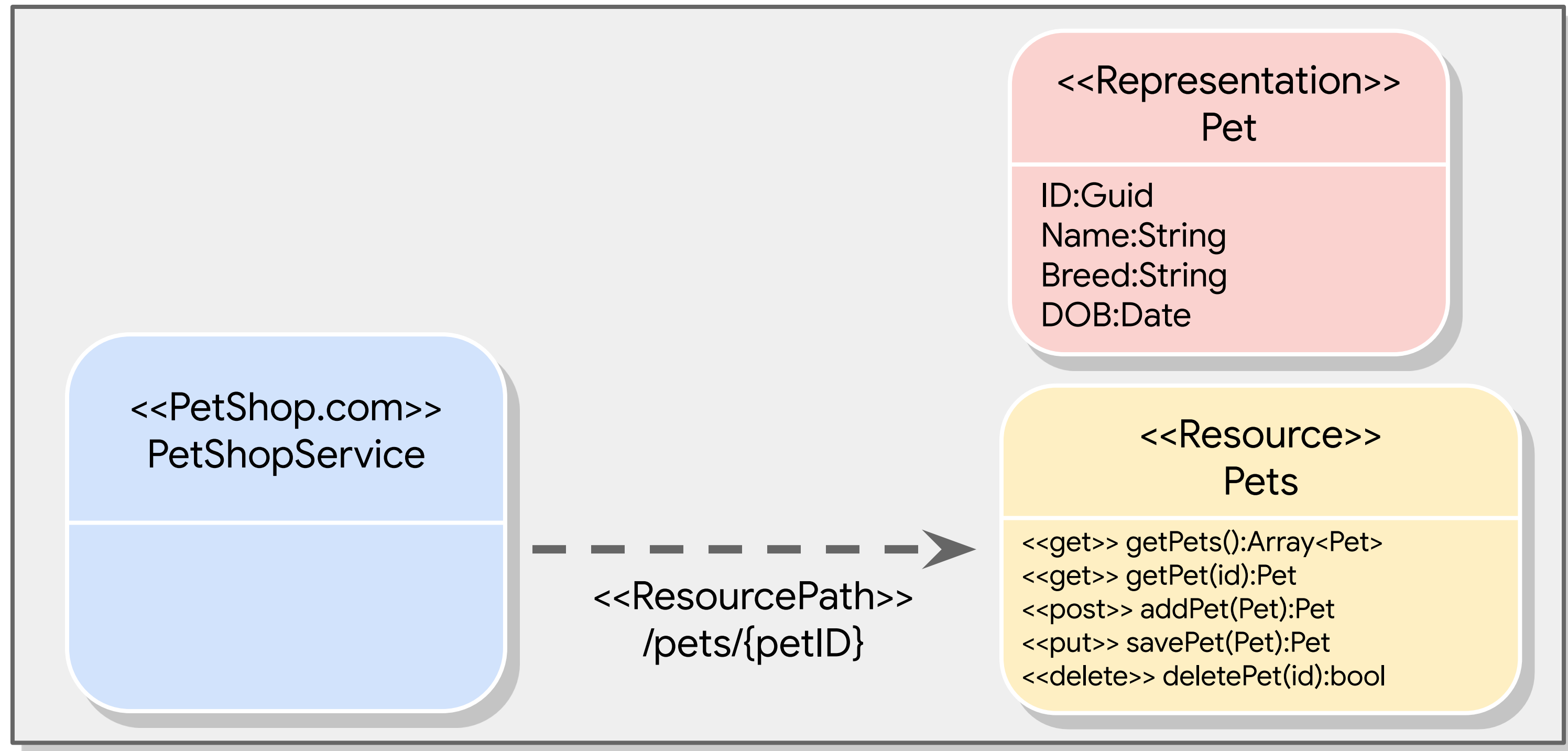| <HTTP Version> | <Response Code> |
|---|---|
| <Response Header> | |
| <Response Body> | |

- Response Code: 3-digit HTTP status code
  - 200 codes for success
  - 400 codes for client errors
  - 500 codes for server errors
- Response Body: contains resource representation
  - JSON, XML, HTML, etc.

Google Cloud

# All services need URIs (Uniform Resource Identifiers)

- Plural nouns for sets (collections)
- Singular nouns for individual resources
- Strive for consistent naming
- URI is case-insensitive
- Don't use verbs to identify a resource
- Include version information

Google Cloud

# Diagramming an example service



<<Representation>>
Pet

ID:Guid
Name:String
Breed:String
DOB:Date

<<PetShop.com>>
PetShopService

<<ResourcePath>>
/pets/{petID}

<<Resource>>
Pets

<<get>> getPets():Array<Pet>
<<get>> getPet(id):Pet
<<post>> addPet(Pet):Pet
<<put>> savePet(Pet):Pet
<<delete>> deletePet(id):bool

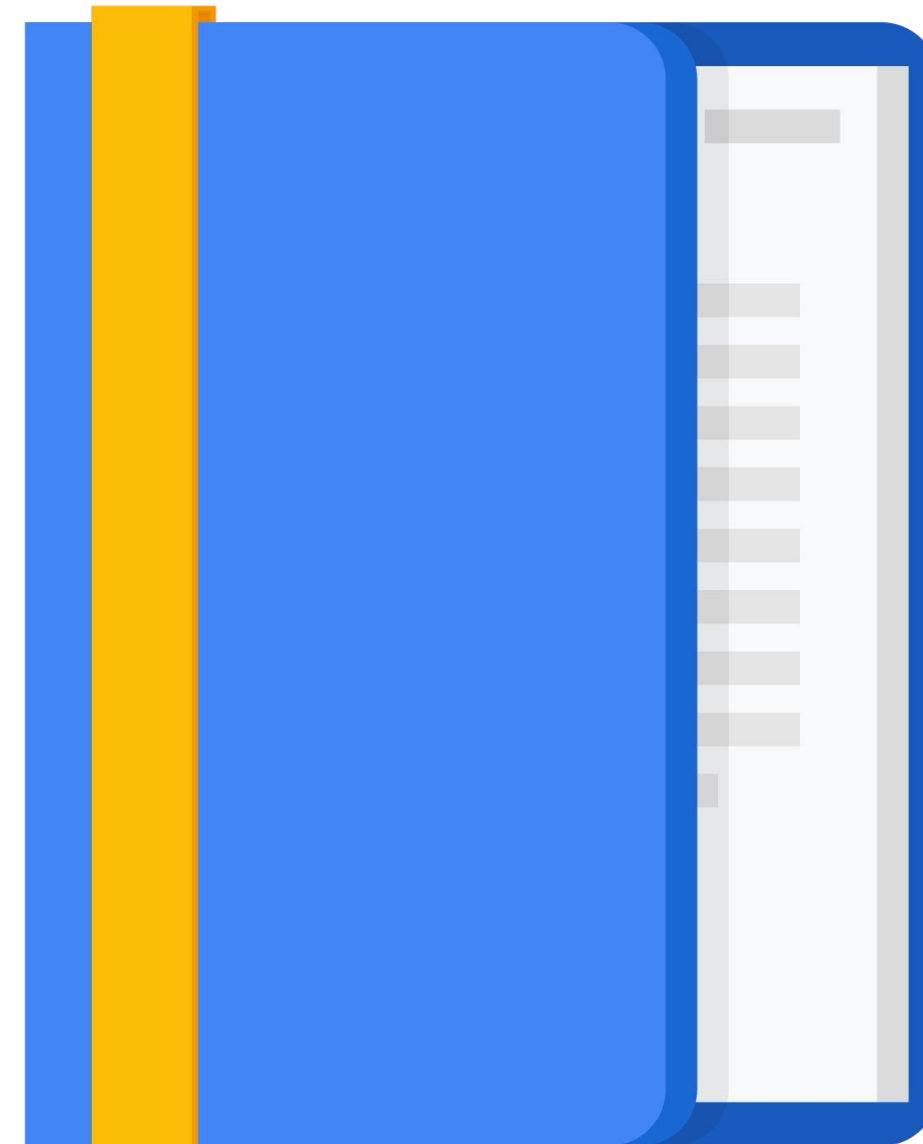Google Cloud

# Agenda

Microservices

Microservice Best Practices

Design Activity #4

REST

APIs

Design Activity #5

Google Cloud

# It's important to design consistent APIs for services

- Each Google Cloud service exposes a REST API
  - Functions are in the form:

    `service.collection.verb`
  - Parameters are passed either in the URL or in the request body in JSON format

- For example, the Compute Engine API has…
  - A service endpoint at: https://compute.googleapis.com
  - Collections include `instances`, `instanceGroups`, `instanceTemplates,` etc.
  - Verbs include insert, list, get, etc.

- So, to see all your instances, make a GET request to:

  `https://compute.googleapis.com/compute/v1/projects/{project}/zones/{zone}/instances`

Google Cloud

# OpenAPI is an industry standard for exposing APIs to clients

- Standard interface description format for REST APIs
  - Language agnostic
  - Open-source (based on Swagger)
- Allows tools and humans to understand how to use a service without needing its source code

```
1   openapi: "3.0.0"
2   info:
3     version: 1.0.0
4     title: Swagger Petstore
5     license:
6       name: MIT
7   servers:
8     - url: http://petstore.swagger.io/v1
9   paths:
10    /pets:
11      get:
12        summary: List all pets
13        operationId: listPets
14        tags:
15          - pets
```

Google Cloud

# gRPC is a lightweight protocol for fast, binary communication between services or devices

- Developed at Google
  - Supports many languages
  - Easy to implement
- gRPC is supported by Google services

  - Global load balancer (HTTP/2)

  - Cloud Endpoints

  - Can expose gRPC services using an Envoy Proxy in GKE

Google Cloud

# Google Cloud provides two tools, Cloud Endpoints and Apigee, for managing APIs

Both provide tools for:
- User authentication
- Monitoring
- Securing APIs
- Etc.

Both support OpenAPI and gRPC

Apigee API
Platform

Cloud
Endpoints

Google Cloud

# Activity 5: Designing REST APIs

Refer to your Design and Process Workbook.

- Design the APIs for your case study microservices.

Google Cloud

# Quiz

List some pros and cons of microservice architectures.

Google Cloud

# Quiz

List some pros and cons of microservice architectures.

| Pros | Cons |
| --- | --- |
| Easier to program and test<br>Scale independently<br>Less risky deployments<br>Easier to add new features<br>Etc. | Communication between services<br>Multiple deployments<br>Latency<br>Versioning<br>Etc. |

Google Cloud

# Quiz

You've re-architected a monolithic web application so state is not stored in memory on the web servers, but in a database instead. This has caused slow performance when retrieving user sessions though. What might be the best way to fix this?

A. Move session state back onto the web servers and use sticky sessions in the load balancer.

B. Use a caching service like Redis or Memorystore.

C. Increase the number of CPUs in the database server.

D. Make sure all web servers are in the same zone as the database.

Google Cloud

# Quiz

You've re-architected a monolithic web application so state is not stored in memory on the web servers, but in a database instead. This has caused slow performance when retrieving user sessions though. What might be the best way to fix this?

A.  Move session state back onto the web servers and use sticky sessions in the load balancer.

B.  Use a caching service like Redis or Memorystore.

C.  Increase the number of CPUs in the database server.

D.  Make sure all web servers are in the same zone as the database.

Google Cloud

# Quiz

Which below would **_violate_** 12-factor app best practices?

A.  Store configuration information in your source repository for easy versioning.

B.  Treat logs as event streams and aggregate logs into a single source.

C.  Keep development, testing, and production as similar as possible.

D.  Explicitly declare and isolate dependencies.

# Quiz

Which below would **_violate_** 12-factor app best practices?

A. Store configuration information in your source repository for easy versioning.

B. Treat logs as event streams and aggregate logs into a single source.

C. Keep development, testing, and production as similar as possible .

D. Explicitly declare and isolate dependencies.

Google Cloud

# Quiz

You're writing a service, and you need to handle a client sending you invalid data in the request. What should you return from the service?

A. An XML exception

B. A 200 error code

C. A 400 error code

D. A 500 error code

Google Cloud

# Quiz

You're writing a service, and you need to handle a client sending you invalid data in the request. What should you return from the service?

A.  An XML exception

B.  A 200 error code

C.  A 400 error code

D.  A 500 error code

Google Cloud

# Quiz

You're building a RESTful microservice. Which would be a valid data format for returning data to the client?

A. JSON

B. XML

C. HTML

D. All of the above

Google Cloud

# Quiz

You're building a RESTful microservice. Which would be a valid data format for returning data to the client?

A. JSON

B. XML

C. HTML

D. All of the above

Google Cloud

# Review

## Microservice Design and Architecture

Google Cloud

# More resources

API Design Guide

https://cloud.google.com/apis/design/

Authenticating service-to-service calls with Google Cloud Endpoints

https://youtu.be/4PgX3yBJEyw

Google Cloud