# Cloud Design Patterns

# Agenda

- ▶ Why Cloud Design Patterns
- ▶ Challenges in cloud development
- ▶ Catalogue of patterns

# Why Cloud Design Patterns

- In enterprise, we need reliable, scalable, secure applications
- Design patterns are useful for building reliable, scalable, secure applications in the cloud.
- They help to migrate existing applications to cloud native applications.

# Challenges in cloud development

**Data Management**

▶ Data management is the key element of cloud applications

▶ It influences most of the quality attributes.

▶ Data is typically hosted in different locations and across multiple servers for reasons such as performance, scalability or availability, and this can present a range of challenges.

▶ For example, data consistency must be maintained, and data will typically need to be synchronized across different locations.

# Challenges in cloud development

**Design and Implementation**

▶ Good design encompasses factors such as consistency and coherence in component design and deployment, maintainability to simplify administration and development, and reusability to allow components and subsystems to be used in other applications and in other scenarios.

▶ Decisions made during the design and implementation phase have a huge impact on the quality and the total cost of ownership of cloud hosted applications and services.

# Challenges in cloud development

**Messaging**

▶ The distributed nature of cloud applications requires a messaging infrastructure that connects the components and services, ideally in a loosely coupled manner in order to maximize scalability

▶ Asynchronous messaging is widely used, and provides many benefits, but also brings challenges such as the ordering of messages, poison message management, idempotency, and more.

# Patterns (Some of them)

- Ambassador
- Anti Corruption Layer
- Backends for Frontends
- Bulkhead
- Gateway Aggregation
- Gateway Offloading
- Gateway Routing
- Strangler
- Side Car

# Ambassador Pattern

*"Create helper services that send network requests on behalf of a consumer service or application. An ambassador service can be thought of as an out-of-process proxy that is co-located with the client."*

- This pattern can be useful for offloading common client connectivity tasks such as monitoring, logging, routing, security (such as TLS), and resiliency patterns in a language agnostic way.

- It is often used with legacy applications, or other applications that are difficult to modify, in order to extend their networking capabilities.

- It can also enable a specialized team to implement those features

# Ambassador Pattern

**Context and problem**

▶ Resilient cloud-based applications require features such as circuit breaking, routing, metering and monitoring, and the ability to make network-related configuration updates.

▶ It may be difficult or impossible to update legacy applications or existing code libraries to add these features, because the code is no longer maintained or can't be easily modified by the development team.
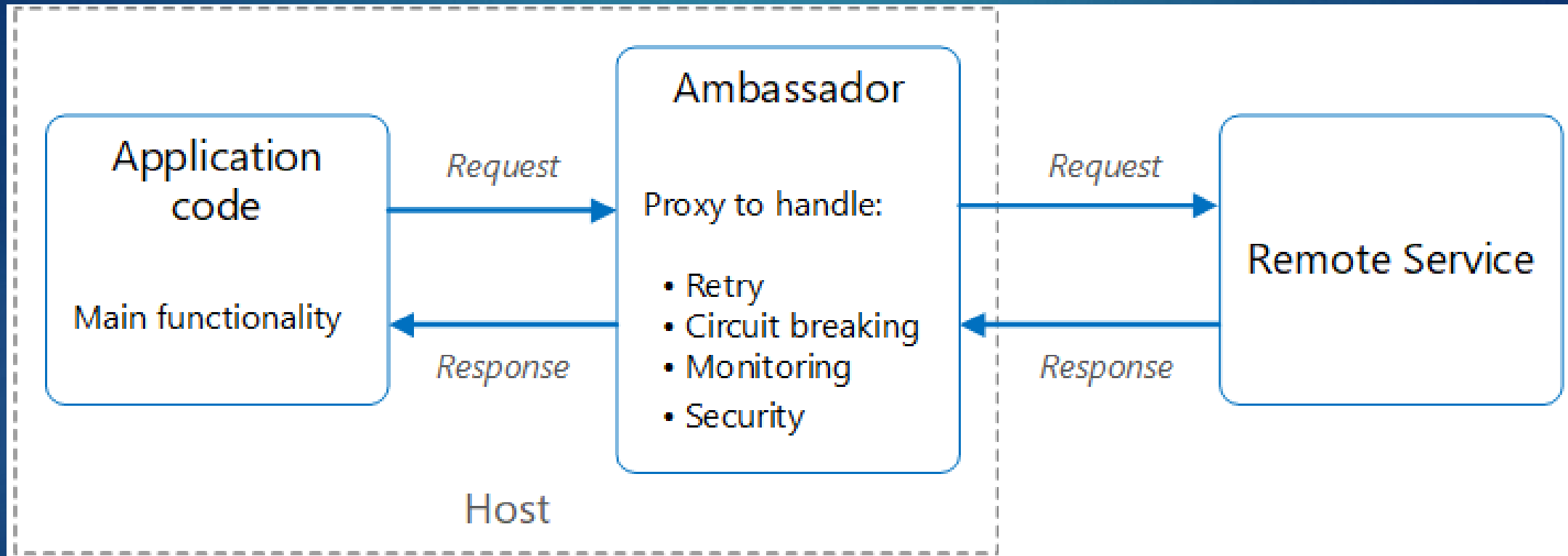
# Ambassador Pattern

**Context and problem**

▶ Network calls may also require substantial configuration for connection, authentication, and authorization.

▶ If these calls are used across multiple applications, built using multiple languages and frameworks, the calls must be configured for each of these instances.

▶ In addition, network and security functionality may need to be managed by a central team within your organization. With a large code base, it can be risky for that team to update application code they aren't familiar with.

# Ambassador Pattern

**Solution**

► Put client frameworks and libraries into an external process that acts as a proxy between your application and external services.

► Deploy the proxy on the same host environment as your application to allow control over routing, resiliency, security features, and to avoid any host-related access restrictions.

► You can also use the ambassador pattern to standardize and extend instrumentation.

► The proxy can monitor performance metrics such as latency or resource usage, and this monitoring happens in the same host environment as the application.

# Ambassador Pattern

# Ambassador Pattern

**Issues and considerations**

- The proxy adds some latency overhead. Consider whether a client library, invoked directly by the application, is a better approach.

- Consider the possible impact of including generalized features in the proxy.

  - For example, the ambassador could handle retries, but that might not be safe unless all operations are idempotent.

- Consider a mechanism to allow the client to pass some context to the proxy, as well as back to the client.

  - For example, include HTTP request headers to opt out of retry or specify the maximum number of times to retry.

- Consider how you will package and deploy the proxy.

- Consider whether to use a single shared instance for all clients or an instance for each client.

# Anti-Corruption Layer pattern

*Implement a façade or adapter layer between different subsystems that don't share the same semantics. This layer translates requests that one subsystem makes to the other subsystem. Use this pattern to ensure that an application's design is not limited by dependencies on outside subsystems.*
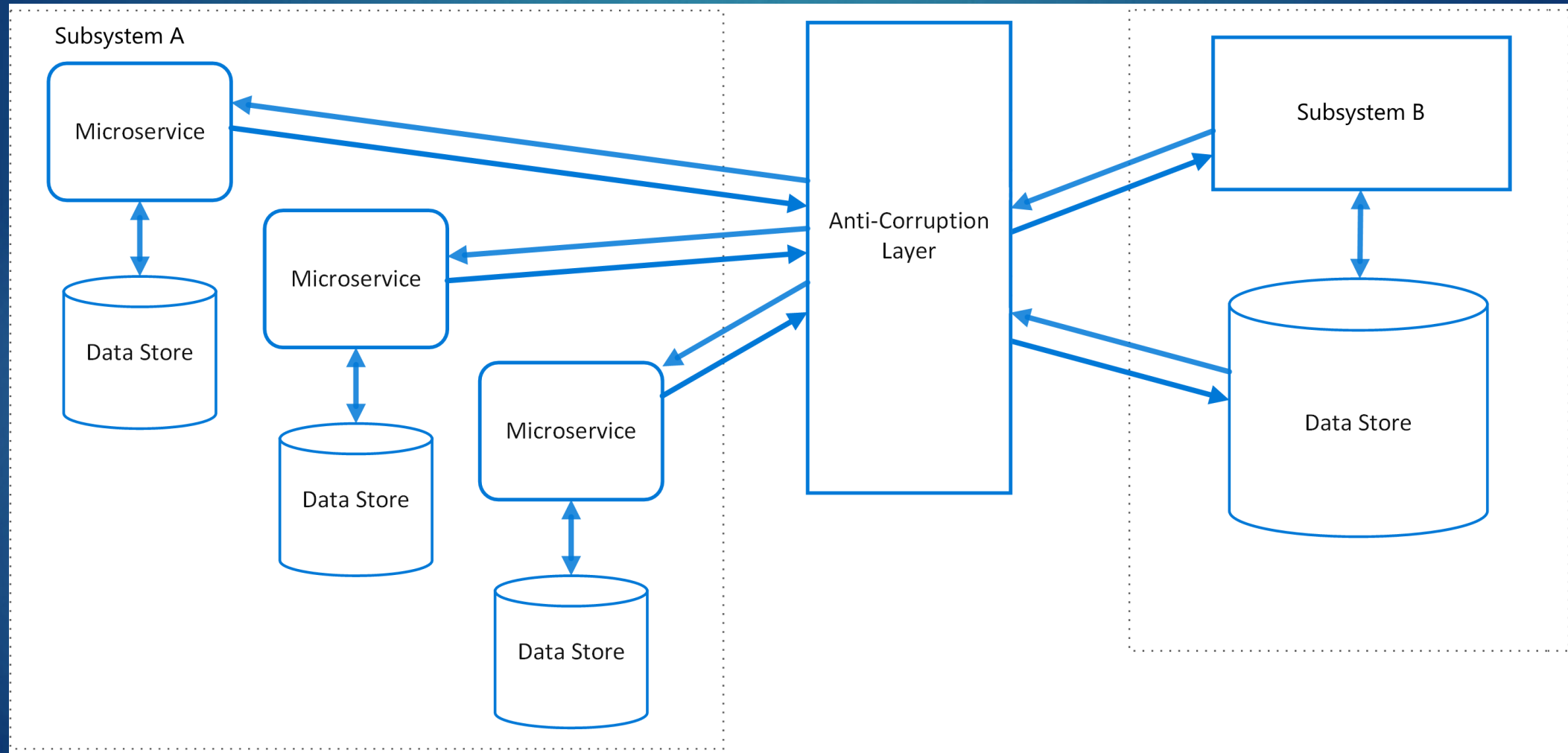
# Anti-Corruption Layer pattern

**Context and problem**

▶ Most applications rely on other systems for some data or functionality

▶ Often these legacy systems suffer from quality issues such as convoluted data schemas or obsolete APIs.

▶ The features and technologies used in legacy systems can vary widely from more modern systems.

▶ To interoperate with the legacy system, the new application may need to support outdated infrastructure, protocols, data models, APIs, or other features that you wouldn't otherwise put into a modern application.

▶ Maintaining access between new and legacy systems can force the new system to adhere to at least some of the legacy system's APIs or other semantics. When these legacy features have quality issues, supporting them "corrupts" what might otherwise be a cleanly designed modern application.

▶ Similar issues can arise with any external system that your development team doesn't control, not just legacy systems.

# Anti-Corruption Layer pattern

**Solution**

▶ Isolate the different subsystems by placing an anti-corruption layer between them.

▶ This layer translates communications between the two systems, allowing one system to remain unchanged while the other can avoid compromising its design and technological approach.

# Anti-Corruption Layer pattern

# Anti-Corruption Layer pattern

▶ **Issues and considerations**

- The anti-corruption layer may add latency to calls made between the two systems.

- The anti-corruption layer adds an additional service that must be managed and maintained.

- Consider how your anti-corruption layer will scale.

- Consider whether you need more than one anti-corruption layer. You may want to decompose functionality into multiple services using different technologies or languages, or there may be other reasons to partition the anti-corruption layer.

- Make sure transaction and data consistency are maintained and can be monitored.

- Consider whether the anti-corruption layer needs to handle all communication between different subsystems, or just a subset of features.

- If the anti-corruption layer is part of an application migration strategy, consider whether it will be permanent, or will be retired after all legacy functionality has been migrated.

# Anti-Corruption Layer pattern

**When to use this pattern**

▶ Use this pattern when:

▶ A migration is planned to happen over multiple stages, but integration between new and legacy systems needs to be maintained.

▶ Two or more subsystems have different semantics, but still need to communicate.

▶ This pattern may not be suitable if there are no significant semantic differences between new and legacy systems.
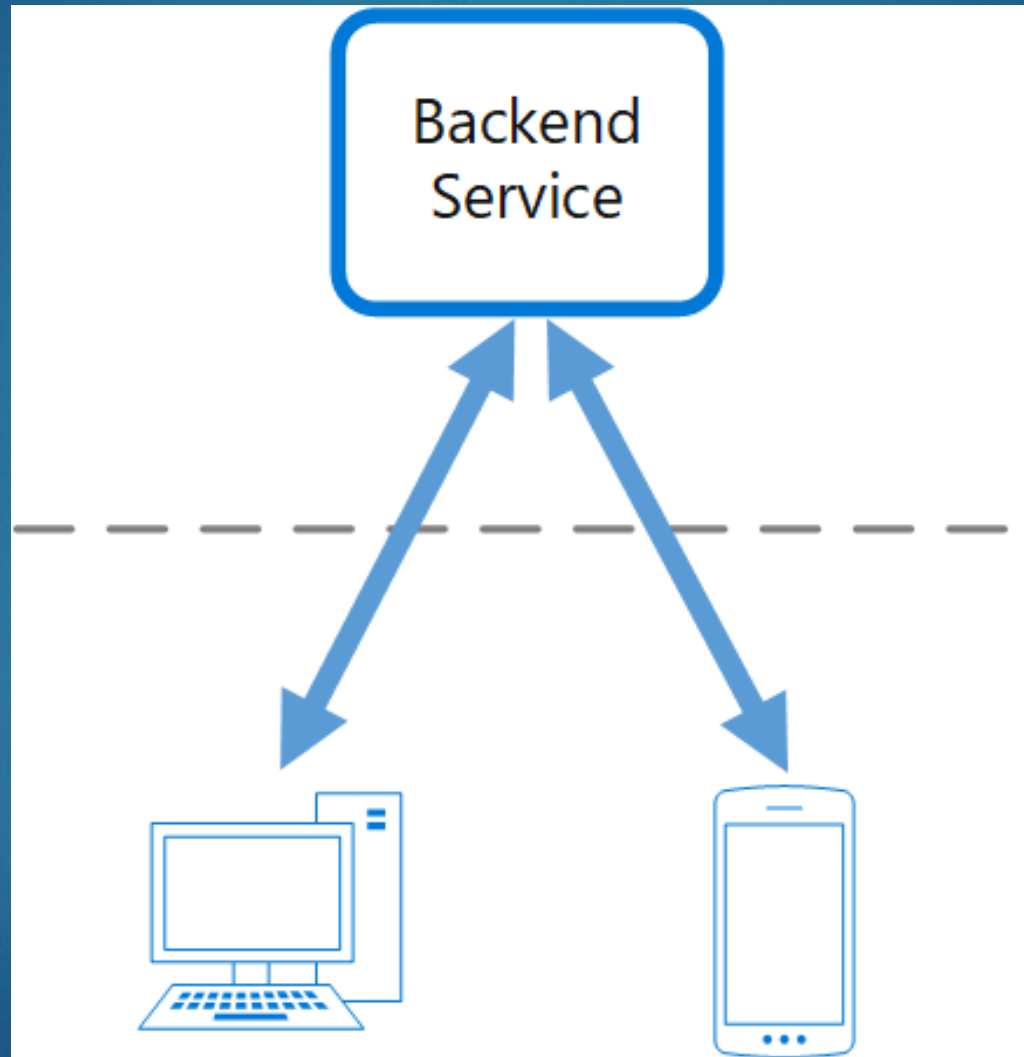
# Backends for Frontends pattern

*Create separate backend services to be consumed by specific frontend applications or interfaces. This pattern is useful when you want to avoid customizing a single backend for multiple interfaces.*

# Backends for Frontends pattern

**Context and problem**

▶ An application may initially be targeted at a desktop web UI. developed in parallel that provides the features needed for that UI. As the application's user base grows, a mobile application is developed that must interact with the same backend. The

▶ Typically, a backend service is backend service becomes a general-purpose backend, serving the requirements of both the desktop and mobile interfaces.

▶ But the capabilities of a mobile device differ significantly from a desktop browser, in terms of screen size, performance, and display limitations. As a result, the requirements for a mobile application backend differ from the desktop web UI.

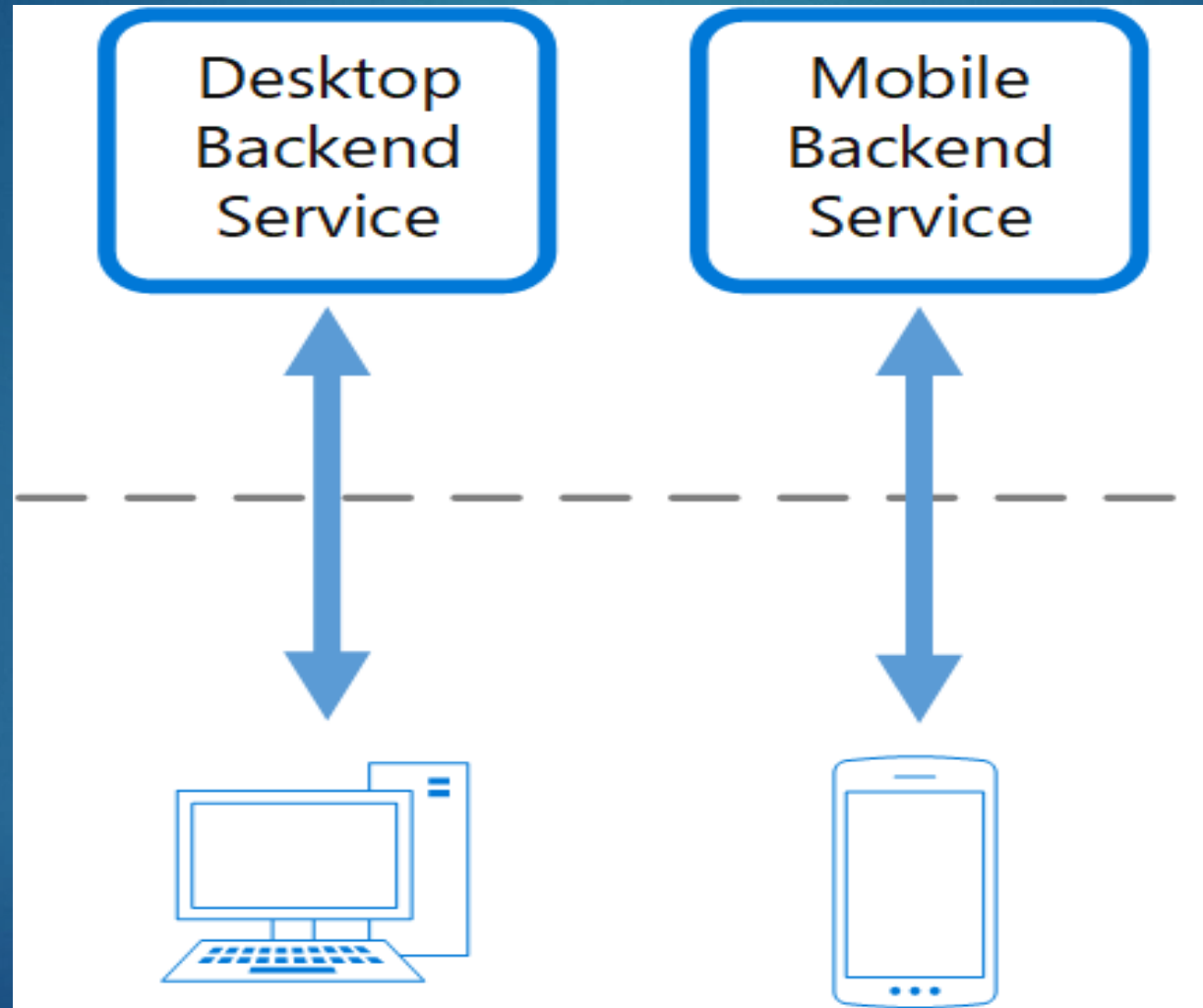# Backends for Frontends pattern

# Backends for Frontends pattern

**Solution**

▶ Create one backend per user interface. Fine-tune the behavior and performance of each backend to best match the needs of the frontend environment, without worrying about affecting other frontend experiences.

▶ Because each backend is specific to one interface, it can be optimized for that interface. As a result, it will be smaller, less complex, and likely faster than a generic backend that tries to satisfy the requirements for all interfaces. Each interface team has autonomy to control their own backend and doesn't rely on a centralized backend development team. This gives the interface team flexibility in language selection, release cadence, prioritization of workload, and feature integration in their backend.

# Backends for Frontends pattern

# Backends for Frontends pattern

**Issues and considerations**

▶ Consider how many backends to deploy.

▶ If different interfaces (such as mobile clients) will make the same requests, consider whether it is necessary to implement a backend for each interface, or if a single backend will suffice.

▶ Code duplication across services is highly likely when implementing this pattern.

▶ Frontend-focused backend services should only contain client-specific logic and behavior. General business logic and other global features should be managed elsewhere in your application.

▶ Think about how this pattern might be reflected in the responsibilities of a development team.

▶ Consider how long it will take to implement this pattern. Will the effort of building the new backends incur technical debt, while you continue to support the existing generic backend?

# Backends for Frontends pattern

**When to use this pattern**

▶ Use this pattern when:

▶ A shared or general purpose backend service must be maintained with significant development overhead.

▶ You want to optimize the backend for the requirements of specific client interfaces.

▶ Customizations are made to a general-purpose backend to accommodate multiple interfaces.

▶ An alternative language is better suited for the backend of a different user interface.

▶ This pattern may not be suitable:

▶ When interfaces make the same or similar requests to the backend.

▶ When only one interface is used to interact with the backend.

# Strangler Fig pattern

*Incrementally migrate a legacy system by gradually replacing specific pieces of functionality with new applications and services. As features from the legacy system are replaced, the new system eventually replaces all of the old system's features, strangling the old system and allowing you to decommission it.*

# Strangler Fig pattern

**Context and problem**

▶ As systems age, the development tools, hosting technology, and even system architectures they were built on can become increasingly obsolete. As new features and functionality are added, the complexity of these applications can increase dramatically, making them harder to maintain or add new features to.

▶ Completely replacing a complex system can be a huge undertaking. Often, you will need a gradual migration to a new system, while keeping the old system to handle features that haven't been migrated yet. However, running two separate versions of an application means that clients have to know where particular features are located. Every time a feature or service is migrated, clients need to be updated to point to the new location.
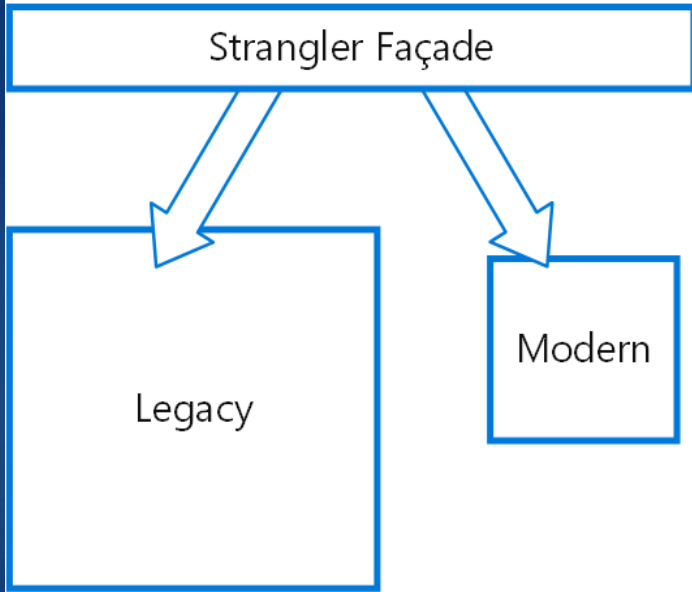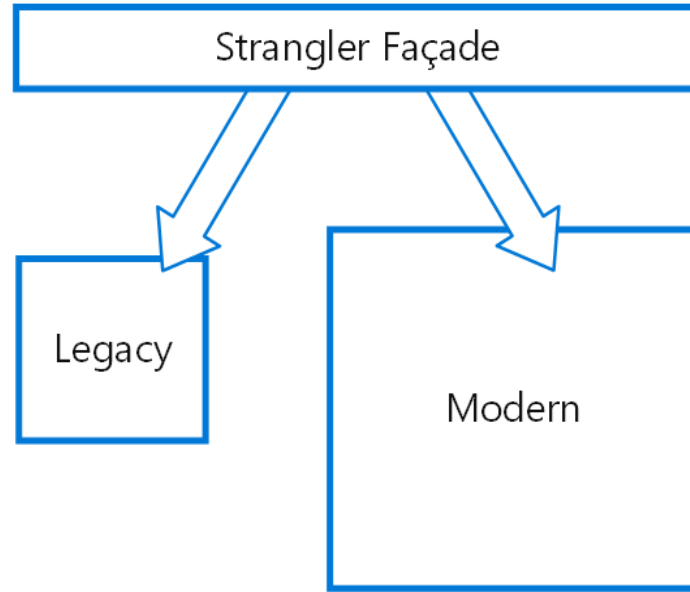
# Strangler Fig pattern

**Solution**

▶ Incrementally replace specific pieces of functionality with new applications and services. Create a façade that intercepts requests going to the backend legacy system.

▶ The façade routes these requests either to the legacy application or the new services.

▶ Existing features can be migrated to the new system gradually, and consumers can continue using the same interface, unaware that any migration has taken place.
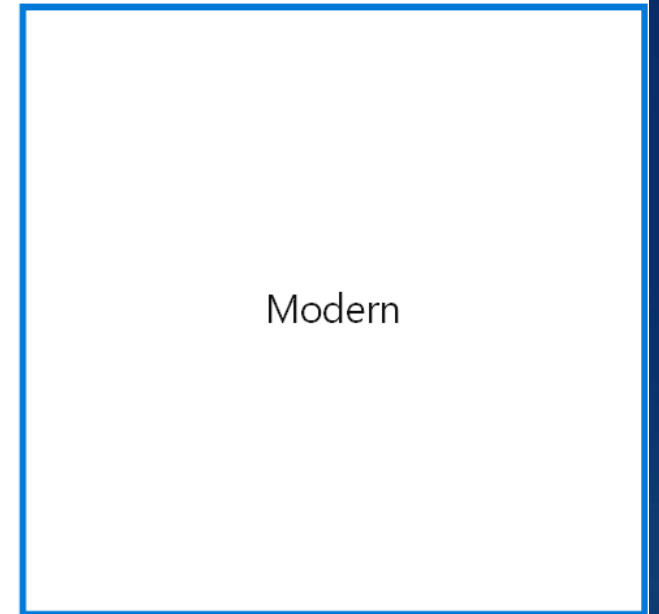
# Strangler Fig pattern

# Strangler Fig pattern

**Issues and considerations**

▶ Consider how to handle services and data stores that are potentially used by both new and legacy systems. Make sure both can access these resources side-by-side.

▶ Structure new applications and services in a way that they can easily be intercepted and replaced in future strangler fig migrations.

▶ At some point, when the migration is complete, the strangler fig façade will either go away or evolve into an adaptor for legacy clients.

▶ Make sure the façade keeps up with the migration.

▶ Make sure the façade doesn't become a single point of failure or a performance bottleneck.

# Strangler Fig pattern

**When to use this pattern**

▶ Use this pattern when gradually migrating a back-end application to a new architecture.

▶ This pattern may not be suitable:

▶ When requests to the back-end system cannot be intercepted.

▶ For smaller systems where the complexity of wholesale replacement is low.

# Sidecar Pattern

*Deploy components of an application into a separate process or container to provide isolation and encapsulation. This pattern can also enable applications to be composed of heterogeneous components and technologies.*

**This pattern is named Sidecar because it resembles a sidecar attached to a motorcycle**

# Sidecar Pattern

**Context and Problem**

▶ Applications and services often require related functionality, such as monitoring, logging, configuration, and networking services. These peripheral tasks can be implemented as separate components or services.

▶ If they are tightly integrated into the application, they can run in the same process as the application, making efficient use of shared resources. However, this also means they are not well isolated, and an outage in one of these components can affect other components or the entire application. Also, they usually need to be implemented using the same language as the parent application. As a result, the component and the application have close interdependence on each other.

▶ If the application is decomposed into services, then each service can be built using different languages and technologies. While this gives more flexibility, it means that each component has its own dependencies and requires language-specific libraries to access the underlying platform and any resources shared with the parent application. In addition, deploying these features as separate services can add latency to the application. Managing the code and dependencies for these language-specific interfaces can also add considerable complexity, especially for hosting, deployment, and management.
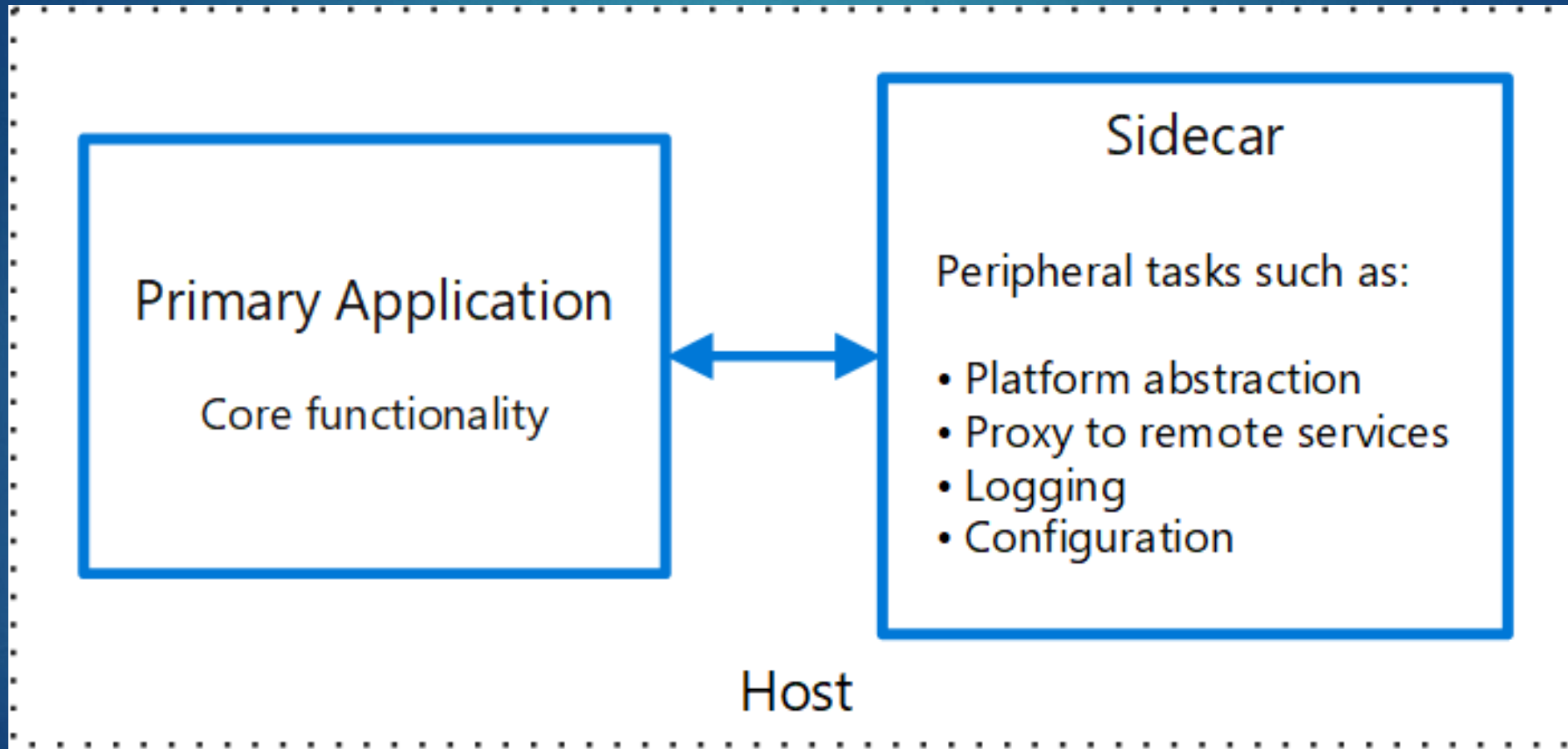
# Sidecar Pattern

**Context and Problem**

Applications and services often require related functionality, such as monitoring, logging, configuration, and networking services. These peripheral tasks can be implemented as separate components or services.

If they are tightly integrated into the application, they can run in the same process as the application, making efficient use of shared resources. However, this also means they are not well isolated, and an outage in one of these components can affect other components or the entire application. Also, they usually need to be implemented using the same language as the parent application. As a result, the component and the application have close interdependence on each other.

If the application is decomposed into services, then each service can be built using different languages and technologies. While this gives more flexibility, it means that each component has its own dependencies and requires language-specific libraries to access the underlying platform and any resources shared with the parent application. In addition, deploying these features as separate services can add latency to the application. Managing the code and dependencies for these language-specific interfaces can also add considerable complexity, especially for hosting, deployment, and management.

# Sidecar Pattern

**Solution**

# Sidecar Pattern

▶ A sidecar service is not necessarily part of the application, but is connected to it. It goes wherever the parent application goes. Sidecars are supporting processes or services that are deployed with the primary application. On a motorcycle, the sidecar is attached to one motorcycle, and each motorcycle can have its own sidecar. In the same way, a sidecar service shares the fate of its parent application. For each instance of the application, an instance of the sidecar is deployed and hosted alongside it.

▶ Advantages of using a sidecar pattern include:

• A sidecar is independent from its primary application in terms of runtime environment and programming language, so you don't need to develop one sidecar per language.

• The sidecar can access the same resources as the primary application. For example, a sidecar can monitor system resources used by both the sidecar and the primary application.

• Because of its proximity to the primary application, there's no significant latency when communicating between them.

• Even for applications that don't provide an extensibility mechanism, you can use a sidecar to extend functionality by attaching it as its own process in the same host or sub-container as the primary application.

▶ The sidecar pattern is often used with containers and referred to as a sidecar container or sidekick container.

# Sidecar Pattern

**Issues and Considerations**

▶ Consider the deployment and packaging format you will use to deploy services, processes, or containers. Containers are particularly well suited to the sidecar pattern.

▶ When designing a sidecar service, carefully decide on the interprocess communication mechanism. Try to use language- or framework-agnostic technologies unless performance requirements make that impractical.

▶ Before putting functionality into a sidecar, consider whether it would work better as a separate service or a more traditional daemon.

▶ Also consider whether the functionality could be implemented as a library or using a traditional extension mechanism. Language-specific libraries may have a deeper level of integration and less network overhead.

# Sidecar Pattern

**When to Use this Pattern**

▶ Use this pattern when:

• Your primary application uses a heterogeneous set of languages and frameworks. A component located in a sidecar service can be consumed by applications written in different languages using different frameworks.

• A component is owned by a remote team or a different organization.

• A component or feature must be co-located on the same host as the application

• You need a service that shares the overall lifecycle of your main application, but can be independently updated.

• You need fine-grained control over resource limits for a particular resource or component. For example, you may want to restrict the amount of memory a specific component uses. You can deploy the component as a sidecar and manage memory usage independently of the main application.

# Sidecar Pattern

**When to Use this Pattern**

▶ This pattern may not be suitable:

- When inter process communication needs to be optimized. Communication between a parent application and sidecar services includes some overhead, notably latency in the calls. This may not be an acceptable trade-off for chatty interfaces.

- For small applications where the resource cost of deploying a sidecar service for each instance is not worth the advantage of isolation.

- When the service needs to scale differently than or independently from the main applications. If so, it may be better to deploy the feature as a separate service.

# Sidecar Pattern

▶ **Example**

▶ The sidecar pattern is applicable to many scenarios. Some common examples:

o Infrastructure API. The infrastructure development team creates a service that's deployed alongside each application, instead of a language-specific client library to access the infrastructure. The service is loaded as a sidecar and provides a common layer for infrastructure services, including logging, environment data, configuration store, discovery, health checks, and watchdog services. The sidecar also monitors the parent application's host environment and process (or container) and logs the information to a centralized service.

o Manage NGINX/HAProxy. Deploy NGINX with a sidecar service that monitors environment state, then updates the NGINX configuration file and recycles the process when a change in state is needed.

o Ambassador sidecar. Deploy an ambassador service as a sidecar. The application calls through the ambassador, which handles request logging, routing, circuit breaking, and other connectivity related features.

o Offload proxy. Place an NGINX proxy in front of a node.js service instance, to handle serving static file content for the service.
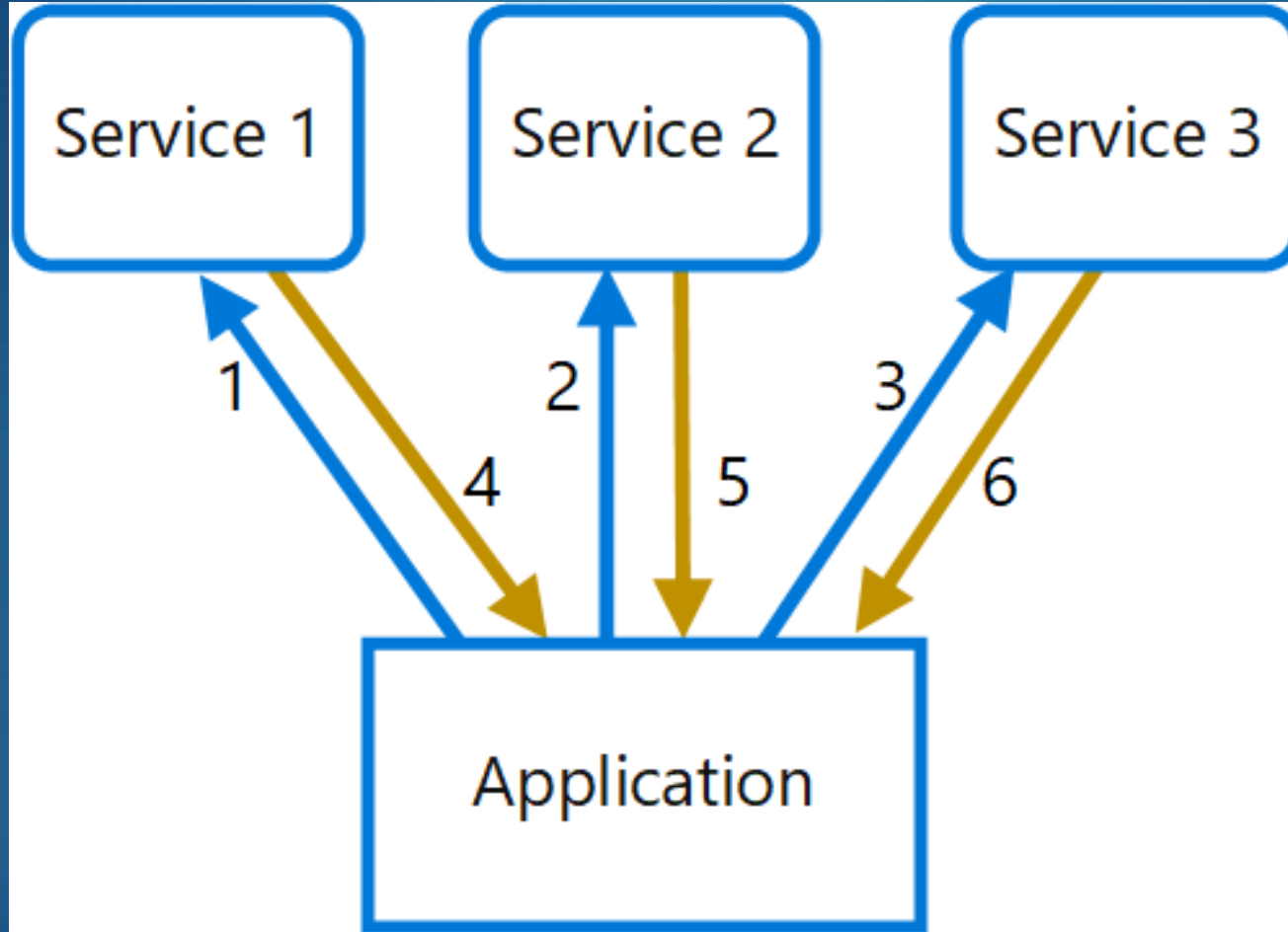
# Gateway Aggregation pattern

*Use a gateway to aggregate multiple individual requests into a single request. This pattern is useful when a client must make multiple calls to different backend systems to perform an operation.*

# Gateway Aggregation pattern

▶ **Context and problem**

▶ To perform a single task, a client may have to make multiple calls to various backend services. An application that relies on many services to perform a task must expend resources on each request. When any new feature or service is added to the application, additional requests are needed, further increasing resource requirements and network calls. This chattiness between a client and a backend can adversely impact the performance and scale of the application. Microservice architectures have made this problem more common, as applications built around many smaller services naturally have a higher amount of cross-service calls.

▶ In the following diagram, the client sends requests to each service (1,2,3). Each service processes the request and sends the response back to the application (4,5,6). Over a cellular network with typically high latency, using individual requests in this manner is inefficient and could result in broken connectivity or incomplete requests. While each request may be done in parallel, the application must send, wait, and process data for each request, all on separate connections, increasing the chance of failure.
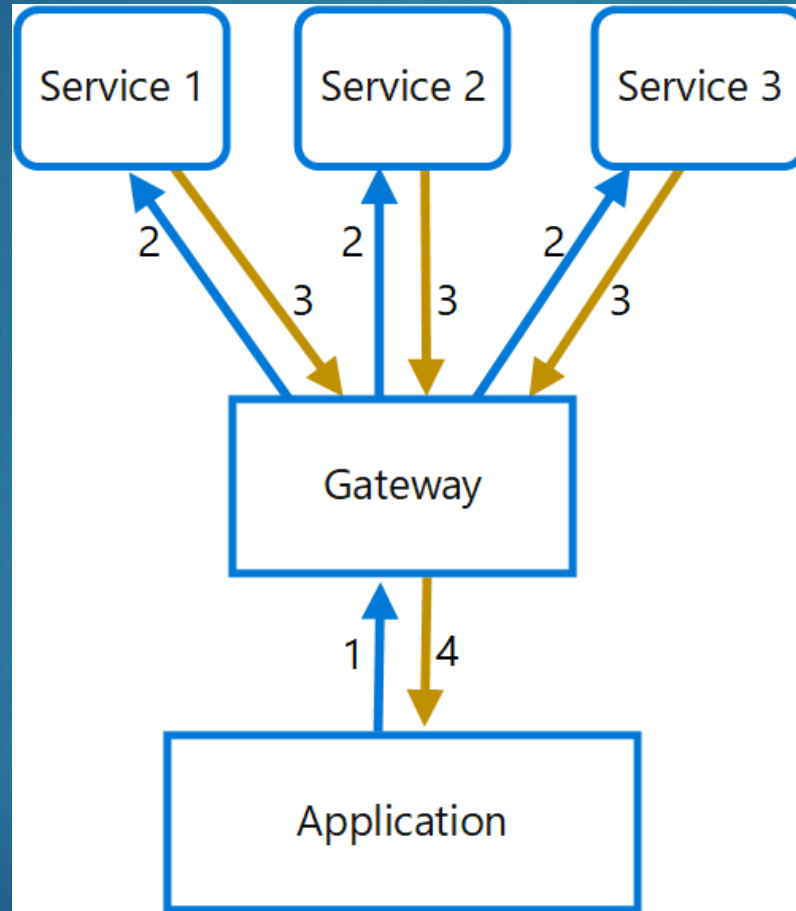
# Gateway Aggregation pattern

# Gateway Aggregation pattern

**Solution**

▶ Use a gateway to reduce chattiness between the client and the services. The gateway receives client requests, dispatches requests to the various backend systems, and then aggregates the results and sends them back to the requesting client.

▶ This pattern can reduce the number of requests that the application makes to backend services, and improve application performance over high-latency networks.

▶ In the following diagram, the application sends a request to the gateway (1). The request contains a package of additional requests. The gateway decomposes these and processes each request by sending it to the relevant service (2). Each service returns a response to the gateway (3). The gateway combines the responses from each service and sends the response to the application (4). The application makes a single request and receives only a single response from the gateway.

# Gateway Aggregation pattern

# Gateway Aggregation pattern

**Issues and considerations**

▶ The gateway should not introduce service coupling across the backend services.

▶ The gateway should be located near the backend services to reduce latency as much as possible.

▶ The gateway service may introduce a single point of failure. Ensure the gateway is properly designed to meet your application's availability requirements.

▶ The gateway may introduce a bottleneck. Ensure the gateway has adequate performance to handle load and can be scaled to meet your anticipated growth.

▶ Perform load testing against the gateway to ensure you don't introduce cascading failures for services.

▶ Implement a resilient design, using techniques such as <u>bulkheads</u>, <u>circuit breaking</u>, <u>retry</u>, and timeouts.

▶ If one or more service calls takes too long, it may be acceptable to timeout and return a partial set of data. Consider how your application will handle this scenario.

▶ Use asynchronous I/O to ensure that a delay at the backend doesn't cause performance issues in the application.

▶ Implement distributed tracing using correlation IDs to track each individual call.

▶ Monitor request metrics and response sizes.

▶ Consider returning cached data as a failover strategy to handle failures.

▶ Instead of building aggregation into the gateway, consider placing an aggregation service behind the gateway. Request aggregation will likely have different resource requirements than other services in the gateway and may impact the gateway's routing and offloading functionality.

# Gateway Aggregation pattern

**When to use this pattern**

▶ Use this pattern when:

• A client needs to communicate with multiple backend services to perform an operation.

• The client may use networks with significant latency, such as cellular networks.

▶ This pattern may not be suitable when:

• You want to reduce the number of calls between a client and a single service across multiple operations. In that scenario, it may be better to add a batch operation to the service.

• The client or application is located near the backend services and latency is not a significant factor.

# Gateway Offloading pattern

*Offload shared or specialized service functionality to a gateway proxy. This pattern can simplify application development by moving shared service functionality, such as the use of SSL certificates, from other parts of the application into the gateway.*
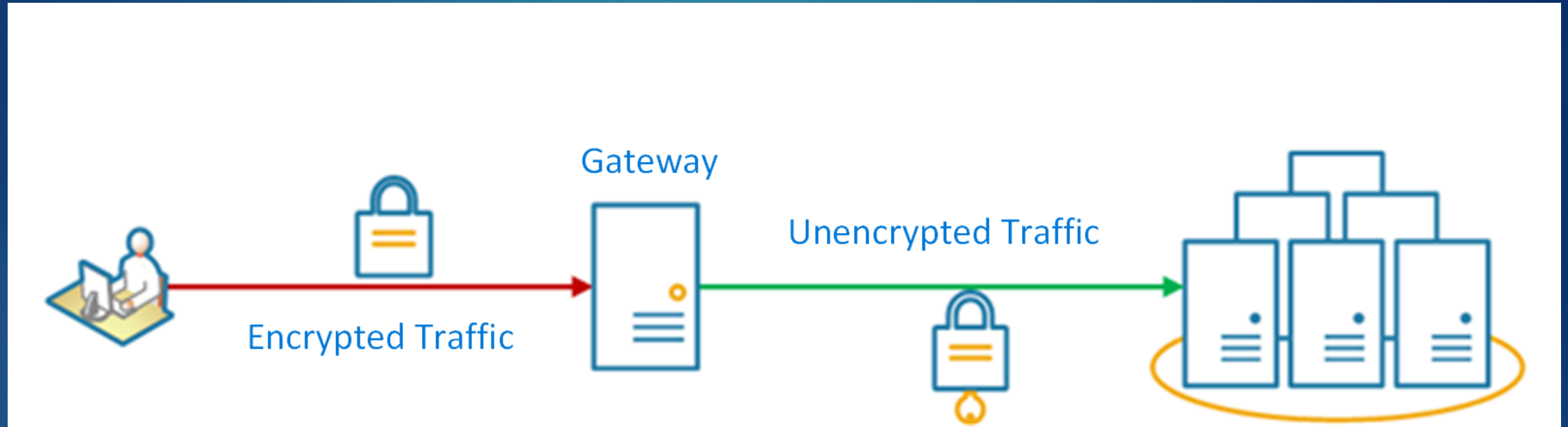
# Gateway Offloading pattern

**Context and problem**

▶ Some features are commonly used across multiple services, and these features require configuration, management, and maintenance. A shared or specialized service that is distributed with every application deployment increases the administrative overhead and increases the likelihood of deployment error. Any updates to a shared feature must be deployed across all services that share that feature.

▶ Properly handling security issues (token validation, encryption, SSL certificate management) and other complex tasks can require team members to have highly specialized skills. For example, a certificate needed by an application must be configured and deployed on all application instances. With each new deployment, the certificate must be managed to ensure that it does not expire. Any common certificate that is due to expire must be updated, tested, and verified on every application deployment.

▶ Other common services such as authentication, authorization, logging, monitoring, or throttling can be difficult to implement and manage across a large number of deployments. It may be better to consolidate this type of functionality, in order to reduce overhead and the chance of errors.

# Gateway Offloading pattern

**Solution**

- ▶ Offload some features into a gateway, particularly cross-cutting concerns such as certificate management, authentication, SSL termination, monitoring, protocol translation, or throttling.

- ▶ The following diagram shows a gateway that terminates inbound SSL connections. It requests data on behalf of the original requestor from any HTTP server upstream of the gateway.

# Gateway Offloading pattern

# Gateway Offloading pattern

Benefits of this pattern include:

- Simplify the development of services by removing the need to distribute and maintain supporting resources, such as web server certificates and configuration for secure websites. Simpler configuration results in easier management and scalability and makes service upgrades simpler.

- Allow dedicated teams to implement features that require specialized expertise, such as security. This allows your core team to focus on the application functionality, leaving these specialized but cross-cutting concerns to the relevant experts.

- Provide some consistency for request and response logging and monitoring. Even if a service is not correctly instrumented, the gateway can be configured to ensure a minimum level of monitoring and logging.

# Gateway Offloading pattern

**Issues and considerations**

- Ensure the gateway is highly available and resilient to failure. Avoid single points of failure by running multiple instances of your gateway.

- Ensure the gateway is designed for the capacity and scaling requirements of your application and endpoints. Make sure the gateway does not become a bottleneck for the application and is sufficiently scalable.

- Only offload features that are used by the entire application, such as security or data transfer.

- Business logic should never be offloaded to the gateway.

- If you need to track transactions, consider generating correlation IDs for logging purposes.

# Gateway Offloading pattern

▶ **When to use this pattern**

▶ Use this pattern when:

- An application deployment has a shared concern such as SSL certificates or encryption.

- A feature that is common across application deployments that may have different resource requirements, such as memory resources, storage capacity or network connections.

- You wish to move the responsibility for issues such as network security, throttling, or other network boundary concerns to a more specialized team.

▶ This pattern may not be suitable if it introduces coupling across services.
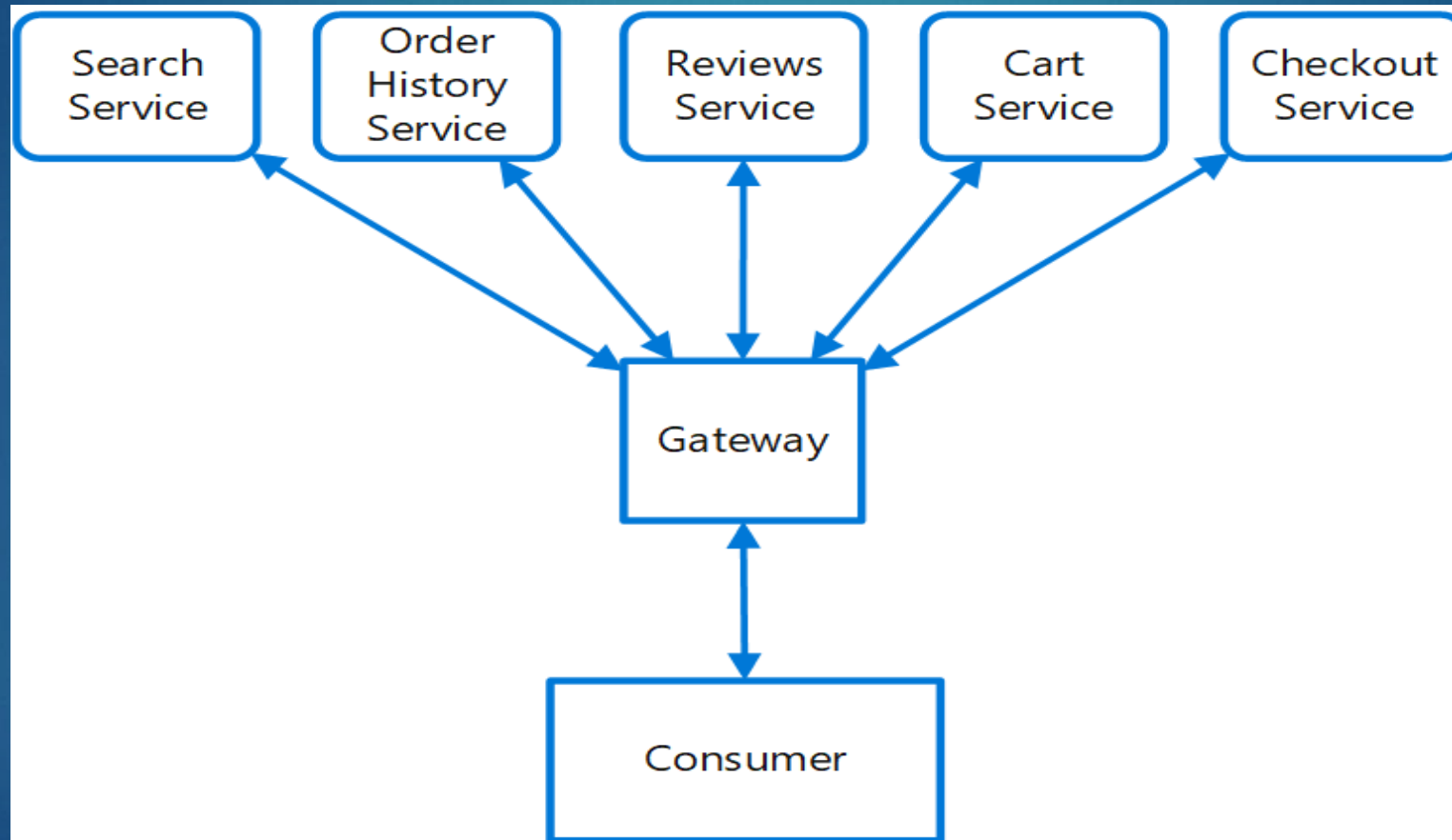
# Gateway Routing pattern

*Route requests to multiple services using a single endpoint. This pattern is useful when you wish to expose multiple services on a single endpoint and route to the appropriate service based on the request.*

# Gateway Routing pattern

▶ **Context and problem**

▶ When a client needs to consume multiple services, setting up a separate endpoint for each service and having the client manage each endpoint can be challenging. For example, an e-commerce application might provide services such as search, reviews, cart, checkout, and order history. Each service has a different API that the client must interact with, and the client must know about each endpoint in order to connect to the services. If an API changes, the client must be updated as well. If you refactor a service into two or more separate services, the code must change in both the service and the client.

# Gateway Routing pattern

# Gateway Routing pattern

▶ Place a gateway in front of a set of applications, services, or deployments. Use application Layer 7 routing to route the request to the appropriate instances.

▶ With this pattern, the client application only needs to know about and communicate with a single endpoint. If a service is consolidated or decomposed, the client does not necessarily require updating. It can continue making requests to the gateway, and only the routing changes.

▶ A gateway also lets you abstract backend services from the clients, allowing you to keep client calls simple while enabling changes in the backend services behind the gateway. Client calls can be routed to whatever service or services need to handle the expected client behavior, allowing you to add, split, and reorganize services behind the gateway without changing the client.

# Gateway Routing pattern

▶ This pattern can also help with deployment, by allowing you to manage how updates are rolled out to users. When a new version of your service is deployed, it can be deployed in parallel with the existing version. Routing lets you control what version of the service is presented to the clients, giving you the flexibility to use various release strategies, whether incremental, parallel, or complete rollouts of updates. Any issues discovered after the new service is deployed can be quickly reverted by making a configuration change at the gateway, without affecting clients.

# Gateway Routing pattern

**Issues and considerations**

- The gateway service may introduce a single point of failure. Ensure it is properly designed to meet your availability requirements. Consider resiliency and fault tolerance capabilities when implementing.

- The gateway service may introduce a bottleneck. Ensure the gateway has adequate performance to handle load and can easily scale in line with your growth expectations.

- Perform load testing against the gateway to ensure you don't introduce cascading failures for services.

- Gateway routing is level 7. It can be based on IP, port, header, or URL.

# Gateway Routing pattern

▶ **When to use this pattern**

▶ Use this pattern when:

- A client needs to consume multiple services that can be accessed behind a gateway.

- You wish to simplify client applications by using a single endpoint.

- You need to route requests from externally addressable endpoints to internal virtual endpoints, such as exposing ports on a VM to cluster virtual IP addresses.

▶ This pattern may not be suitable when you have a simple application that uses only one or two services.

# Bulkhead pattern

The Bulkhead pattern is a type of application design that is tolerant of failure. In a bulkhead architecture, elements of an application are isolated into pools so that if one fails, the others will continue to function.
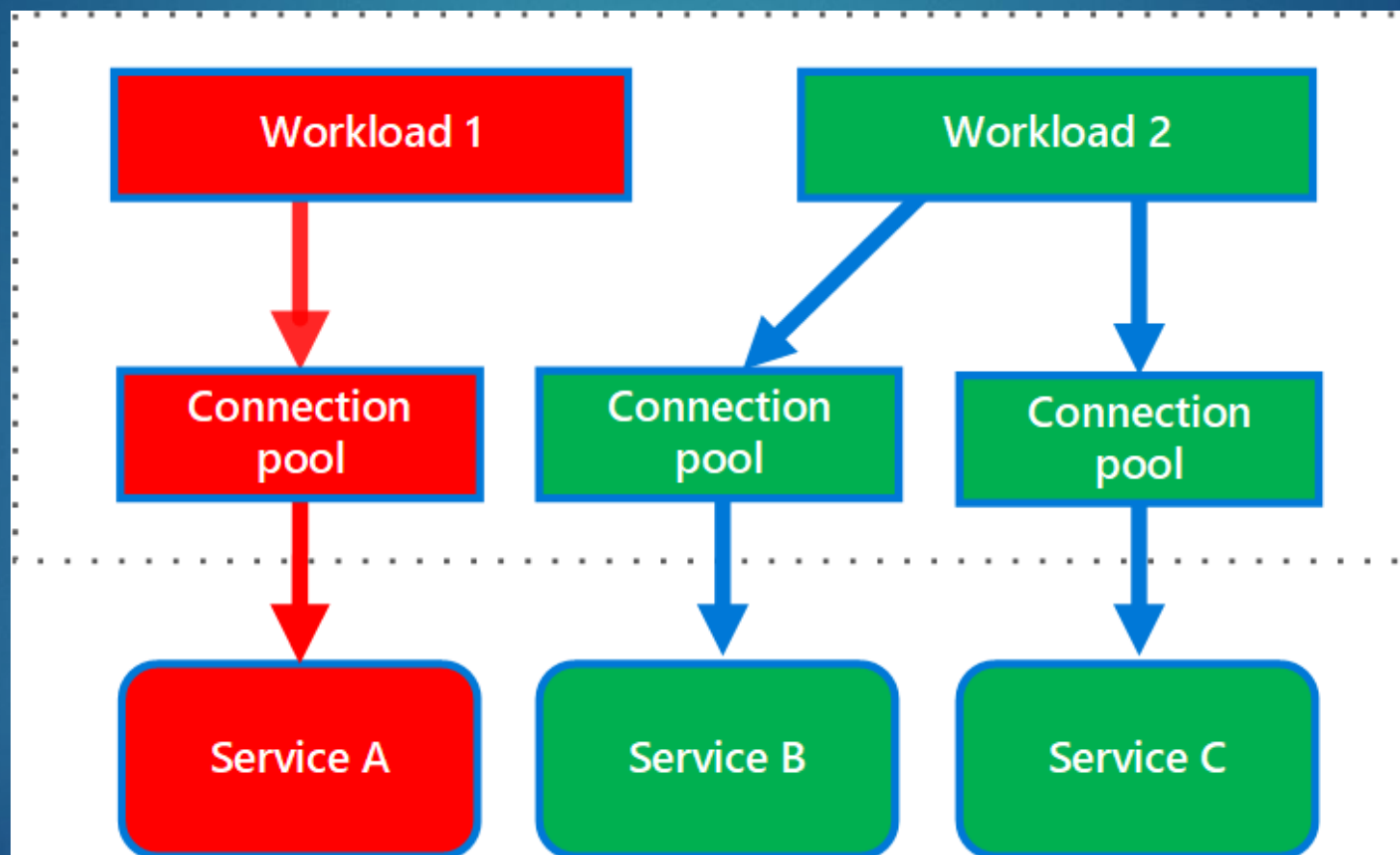
It's named after the sectioned partitions (bulkheads) of a ship's hull. If the hull of a ship is compromised, only the damaged section fills with water, which prevents the ship from sinking.

# Bulkhead pattern

**Context and problem**

▶ A cloud-based application may include multiple services, with each service having one or more consumers. Excessive load or failure in a service will impact all consumers of the service.

▶ Moreover, a consumer may send requests to multiple services simultaneously, using resources for each request. When the consumer sends a request to a service that is misconfigured or not responding, the resources used by the client's request may not be freed in a timely manner. As requests to the service continue, those resources may be exhausted. For example, the client's connection pool may be exhausted. At that point, requests by the consumer to other services are affected. Eventually the consumer can no longer send requests to other services, not just the original unresponsive service.

▶ The same issue of resource exhaustion affects services with multiple consumers. A large number of requests originating from one client may exhaust available resources in the service. Other consumers are no longer able to consume the service, causing a cascading failure effect.
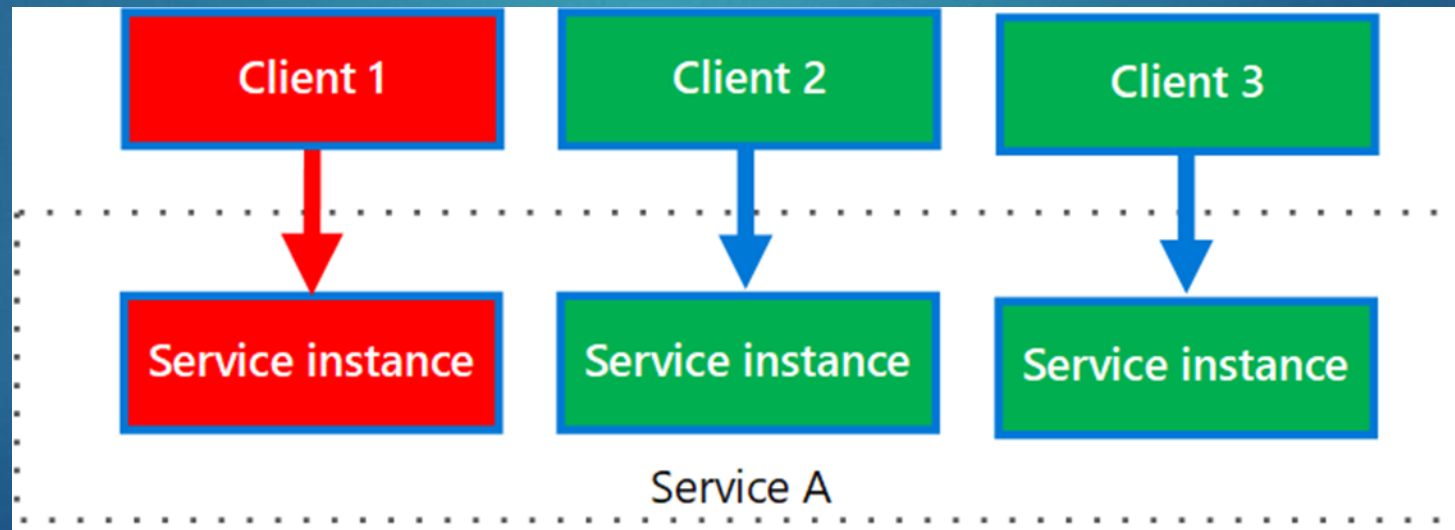
# Bulkhead pattern

- **Solution**

- Partition service instances into different groups, based on consumer load and availability requirements. This design helps to isolate failures, and allows you to sustain service functionality for some consumers, even during a failure.

- A consumer can also partition resources, to ensure that resources used to call one service don't affect the resources used to call another service. For example, a consumer that calls multiple services may be assigned a connection pool for each service. If a service begins to fail, it only affects the connection pool assigned for that service, allowing the consumer to continue using the other services.

- The benefits of this pattern include:

  - Isolates consumers and services from cascading failures. An issue affecting a consumer or service can be isolated within its own bulkhead, preventing the entire solution from failing.

  - Allows you to preserve some functionality in the event of a service failure. Other services and features of the application will continue to work.

  - Allows you to deploy services that offer a different quality of service for consuming applications. A high-priority consumer pool can be configured to use high-priority services.

- The following diagram shows bulkheads structured around connection pools that call individual services. If Service A fails or causes some other issue, the connection pool is isolated, so only workloads using the thread pool assigned to Service A are affected. Workloads that use Service B and C are not affected and can continue working without interruption.

# Bulkhead pattern

▶ The next diagram shows multiple clients calling a single service. Each client is assigned a separate service instance. Client 1 has made too many requests and overwhelmed its instance. Because each service instance is isolated from the others, the other clients can continue making calls.

# Bulkhead pattern

**Issues and considerations**

- Define partitions around the business and technical requirements of the application.

- When partitioning services or consumers into bulkheads, consider the level of isolation offered by the technology as well as the overhead in terms of cost, performance and manageability.

- Consider combining bulkheads with retry, circuit breaker, and throttling patterns to provide more sophisticated fault handling.

- When partitioning consumers into bulkheads, consider using processes, thread pools, and semaphores. Projects like resilience4j and Polly offer a framework for creating consumer bulkheads.

- When partitioning services into bulkheads, consider deploying them into separate virtual machines, containers, or processes. Containers offer a good balance of resource isolation with fairly low overhead.

- Services that communicate using asynchronous messages can be isolated through different sets of queues. Each queue can have a dedicated set of instances processing messages on the queue, or a single group of instances using an algorithm to dequeue and dispatch processing.

- Determine the level of granularity for the bulkheads. For example, if you want to distribute tenants across partitions, you could place each tenant into a separate partition, or put several tenants into one partition.

- Monitor each partition's performance and SLA.

# Bulkhead pattern

**When to use this pattern**

▶ Use this pattern to:

• Isolate resources used to consume a set of backend services, especially if the application can provide some level of functionality even when one of the services is not responding.

• Isolate critical consumers from standard consumers.

• Protect the application from cascading failures.

▶ This pattern may not be suitable when:

• Less efficient use of resources may not be acceptable in the project.

• The added complexity is not necessary

# References:

https://docs.microsoft.com/en-us/azure/architecture/patterns/

Thank You