

Microservice Scaling Cube

X-axis scaling

X-axis scaling consists of running multiple copies of an application behind a load balancer. If there are N copies then each copy handles 1/N of the load. This is a simple, commonly used approach of scaling an application. One drawback of this approach is that because each copy potentially accesses all of the data, caches require more memory to be effective. Another problem with this approach is that it does not tackle the problems of increasing development and application complexity.

Y-axis scaling

Unlike X-axis and Z-axis, which consist of running multiple, identical copies of the application, Y-axis axis scaling splits the application into multiple, different services. Each service is responsible for one or more closely related functions. There are a couple of different ways of decomposing the application into services. One approach is to use verb-based decomposition and define services that implement a single use case such as checkout. The other option is to decompose the application by noun and create services responsible for all operations related to a particular entity such as customer management. An application might use a combination of verb-based and noun-based decomposition.

Z-axis scaling

When using Z-axis scaling each server runs an identical copy of the code. In this respect, it's similar to X-axis scaling. The big difference is that each server is responsible for only a subset of the data. Some component of the system is responsible for routing each request to the appropriate server. One commonly used routing criteria is an attribute of the request such as the primary key of the entity being accessed. Another common routing criteria is the customer type. For example, an application might provide paying customers with a higher SLA than free customers by routing their requests to a different set of servers with more capacity. Data partition across servers.

External vs Internal Quality Attributes

Internal

- Not directly observed when software is executing
- Perceived by the developers and maintainers
- Encompasses aspects of design that may impact external attributes

Availability	Performance	Efficiency	Scalability	Robustness
Safety	Security	Reliability	Integrity	Verifiability
Deployability	Compatibility	Installability	Portability	Maintainability
Usability	Testability	Modifiability	Reusability	Interoperability

Synchronous communication :

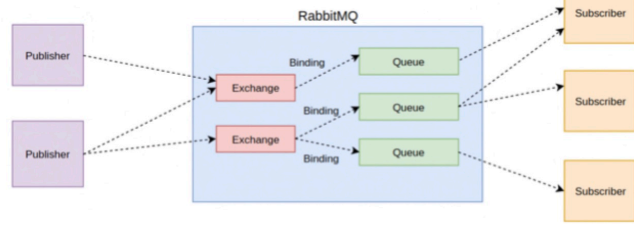
- the caller sends a message and waits for the receiver to respond eg in login action in which the caller must have a reply.
- sender sends a request to the receiver and receiver sends a response

e.g. HTTP/S protocol for synchronous communication

Asynchronous communication :

- a communication that has a lag between when a message is sent and when the receiver receives/interprets/responds to it.
- the caller skips the wait and continues executing whatever code is necessary.
- It enables independent functioning of sender and receiver
- It also enables one-to-many communication where sender can send request to multiple receivers at once.

e.g. AMQP for asynchronous communication.
AMQP: Advanced Message Queue Protocol



Creational Patterns

- Creational patterns help designers handling issues with creation of objects.
- e.g. Factory pattern

Structural Patterns

- Structural Patterns are used to provide a structure to the relationship between objects.
- e.g. Façade pattern

Behavioral Patterns

- Behavioral Patterns help define how objects interact with each other to deliver a task.
- e.g. Observer pattern

RabbitMQ: RabbitMQ is a message broker that implements the Advanced Message Queuing Protocol (AMQP). It excels at asynchronous message communication and supports various messaging patterns.

ActiveMQ: ActiveMQ is also a message broker, but it implements the Java Message Service (JMS) API. It is well-suited for Java applications.

ActiveMQ provides support for both point-to-point and publish-subscribe messaging patterns.

ActiveMQ also supports message persistence to disk, which ensures that messages are not lost even in the event of a system failure.

Redis: Redis is primarily an in-memory data store, but it can also be used for message queuing via its pub-sub (publish-subscribe) mechanism.

While Redis is fast and lightweight, it may not be as durable as RabbitMQ or ActiveMQ since it relies on memory for message storage.

Since Redis is primarily relies on in-memory storage, it makes it extremely fast but less durable. It can be configured to periodically save data to disk, but this introduces latency and potential data loss in case of a crash.

Redis is simple and easy to use due

Feature	RabbitMQ	ActiveMQ	Redis
Messaging Patterns	Point-to-point, Publish/Subscribe, Request/Reply	Point-to-point, Publish/Subscribe, Request/Reply	Publish/Subscribe
Message Durability	Supports persistent and non-persistent messages	Supports both persistent and non-persistent messages	Not persistent by default, but can be configured
Clustering	Yes	Yes	Yes
High Availability	Yes	Yes	Yes
Message Acknowledgment	Supported	Supported	Not applicable (messages are received immediately)
Message Routing	Routing based on message attributes	Flexible routing options, including virtual topics	Messages sent to specific channels/subscribers
Supported Protocols	AMQP, MQTT, STOMP, HTTP, and more	JMS, MQTT, STOMP, AMQP, OpenWire	Pub/Sub via Redis protocol
Language Support	Multiple client libraries for various languages	Primarily Java clients, but supports other languages	Multiple client libraries for various languages
Performance	High throughput and low latency	High performance, optimized for Java applications	Extremely fast due to in-memory processing
Scalability	Highly scalable, supports large-scale deployments	Scalable, suitable for both small and large setups	Scalable, suitable for real-time applications
Speed	~10k msg/sec	~10k msg/sec	~10m msg/sec

Persistent Asynchronous

- Sender keeps executing without blocking.
- The message may take an arbitrary amount of time to reach the receiver.
- Sender may or may not be running by the time the message reaches receiver.
- Disk or multiple memory queues could be used for persistence at receiver's side.
- There is a guarantee that the message will eventually reach the receiver.

Persistent Synchronous

- Sender is blocked until an ack for receipt(not for delivery/response) is received.
- The message persists -stays in receiver's queue (or in any router along the way) for an arbitrary amount of time.

Transient Asynchronous

- Sender continues execution (nonblocking) after sending a message.
- Receiver has to be running, because if it is not running the message will be discarded.
- Even if any router along the way is down, the message will be discarded.

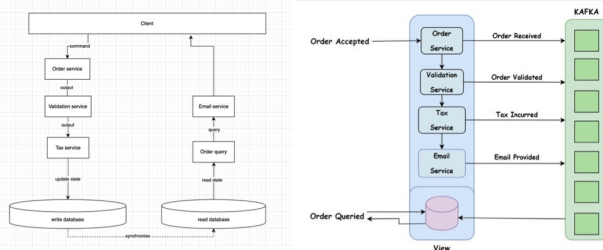
Transient Synchronous

- Receipt-based**
 - Sender blocks (synchronous) until an ack is received. This ack is simply a receipt, it does not tell us anything about whether the process has started on receiver end.
- Delivery-based**
 - Sender will block until receiver takes the delivery of the message.
 - The ack comes a little bit later than that in receipt-based.
 - Essentially an asynchronous RPC.
- Response-based**
 - Sender blocks until it receives a response.
 - This is traditional RPC. The sender is blocked for the entire duration the reply comes

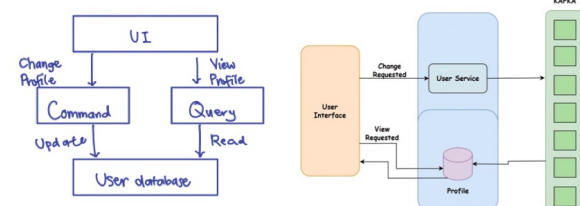
PART 1: Q1. same as Tutorial 6 Question 4: Event-driven service communication and CQRS:

- Illustrate (using a simple diagram) the service communication and use of CQRS in the following scenario:

- An order is accepted in the order service. This is picked up by the validation service, where it is validated. Sales tax is added (Tax service). Information for an email is made available(provided) to an Email service. The updated order goes back to the orders service, where it can be queried via the orders view (Hint: CQRS). After being sent a confirmation email, the user can click through to the order.



Illustrate use of CQRS pattern for the PeerPrep User Service for use cases of the user making a change in their profile at the UI and the user viewing their profile in the UI.



General design principles --

SoC, SRP, High Cohesion, Loose Coupling, Abstraction, Encapsulation, Information Hiding, Interface Segregation principle(ISP), Dependency Inversion, Open-close principle, Design for reuse

Principles of microservice design - Design with ubiquitous language, loose coupling, domain specific requirements/ around business capabilities, small team structures, containerised deployments, scaling components at different pace, decentralized development and deployment

Principles in Event driven and messaging systems - asynchronous communication, decoupling(or loose coupling)

Messaging Patterns

Message Construction. A message contains:

- Header:** Contains metadata. Can contain message intent command message, which tells the receiver what to do; document message, which sends data but does not really tell the receiver what to do; event message, which simply notifies the receiver about a change.
- Properties:** Optional. For message selection and filtering. Three kinds — application-related, providerrelated and standard properties.
- Payload:** Body, data structures.

Message Channels. Connect collaborating senders and receivers.

A channel transmits one way. Two-way messages need two channels. Some concepts:

- Return Address:** Request contains an address to tell the replier where to send reply to.
- Correlation ID:** Specifies which request the reply is for. Can be chained.
- Message Sequence:** Basically break the message down into smaller segments.
- Point to Point (P2P):** Request processed by single consumer.
- Publish-Subscribe Channel:** Request broadcasted to all interested parties via topics.
- Invalid Messages:** Queue to move a message to when cannot be interpreted.
- Dead Letter:** Queue to move a message to when it cannot be delivered.
- Datatype Channel:** Separate channel for each type of data, e.g. XML, byte array, etc.

Message Routing. Consumes messages from one channel and reinserts them into different channels based on condition

Some concepts:

- Content-Based:** Examines message content and routes. Message filter is a special kind of content-based router that discards messages based on content.
- Context-Based:** Decides destination based on context, e.g. load-balancing.
- Message Splitter:** Splits a single message into multiple.
- Message Aggregator:** Aggregates correlated messages into a single message.
- Scatter-Gather:** Broadcast to multiple participants and aggregates replies into a single message.

Message Transformation. Transform application-layer data structures, data types of fields, data representations

e.g. ASCII to Unicode, transport protocol, e.g. TCP/IP to sockets. Also called Message Translators or Channel Adapters. A Canonical Data Model may also be used, which is an "adapter" superset model.

Message Endpoints. Interface between app and messaging system. Channel-specific, one instance handles either send or receive. Can be synchronous, i.e. polling consumer, or can be asynchronous, i.e. event-driven receiver.

Design principles provide high level guidelines to design better software applications.

- They do not provide implementation guidelines and are not bound to any programming language.

-Modularity and the SOLID (SRP, OCP, LSP, ISP, DIP) principles are general set of design principles. There are specific principles for specific design style eg for Microservices or Event driven design approach.

Design Pattern provides low-level solutions related to implementation, of commonly occurring problems.

- Design pattern suggests a specific implementation for a specific problem.

- Eg. Creating a class that can only have 1 object at a time => singleton design pattern