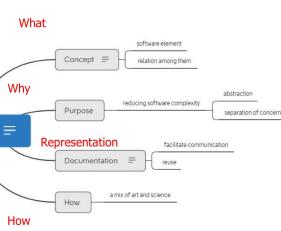
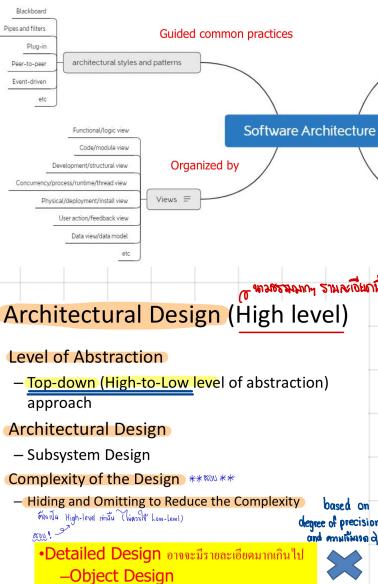


Week 1

Introduction

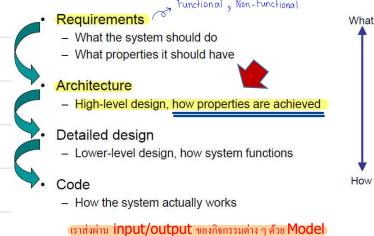


Software Architecture
as
an Architectural Design

Related to Abstraction!



Where Architecture Fits



Architectural Design (High level)

- Level of Abstraction**
 - Top-down (High-to-Low level of abstraction) approach
- Architectural Design**
 - Subsystem Design
- Complexity of the Design** ★★★★★
 - Hiding and Omitting to Reduce the Complexity (High level info / Low-level details)
 - Detailed Design (Object Design)
 - Object Design

Model គឺគេងនៅថ្មីថាបានការស្ថិតិស្សន៍
ពីទៅរួមចំណែកអនុវត្តសារការណា
(រួចរាល់ខ្លួនទាំងឡាយនៃសេដក)

ការណាំនៃការគូនឲ្យក្នុង

- Model គឺខ្លួន → ពីការណាំនៃការគូនឲ្យក្នុង
- Model សាកស្អួលបាន → more model ជាពីរនាក់
- Informal model → Semi formal model
 - Formal model គឺខ្លួន
 - Regular Expression (ទិន្នន័យការណាំនៃការគូនឲ្យក្នុង)

និរន្តរភាព Software Architecture (ស្រើ!)

IEEE:

Architecture is defined by the recommended practice as the fundamental organization of a system, embodied in its components, their relationships to each other and the environment, and the principles governing its design and evolution.

Bass, Clements , Kazman :

The software architecture of a program or computing system is the structure or structures of the system, which comprise software elements, the externally visible properties of those elements, and the relationships among them.

About Software Architecture

- The extension of software engineering discipline
- Software Architecture presents a view of software system as components and connectors
 - Components encapsulate some coherent set of functionality
 - Connectors realize the runtime interaction between components

Architecture នៃ Abstraction level

- ផ្ទាល់ខ្លួន Software Architecture ទៅជា model
 - UML Subsystem diagram (design-time)
 - UML Component diagram (run-time)

SW Arch គាយកិច diagram?

គឺជា diagram សម្រាប់យុទ្ធសាស្ត្រ View នៃកែតាមរបាយ នៃកំណើនការងារ

In summary, a model is an abstract representation of a system, and a diagram is a specific visual representation of a model. Diagrams serve as tools for communicating and visualizing models, making them more accessible to various stakeholders in the software development process.

- Software architecture descriptions are organized into Views.
- Each view addresses a set of system concerns
- Each view is drawn following the conventions of its viewpoint.
- A viewpoint is a specification that describes the notations, modeling techniques to be used in a view to express the architecture.
- Viewpoints are generic templates, while views are what you actually see.

ADL (Formal Model)

- An architecture description language (ADL) is any means of expression used to describe a software architecture (ISO/IEC/IEEE 42010).
- Many special-purpose ADLs have been developed since the 1990s, including
 - AADL (SAE standard),
 - Wright (developed by Carnegie Mellon),
 - Acme (developed by Carnegie Mellon),
 - xADL (developed by UCI),
 - Darwin (developed by Imperial College London),
 - ArchC
- UML is an alternative ADL

- Architecture addresses non-functional requirements (NFR)

Use case diagram ឬ Non-functional req. នៅណែនាំ? Class diagram ហើយណែនាំ?

Non-functional មានចំណាំនៃការងារ

Martin Fowler ឈ្មោះ Internal quality

រូបរាងស្ថិតិស្សន៍

- Architectural styles are high-level design
- How components and connectors may be combined
- Architectural styles = Components + Connectors + Constraints

គិតគុណ

- Software Architecture
 - Model, Level of Abstraction
 - Architectural Views
 - Architectural Viewpoints
 - Architectural Styles (Patterns)
 - Architectural Representation/Documentation
 - Architectural Decision Record (ADR)

Quality Factors (Some NFR)

- Scalability
- Security
- Performance
- Modifiability
- Availability
- Integration

High Available នូវ Performance

* Quality Factors ឱ្យអាជីវកម្ម

* រាយការណ៍ Cost នូវការទិន្នន័យ Building និងការបារិករាយ maintain តួនាទី

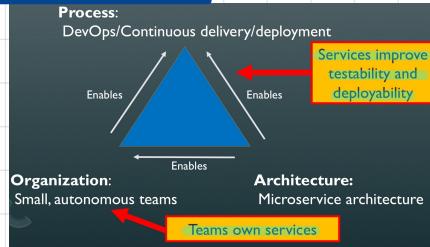
សម្រាប់ការកំណត់គោរគោរ៉ែនកំណត់ model នៅក្នុង

- Constraints
- Risks
- Non-functional requirements
 - ឯកតាម: Online shopping ការបង់ប្រាក់ និងការបង់ប្រាក់
 - Scalability
 - Security
 - Cost constraints
 - Connecting to the Legacy systems
 - 24x7 Availability
 - ...

③ The essential characteristics of the microservice architecture

ilities of microservice architecture

- Maintainability
- Testability
- Deployability
- ...



Microservices ⇒

Evolve the technology stack

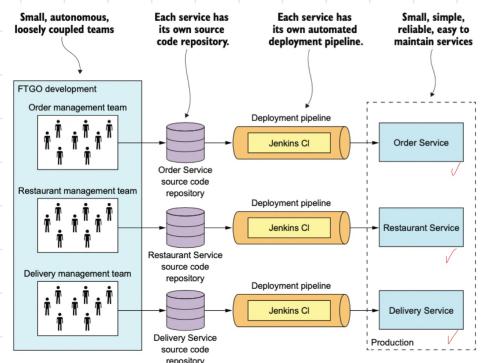
- Pick a new technology when
 - Writing a new service
 - Make major changes to an existing service
- Let you experiment and fail safely

④ Benefits and drawbacks of microservice & How it enables DevOps

- Benefits:**
 - Enable CI/CD of large, complex applications part of DevOps in microservice & testability, deployability
 - Services are small and easily maintained
 - Services are independently deployable (each service can be scaled independently)
 - It has better fault isolation
 - It allows easy experimenting and adoption of new technologies

- Drawbacks:**
 - Finding the right set of services is challenging (transition from monolithic distributed monolith)
 - Distributed systems are complex → Requires the use of unfamiliar techniques
 - Deploying features that span multiple services requires careful coordination
 - Need to create a rollout plan that orders service deployments based on the dependencies between services
 - Deciding when to adopt the microservice architecture is difficult
 - A startup should begin with a monolithic application
 - When the problem is how to handle complexity
 - It makes sense to decompose the application into microservices

required by CI/CD



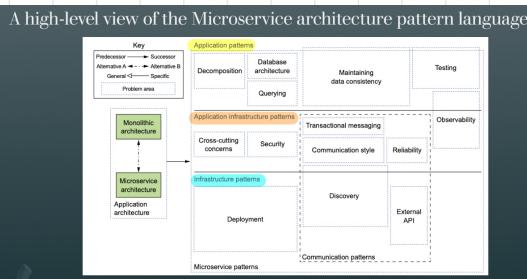
⑤ The microservice architecture pattern language and why you should use it

A pattern is a reusable solution to a problem that occurs in a particular context.

The pattern language guides you when developing an architecture.

What architectural decisions you must make for each decision:

- Available options
- Trade-offs of each option

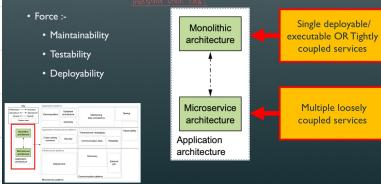


Application patterns — These solve problems faced by developers.

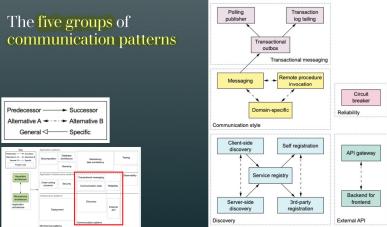
Application infrastructure patterns — These are for infrastructure issues that also impact development.

Infrastructure patterns — These solve problems that are mostly infrastructure issues outside of development.

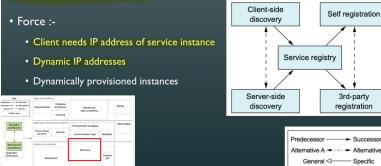
Issue: What's the deployment architecture?



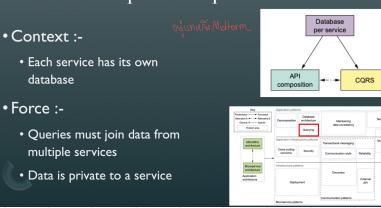
The five groups of communication patterns



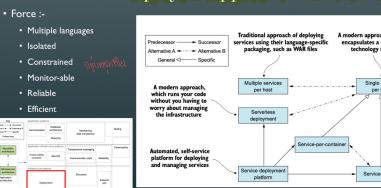
Issue: How to discover a service instance's network location?



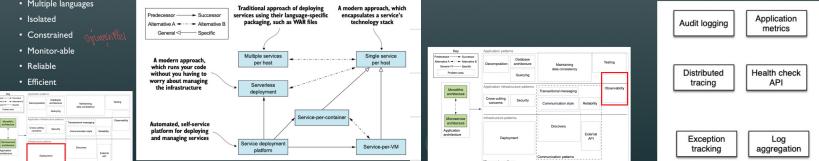
Issue: How to perform queries?



Issue: How to deploy an application's services?



Issue: How to handle cross cutting concerns?

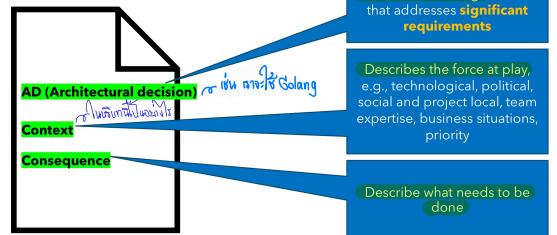


Summary

- The goal of architecture is to satisfy non-functional requirements
- For continuous delivery/deployment use the appropriate architectural style
 - Small applications → Monolithic architecture
 - Complex applications → Microservice architecture
- Use the pattern language to guide your decision making

Architectural Decision Record (ADR)

ADR (Architectural Design Record)



How to start using ADRs?

1. Decision Identification → ដែល? / តើអ្វី? / នៅក្នុង?
2. Decision Making → នូវរាយ Dialog Mapping / Group decision making
3. Decision Enactment and enforcement → ផ្តល់ព័ត៌មានសម្រាប់គ្មានទាំងអស់នៃការងារ Stakeholders via systems / code review នូវការណ៍ architecture
4. Decision Sharing (optional) → Many ADRs recur across project
5. Decision Documentation → នូវ M. Nygard's ADRs.

How to write a good ADR?

Characteristics of a good ADR

Point in Time - Identify when the AD was made

Rationality - Explain why when making the particular AD

Immutable record - The previously published AD should not be altered

Specificity - Each ADR should be about a single AD

↳ includes AD per specific requirement

Decision record template by Jeff Tyree and Art Akerman

Issue	Why do we address this?
Decision	The architecture's direction; the position we selected.
Status	The decision's status; <u>pending</u> , <u>decided</u> , <u>approved</u>
Group	Assigned category for the decision; help organize the set of decisions
Assumptions	Assumptions when the decision is made; <u>schedule</u> , <u>cost</u> , <u>tech</u> ...
Constraints	Capture any additional constraints that the decision might pose.
Positions	Other options / <u>alternatives</u> ; to avoid question "Did you think about..?"
Argument	Outlines why you selected a position; implementation cost, total ownership cost, time to market, and required development resources' availability.
Implication	A decision may come with many implications; creating new requirements; modifying existing requirements, pose additional constraints, require negotiating scope or schedule with customer, require additional staff training

Application in 12-factor app

I. Codebase	One codebase tracked in revision control, many deploys
II. Dependencies	Explicitly declare and isolate dependencies
III. Config	Store config in the environment
IV. Backing services	治本សេវាដែលបានចូលរួមនៅក្នុងការងារ ដែលត្រូវបានគ្រប់គ្រងជាអត្ថបន្ទីរ និងគ្រប់គ្រងជាសេវាដែលបានចូលរួមនៅក្នុងការងារ
V. Build, release, run	Strictly separate build and run stages
VI. Processes	Execute the app as one or more stateless processes
VII. Port binding	Export services via port binding
VIII. Concurrency	Scale out via the process model
IX. Disposability	Maximize robustness with fast startup and graceful shutdown
X. Dev/prod parity	Keep development, staging, and production as similar as possible
XI. Logs	Treat logs as event streams
XII. Admin processes	Run admin/management tasks as one-off processes

The Twelve-Factor App

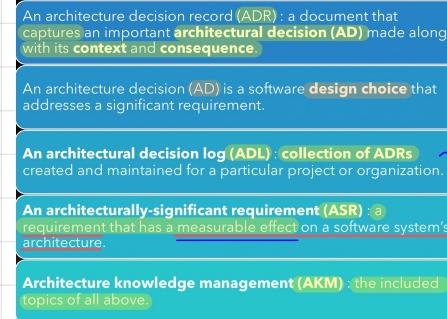
The 12-Factor App Framework was created in 2012 by developers at early cloud pioneer Heroku as a set of rules and guidelines for organizations building modern web applications that run "as a service." The framework was inspired by Martin Fowler's work on application architecture and what he saw as suboptimal application development processes.

Telemetry data refers to the logs, metrics and traces in observability – what is sometimes called "the three pillars of observability". Logs are text records of events that include a timestamp (when the event happened) and a payload that offers context.

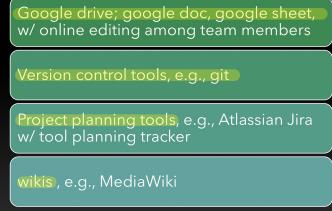
Then, what is the Fifteen-Factor App??

XIII. API First ↗ monitor infrastructure Software as code
XIV. Telemetry
XV. Authentication and Authorization by Hoffman

<https://github.com/ciudd/15-factor-app-workshop>



How to start ADRs with tools



Characteristics of a good "Context" section in ADR

Explain your organization's **situation** and business **priorities**

Include rationale and considerations based on social and skill makeups of your teams

ឧបករណ៍របស់យើង និងការងាររបស់យើង តាមបច្ចេកទេសការងារ និងការងាររបស់យើង ដែលមានការងារខ្លួន និងការងាររបស់យើង

Characteristics of a good "Consequences" section in ADR

Right approach - "We need to start doing X instead of Y"

Wrong approach - Do not explain the AD in terms of "Pros" and "Cons"

14 slide!!!

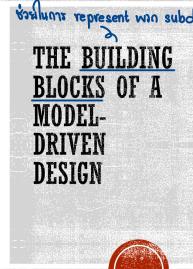
→ ការងារដែលបានរួមចំណាំត្រូវបានរួមចំណាំ

Decision record template by Michael Nygard

Title	Context
	Describes the forces at play, including technological, political, social, and project local.
Decision	It is simply describing facts.
	Describes our responses to these forces. State it in full sentences, with active voices. "We will..."
Status	<p>proposed: if the project stakeholders haven't agreed with it yet.</p> <p>accepted: if it is agreed.</p> <p>deprecated or superseded: if it is replaced by a new ADR.</p>
Consequences	describes the resulting context, after applying the decision. All consequences: positives.

Domain Driven Design (DDD)

- DDD focus on modelling the software that match the business's core domain.
- Software needs to incorporate the core concepts and elements of the domain, and to precisely realize the relationship between them.
- DDD emphasizes collaboration between domain experts and software developer.
- **Ubiquitous Language**: A common language between developers and domain experts (ภาษาที่ใช้กันในทุกๆ ที่)
- **Model-Driven Design (MDD)**: emphasizes the use of models to describe various aspects of a software system, including its structure, behavior, and interactions.



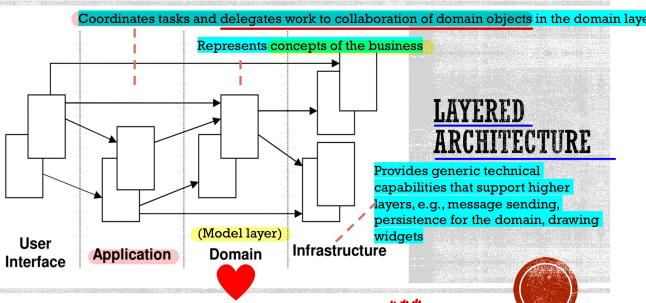
THE BUILDING BLOCKS OF A MODEL-DRIVEN DESIGN

— MDD is an approach to software development where extensive models are created before source code is written.

— MDD ~~integrates~~ ^{model} ~~functional~~ ^{non-functional} aspects into the software ~~for all its domain~~

^{incorporate concepts from DDD}

- **Layered Architecture**: A common architectural solution for DDD contain 4 conceptual layers



LAYERED ARCHITECTURE

- **Entities** :
 - A category of objects which seem to have an identity, which remains the same throughout the states of the software.
 - For these objects it is not attributes which matters but a thread of continuity and identity, which spans the life of a system.

* Entities represents persistent information tracked by the system.

- **Value Objects** :
 - An object used to describe certain aspects of a domain, and which does not have identity.
 - It is recommended to select as entities only those objects conform to the entity definition and make the rest of the objects to be Value Objects.
 - It is highly recommended that value objects are immutable.
 - Examples: Air travel booking system create object for each flight, the flight code can be shared among the customers, who book the same flight.

- **Aggregates** :
 - A group of associated objects which are considered as one unit regarding data changes.
 - The Aggregate is demarcated by a boundary which separates the objects inside from outside.
 - Each Aggregate has one root, an Entity, the only accessible object from outside.

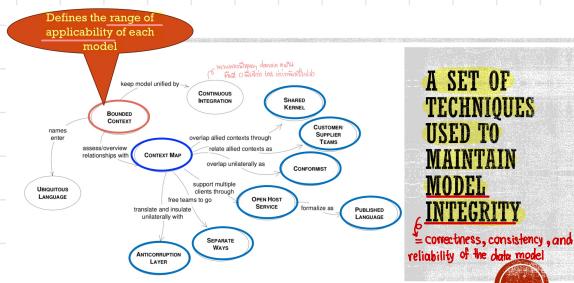
- **Services** :
 - Service is used to represent behavior within a domain.
 - It can group related functionalities which serves the Entities and Value objects.
 - It is much better to declare service explicitly.
 - Services act as interfaces which provide operations, which should be stateless.
 - Services provide a point of connection for many objects.

- **Factories** :
 - Factories are used to encapsulate the knowledge necessary for object creation, and they are especially useful to create Aggregates.
 - Factories is responsible to create other objects.
 - For immutable Value Objects, all attributes are initialized to their valid state.

- **Repositories** : minimum database layer to encapsulate logic

- **Domain Events** :
 - A domain event is representation of something that happened in the domain.
 - Something happened that domain experts care about.
 - Each event is represented as a domain object.

- **Preserving model integrity**



A SET OF TECHNIQUES USED TO MAINTAIN MODEL INTEGRITY

= correctness, consistency, and reliability of the data model

PROBLEM IN A LARGE PROJECT

- Multiple models are in play on any large project.
- When code based on distinct models is combined, software becomes buggy, unreliable, and difficult to understand.
- Communication among team members becomes confused.
- It is often unclear in what context a model should not be applied.

TO MANAGE A LARGE PROJECT

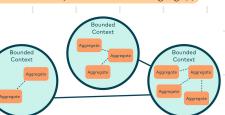
- Explicitly define the context within which a model applies.
- Identify contextual boundaries in terms of team organization, usage within specific parts of the application, and physical manifestations such as code bases and database schemas.
- Keep the model strictly consistent within these boundaries, but don't be distracted or confused by issues outside.

Boundary vs domain

- Bounded Context** : A Bounded Context defines the range of applicability of each model. A Bounded Context provides the logical frame inside of which the model evolves.
- Continuous Integration** : All work within the context is being merged and made consistent frequently enough that when splinters happen they are caught and corrected quickly.

- Continuous Integration separates at two levels:
 - the integration of model concepts
 - the integration of the implementation.

- Aggregates represent business entities, and thus are smaller in scope than Bounded Contexts, which represent entire domains
- Aggregates encapsulate data, whereas Bounded Contexts encapsulate services, team members, and their shared language(s)



"Address" ~~entity~~ ^{value object} ~~model~~ ^{domain}

Is "Address" a VALUE OBJECT? Who's Asking?

In software for a mail-order company, an address is needed to confirm the credit card, and to address the parcel. But if it is normally also orders from the same company, is it not important to realize they are in the same context?

In software for the postal service, intended to organize delivery routes, the country could be formed into a hierarchy of regions, cities, postal zones, and blocks, terminating in individual addresses. These address objects are likely to be part of a larger context, and therefore need to be designed to be reusable across contexts.

In software for an electric utility company, an address corresponds to a distribution for the company's lines and services. In this context, the address is a key part of the domain model, and therefore needs to be an ENTITY. Alternatively, the model could associate utility service with a "swelling", an ENTITY with an attribute of address. Then Address would be a VALUE OBJECT.

Figure 4.1. Objects carry out responsibilities consistent with their layer and are more coupled to other objects in their layer.



"Address" ~~entity~~ ^{value object} ~~model~~ ^{domain}

Is "Address" a VALUE OBJECT? Who's Asking?

In software for a mail-order company, an address is needed to confirm the credit card, and to address the parcel. But if it is normally also orders from the same company, is it not important to realize they are in the same context?

In software for the postal service, intended to organize delivery routes, the country could be formed into a hierarchy of regions, cities, postal zones, and blocks, terminating in individual addresses. These address objects are likely to be part of a larger context, and therefore need to be reusable across contexts.

In software for an electric utility company, an address corresponds to a distribution for the company's lines and services. In this context, the address is a key part of the domain model, and therefore needs to be an ENTITY. Alternatively, the model could associate utility service with a "swelling", an ENTITY with an attribute of address. Then Address would be a VALUE OBJECT.

Customer customerID name street city state



ENTITY CAN CONTAIN VALUE OBJECTS

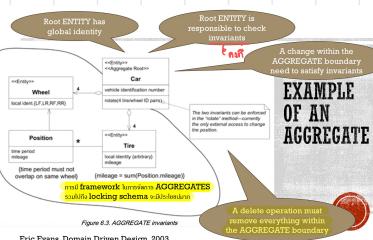


Figure 6.3. AGGREGATE invariants

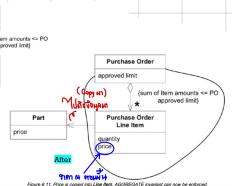


Figure 6.4. Purchase Order LINE ITEM AGGREGATE invariant can now be enforced.

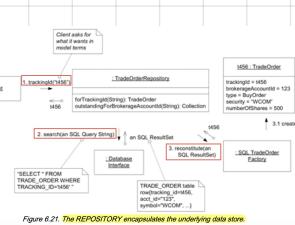
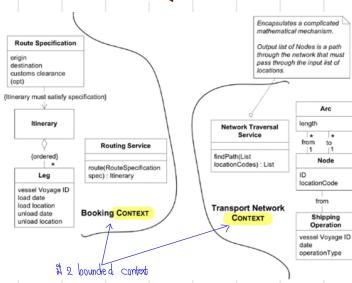


Figure 6.21. THE REPOSITORY encapsulates the underlying data store.

- Context Map = Document which outlines the different Bounded Context and relationships between them.

Code reuse between BOUNDED CONTEXTS is a hazard to be avoided.

- o 1. **unbounded** Context នៅរាល់ business logic នៃ rule និងការសម្រាប់គ្រប់ code នឹងរាយការណ៍នា work នៃ context វាន់ទៀត
2. **multiBounded Context** នឹងរាយការណ៍នាមពេល ដើម្បី code reuse និងការ coupling នៃការសម្រាប់គ្រប់ context (នឹងការសម្រាប់គ្រប់នៅក្នុងគ្រប់)

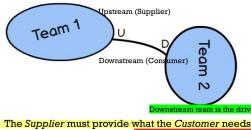


1. Shared Kernel: 2 or more bounded contexts can share a common model
 (ex shared libs, shared)



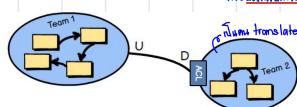
The teams must agree on what model elements they are to share.

2. Customer-Supplier : supplier នឹងរាយក្រឹង Consumer និង Consumer ផ្លូវ



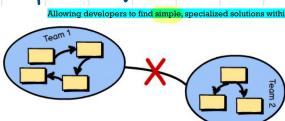
Downstream team is the driver

4. Anticorruption layer (ACL) : ชั้น ACL หรือ downstream
ที่มีหน้าที่ตรวจสอบ corruption



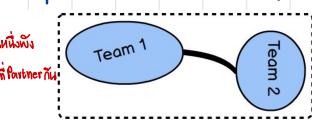
An Anticorruption Layer is the **most defensive** Context Mapping relationship. The layer isolates the downstream model from the upstream model and translates between the two.

- #### 7. Separate Ways: ດີນເວັ້ນທີ່ມີການສ່ວນຮ່າງ



In this case produce your own specialized solution in your Bounded Context and forget integrating for this special case.

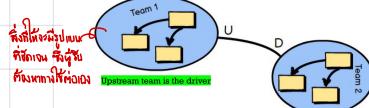
8. Partnership : ≥ 2 bounded context have a strategic collaboration.



Each team is responsible for one Bounded Context.

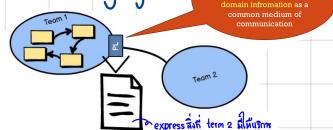
3. Co

3. Conformist: One bounded context makes a decision to conform to another bounded context's model, rules, standards.



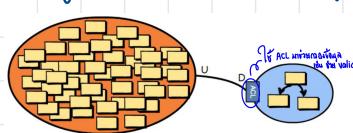
For various reasons the downstream team cannot sustain an effort to translate the Ubiquitous Language of the upstream model to fit its specific needs, so the team conforms to the upstream model as is.

- ## 6. Publish language



Such a Published Language can be defined with XML Schema, JSON Schema, or a more optimal wire format, such as Protobuf or Avro.

- ### 9. Big Ball Of Mud : នាយកដែលមិនអាចរំលែកបាន



Avoid to create big ball of mud, and if need to use it, try ACL

TAKE HOME MESSAGE

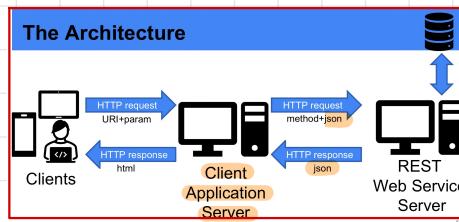
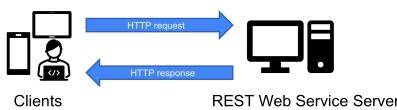
- **“DDD is less about software design patterns and more about problem solving through collaboration.**
 - Evans presents techniques to use software design patterns to enable models created by the development team and business experts to be implemented using the UL.
 - However, without the practices of analysis, and collaboration, the coding implementation really means very little on its own
 - **DDD is not code centric;** its purpose is not to make elegant code. Software is merely an artifact of DDD.
 - **Strategic design** focuses on **defining the bounded contexts**, ubiquitous languages, and context maps
 - **Tactical design** is a set of **technical resources used in constructing domain model in a bounded context.**

Week 5 & 6

REST API

Architecture

The Client-Server Architecture



HTTP Request Methods

GET	Retrieve Resource
POST	Submit Resource Request body (Request Body)
PUT/PATCH	Update Resource
DELETE	Delete/Destroy Resource

API Security

- Logout to clear token cookie
- Prevent NoSQL Injection & Sanitize Data → នៅក្នុងនឹង express-mongo-sanitize



- Security Headers → នៅក្នុងនឹង helmet



Creating OpenAPI Document with Swagger

OpenAPI = A specification language for REST APIs that defines structure and syntax

■ A standard, language-agnostic interface to REST APIs which allows both humans and computers to discover and understand the capabilities of the service without access to source code documentation or through network traffic inspection.

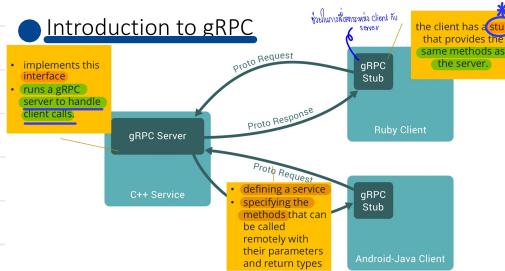
Support to many languages!
= Framework និងក្រុមហ៊ុន distributed system

gRPC Remote Procedure Calls

• gRPC is a modern, open source remote procedure call (RPC) framework that can run anywhere. It enables client and server applications to communicate transparently, and makes it easier to build connected systems.
• It has been used by Google, Square, Netflix, CoreOS, Docker, CockroachDB, Cisco, Juniper Networks and many other organizations

OpenAPI Specification

- Allow machines/tools to integrate with API
- Allow humans to implement API code
- Allow humans to read and generate API documentation and test case
- OpenAPI = Specification
- Swagger = tools



4 kinds of gRPC Service Method

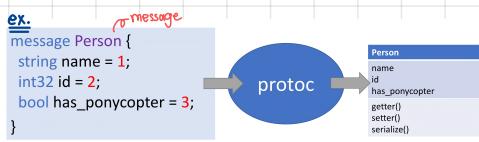
#	Kind	Client's Request	Server's Response
1	Unary RPCs	Single	Single
2	Server streaming RPCs	Stream	Stream
3	Client streaming RPCs	Stream	Single
4	Bidirectional streaming RPCs	Stream	Stream

*gRPC guarantee message order

```
// The greeter service definition.
service Greeter {
    rpc SayHello(HelloRequest) returns (HelloResponse) {}
}

// The request message containing the user's name.
message HelloRequest {
    string greeting = 1;
}

// The response message containing the greetings.
message HelloResponse {
    string reply = 1;
}
```



Protocol Buffer (.proto file)

: Interface Definition Language (IDL) for describing both the service interface and the structure of the payload messages

• required [1]
• optional [0..1]
• repeated [0..*]

• double
• float
• bool
• String
• bytes
• int32/int64
• uint32/uint64
• int32/int64
• uint32/uint64

↳ unique number
↳ From 1 to 2^31 - 1, or
↳ 536,870,911
↳ Except 19999 through 19999

↳ unique field & the Protobuf Run instructions
↳ This is the same as Java
↳ Compatible to Ver 3

Week 7

(IPC) Interprocess communication

Overview

- IPC
 - Synchronous: request/response-based communication mechanism (via REST, gRPC)
 - HTTP based
 - Asynchronous: message-based communication mechanism (via Advanced Message Queuing Protocol (AMQP), STOMP)
- Different messages formats
 - Human-readable: text-based formats (via JSON, XML)
 - Binary format: via Avro, Protocol Buffers

- Interaction styles
 - First Dimension
 - One-to-one → Each client request is processed by exactly one service.
 - One-to-many → Each request is processed by multiple services.
 - Second Dimension
 - Synchronous → The client expects a timely response from service and might even block while it waits.
 - Asynchronous → The client doesn't block, and the response, if any, isn't necessarily sent immediately.

	one-to-one	one-to-many
Synchronous	Request/response	—
Asynchronous	Asynchronous request/response One-way notifications	Publish/subscribe Publish/async responses

1. One-to-one interactions

- Request/response: waits for response + might even block while waiting
- Asynchronous request/response: NOT block while waiting
- One-way notifications: A service client sends a request to a service, but no reply is expected or sent.

2. One-to-many interactions

- Publish/Subscribe: A client publishes a notification message + consumed by ≥ 0 interested service
- Publish/async responses: A client publishes a request message + waits for a certain amount of time for responses from interested services.

Defining APIs in microservice

- A well-designed interface exposes useful functionality while hiding the implementation
 - Enables the implementation to change without impacting clients
 - Define a service's API using Interface Definition Language (IDL)

↳ IDL was REST transitioning to OpenAPI Specification (REST did not originally have an IDL)
 ↳ gRPC is Protocol Buffers-based IDL

Communicating using the Synchronous Remote procedure invocation pattern

- Using a remote procedure invocation-based IPC mechanism: client → service → response return + client doesn't block until response + Client mustn't receive response "arrive in a timely fashion"
- REST vs. gRPC

REST	gRPC
Benefits <ul style="list-style-type: none"> simple & familiar Directly supports request/response style communication Firewall friendly Doesn't require an intermediate broker (Simplifies the system's architecture) 	<ul style="list-style-type: none"> Straightforward to design an API with a rich set of update operations Efficient especially when exchanging large messages Bidirectional Streaming enables both API and messaging styles of communication Enables interoperability between clients and services written in a wide range of languages. ↳ gRPC has many implementations
Drawback <ul style="list-style-type: none"> Only supports request/response style communication Reduced availability <ul style="list-style-type: none"> Client and service communicate directly without an intermediary to buffer messages. ↳ They must both be running for the duration of the exchange ↳ Using "service discovery" mechanism Client must know the location (URLs) of the service instance(s) Fetching multiple resources in a single request is challenging 	<ul style="list-style-type: none"> More work for javascript clients to consume gRPC-based API Older firewalls might not support HTTP/2

1 2 3 → ms = partial failure

* ms = synchronous request → full automation partial failure

Handling partial failure

① using the Circuit breaker pattern

② 2 solution

1. Design RPC proxies to handle unresponsive remote services.

1. Network timeouts
 - Never block indefinitely and always use timeouts when waiting for a response.
2. Limiting the number of outstanding requests from a client to a service
 - Impose an upper bound on the number of outstanding requests that a client can make to a particular service.
3. Circuit breaker pattern
 - Track the number of successful and failed requests
 - If the error rate exceeds some threshold, trip the circuit breaker so that further attempts fail immediately.

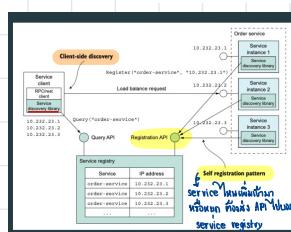
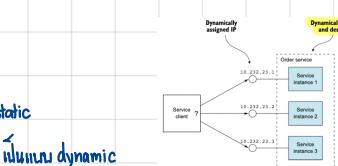
2. Design how to recover from a failed remote service

- We must also decide how your services should recover from an unresponsive remote service
- Option 1 - Return an error to its client
- Option 2 - Returning a fallback value

② using Service Discovery

- សំគាល់ → Network location has service instance និងមុន static
- សំគាល់ → អាមេរិក cloud-based និង Network location និង dynamic

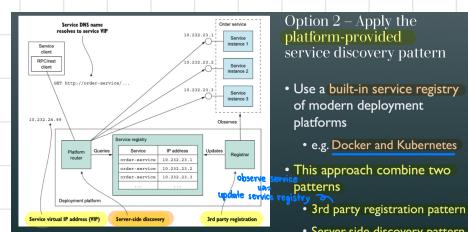
- Service discovery key component is a **service registry** → ពីរ service និង ip address ទៅ
 - A database of the network locations of an application's service instances.
- The service discovery mechanism updates the **service registry** when service instances start and stop.
- There are two main ways to implement service discovery:
1. The services and their clients interact directly with the **service registry**.
 2. The deployment infrastructure handles service discovery.



Option 1 – Apply the application-level service discovery pattern

- Implement service discovery for the application's services and their clients
- This approach combine two patterns
 - Self registration pattern
 - Client-side discovery pattern

Client knows discover service



Option 2 – Apply the platform-provided service discovery pattern

- Use a built-in service registry of modern deployment platforms
 - e.g. Docker and Kubernetes
- This approach combine two patterns
 - 3rd party registration pattern
 - Server-side discovery pattern

return network addr was service VIP client

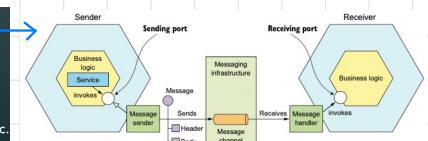
The client doesn't block waiting for a reply + assuming that the reply won't be received immediately.

Communicating using the Asynchronous messaging pattern

Overview

- Messages are exchanged over message channels.
- A **sender** (an application or service) writes a message to a channel.
- A **receiver** (an application or service) reads messages from a channel.
- There are several different kinds of messages:
 - Document – A generic message that contains only data.
 - Command – A message that's the equivalent of an RPC request.
 - Event – A message indicating that something notable has occurred in the sender.

1. The business logic in the sender invokes a **sending port interface**.
2. The **message sender** sends a message to a receiver via a **message channel**.
3. The **message channel** is an abstraction of messaging infrastructure.
4. A **message handler** adapter in the receiver is invoked to handle the message.
5. It invokes the **receiving port** interface implemented by the receiver's business logic.



– Channel និង វិធាន នូវ 1. A point-to-point channel : នឹង message នូវ 1 consumer (និង command message)

2. A publish-subscribe channel : នូវ នូវយុទ្ធសាស្ត្រ consumer (និង event message)

Implementing the interaction style using messaging

1. Request/response and Asynchronous Request/response

- Messaging is inherently asynchronous, so only provides asynchronous request/response.
- But a client could block until a reply is received. → Synchronous

3. Publish/subscribe

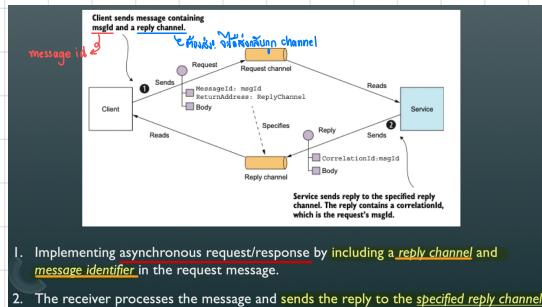
- A client publishes a message to a publish-subscribe channel that is read by multiple consumers
- A service that's interested in a particular domain object's events only has to subscribe to the appropriate channel.

2. One-way notification

- The service subscribes to the channel and processes the message.
- But it doesn't send back a reply.

4. Publish/Asynchronous subscribe

- Combining elements of publish/subscribe and request/response



1. Implementing asynchronous request/response by including a **reply channel** and **message identifier** in the request message.

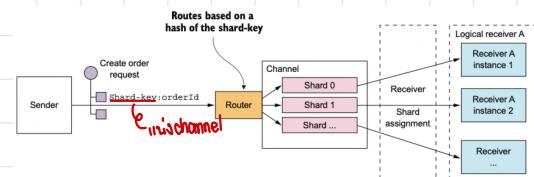
2. The receiver processes the message and sends the reply to the specified **reply channel**.

Brokerless vs. Broker-based architectures

	Brokerless	Broker-based
Benefits	<ul style="list-style-type: none"> • Allows lighter network traffic and better latency. • Eliminates the possibility a performance bottleneck or a single point of failure. • Features less operational complexity 	<ul style="list-style-type: none"> • Loose coupling – The client doesn't need to use a discovery mechanism • Message buffering (from client to broker) • Flexible communication – Messaging supports all the interaction styles • Explicit interprocess communication
Drawbacks	<ul style="list-style-type: none"> • Services must use one of the discovery mechanisms • Availability is reduced, because both the sender and receiver must be available while the message is being exchanged. 	<ul style="list-style-type: none"> • Potential performance bottleneck • Potential single point of failure • Additional operational complexity

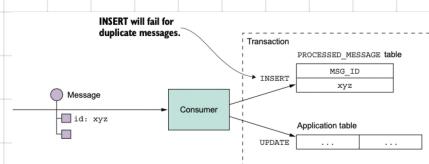
Competing receivers and message ordering

- A common solution, used by modern message brokers like Apache Kafka and AWS Kinesis, is to use sharded (partitioned) channels
- Each Order event message has the **orderID** as its shard key → និង នូវ instance នូវយុទ្ធសាស្ត្រ
- Each event is published to the same shard, which is read by a single consumer instance



Handling duplicate messages

- A message consumer has to discard any duplicates message → Tracking by the message id
- A consumer record the IDs of processed messages in a database table.
- If a message has been processed before, the **INSERT** into the **PROCESSED_MESSAGES** table will fail.



Summary

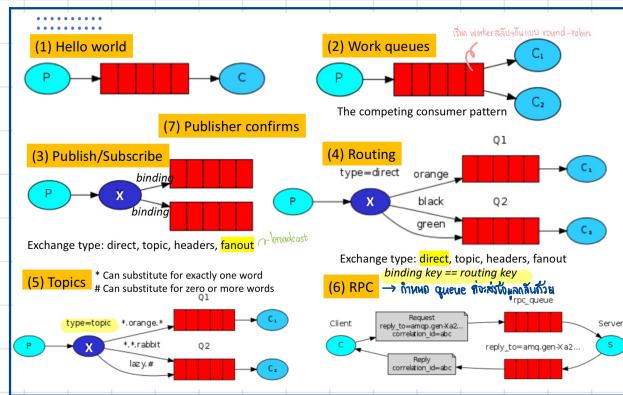
- The microservice architecture is a distributed architecture, so **interprocess communication** plays a key role.
- It's essential to **carefully manage** the evolution of a service's API.
- There are **numerous IPC technologies**, each with different trade-offs.
- To prevent failures, a service client that **uses a synchronous protocol** must be designed to **handle partial failures**.

- An architecture that uses **synchronous protocols** must include a **service discovery mechanism** to determine the network location of a service instance.
- A good way to design a messaging-based architecture is to **use the messages and channels model, typically message broker-based**.
- A key challenge when using messaging is atomically updating the database and publishing a message and a good solution is to use the Transactional outbox pattern

RabbitMQ

- RabbitMQ is an open-source distributed message broker.
- It is built on AMQP (Advanced Messaging Queuing Protocol).
- AMQP is an open standard application layer protocol.
- RabbitMQ accepts and forwards messages.
- It accepts, stores, and forwards binary blobs of data – *messages*.

Producing A program which sends message is a <i>producer</i> .	P
Queue : a post box living inside RabbitMQ.	queue_name
Consuming A consumer is a program that mostly waits to receive messages.	C



round-robin → လုပ်မှုများအတွက် အသေစိတ်ပါ။
 fan-out-dispatch → အနေဖြင့် အပေါ် ပေါ်မှုများ

RabbitMQ & Apache Kafka

RabbitMQ	Kafka
Direct messaging: users can set sophisticated rules for message delivery.	Kafka is ideal for handling large amounts of homogeneous messages, such as logs or metric, and it is the right choice for instances with high throughput.
RabbitMQ will be usually used with Cassandra for message playback.	Replay messages
Multiprotocol supports: AMQP, STOMP, MQTT, Web sockets and others.	Big data consideration with e.g., Elastic search, Hadoop
Flexibility: varied point-to-point, request/reply, publish/subscribe	Scaling capability; topics can be split into partitions
Communication: supports both async and sync	Batches: Kafka works best when messages are batched.
Security	Security
Mature platform	Mature platform
Slower than Kafka	Don't come with user-friendly GUI but can be monitored via Kibana

Kafka: Best used for basic streaming without complex routing with maximum throughput, ideal for event-sourcing, stream processing, multi-stage pipelines, routinely audited systems, real-time processing and analyzing data

RabbitMQ: Used for high throughput and reliable background jobs, integration and intercommunication between and within applications, perform complex routing to consumers, integrate multiple applications and services with non-trivial routing logic.

Benefits of an API Gateway

- Encapsulates internal structure of the application.
- Reduces the number of round-trips between the client and application.
- Simplifies the client code.

-API Gateway Design Issues

1. Performance and scalability

"API Gateway mit synchronous vs asynchronous I/O ?"

- Synchronous I/O model → program's execution is blocked, or synchronized while waiting
- Each network connection is handled by a dedicated thread → Special thread for I/O operation to finish
 - This is a simple programming model and works reasonably well for specific task
 - Limitation is that operating system threads are heavyweight.
↳ Many threads need many dedicated threads
 - So there is a limit on the number of threads, and hence concurrent connections, that an API gateway can have.

Thread with Heavyweight threads
Dedicated thread with few threads
Few threads lots

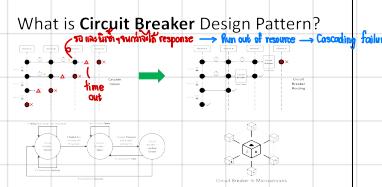
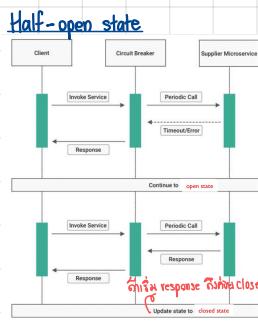
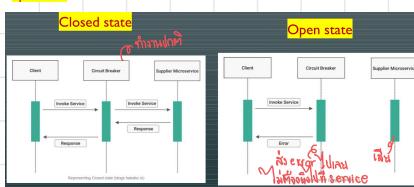
Asynchronous I/O model

- A single event loop thread dispatches I/O requests to event handlers.
- Much more scalable because it doesn't have the overhead of using multiple threads.
- The drawback is the callback-based programming model is much more complex.
- Event handlers must return quickly to avoid blocking the event loop thread.

3. Handling partial failure

- API gateway must be reliable.
- One way is to run multiple instances of the gateway behind a load balancer.
- Another way is to properly handle failed requests and requests that have unacceptably high latency.
- Using the Circuit breaker pattern when invoking services.

Circuit breaker pattern (CBP) = Allow developers to prevent cascading failure
3 states - closed, open, half-open



Implementing an API Gateway

1. Using an off-the-shelf API gateway product/service → Requires little or no development but is the least flexible
2. Developing your own API gateway → with framework like Traefik → Flexible but it requires some development effort.

Implementation / Tools

1. AWS Application Load Balancer

- A load balancer for HTTP, HTTPS, WebSocket, and HTTP/2
- You define routing rules that route requests to backend services
 - Drawback:
 - It does not implement HTTP method-based routing.
 - Nor does it implement API composition or authentication.

2. Kong with Traefik

- Kong is based on the NGINX HTTP server
- Traefik is written in GoLang
- Support edge functions and powerful routing capabilities
- Drawback:
 - They don't support API composition

Summary (I)

- Your application's external clients usually access the application's services via an API gateway.
- Your application can have a single API gateway or it can use the Backends for frontends pattern.
- The main advantage of the Backends for frontends pattern is that it gives the client teams greater autonomy because they develop, deploy, and operate their own API gateway.
- There are numerous technologies you can use to implement an API gateway, including off-the-shelf API gateway products or develop your own API gateway using a framework.

API gateway คืออะไร คืออะไร
Backend for Frontend pattern

-Drawbacks of an API Gateway

- Another highly available component that must be developed, deployed, and managed.
- Also a risk that the API gateway becomes a development bottleneck.
- Developers must update the API gateway in order to expose their service's API.
- The process for updating the API gateway be as lightweight as possible.
→ Otherwise, developers will be forced to wait in line in order to update the gateway.

Lightweight API gateway → One process per microservice
Sequential calls

2. Writing maintainable code by using reactive programming abstractions

Illustrate that service call Sequential → ที่สุด: Response Time = Sum(Service Response Time)

↳ วิธี minimize วนการคำนวณ service ให้ concurrent

- ไม่ต้องห่วงเรื่อง queue หรือ callback

• The problem is this method returns a Future, which has a blocking API.

• The callback accumulates results.

• Once all of them have been received, it sends back the response to the client.

Can lead to callback hell
(Callback ลากยาวที่สุด)

Listing 8.1 Fetching the order details by calling the backend services sequentially

```
@RestController
public class OrderDetailsController {
    @GetMapping("/order/{orderId}")
    public OrderDetails getOrderDetails(@PathVariable String orderId) {
        OrderInfo orderInfo = orderService.findOrderById(orderId);
        TicketInfo ticketInfo = ticketService.findTicketByOrderId(orderId);
        DeliveryInfo deliveryInfo = deliveryService.findDeliveryByOrderId(orderId);
        BillInfo billInfo = accountingService.findBillByOrderId(orderId);
        OrderDetails orderDetails = OrderDetails.newBuilder()
            .setOrderInfo(orderInfo)
            .setTicketInfo(ticketInfo)
            .setDeliveryInfo(deliveryInfo)
            .setBillInfo(billInfo)
            .build();
        return orderDetails;
    }
    ...
}
```

Sequential calls cause the long response time ... concurrently calls are the solution but difficult to implement.

- mention the Write API composition Code in a declarative style
using a reactive approach
✓ Code simple and easy to understand

Examples of reactive abstractions for the JVM include the following

- Java 8 CompletableFuture
- Project Reactor Monos
- RxJava (Reactive Extensions for Java) by Netflix
- Scala Futures

Week 9

Implementing Queries

Overview of querying data in a microservice architecture

- An Application implements a variety of query operations → **Querying in Monolithic** → Straightforward (Because it has a single database)
- **Querying in Microservice** → Queries often need to retrieve data scattered among databases owned by multiple services.

- 2 Pattern for Implement query operations
 1. **The API composition pattern**
 - The simplest approach and should be whenever possible.
 2. **The Command query responsibility segregation (CQRS) pattern**
 - More powerful than the API composition pattern, but also more complex.

↳ Queries must join data from multiple services (Having data from private information service!?)

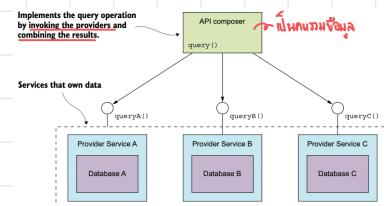
Querying using the API composition pattern

- This pattern implements a query operation by
 1. Invoking the services that own the data
 2. Combining the results
- The structure of this pattern has 2 types of participants
 1. Provider Service = A service that owns some of the data that the query returns
 2. API composer = Implements the query operation by querying the provider services

↳ might be

 - A client, such as a web application
 - A service, such as an API gateway

↳ might need to perform **in-memory join of large datasets**.
- To implement a particular query operation depends on
 - How the data is partitioned
 - The capabilities of the APIs exposed by the services that own the data
 - The capabilities of the databases used by the services



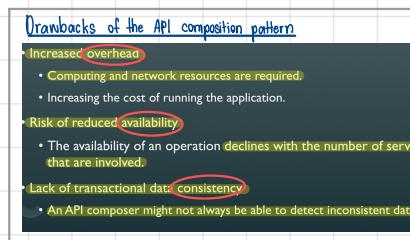
* Design issues

- 1. Which component in architecture is the query operation's API composer?
 - 1) Using a Client to be the API Composer

Drawbacks: Probably not practical for clients that are outside of the firewall and access services via a slower network + unresponsive clients!
 - 2) Using an API Gateway to play the role of an API Composer

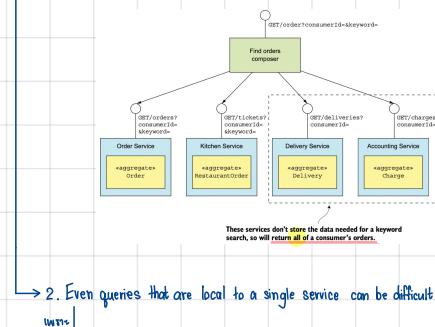
Drawbacks: Not practical if the query operation is part of the application's external API.
 - 3) Implement an API composer as a stand-alone service

Drawbacks: Enables a client that is running outside of the firewall to efficiently retrieve data from numerous services with a single API call.



Querying using the CQRS pattern

- Motivations for using CQRS → 1. Not all services store the attributes that are used for filtering or sorting



Solution I : Do an in-memory join (Retrieves all orders from Delivery to Accounting Service → performs a join with orders from Order to Kitchen Service)

Drawbacks : API composer must retrieve and join large datasets, which is inefficient.

Solution II : Matching then request by fetching ID (Matching orders from Order to Kitchen Service → then request order from service Delivery to Kitchen Service)

Drawbacks : Only practical if those services have a bulk fetch API
Inefficient because of excessive network traffic

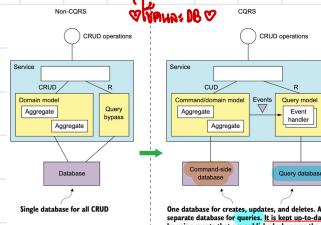
Command Query Responsibility Segregation

Overview of CQRS

- Three problems commonly encountered when implementing queries in a microservice architecture :
 1. Using the API composition pattern to retrieve data scattered across multiple services results in expensive, inefficient in-memory joins.
 2. The service that owns the data stores the data in a form or in a database that doesn't efficiently support the required query.
 3. The need to separate concerns means that the service that owns the data isn't the service that should implement the query operations.

↳ real use with microservices
- The solution to all these problems is to use the **CQRS pattern**!!

- CQRS separates command from queries
- Restructures a service into **command-side** and **query-side** modules, which have separate databases.



We can use CQRS pattern to define query services.

query-only service (no command operation)

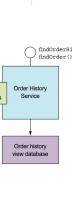
→ Keep it up-to-date by subscribing to events published by > 1 other services

→ This kind of view doesn't belong to any particular service So, it makes sense to implement it as a standalone service

One database for creates, updates, and deletes. A separate database for queries. It is kept up-to-date by using events that are published whenever the command-side database changes.

- Example

- It implements the `findOrderHistory()` query operation by querying a database, which it maintains by **subscribing to events published by multiple other services**.
- The popular approach :-
 - Use RDBMS as the system of record
 - Use a text search engine, such as Elasticsearch, to handle text queries



- Benefits

- Enables the **efficient implementation of queries** in a microservice architecture.
- Enables the **efficient implementation of diverse queries**.
- Makes querying possible in an **event sourcing-based application**.
 - Capture all changes to an application state as a sequence of events.
- Improves separation of concerns**.

- Drawbacks

- More **complex architecture**
 - Developers must write the query-side services that update and query the views.
 - An application might use different types of databases.
- Dealing with the **replication lag**
 - There's delay between the command side publishes an event and the query side processes the view updated.

Solution to the "Replication Log"

- Solution 1**: Supply the client with **version information**.
 - Enables the client to tell that the query side is out-of-date.
 - A client can poll the query-side view until it's up-to-date.
- Solution 2**: **Updating its local model** once the command is successful
 - Update its model using data returned by the command.
 - Drawback - The UI code may need to duplicate server-side code in order to update its model.

● Designing CQRS views

- CQRS view module = 1 view database + 3 submodules

- Important decisions when developing a view module

① Choose a database and design the schema

Issue 1: SQL vs. NoSQL

- NoSQL database is often a good choice for a CQRS view
 - More flexible data model + Better performance and scalability

Query-side view stores

If you need	Use	Example
PK-based lookup of JSON objects	A document store such as MongoDB or DynamoDB, or a key value store such as Redis	Implement order history by maintaining a MongoDB document containing the per-consumer.
Query-based lookup of JSON objects	A document store such as MongoDB or DynamoDB	Implement customer view using MongoDB or DynamoDB.
Text queries	A text search engine such as Elasticsearch	Implement text search for orders by maintaining a per-order Elasticsearch document.
Graph queries	A graph database such as Neo4j	Implement fraud detection by maintaining a graph of customers, orders, and other data.
Traditional SQL reporting/BI	An RDBMS	Standard business reports and analytics.

Issue 2: Supporting update Operations

- Implement the **efficient update operations** executed by the event handlers

Using the primary key

Sometimes, using foreign key

- Some types of databases efficiently support foreign-key-based update operations.

RDBMS or MongoDB can create an index on the necessary columns

→ supports new query

③ When implementing a new view or changing the schema → Must implement a mechanism to efficiently build or rebuild the view
(*enables re-create view from initial state* Schema change into *allowing bug uninitialization View*)

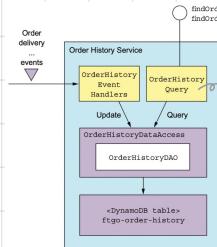
- To create a new view,
 - Develop the query-side module → set up the datastore, → deploy the service
- To update an existing view
 - Change the event handlers → rebuild the view from scratch or update incrementally.

④ Decide how to enable a client of the view to cope with the replication lag

NoSQL

● Implementing a CQRS view with AWS DynamoDB

Architecture Implement CQRS First Design Issues



Design issues

1. Design Ago-order-history table

Primary key : consumerId, orderCreationTime

return a consumer's orders sorted by increasing age

2. Define an index for the findOrderHistory query

Partition key → Primary key that has the consumer
Sort key → orderCreationDate
Index → AWS Lambda function that triggers on order creation
Index key → consumerId, orderCreationTime, global secondary index

5. Update order

- DynamoDB supports two operations for adding and updating items: `PutItem()` and `UpdateItem()`.
 - The `PutItem()` operation creates or replaces an entire item by its primary key.
 - The `UpdateItem()` operation updates individual attributes of the item, creating the item if necessary.
- The `UpdateItem()` operation is more efficient because there's no need to first retrieve the order from the table.

- key condition expression supports a range restriction on the sort key.
- The other filter criteria (Boolean) can be implemented using a filter expression.

has a filter parameter that specifies the search criteria

3. Implementing the findOrderHistory query

- The `OrderHistoryDaoDynamoDb` also enables the keyword search.
 - By tokenizing the restaurant name and menu items.
 - Storing the set of keywords in a se-valued attribute called `keywords`.

4. Paginate the query result

- The `DynamoDB Query` operation has an operation `pageSize` parameter.
 - specifies the maximum number of items to return
 - If there are more items, the result of the query has a non-null `LastEvaluatedKey` attribute

6. Detect duplicate events

- `OrderHistoryDaoDynamoDb` can detect duplicate events by recording in each item the events that have caused it to be updated.
- It can then use the `UpdateItem()` operation's conditional update mechanism to only update an item if an event isn't a duplicate.
- A conditional update is only performed if a condition expression is true.

Summary

- Implementing queries that **retrieve data from multiple services** is challenging because each service's data is private.
- There are two ways to implement these kinds of query:
 - The API composition pattern
 - Command query responsibility segregation (CQRS) pattern
- The API composition pattern is the simplest way to implement queries and should be used whenever possible.
- A limitation of the API composition pattern is that some complex queries require inefficient in-memory joins of large datasets.

- The **CQRS pattern**, which implements queries using **view databases**, is more powerful but more **complex to implement**.
- A **CQRS view module** must handle concurrent updates as well as detect and discard duplicate events.
- CQRS improves separation of concerns by enabling a service to implement a query that returns data owned by a different service.
- Clients must handle the eventual consistency of CQRS views.

Week 10

Deploying MicroService

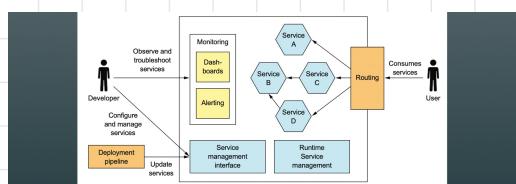
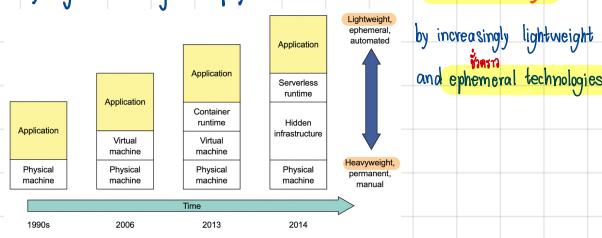
An overview of Deploying Microservice

Deployment = combination of two interrelated concepts \Rightarrow process + architecture

- Process** consists of the steps that must be performed by people—developers and operations—in order to get software into production.
- Architecture** defines the structure of the environment in which that software runs.

* Both aspects of deployment have changed radically since 1990s.

Heavyweight and long-lived physical machines has been abstracted away



A simplified view of the production environment.

1. Service management interface enables developers to deploy and manage their services
2. Runtime service management ensures that the services are running
3. Monitoring visualizes service behavior and generates alerts
4. Request routing routes requests from users to the services

Deploying services using the Language-specific packaging format pattern

- Deploy 服务到 Bare Machine จัดการผ่าน JAR file หรือ战部署

- Benefits**
 - Fast deployment
 - Efficient resource utilization

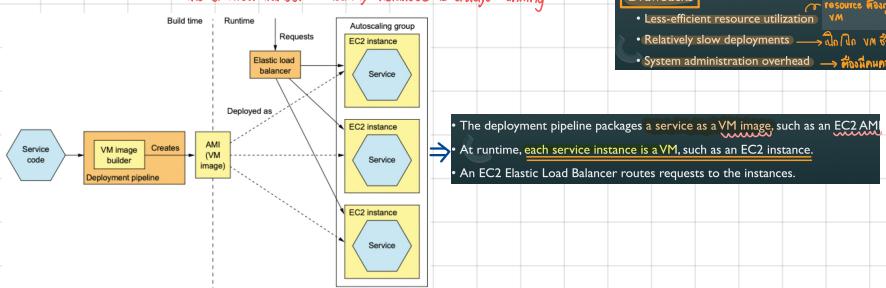
- Drawbacks**
 - Lack of encapsulation of the technology stack
 - No ability to constrain the resources consumed by a service instance
 - Lack of isolation when running multiple service instances on the same machine
 - Automatically determining where to place service instances is challenging

Encapsulation in the context of software deployment typically involves packaging applications and their dependencies in a self-contained unit, ensuring that they can run independently without interfering with other services or applications on the same machine. This is a common feature provided by virtualization technologies, containers, and other deployment solutions.

Deploying services using the Service as a virtual machine pattern

- 用 AWS Config EC2 instance 之間的 Executable File 進行部署

- Created from Amazon Machine Image (AMI) < Better Option
- Managed by an AWS Auto Scaling Group
- The desired number of healthy instances is always running

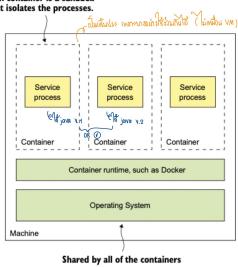


- Benefits**
 - The VM image encapsulates the technology stack
 - Service instances are isolated
 - Uses mature cloud infrastructure \rightarrow cloud VM is stable

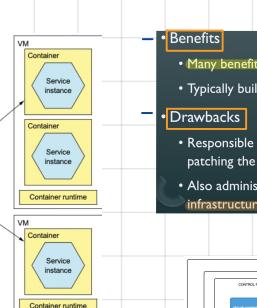
- Drawbacks**
 - Less-efficient resource utilization \rightarrow resource 消耗高
 - Relatively slow deployments \rightarrow 需要更多时间 + Set up environment
 - System administration overhead \rightarrow 管理复杂

Deploying services using the Service as a container pattern

Each container is a sandbox that isolates the processes.

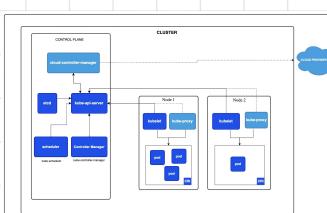


- A service is packaged as a container image
- At runtime the service consists of multiple containers instantiated from that image.
- A single VM will usually run multiple containers.



- Benefits**
 - Many benefits of VM but more lightweight
 - Typically build faster than VM

- Drawbacks**
 - Responsible for the heavy lifting of administering the container images, patching the operating system and runtime.
 - Also administering the container infrastructure and possibly the VM infrastructure it runs on.



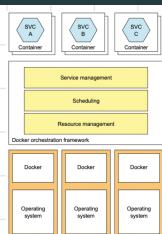
* Placing containers into pods to run on Nodes

Deploying services using Kubernetes

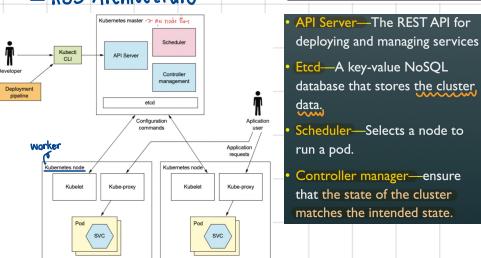
- Kubernetes is a Docker orchestration framework.
 - A layer of software on top of Docker that turns a set of machines into a single pool of resources for running services.
- The desired number of instances of each service keeps running at all times. \heartsuit Self-Healing
 - Even when service instances or machines crash.

3 Main Functions

- Resource management**—Treats a cluster of machines as a pool of CPU, memory, and storage volumes
- Scheduling**—Selects the machine to run your container
- Service management**—Implements the concept of named and versioned services that map directly to services in the microservice architecture.



- K8s Architecture



- API Server**—The REST API for deploying and managing services
- Etc**—A key-value NoSQL database that stores the cluster data.
- Scheduler**—Selects a node to run a pod.
- Controller manager**—Ensures that the state of the cluster matches the intended state.
- Kubelet**—Creates and manages the pods running on the node
- Kube-proxy**—Manages networking, including load balancing across pods
- Pods**—The application services

♥ K8s has Service discovery mechanism built-in

♥ Service = K8s object that provides the clients of one or more pods with a stable endpoint \heartsuit DNS Name

- Deploying the API Gateway (Route traffic from the outside world to the service)

ClusterIP service ⇒ A type of service that exposes an application or set of pods **within the cluster** to other pods within the same cluster (Traffic default) → for internal communication

NodePort service ⇒ A type of service that exposes an application or set of pods **to the external world**, allowing external traffic to reach the service on a specific port (using nodeport)

LoadBalancer ⇒ It provisions an **external load balancer** from a cloud provider to distribute traffic across the nodes in your cluster.

- Using a service mesh to separate deployment from release

- Deployment—Running in the production environment
- Releasing a service—Making it available to end-users

Steps :

- Deploy the new version into production without routing any end-user requests to it.
- Test it in production.
- Release it to a small number of end-users.
- Incrementally release it to an increasingly larger number of users until it's handling all the production traffic.
- If at any point there's an issue, revert to the old version.

การปล่อย release ที่ดีที่สุดคือ
ให้ผู้ใช้เข้าถึงเวอร์ชันใหม่ อย่างช้าๆ
ตั้งแต่ 20% ไปจน
ถึง 100% ของจำนวนผู้ใช้

using "service mesh" to make this style of deployment is a lot of easier

- A service mesh is a networking infrastructure that mediates all communication between a service and other services and external applications.
- A service mesh provides rule-based load balancing and traffic routing
 - let you safely run multiple versions of your services simultaneously. *
 - *

⇒ โครงสร้าง Istio :

Istio as an "An open platform to connect, manage, and secure microservices" (<https://istio.io>)

- Istio has a rich set of features organized into four main categories:
 - Traffic management—Includes service discovery, load balancing, routing rules, and circuit breakers
 - Security—Secures interservice communication using Transport Layer Security (TLS)
 - Telemetry—Captures metrics about network traffic and implements distributed tracing
 - Policy enforcement—Enforces quotas and rate limits

Pilot = ผู้จัดการ data plane ที่อยู่ใน Cloud
Mixer = โมดูล monitoring infra ที่อยู่ใน Cloud

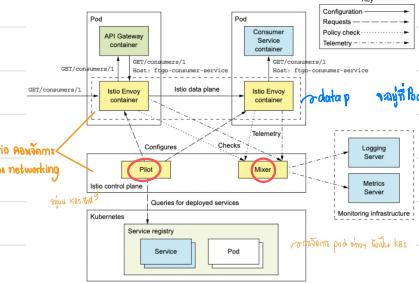
- It consists of a control plane and a data plane.
- The control plane implements management functions, including configuring the data plane to route traffic.
- The data plane consists of Envoy proxies, one per service instance.

The two main components of the control plane are the Pilot and the Mixer.

- The Pilot extracts information about deployed services from the underlying infrastructure.
- The Mixer collects telemetry from the Envoy proxies and enforces policies.

Proxy server ที่อยู่ใน Cloud

A proxy server is an intermediary server that acts as a gateway between client devices (such as computers or smartphones) and other servers (typically web servers or other online services). It sits between the client and the destination server, receiving client requests to the server and returning the server's responses back to the client.



● Deploying services using the Serverless deployment pattern

- Undifferentiated "Heavy Lifting"

Heavy Lifting = ไม่ต้องห่วงเรื่อง hardware (VM, Container, Physical Machine) หรือการติดต่อเชื่อมต่อเครือข่าย
ให้ AWS ดูแลให้หมด

"Serverless"

- AWS Lambda = example of serverless deployment technology

- AWS Lambda was initially for deploying event-driven services.
- AWS Lambda automatically runs enough instances of your microservice to handle incoming requests.
- Billed for each request based on the time taken and the memory consumed.
- No need to worry about any aspect of servers, virtual machines, or containers!!

- Benefits of AWS Lambda

- Integrated with many AWS services
- Eliminates many system administration tasks
- Elasticity—No need to predict the capacity of VMs or containers.
- Usage-based pricing

Drawbacks of AWS Lambda

- Long-tail latency
 - Because AWS Lambda dynamically runs your code.
 - Some requests have high latency because of the time it takes for AWS to provision an instance of your application to start.
- Limited event/request-based programming model
 - AWS Lambda isn't intended to be used to deploy long-running services.
 - Such as a service that consumes messages from a third-party message broker.

Summary

You should choose the **most lightweight deployment** pattern that supports your service's requirements.

Evaluate the options in the following order:

- serverless
- containers
- virtual machines
- language-specific packages

Deploying your service as a **virtual machine** is a **heavyweight** deployment option and will most likely use more resources than Docker containers.

Deploying your services as **language-specific packages** is generally best avoided unless you only have a small number of services.

A serverless deployment eliminates the need to administer operating systems and runtimes.

It provides automated elastic provisioning and request-based pricing.

A serverless deployment isn't a good fit for every service.

Long-tail latencies

Required to use an event/request-based programming model.

Docker containers, a lightweight, OS-level virtualization technology, are more flexible than serverless deployment and have more predictable latency.

It's best to use a Docker orchestration framework such as Kubernetes.

The drawback is that you must administer the operating systems and runtimes.

Week 11

Testing in Microservices Part I

● Overview of testing microservice

- Testing is primarily an activity that happens after development
- Why we need automated testing → ① Manual testing is extremely inefficient ② Manual testing is done far too late in the delivery process
- Workflow should be: edit code → run tests (ideally with a single keystroke) → repeat

● Testing strategies for microservice architectures

- Goal: verify the behavior of the System Under Test (SUT)

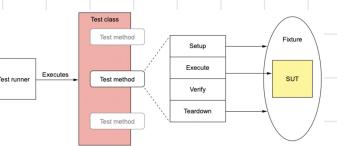
might be as small as a class
or as large as an entire application

— Writing Automated Tests

- An automated test typically consists of four phases:
 1. **Setup** — Initialize the test fixture, which consists of the SUT and its dependencies, to the desired initial state.
 2. **Exercise** — Invoke the SUT, for example, invoke a method on the class under test.
 3. **Verify** — Make assertions about the invocation's outcome and the state of the SUT.
 - (optional) 4. **Teardown** — Clean up the test fixture, if necessary. Many tests omit this phase, but some types of database test will, for example, roll back a transaction initiated by the setup phase.

— The different type of tests

(We focus on automated tests that verify the functional aspects of the application or service)



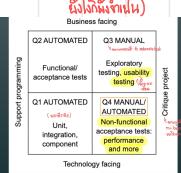
Test suite = A collection of related tests

— Testing using mocks and stub

- Replace the SUT's dependencies with **Test doubles** *
- SUT can be tested in isolation
- A **stub** is a test double that **returns values to the SUT**.
- A **mock** is a test double that a test uses to verify that the SUT correctly invokes a dependency. Also, a mock is often a stub.

(unit test vs mock library)

— Brian Marick's test quadrant (manual test vs automated)



— Test Pyramid

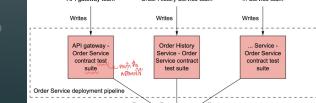


● The challenge of testing microservices

→ Layen: Many communication!?

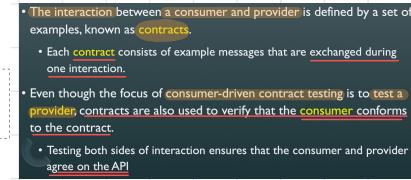
— Customer-Driven Contract Testing

- A **consumer contract test** focuses on verifying that the "shape" of a provider's API meets the consumer's expectations.
- For a REST endpoint, a contract test verifies that the provider implements an endpoint (for the consumer) that
 - Has the expected HTTP method and path
 - Accepts the expected headers, if any
 - Accepts a request body, if any
 - Returns a response with the expected status code, headers, and body



- The interaction between a consumer and provider is defined by a set of examples, known as **contracts**.

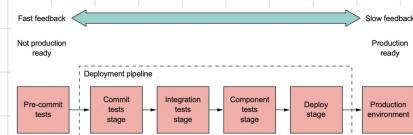
- Each **contract** consists of example messages that are exchanged during one interaction.
- Even though the focus of **consumer-driven contract testing** is to **test a provider**, contracts are also used to verify that the **consumer** conforms to the contract.
- Testing both sides of interaction ensures that the consumer and provider agree on the API



● The deployment pipeline

- Every service has a **deployment pipeline**.
- **Continuous Delivery (CD)** describes a deployment pipeline as the automated process of getting code from the developer's desktop into production.
 - Ideally, it's fully automated, but it might contain manual steps.
- A deployment pipeline is often implemented using a **Continuous Integration (CI)** server, such as **Jenkins**.

1. **Pre-commit tests stage** — Runs the unit tests. This is executed by the developer before committing their changes.
2. **Commit tests stage** — Compiles the service, runs the unit tests, and performs static code analysis.
3. **Integration tests stage** — Runs the integration tests.
4. **Component tests stage** — Runs the component tests for the service.
5. **Deploy stage** — Deploys the service into production.



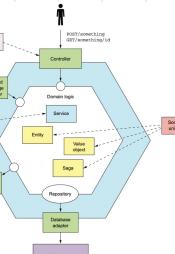
● Writing unit tests for a service

- **Unit tests** are the **lowest level** of the test pyramid.
- They're **technology-facing** tests that support development.
- A unit test verifies that a unit, which is a very small part of a service, works correctly.
- A unit is typically a class, so the goal of unit testing is to verify that it behaves as expected.
- Entities, such as Order, are objects with persistent identity → **sociable** unit tests.
- Value objects, such as Money, are objects that are collections of values → **sociable** unit tests.
- Sagas, such as CreateOrderSaga, maintain data consistency across services → **sociable** unit tests.
- Domain services, such as OrderService, are classes that implement business logic that doesn't belong in entities or value objects → **solitary** unit tests.
- Controllers, such as OrderController, which handle HTTP requests → **solitary** unit tests.
- Inbound and outbound messaging gateways → **solitary** unit tests.

- There are two types of **unit tests**:
 - **Solitary unit test** — Tests a class in isolation using mock objects for the class's dependencies
 - **Sociable unit test** — Tests a class and its dependencies



- The responsibilities of the class and its role determine which type of test to use.
 - **Controller** and **service** classes are often tested using **solitary** unit tests.
 - **Domain objects**, such as entities and value objects, are typically tested using **sociable** unit tests.

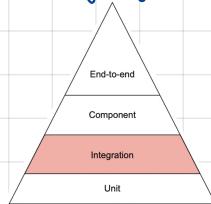


Summary

- Automated testing is the **key foundation** of rapid, safe delivery of software
- The purpose of a test is to **verify** the behavior of the system under test (SUT).
- A **good way** to simplify and speed up a test is to use test doubles.
- Use the **test pyramid** to determine where to focus your testing efforts for your services.

Testing in Microservices Part II

● Writing integration tests



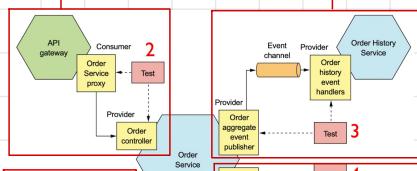
- They verify that a service can communicate with its dependencies and clients.

(Includes infrastructure services such as the database, application services.)

- ~~Testing whole services~~, The strategy is to test the individual adapter classes that implement the communication

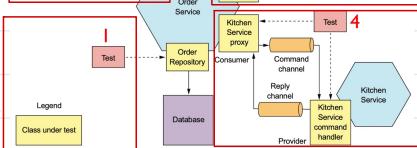
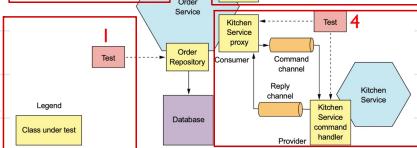
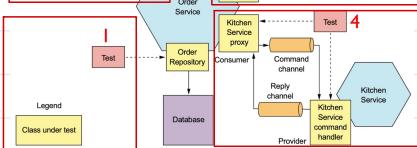
2. Integration testing REST-based request/response style interactions

- The contracts are used to verify that the adapter classes on both sides.
- The consumer-side tests verify that OrderServiceProxy invokes Order Service correctly.
- The provider-side tests verify that OrderController implements the REST API endpoints correctly.



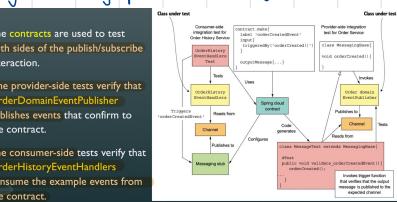
1. Persistence Integration tests

- Verifying that a service's database access logic works as expected.
- Each phase of a persistence integration test behaves as follows:
 - Setup** — Set up the database by creating the database schema and initializing it to a known state. It might also begin a database transaction.
 - Execute** — Perform a database operation.
 - Verify** — Make assertions about the state of the database and objects retrieved from the database.
 - Teardown** — An optional phase that might undo the changes made to the database.
 - for example, rolling back the transaction that was started by the setup phase.



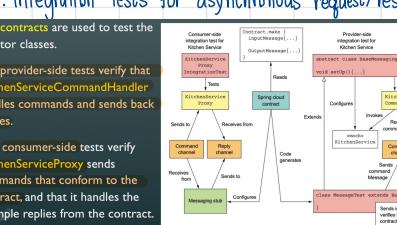
3. Integration testing publish/subscribe-style interactions

- The contracts are used to test both sides of the publish/subscribe interaction.
- The provider-side tests verify that OrderDomainEventPublisher publishes events that conform to the contract.
- The consumer-side tests verify that OrderHistoryEventHandlers consume the example events from the contract.



4. Integration tests for asynchronous request/response interactions

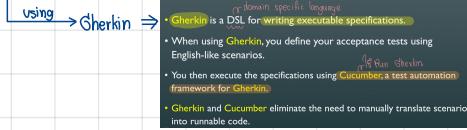
- The contracts are used to test the adaptor classes.
- The provider-side tests verify that KitchenServiceCommandHandler handles commands and sends back replies.
- The consumer-side tests verify KitchenServiceProxy sends commands to conform to the contract, and that it handles the example replies from the contract.



● Developing component tests

- We want to write the service's acceptance tests, which treat it as a black box and verify its behavior through its API.
- Component testing verifies the behavior of a service in isolation.
- It replaces a service's dependencies with stubs that simulate their behavior.

- Acceptance Tests = Business-facing tests for a software development



- They describe the desired externally visible behavior from the perspective of the component's clients rather than in terms of the internal implementation.
- These tests are derived from user stories or use cases.

● Writing end-to-end tests

- End-to-end tests are business-facing tests.
- You can write the end-to-end tests using Gherkin and execute them using Cucumber.
- The main difference is that rather than a single "Then", this test has multiple actions.

- End-to-end tests must run the entire application, including any required infrastructure services.
- The implementation of the end-to-end test is quite similar to the implementation of the component tests.
- These tests are written using Gherkin and executed using Cucumber.

Summary (1)

- Use contracts to drive the testing of interactions between services.
- Write tests that verify that the adapters of both services conform to the contracts.
- Write component tests to verify the behavior of a service via its API.
- You should simplify and speed up component tests by testing a service in isolation, using stubs for its dependencies.

Summary (2)

- Write user journey tests to minimize the number of end-to-end tests, which are slow, brittle, and time consuming.
- A user journey test simulates a user's journey through the application and verifies high-level behavior of a relatively large slice of the application's functionality.

Story

As a consumer of the Order Service I should be able to place an order

Map to execute phase

