
AgileDev
v. 0.1.9-SNAPSHOT.20150728
Guidelines

Table of Contents

1. Table of Contents	i
2. Branching Models	1
3. Javadoc guidelines	6
4. Release notes	10
5. Software Designing Principles	14
5.1. Behavior Driven Development	15
5.2. Code Reviews	16
5.3. Defensive Programming	18
5.4. Design by Contract	19
5.5. Intention revealing interfaces	20
5.6. Reactive System	21
5.7. Semantic versioning	22
5.8. Side effect free functions	23
6. Patterns	24
6.1. API gateway	26
6.2. Automated testing	28
6.2.1. Consumer-based testing	30
6.2.2. Unit testing	32
6.3. Builder	33
6.4. Bulkheads	36
6.5. Caching	37
6.6. Circuit Breaker	39
6.7. Client Side Discovery	40
6.8. Command Query Responsibility Separation	41
6.9. Deployment microservices	43
6.10. Domain Driven Design	45
6.10.1. Domain Model	47
6.10.2. Building blocks DDD	49
6.10.2.1. Aggregates	50
6.10.2.2. Entities	51
6.10.2.3. Factories	52
6.10.2.4. Layered Architecture	53
6.10.2.5. Modules	55
6.10.2.6. Repositories	56
6.10.2.7. Services	57
6.10.2.8. Value Objects	58
6.10.3. Context mapping	59
6.11. Enterprise Integration Pattern	60

6..11..1. Correlation ID	61
6..11..2. Exclusive Consumer	62
6..11..3. Message Router	63
6..12. Exception Handling	64
6..13. Microservice Architecture	66
6..14. Mocks	72
6..15. Monitoring	73
6..15..1. Canary endpoint monitoring	74
6..15..2. Log aggregation	76
6..15..3. Synthetic monitoring	77
6..16. Monolithic Architecture	78
6..17. Self registration	80
6..18. Server Side Discovery	81
6..19. Service Connector	82
6..20. Service Registry	83
6..21. Service Statelessness	84
6..22. Single Responsibility Principle	85
6..23. System of Record	86
6..24. Timeouts	88
6..25. Try-Cancel/Confirm	89
6..26. UUID	90
6..26. 3rd party registration	91
7. Glossary	
8. Abbreviations	92

1 Branching Models

Branching Models

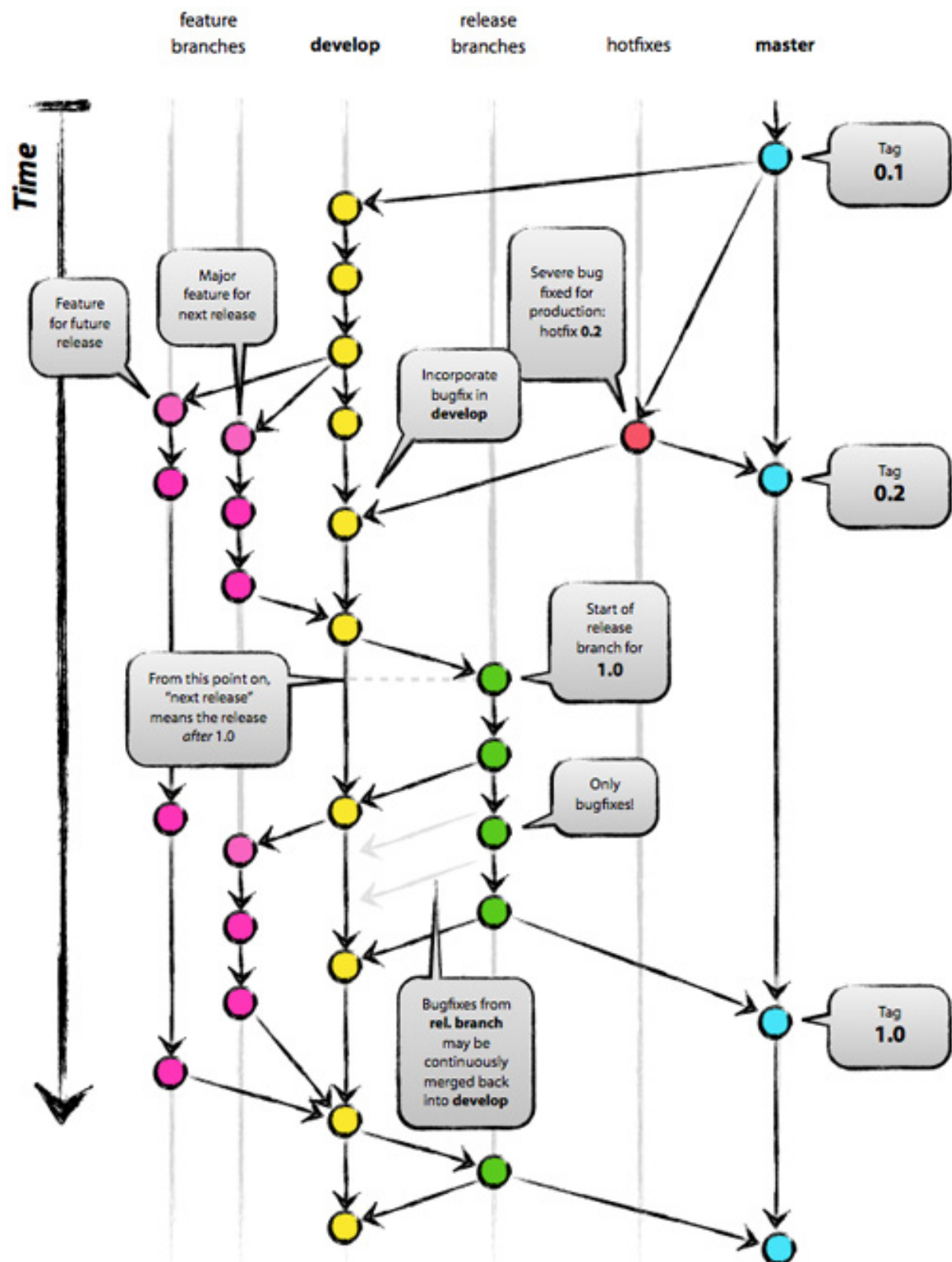
Within the following paragraphs will discuss some of the popular branching models which are around these days. After reading this you should be able to decide which one fits your project the most.

1.1 Git-Flow

This workflow has been published by [Vincent Driesen](#) as a successful branching model for git and covers most of the standard needs for a 'classical' development project.

1.1.1 How it works

Git-Flow has the following branching model:



There are 2 primary branches:

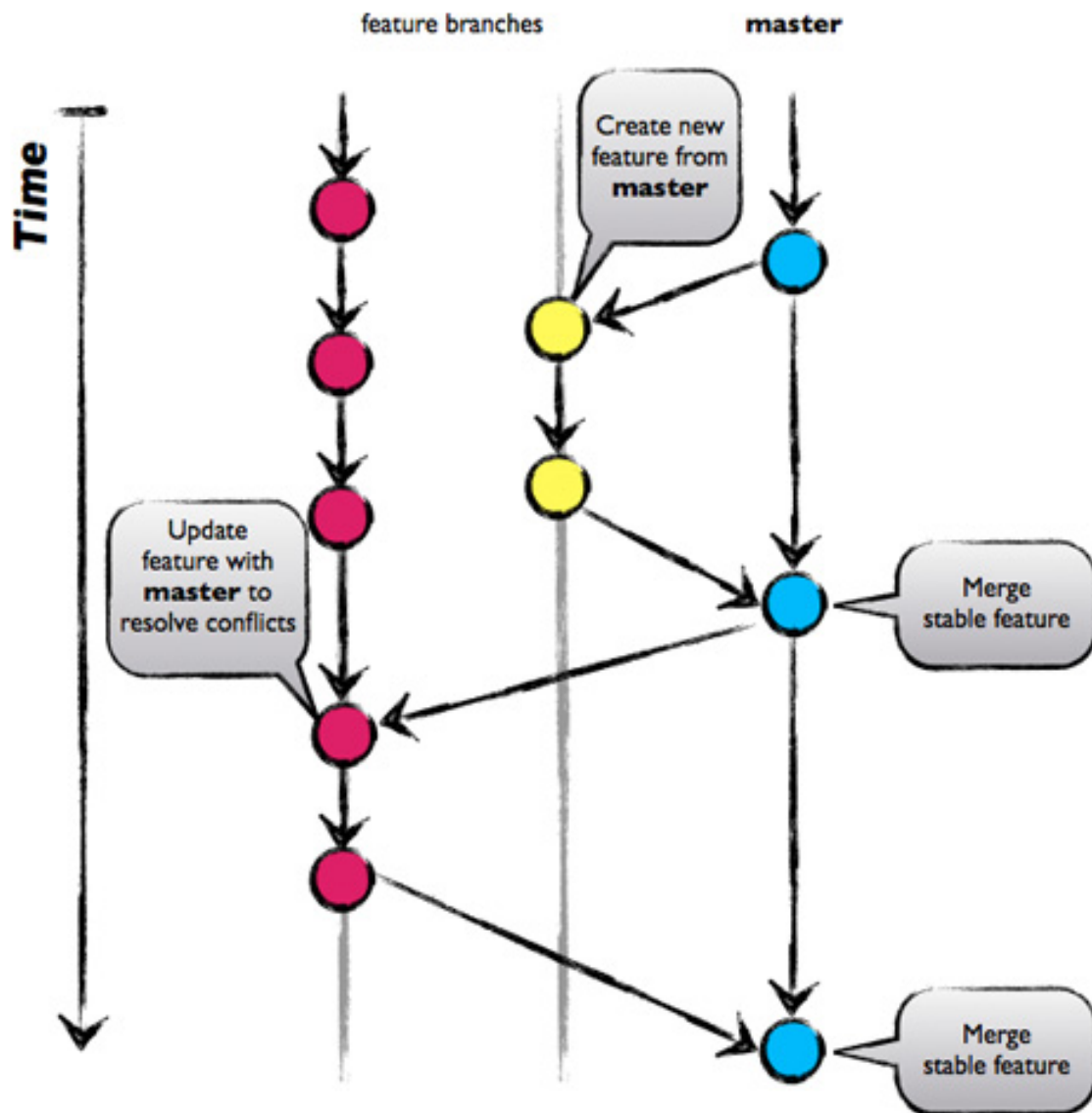
- `master` is the main branch where everything is stable. Each commit is a stable (fully tested) version of the project, (i.e. a release) which could be deployed to production and tagged accordingly.
- `develop` is the main branch where development is done. It will contain prepared changes for the next release in `master`.

And secondary branches which are flexible over time:

- `feature` starts from `develop` and merge into `develop`. When you are working on a specific feature, you create a `feature/xxx` branch and once it's done, you merge it back into `develop` to add the stable feature to the scope for the next release. The essence of a feature branch is that it exists as long as the feature is in development, but will eventually be merged back into `develop` or discarded.
- `release` starts from `develop` and merge into `master` and/or `develop`. When `develop` is reflecting the desired state of the feature release, (all features for the release haven been merged), you create a `'release/xxx'` branch. By doing this, you can prepare the next release, correct eventual bugs and continue development in parallel. All features targeted at future releases may not be merged into the `develop` branch until the `release` branch is branched off. It is exactly at the start of a release branch that the upcoming release gets a version number assigned.
- `hotfix` starts from `master` and merge into `master` and/or `develop/ release`. When you want quickly to resolve critical bugs in production you create a `hotfix/xxx` branch. When the hotfix is developed, you merge it back into `master` with the appropriate version number, and into `develop` and/or `release` branch to update it with the modifications you made.

1.2 GitHub

GitHub has the following branching model:



As you can see in the picture above, there is only a single `master` branch and the following 6 basic rules should be followed:

1. **Everything that is in `master` could be deployed in production** - The `master` branch is the only meaningful branch of the project and it should stay stable in any circumstance so you can base your work upon and deploy it to production at anytime.
2. **Create descriptive feature branches from `master`** - When you want to develop a feature or a hotfix, you just create your branch from `master` with an explicit name that describe your work.
3. **Regularly push to `origin`** - In contrary of the Git-Flow where developers doesn't have to push its local feature branch to the main repo, you have to do that regularly within the GitHub model.
4. **Open a pull-request at anytime**
5. **Only merge after a pull request review** - This is more an advice than an absolute rule. It's a best practice that another developer should review the pull request and confirm that the branch is stable. From here you can, merge the branch back into `master` and delete the merge branch.

6. **Immediately deploy after you merge into master** - Once your branch is merged into `master` the whole thing is deployed to production. Doing so, you will stress on the necessity to keep `master` stable. Developers don't want to break everything because of its modifications were deployed, so they are more likely to pay attention to code stability before merge.

The GitHub model perfectly fits projects that don't have releases nor versions.

You do continuously integrate into `master` and you deploy the stable project to production regularly; sometimes several times a day.

Due to this it's very unlikely to add series of big bugs. If problems appear, they are quickly fixed on the go. There is no difference between a big feature and a small hotfix in terms of process.

1.3 Final words

To choose the best model depends on the nature and needs of your project. It's up to you to consider which of these models fit more or less to your environment. If necessary adapt the framework to your needs.

2 Javadoc guidelines

JavaDoc standards

Javadoc is a key part of coding in Java. Within this chapter best practices regarding javadoc are shared:

2.1 Document `public` and `protected` methods

All `public` and `protected` methods **MUST** be fully defined with javadoc. `package` and `'private'` methods do not have to be, but may benefit from it.

If a method is overridden in a subclass, javadoc should only say something about what it distinct from the original definition of the method. The `@Override` annotation should be used to indicate to source code readers that teh javadoc is inherited in addition to its normal meaning.

2.2 Use simple HTML tags, not valid XHTML

Javadoc uses HTML tags to identify paragraphs and other elements. Many developers get drawn to the thought that XHTML is better and ensure that all tags open and close correctly. This is a mistake. XHTML adds many extra tags that make the javadoc harder to read within the source code. The javadoc parser will interpret the incomplete HTML tags just fine.

2.3 Use a single `<p>` tag between paragraphs

Longer javadoc always need multiple paragraphs. To seperate these from each other place a single `<p>` tag on the blank line between the paragraphs.

```
/**
 * First paragraph.
 * <p>
 * Second paragraph.
 * May be on multiple lines.
 * <p>
 * Third paragraph.
 */
public ...
```

2.4 Use a single `` tag for items in a list

Lists are useful in javadoc when explaining a set of options, choices or issues. Place in front of every item `` at the start of the line and no closing tag. In order to get correct paragraph formatting, extra paragraph tags are required.

```
/**
 * First paragraph.
 * <p><ul>
 * <li>the first item
 * <li>the second item
 * <li>the third item
 * </ul><p>
 * Second paragraph.
 */
public ...
```

2.5 Define a punchy first sentence

The first sentence, ended by a dot, is used in the next level higher javadoc. As such it has the responsibility of summing up the method or class. To achieve this the, first sentence should be clear and punchy, and generally short.

While not required, it is recommended that the first sentence is a paragraph itself. This helps retain the punchiness for the readers of the code.

2.6 Use 'this' to refer to an instance of the class

When referring to an instance of the class being documented, use 'this' to reference it. For example Returns a copy of this.

2.7 Aim for short single line sentences

Whenever possible, make javadoc sentences fit on a single line (80 till 120 characters). In most cases, each new sentence should start on a new line. This aids readability as source code, and simplifies refactoring re-writes of complex javadoc.

```
/**
 * This is the first paragraph, on one line.
 * <p>
 * This is the first sentence of the second paragraph, on one line.
 * This is the second sentence of the second paragraph, on one line.
 * This is the third sentence of the second paragraph which is a bit longer so has
 * split onto a second line, as that makes sense.
 * This is the fourth sentence, which starts a new line, even though there is space
 */
public ...
```

2.8 Use @link and @code

Many javadoc descriptions reference other methods and classes. This can be achieved most effectively using the @link and @code features.

The @link feature creates a visible hyperlink in generated Javadoc to the target. The @link target is one of the following forms:

```

/**
 * First paragraph.
 * <p>
 * Link to a class named 'Foo': {@link Foo}.
 * Link to a method 'bar' on a class named 'Foo': {@link Foo#bar}.
 * Link to a method 'baz' on this class: {@link #baz}.
 * Link specifying text of the hyperlink after a space: {@link Foo the Foo class}.
 * Link to a method handling method overload {@link Foo#bar(String,int)}.
 */
public ...

```

The `@code` feature provides a section of fixed-width font, ideal for references to methods and class names. While `@link` references are checked by the Javadoc compiler, `@code` references are not.

Only use `@link` on the first reference to a specific class or method. Use `@code` for subsequent references. This avoids excessive hyperlinks cluttering up the Javadoc.

2.9 Never use `@link` in the first sentence

The first sentence is used in the higher level Javadoc. Adding a hyperlink in that first sentence makes the higher level documentation more confusing. Always use `@code` in the first sentence if necessary. `@link` can be used from the second sentence/paragraph onwards.

2.10 Do not use `@code` for `null`, `true` or `false`

The concepts of `null`, `true` and `false` are very common in Javadoc. Adding `@code` for every occurrence is a burden to both the reader and writer of the Javadoc and adds no real value.

2.11 Use `@param`, `@return` and `@throws`

Almost all methods take in a parameter, return a result or both. The `@param` and `@return` features specify those inputs and outputs. The `@throws` feature specifies the thrown exceptions.

The `@param` entries should be specified in the same order as the parameters. The `@return` should be after the `@param` entries, followed by `@throws`.

2.12 Use `@param` for generics

If a class or method has generic type parameters, then these should be documented. The correct approach is an `@param` tag with the parameter name of `<T>` where T is the type parameter name.

2.13 Use one blank line before `@param`

There should be one blank line between the Javadoc text and the first `@param` or `@return`. This aids readability in source code.

2.14 Treat `@throws` as an if clause

The `@throws` feature should normally be followed by “if” and the rest of the phrase describing the condition. For example, `@throws IllegalArgumentException if the file could not be found`. This aids readability in source code and when generated.

2.15 Define null-handling for all parameters and return types

Whether a method accepts null on input, or can return null is critical information for building large systems. All non-primitive methods should define their null-tolerance in the @param or @return. Some standard forms expressing this should be used wherever possible:

- “not null” means that null is not accepted and passing in null will probably throw an exception , typically `NullPointerException`
- “may be null” means that null may be passed in. In general the behaviour of the passed in null should be defined
- “null treated as xxx” means that a null input is equivalent to the specified value
- “null returns xxx” means that a null input always returns the specified value

Other simple constraints may be added as well if applicable, for example “not empty, not null”. Primitive values might specify their bounds, for example “from 1 to 5”, or “not negative”.

```
/**
 * Javadoc text.
 *
 * @param foo the foo parameter, not null
 * @param bar the bar parameter, null returns null
 * @return the baz content, null if not processed
 */
public String process(String foo, String bar) {...}
```

2.16 Avoid @author

The @author feature can be used to record the authors of the class. This should be avoided, as it is usually out of date, and it can promote code ownership by an individual. The source control system is in a much better position to record authors.

3 Release notes

Release Notes

Software release notes convey two important things:

1. They explain what changes, deficiencies and defects were addressed in the release.
2. They inform the reader about installation and deployment.

The audience for release notes can be testers, support staff and management.

Public release notes should contain at least:

- Product Name
- Version Number (Be consistent with major.minor.revision numbering in your releases, see [semantic versioning](#) for more details).
- release, buildnumber
- all fixed public bugs
- all added public features
- Release Date (name, version, and release date often repeated in the document footer)
- Type of Release (e.g. “This is a developer release for internal evaluation only.” or “This is a public production release.”)
- Major announcement (Not often used, but important when you need to note “This release corrects a major flaw in XYZ from the previous release. Please upgrade as soon as possible.”)
- What’s New (A List of major enhancements visible to the user, in user-friendly terms, e.g. “increased response time” versus “manipulated line 9 of widget #745 to change algorithm”)
- Installation Instructions (Often is a list of system requirements with a link to a full installation document, but if there are compatibility issues or deviations from previous version instructions, note them.)
- Headings for: Additions, Removals, Changes, Bugfixes (The first three mean additions, removals, changes to functionality; the latter means fixes of anything. The info under each heading can come right out of your bug-tracking system and should be simple bullet points like “FIX ### Title_of_Bug”. This assumes tickets are given meaningful titles. Bonus if you can link users to the bug-tracking system so they can see the entire ticket and commentary.)
- Known Issues (List of open bugs slated for next release, or prose describing known issues; this depends on your audience and what the company wants to make transparent, and how.)
- Additional Documentation (links to user guide, administrator’s guide, other relevant documents)

The release notes should be published in plain text or at the very least html.

3.1 Generate release notes from JIRA

With the `maven changes plugin` one can generate release notes from JIRA tasks and also mail it to a list of email addresses.

To generate the release notes with the `maven changes plugin` you should include the following snippet into your `pom.xml`

```

<plugin>
  <groupId>org.apache.maven.plugins</groupId>
  <artifactId>maven-changes-plugin</artifactId>
  <version>${maven-changes-plugin.version}</version>
  <executions>
    <execution>
      <id>generate-release-notes</id>
      <phase>generate-resources</phase>
      <goals>
        <goal>announcement-generate</goal>
      </goals>
      <configuration>
        <!-- This will generate release-notes.txt file under META-INF folder within
        <announcementFile>release-notes.txt</announcementFile>
        <announcementDirectory>${project.build.outputDirectory}/META-INF</announcementDirectory>
      </configuration>
    </execution>
  </executions>
  <configuration>
    <issueManagementSystems>
      <!-- generate announcement based on both `changes.xml` and JIRA -->
      <issueManagementSystem>changes.xml</issueManagementSystem>
      <issueManagementSystem>JIRA</issueManagementSystem>
    </issueManagementSystems>

    <smtpHost>mail.yourhost.com</smtpHost>
    <smtpPort implementation="java.lang.Integer">25</smtpPort>

    <!-- Specifies the sender of release notes -->
    <mailSender>
      <name>Release Notification</name>
      <email>build@example.com</email>
    </mailSender>

    <!-- Recipient email addresses -->
    <toAddresses>
      <toAddress implementation="java.lang.String">to@example.com</toAddress>
    </toAddresses>

    <!-- Use the JIRA query language -->
    <useJql>true</useJql>

    <!--
    Take the version from the project's POM and match it against the
    Fix-For version of the JIRA issues. The names of your versions in JIRA
    must match the ones you use in your POM. The -SNAPSHOT part of the
    version in your POM is handled automatically by the plugin, so you
    don't need to include -SNAPSHOT in the names of your versions in JIRA.
    <onlyCurrentVersion>true</onlyCurrentVersion>
    -->
    <!--
    IDs of version(s) you want to include in the announcements
    These are JIRA's internal version IDs, NOT the human readable display
    ones. Multiple version(s) can be separated by commas. If this is set
    to empty, that means all fix versions will be included.
    -->
    <fixVersionIds>11412</fixVersionIds>

    <!-- generate announcements only at the top of module tree -->
    <runOnlyAtExecutionRoot>true</runOnlyAtExecutionRoot>

    <!-- only include JIRA issues with resolution IDs Fixed and Done -->
  
```


To include links to the issues in your issue management system you have to enter the type of issue management system (see [Usage maven changes plugin](#) for list of pre-configured issue management systems) and the URL to it:

```
<project>
  ...
  <issueManagement>
    <system>JIRA</system>
    <url>http://jira.company.com/</url>
  </issueManagement>
  ...
</project>
```

Note: Make sure that your `<issueManagement>/<url>` is correct. In particular, make sure that it has a trailing slash if it needs one.

To generate the release notes run this command:

```
mvn changes:jira-report
```

3.2 Related patterns

- [Semantic versioning](#)

See also:

- [Release notes from JIRA](#)

4 Software Designing Principles

Software Design Principles

Software design principles represent a set of guidelines at the highest level. The design principles define the overall generic shape and structure of software applications. A level lower are the (architectural) patterns which focus on a specific problem.

There are 4 primary symptoms that indicate that the design is rotten:

- *Rigidity* - Rigidity is the tendency for software to be difficult to change even in case of simple changes. Every change causes a cascade of subsequent changes in dependent modules.
- *Fragility* - Closely related to rigidity is fragility. Fragility is the tendency of the software to break (in many places) every time it is changed. Often the breakage occur in places that have no conceptual relationship with the area that was changed.
- *Immobility* - Immobility is the inability to reuse software from other projects or from parts of the same project.
- *Viscosity* - Viscosity comes in 2 forms: viscosity of the design, and viscosity of the environment. When changes in which the design is preserved are harder to employ than the hacks, then the viscosity of the design is high. Viscosity of environment comes about when the development environment is slow and inefficient. In this case developers are tempted to make changes that prevent a lot of changes, which are usually not optimal from a design point.

These four symptoms are the signs of poor architecture. Within this part we will focus on software design principles by which we can prevent this.

5 Behavior Driven Development

Behavior Driven Development

Behavior Driven Development is a software development process that emerged from [Test Driven Development](#). Behavior Driven Development combines the general techniques and principles of TDD with ideas from [Domain Driven Design](#). TDD states that for each unit of software (class in case of OO), a software developer must:

1. Define a test set for the unit *first*.
2. Then implement the unit.
3. Verify that the implementation of the unit makes the tests succeed.

Extension to this is that within BDD the tests should be defined in terms that has business value; e.g. the tests should be written in a form that is using business terminology (defined within the [ubiquitous language](#)).

On top of that BDD defines how the desired behavior should be specified. Business analysts and developers should work together and specify the desired behavior in the form of user stories that have the following structure:

```
Title: Explicit and clear title

Narrative: Short introduction that specifies
  * Who is the driver or primary stakeholder of the story, (business actor).
  * What effect the story would have.
  * What benefits the stakeholder will derive from this effect.

Scenarios: Description of each specific case according to the following structure

  * Given: Specifies the initial situation. This may consist of
           single clause, or several (combined by AND).
  * When : States which event triggers the start of the scenario.
  * Then : States the expected outcome in one or more clauses.
```

5.1 Related Patterns

- [Test Driven Development](#) - Behavior Driven Development emerged from Test Driven Development.
- [Domain Driven Design](#) - Behavior Driven Development borrows the concept of the ubiquitous language from DDD to define the tests in business terminology.

5.2 See also

- [Cucumber](#) - BDD testing framework for several languages, (Ruby, Java, ...).
- [An Introduction to BDD Test Automation with Serenity and JUnit](#)

6 Code Reviews

Code Reviews

Code reviews are a systematic examination (also known as peer reviews) of source code. It is intended to find and fix mistakes overlooked in the initial development phase, improving both the overall quality of software and the developers skills.

6.1 Best practices code reviews

- *Review fewer than 200-400 lines of code at a time* - For optimal effectiveness, developers should review fewer than 200-400 lines of code (loc) at a time.
- *Aim for an inspection rate of fewer than 300-500 lines of code per hour* - Take your time with code review. Research has shown that you will achieve optimal results at an inspection rate of less than 300-500 lines of code per hour.
- *Never review code for more than 90 minutes at a stretch* - After about 60 minutes reviewers simply get tired and stop finding additional defects. It's generally known that when people engage in any activity requiring concentrated effort, performance starts dropping off after 60-90 minutes. Given these human limitations, a reviewer will probably not be able to review more than 300-600 lines of code before performance drop. On the flip side, you should always spend at least 5 minutes reviewing code, even if it's just one line. Often, a single line can have consequences throughout the system, and it's worth the five minutes to think through the possible effects that a change could have.
- *Have authors annotate source code before the review begins* - Annotations guide the reviewers through the changes, showing which files to look at first and defending the reason and methods for each code modification. These notes are not comments in the code, but rather comments given to other reviewers. By annotating the source code the developer is required to double-check his work which will uncover many of the defects before the review even begins.
- *Establish quantifiable goals for code reviews, and capture metrics so you can improve your process* - When you have defined specific goals, you will be able to judge whether peer review is truly achieving the results that you require. It's best to start with external metrics such as "reduce support calls by 20%", or "halve the percentage of defects injected by development". However, it can take a while before external metrics show results.
- *Use checklists* - Checklists are a highly recommended way to check for the things that you might forget to do. Omissions are the hardest defects to find. as soon as you start recording your defects in a checklist, you will start making fewer of these errors.
- *Verify that the defects are actually fixed* - Many teams that review code don't have a good way of tracking defects found during review and ensuring that bugs are actually fixed before the review is complete.
- *Foster a good review culture in which finding defects is viewed positively* - It's easy to see defects as a bad thing, but fostering a negative attitude toward defects found can sour a whole team. It must be promoted that defects are positive. After all, each one is an opportunity to improve the code, and the goal of the code review process is to make the code as good as possible. Every defect found and fixed in peer review is a defect that a customer never sees and another problem QA doesn't have to spend time tracking down. Reviews present opportunities for all developers to correct bad habits, learn new tricks, and expand their capabilities.
- *Beware of the big brother effect* - Metrics are vital for process measurement, which form the basis for process improvement. If developers believe that metrics will be used against them, not only will they be hostile to the process, but they will probably focus on improving their metrics rather than truly writing better code and being more productive.
- *Review at least part of the code to benefit from the ego effect* - The ego effect drives developers to write better code because they know that others will be looking at their code and metrics.

No one wants to be known as the guy who makes all those mistakes. The ego effect drives developers to review their own work carefully before passing it on to others.

7 Defensive Programming

Defensive Programming

Defensive Programming is the practice of anticipating all possible ways that an end user could misuse an software system, and designing the system in such way that this is impossible, or to minimise the negative consequences. Goal of defensive programming is:

- Reducing the number of software bugs.
- Making software understandable; the source code should be readable and understandable and verified during code audit.
- Making the software behave in a predictable manner despite unexpected inputs or user actions.

Overly defensive programming could introduce code to prevent errors that can't happen, but needs to be executed on runtime and to be maintained by the developers. There is also the risk that the code catches or prevents too many exceptions. In those cases, the error would be suppressed and go unnoticed while the result would still be wrong.

Some defensive programming techniques are:

- Code reuse - Existing code is tested and known to work, reusing it may reduce the change of bugs being introduced.
- Legacy problems - Before reusing old libraries, APIs, and so forth, it must be validated whether the old work is valid for reuse. The library might for example have a much lower quality than the newly designed system.
- Low tolerance against 'potential' bugs - Assume that code constructs that appear to be problem prone (for example reported by a source code analyzer) are bugs and potentially security flaws.
- Encrypt/authenticate all data transmitted over networks. Do not implement your own encryption scheme, but use a proven one instead.
- [Design by Contract](#) - Use Design by Contract methodology to ensure that provided data (and the state of the program as a whole) is verified.
- Error codes - Prefer exceptions to return error codes, see [error handling](#) for more details.

7.1 See also:

- [PMD](#) - 'Free' source code analyzer to detect potential bugs.

8 Design by Contract

Design By Contract

Design by Contract was first explored by Bertrant Meyer. He has invented a language named Eiffel in which contracts are explicitly stated for each method, and checked at each invocation. The contracts define preconditions, postconditions and invariants.

- Preconditions are certain expectations that need to be guaranteed before a client is allowed to perform certain functionality.
- Postconditions are guarantees on exit of certain functionality.
- An invariant is a certain constraint that must be true during the whole lifecycle of the entity.

Design by Contract advocates writing the conditions and invariants first. Constraints can be written by annotations and enforced by a test suite.

The contract will be applied on the method and will normally contain the following pieces of information:

- Acceptable and unacceptable input values/types and their meanings.
- Return values/types and their meaning.
- Exceptions that can occur during execution and their meaning.
- Side effects (for example notifications that are published, etc).
- Preconditions
- Postconditions
- Invariants

In rare cases performance guarantees are also part of the contract.

When applying design by contract a client should not try to verify that the contract conditions are satisfied. Instead the service will throw an exception and the client should be able to cope with that exception.

Design by contract will facilitate code reuse since the contract documents each piece of the code fully; e.g. the contract for a method can be regarded as a form of software documentation for the behavior of that module.

8.1 Related Patterns

- [Defensive Programming](#) - In case of multi-channel client server or distributed computing the defensive design approach should be taken, meaning that a server component tests (before or while processing a client's request) that all relevant preconditions hold true.
- [Behavior Driven Development](#) - The enforcing of constraints can be applied by following the Behavior Driven Development methodology.

TODO: http://en.wikipedia.org/wiki/Design_by_contract

9 Intention revealing interfaces

Intention Revealing Interfaces.

If a developer must consider the implementation of a component in order to use it, the value of encapsulation is lost. If someone other than the original developer must infer the purpose of an object or operation based on its implementation, that new developer may infer a purpose that the operation or class fulfills only by chance. If that was not the intent, the code may work for that moment, but the conceptual basis of the design will have been corrupted, and the 2 developers will be working at cross purpose.

Therefore, name classes and operations to describe their effect and purpose without reference to the means by which they do what they promise. This relieves the client developer of the need to understand the internals. These names should conform to the ubiquitous language so that team members can quickly infer their meaning. Write a test for a behavior before creating it, to force your thinking into client developer mode.

10 Reactive System

Reactive System

Systems built as Reactive Systems are more flexible, loosely-coupled and scalable. This makes them easier to develop and amenable to change. They are significantly more tolerant of failure and when failure does occur they meet it with elegance rather than disaster. Reactive Systems are highly responsive, giving users effective interactive feedback.

10.1 What are reactive systems

- *Responsive*: The system responds in a timely manner if at all possible. Responsiveness is the cornerstone of usability and utility, but more than that, responsiveness means that problems may be detected quickly and dealt with effectively. Responsive systems focus on providing rapid and consistent response times, establishing reliable upper bounds so they deliver a consistent quality of service. This consistent behaviour in turn simplifies error handling, builds end user confidence, and encourages further interaction.
- *Resilient*: The system stays responsive in the face of failure. This applies not only to highly-available, mission critical systems — any system that is not resilient will be unresponsive after a failure. Resilience is achieved by replication, containment, isolation and delegation. Failures are contained within each component, isolating components from each other and thereby ensuring that parts of the system can fail and recover without compromising the system as a whole. Recovery of each component is delegated to another (external) component and high-availability is ensured by replication where necessary. The client of a component is not burdened with handling its failures.
- *Elastic*: The system stays responsive under varying workload. Reactive Systems can react to changes in the input rate by increasing or decreasing the resources allocated to service these inputs. This implies designs that have no contention points or central bottlenecks, resulting in the ability to shard or replicate components and distribute inputs among them. Reactive Systems support predictive, as well as Reactive, scaling algorithms by providing relevant live performance measures. They achieve elasticity in a cost-effective way on commodity hardware and software platforms.
- *Message Driven*: Reactive Systems rely on asynchronous message-passing to establish a boundary between components that ensures loose coupling, isolation, location transparency, and provides the means to delegate errors as messages. Employing explicit message-passing enables load management, elasticity, and flow control by shaping and monitoring the message queues in the system and applying back-pressure when necessary. Location transparent messaging as a means of communication makes it possible for the management of failure to work with the same constructs and semantics across a cluster or within a single host. Non-blocking communication allows recipients to only consume resources while active, leading to less system overhead.

Large systems are composed of smaller ones and therefore depend on the Reactive properties of their constituents. This means that Reactive Systems apply design principles so these properties apply at all levels of scale, making them composable.

11 Semantic versioning

Semantic versioning

Semantic versioning are providing strict guidelines regarding version numbers.

Given a version number `MAJOR.MINOR.PATCH`, increment the:

- **MAJOR** version when you make incompatible API changes.
- **MINOR** version when you add/change functionality which is backward compatible.
- **PATCH** version when you make backward compatible bug fixes.

The numbers that are used for `MAJOR`, `MINOR`, and `PATCH` are non negative integers, and don't contain leading zeroes. A `MAJOR` version 0 is for initial development and anything may change at any time. Version `1.0.0` defines the first public stable version.

A pre-release version may be denoted by appending a hyphen and a series of dot separated identifiers immediately following the `PATCH` number. The pre-release version identifier must comprise only ASCII alphanumerics (`[0-9A-Za-z]`). A pre-release version indicates that the version is unstable and might not satisfy the intended compatibility requirements as denoted by its associated normal version. An example of a pre-release version is `1.0.0-alpha`.

Build metadata may be denoted by appending a plus sign and a series of dot separated identifiers immediately following the `PATCH` or pre-release version number. The build metadata must comprise only ASCII alphanumerics (`[0-9A-Za-z]`). An example of a version number with build metadata is `1.0.0+201506221617`

Once a versioned package is released, the content of that package **MUST NOT** be modified. Any modifications **MUST** be released as a new version.

11.1 See also

- [Semantic versioning](#)
- [Best practices for Artifact versioning in Service Oriented Systems](#)
- <http://www.sitepoint.com/semantic-versioning-why-you-should-using/>

12 Side effect free functions

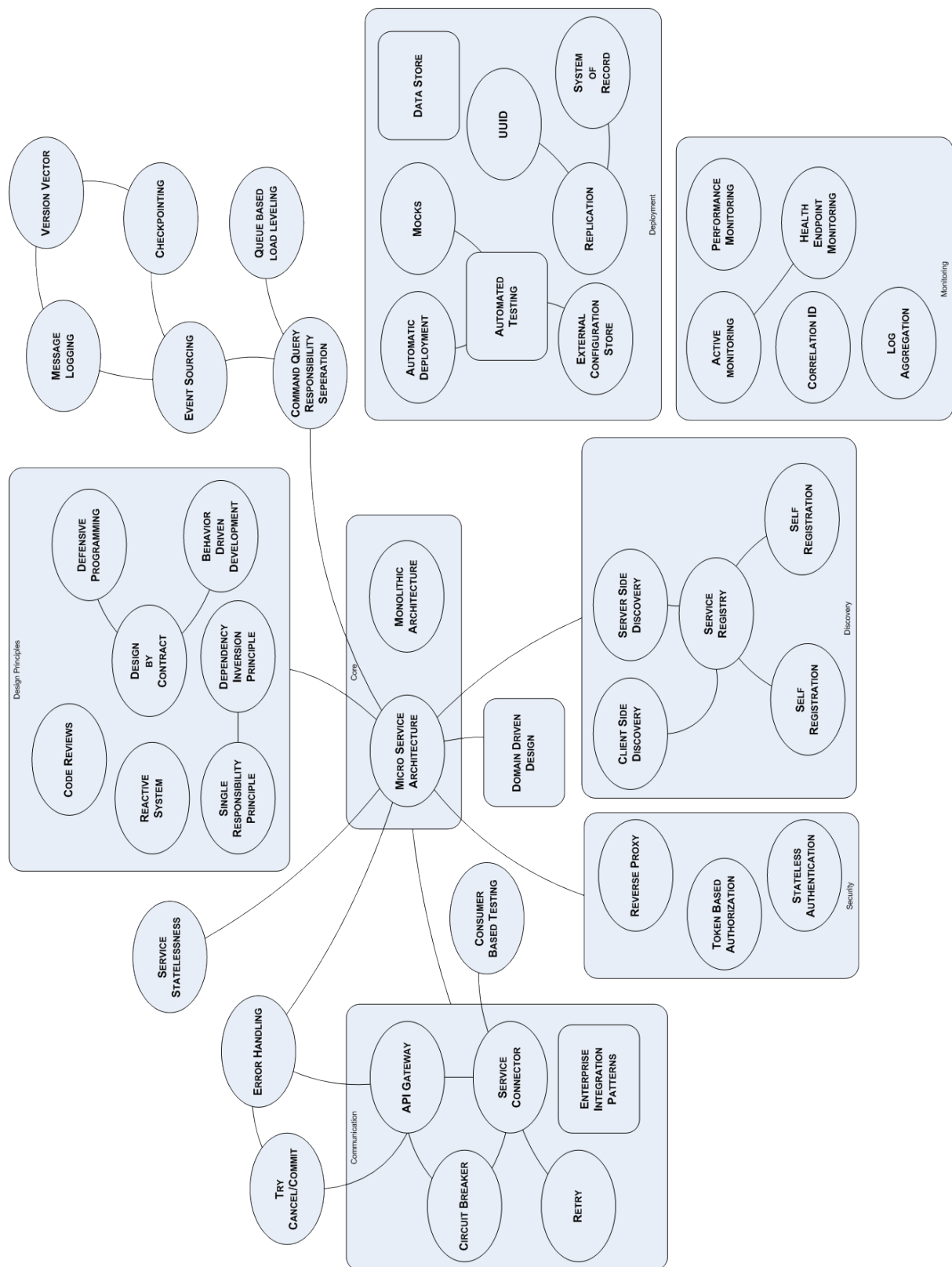
Side effect free functions.

Interactions of multiple rules or compositions of calculations become extremely difficult to predict. The developer calling an operation must understand its implementation and the implementation of all its delegations in order to anticipate the result. The usefulness of any abstraction of interfaces is limited if the developers are forced to pierce the veil.

Therefore, place as much as possible of the program into functions/operations that return results with no observable side effects. Strictly segregate method which result in modifications to observable state into very simple operations that do not return domain information. Further control side effects by moving complex logic into `Value Objects` when a concept fitting the responsibility present itself.

13 Patterns

.....
Patterns

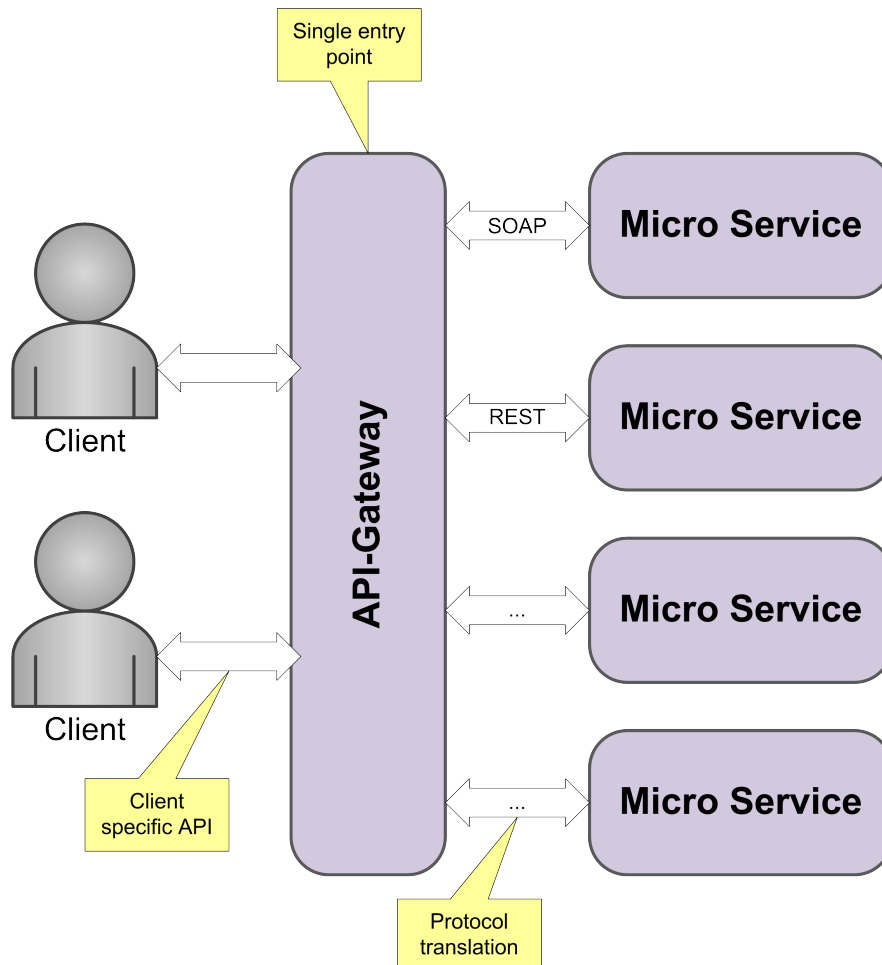


14 API gateway

API Gateway

Your system must expose to clients functionality that they can use. However the granularity of the APIs provided by the microservices is often different from what the client needs. Microservice APIs typically provide fine-grained APIs, which means that clients need to interact with multiple services. Also the number of service instances and their locations (hoss, port, ...) changes dynamically, and the partitioning of the services can change over time. All this should be hidden from the clients.

To solve this we expose an API gateway that is the single entry point for all clients.



The API gateway handles the requests in the following 2 ways:

1. Requests are simply proxied/routed to the appropriate service.
2. Requests are fanning out to multiple services

Rather than exposing a one size fits all style API, the API gateway can expose a different API for each client. The API Gateway might also verify that the client is authorized to perform the request.

Related Patterns

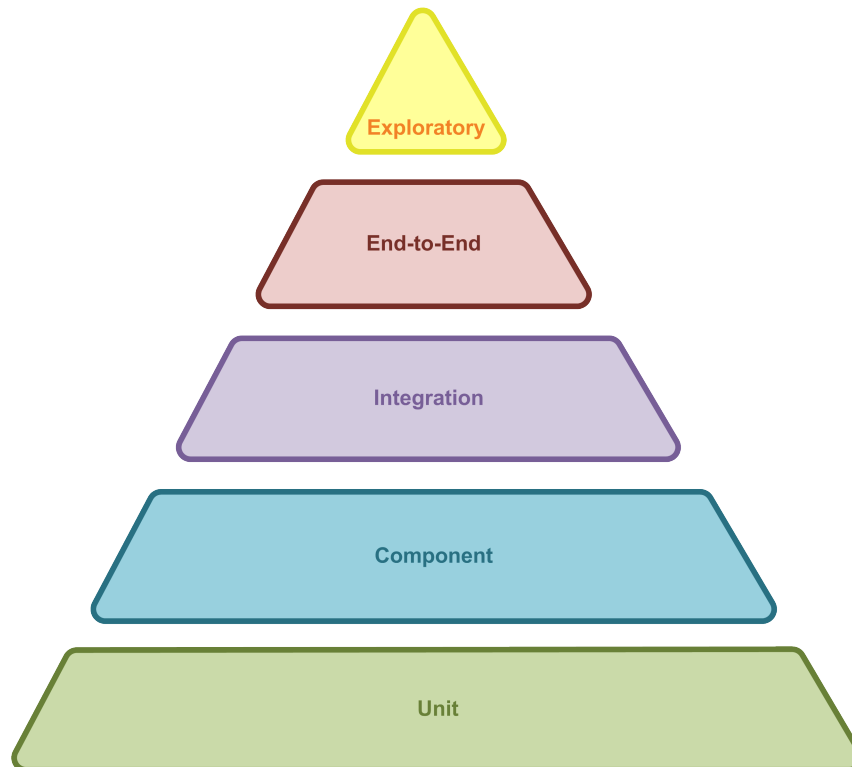
- [Service Connector](#) - Provide high level interface that hides implementation details regarding communication, thereby making the use of the microservice easier.

- [Try-Cancel/Commit](#) - By this pattern we can realize distributed transactions for microservices.

15 Automated testing

Automated testing

TODO



The testing pyramid essentially points out that you should have many more low level unit tests than high level end-to-end tests running. The layers we can distinct are:

- Exploratory -
- End-to-End
- Integration - By integration tests the interaction between two or more services is explicitly tested.
- Component -
- Unit - In this context, a 'unit' is often a function or a method of a class instance. See [unit testing](#) for more details.

15.1 Test Data

One of the harder challenges within automated testing is generating valuable test data. Hard coding assumptions about data availability can be a fragile approach as there are no guarantees that the data continue to exist. Furthermore, some tests may require data consistency across multiple services.

One of the robust strategies is to create the test data during the test run, as you guarantee the data exists before using it. This requires however that each service allows creating new resources, which isn't always the case. A solution for this could be to expose in the test environments test-only endpoints to facilitate test data creation. These end-points are of course not exposed in the production environment.

An alternative strategy is to have each service publish a cohesive set of test data that is guaranteed to be stable.

15.2 See also

- [Consumer-based testing] - To verify the integration between multiple services we could rely on consumer based testing.

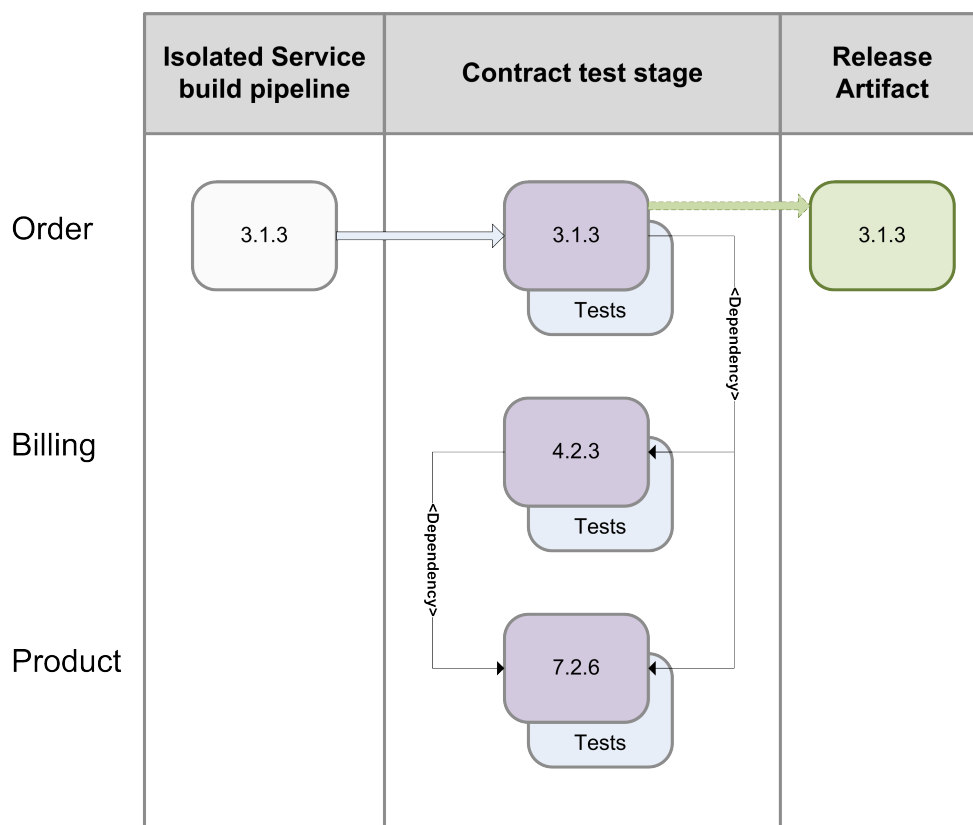
16 Consumer-based testing

Consumer-based testing

Consumer based testing is counter intuitive as it relies on the consumer writing tests for the producer. When writing contract tests, a consumer writes a test against a service it uses to confirm that the service contract satisfies the consumer needs.

By doing this we enable a neat trick in the deployment pipeline. After a service passes its internal build and QA process, all services and consumers go through a unified integration test stage. It could be triggered by the consumer changing tests, or the producer committing changes to the service. During this stage each consumer runs its tests against the new version of the changed service. Any failure prevents the new version from progressing in the pipeline.

Now for example, if we would have a `Order` Service that depends on the `Product` and `Billing` service, a new build of `Order` would trigger the execution of its consumer based tests against the latest version of `Product` and `Billing` service to pass the integration test stage.



Triggering only tests that are associated with a particular change can get tricky, however you can go a long way by simply running all contract tests each time a new service is deployed to the integration test stage within the pipeline.

This would mean that the consumer based tests of billing would be included (gray arrow), which aren't relevant to the change that was introduced.

Assuming that all the tests pass, we have a set of services that have been proven to work together. We can record the set of them working together by creating a Deployable Artifact Set (DAS). This DAS

can become the single deployable artifact for higher stages within the deployment pipeline, or it can become a compatibility reference.



The DAS could become the input for creating a container that could be deployed to other test environments.

16.1 See also

- [Service connector] - Since all services are accessed through a service connector the consumer based tests to verify whether the service satisfies the needs of the consumer will be included in the service connector.

17 Unit testing

Unit Testing

Unit testing exercises the smallest piece of testable software in the application to determine whether it behaves as expected. In this context, a ‘unit’ is often a function or a method of a class instance. Often, difficulty in writing a unit test can highlight when a module should be broken down into independent more coherent pieces and tested individually. Thus alongside being a useful testing strategy, unit testing is also a powerful design tool, especially when combined with [behavior driven development](#).

Within unit testing we can make a distinction based on whether or not the unit under test is isolated from its collaborators:

- **Sociable unit testing** - focuses on testing the behaviour of units by observing changes in their state. This treats the unit under test as a black box tested entirely through its interface.
- **Solarity unit testing** - looks at the interactions and collaborations between an object and its dependencies, which are replaced by [mocks](#).

These styles are not competing and are frequently used in the same codebase to solve different testing problems.

The domain logic within a microservice often manifests as complex complex calculations and a collection of state changes. Since this kind of logic is highly state based there is often little value in trying to isolate the units. This means that we should apply solarity unit testing to the [domain model](#).

With plumbing code like [repositories](#), and [service connectors](#) the purpose of the unit tests is to verify the logic used to produce requests or map responses from external dependencies. As such using [mocks](#) provides a way to control the request-response cycle in a reliable and repeatable manner. Advantage of testing these components at this level is that they provide faster feedback than integration tests and we can force error conditions by having the mocks replicate error conditions.

Coordination logic such as [services](#) care more about the messages passed between the modules and the domain than any complex logic. Using mocks allows the changes to the domain to be verified. If a service requires too many mocks, it is usually a good indicator that some concepts should be extracted and tested in isolation.

17.1 See also

- [ContiPerf] - A Java library for measuring performance by running JUnit tests. This could provide a early hint regarding performance bottle necks.

18 Builder

Builder

Purpose of the builder pattern is to reduce the number of parameters required for the construction of an instance.

Within the builder pattern we define a static inner class `Builder`, whose job is to collect the parameters and then construct the object in one fell swoop

```
public class Pizza {
    private int size;
    private boolean cheese;
    private boolean pepperoni;
    private boolean bacon;

    public static class Builder {
        //required
        private final int size;

        //optional
        private boolean cheese = false;
        private boolean pepperoni = false;
        private boolean bacon = false;

        public Builder(int size) {
            this.size = size;
        }

        public Builder cheese(boolean value) {
            cheese = value;
            return this;
        }

        public Builder pepperoni(boolean value) {
            pepperoni = value;
            return this;
        }

        public Builder bacon(boolean value) {
            bacon = value;
            return this;
        }

        public Pizza build() {
            return new Pizza(this);
        }
    }

    private Pizza(Builder builder) {
        size = builder.size;
        cheese = builder.cheese;
        pepperoni = builder.pepperoni;
        bacon = builder.bacon;
    }
}
```

Note that `Pizza` is immutable and that parameter values are all in a single location. Because the `Builder`'s setter methods return the `Builder` object they can be chained:

```
Pizza pizza = new Pizza.Builder(12)
                .cheese(true)
                .pepperoni(true)
                .bacon(true)
                .build();
```

This results in code that is easy to write and easy to understand. The `build()` method could be modified to check parameters after the `Pizza` instance has been instantiated and throw an `IllegalStateException` if an invalid value has been supplied.

18.1 When to use the builder pattern

We all have seen constructors where each version adds a new optional parameter:

```
Pizza(int size) { ... }
Pizza(int size, boolean cheese) { ... }
Pizza(int size, boolean cheese, boolean pepperoni) { ... }
Pizza(int size, boolean cheese, boolean pepperoni, boolean bacon) { ... }
```

This is called the telescoping constructor pattern. The problem with this pattern is that once constructors are 4 or 5 parameters long it becomes difficult to remember the required order of the parameters, as well as what particular constructor you might want in a given situation.

An alternative you have to the telescoping constructor pattern is the JavaBean pattern where you call a constructor with the mandatory parameters and then call any optional setter after:

```
Pizza pizza = new Pizza(12);
pizza.setCheese(true);
pizza.setPepperoni(true);
pizza.setBacon(true);
```

The problem here is that, because the object is created over several calls, it may be in an inconsistent state partway through its construction. This also requires a lot of extra effort to ensure thread safety.

19 Bulkheads

Bulkheads

Bulkheads are a way to isolate microservices from failures. There are a lots of different bulkheads we can consider. Separation of concerns can be a way to implement bulkheads. By separating functionality into separate microservices we reduce the the change of an outage in one area affecting the another. Another form of bulkhead is using different connection pools for each downstream connection. That way, if one connection pool gets exhausted, the other connection pools aren't impacted. This would ensure that, if a downstream service start behaving slowly, only that one connection pool would be detected, allowing other calls to proceed as normal.

We can regard a [circuit breaker](#) as an automatic mechanism to enable a bulkhead. This way we not only protect the consumer from the downstream problem, but also potentially protect the downstream service from more calls that may be having an adverse impact.

20 Caching

Cache

Repetitious access to remote resources or global [value objects](#) form a bottleneck for many services. Caching is a technique that can drastically improve the performance. For example by avoiding multiple read operations for the same data.

However there is a price, caching data that the application is accessing will increase the memory usage. Therefore it is very important to obtain a proper balance between the retrieval of the data and the memory usage. The quantity of data being cached and the moment when to load, either in the beginning when the application initializes or whenever it is required for the first time, depends on the requirements of the application.

The cache will identify the buffered resources using unique identifiers. When the resources stored in the cache are no longer required they could be released in order to lower the memory consumption.

Basically there are 2 main caching strategies: Primed and Demand cache.

A cache that is initialized from the beginning with default values is primed cache. A primed cache should be considered whenever it is possible to predict a subset or the entire set of data that the client will request, and to put it to the cache. In case the client request data with a key that matches one of the primed keys having no corresponding data in the cache it is assumed that there is no data in the datasource for that key as well. Disadvantage of primed cache is that it takes longer for the application to start-up and become functional. There is however a hidden disadvantage and that is that developers assume that the cache will always be populated and that it will always be populated with the entire data set. This assumption could lead to a disaster for your application because:

- the logic to detect cache misses and fetch individual items on a miss never gets written or is written badly.
- Nothing can ever be evicted from the cache.
- Operations has no options to cleanup when faced with memory or data consistency issues in production.

Draw back of not using a primed cache is that when the application is cold and the cache is empty grabbing data from the datastore is inefficient. A solution to bypass the problems as described above is to provide cache miss logic to operate on sets of identifiers beside the single identifier logic. Now, when the system starts, spawn a couple of workers that fetch ids from the resource (in configurable batch sizes) and insert those in the cache. This way you:

- Ensure that code no longer make assumptions about all data being in cache at all times.
- Caches still get filled close to application start, but clients aren't held back by it.
- Cache warming code and normal get-miss-fetch scenarios share the same code.

Demand cache loads information and stores it in cache whenever the information is requested by the system. A demand cache implementation will improve performance while running. A demand cache should be considered whenever populating the cache is either unfeasible or unnecessary.

Object which are applicable for caching are:

- Results of database queries
- XML message
- Results of I/O

- Any other object that is expensive to get.
- Any object that is mostly read.

A cache requires: * max cache size - Defines how many elements a cache can hold * eviction policies - Defines what to do when the number of elements in cache exceeds the max cache size. The Least Recently Used (LRU) works best. Policies like MRU, LFU, etc, are usually not applicable in most practical situations and are expensive from performance point of view. * Time to live - Defines time after that a cache key should be removed from the cache (expired). * statistics Professional caches like [EHCACHE](#) or [Guava Cache](#) provide this kind of functionality out of the box and are constantly tested.

20.1 See also

- [Distributed cache](#) - a distributed cache is an extension of the traditional concept of cache used in a single locale. A distributed cache may span multiple servers so that it can grow in size and in transactional capacity. Distributed cache might be interesting in case of demand cache since it will fill up quicker.
- [Spring Cache Abstraction](#)
- [Hibernate 2nd level cache](#)

21 Circuit Breaker

Circuit Breaker

One of the properties of a remote call is that it can fail, or hang without a response until some timeout limit is reached. If there are many callers of this unresponsive service this can result in the caller run out of critical resources leading to cascading failures across multiple systems. To prevent this kind of catastrophic cascade we can use a circuit breaker.

The basic idea behind the circuit breaker is that you wrap external calls in a circuit breaker object. The circuit breaker object monitors for failures and once the failures reach a certain threshold, the circuit breaker trips, and all further calls via the circuit breaker return an error without the call being made at all. Usually there is also some kind of monitoring whether the circuit breaker has tripped.

This simple form of circuit breaker would need an external intervention to reset it when things are well again. An improvement would be to have the breaker itself detect if the underlying cells are working again.

21.1 Related patterns

- [Active Monitoring](#) - Circuit breakers are a valuable place for monitoring and any change in the breaker state should be logged. Breaker behavior is often a good source of warnings about deeper troubles in the environment.
- [API Gateway](#) - An API gateway orchestrates the calls between the diverse number of microservices. Each of these calls should be via a orchestrator.
- [Bulkheads](#) - Bulkheads are a way to isolate microservices from failures. By a circuit breaker we can enable bulkheads.

21.2 See also

- [Hystrix](#) - A sophisticated tool for dealing with latency and fault tolerance for distributed systems. It includes an implementation of the circuit breaker pattern
- [JRugged](#) - A Java library of robustness design patterns.

22 Client Side Discovery

Client Side Discovery

Services typically need to call one another. In a [monolithic](#), services invoke one another through direct calls. In a traditional distributed system deployment, services run at fixed, well known locations (host, port, etc) and can easily call one another using HTTP/REST or some RPC mechanism. A modern microservice based application typically runs in a virtualized or containerized environment where the number of instances of a service and their locations change dynamically. As a consequence of this, you must implement a mechanism that enables clients of services to make requests to a dynamically changing set of service instances.

To overcome this the client obtains the location of a service instance by querying a [Service Registry](#) which knows the locations of all service instances.

Disadvantage of this approach is that it couples the client to the service registry and client side service discovery logic needs to be implemented for each language/framework used within the client.

22.1 Related patterns

- [Service Connector](#) - Logic that is required for obtaining the location of a service can be implemented within the service connector.
- [Service Registry](#) - queried by client, containing locations of all service instances.
- [Server Side Discovery](#) is an alternative solution for the lookup of services.
- [System of record](#) - To assure that changes for a record are only happening on a single site.

22.2 See also

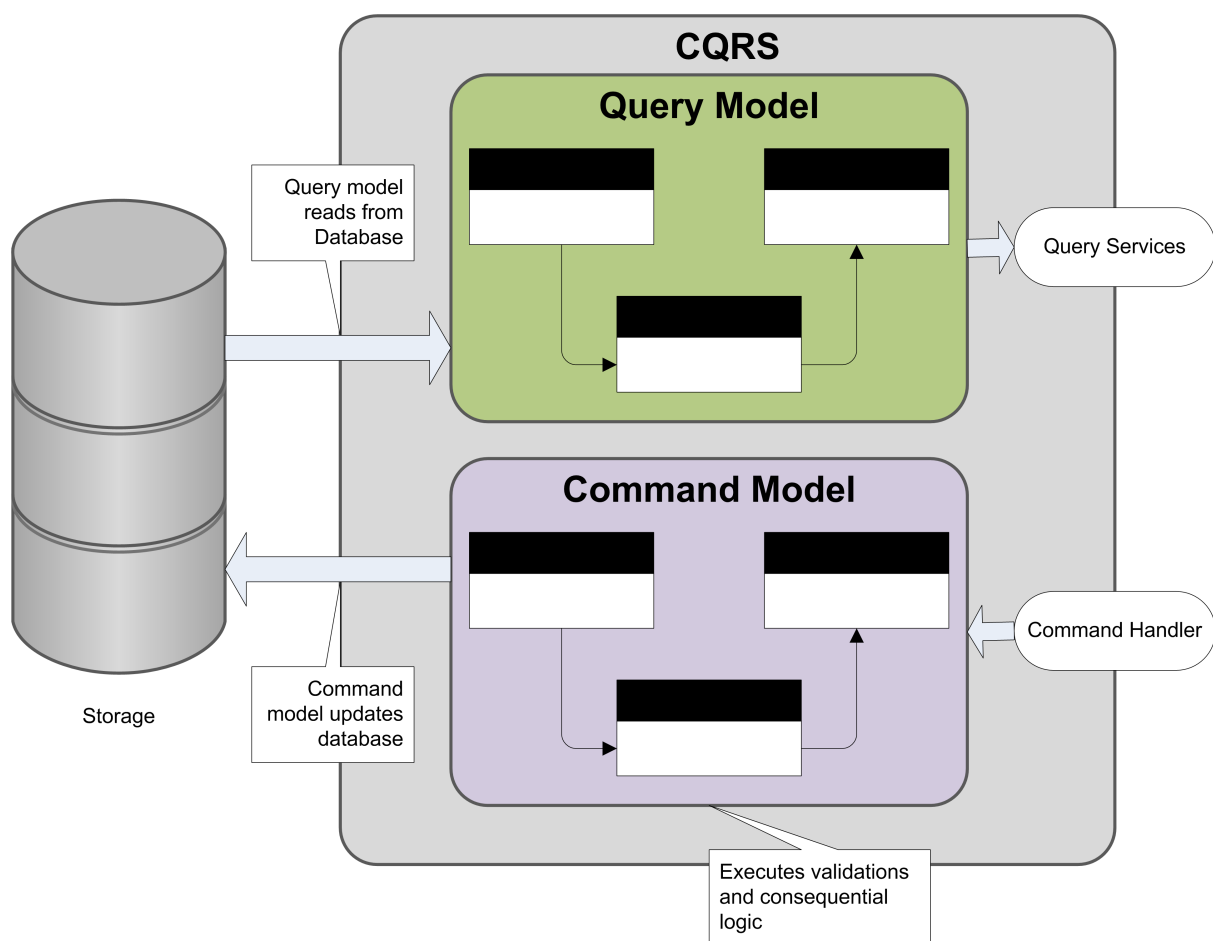
- [Ribbon Client](#) is an HTTP client that queries Eureka to route HTTP requests to an available service instance.

23 Command Query Responsibility Separation

Command Query Responsibility Segregation

CQRS stands for Command Query Responsibility Segregation and helps developers to develop scalable, extensible and maintainable applications. At its heart the pattern makes a distinction between the model in which information is updated and the model from which information is read. Following the vocabulary of *Command Query Separation* the models are called Command and Query. The rationale behind this separation is that for many problems, particularly in complicated domains, having the same conceptual model for updating information and reading information leads to a complex model that does neither well.

The models may share the same database, in which case the database acts as the communication between the two models, but they may also use separate databases, making the query side database a real-time reporting database. In later case there needs to be some communication mechanism between the two models, or their databases. Commonly this communication is realized by events.



23.1 Querying

A query returns data and does not alter the state of the model. Once the data has been retrieved by an actor, that same data may have been changed by another actor, in other words it is stale. If the data we are going to return to actors is stale anyway, is it really necessary to go to the master database and get it from there? Why transform the persisted data into entities if we just want data, not any rule preserving behavior? Why transform those entities to DTOs to transfer them across.

In short, it looks like we are doing a lot of unnecessary work. So why not creating a additional model, (whose data can be a bit out of sync), that matches the DTOs as they are expected by the requester. As data store you can use a regular database but that is not mandatory.

23.2 Commands

A command changes the state of an entity within the command model. A command may be accepted or rejected. An accepted command leads to zero or more events being emitted to incorporate new fact into the system. A rejected command leads to some kind of exception.

One should regard a command as a request to perform a unit of work which is not depending on anything else.

23.3 Domain event

In CQRS, domain events are the source all changes in the query model. When a command is executed, it will change state of one or more entities within the command model. As a result of these changes, one or more events are dispatched. The events are picked up by the event handlers of the query model and those update the query model.

23.4 When to use CQRS.

Like any pattern, CQRS is useful in some places and in others not. Many systems do fit a CRUD model, and should be done in that style. In particular CQRS should only be used on specific portions of a system (Bounded context in DDD lingo) and not to the system as whole; e.g. each bounded context needs its own decision on how it should be modeled.

So far there are the following benefits:

- Handling complexity - a complex domain may be easier to tackle by using CQRS.
- Handling high performance applications - CQRS allows you to separate the load from reads and writes allowing you to scale each independently.

24 Deployment microservices

Deployment of microservices

So you have applied the different patterns as described within these guidelines and architected your application as a set of services. Each service should be deployed as a set of service instances for throughput and availability.

Forces that (could) influence the decision are:

- Services are written using a variety of languages, frameworks and framework versions.
- Services are independently deployable and scalable.
- Services need to be isolated from each other.
- Need to be able to quickly build and deploy service
- Need to be able to restraint the resources (CPU/Memory) consumed by a service.
- Need to monitor behavior for each service instance.
- Deployment of application must be cost-effective

24.1 Multiple service instances per host

Multiple instances of different services are running on a single host (physical or virtual machine). There are various ways of deploying:

- Deploy each service instance as a JVM process, (for example a Tomcat instance per service instance).
- Deploy multiple service instances in the same JVM, (for example as OSGI bundles)

The resource utilization within this approach is more efficient than the service instance per host. The drawbacks of this approach include:

- Risk of conflicting resource requirements.
- Risk of conflicting dependency versions.
- Difficult to limit the resources consumed by a service instance.
- Difficult to monitor the resource consumption of each service.
- Impossible to isolate each service instance.

24.2 Single service instance per host

Within this approach we deploy each service instance on it's own host. The benefits of this approach include:

- Service instances are isolated from one other.
- There is no possible conflict regarding resource requirements or dependency versions.
- A service instance can consume at most the resources of a single host.
- It is straight forward to monitor, manage, and to redeploy a service instance.

The drawback is less efficient resource utilization compared to multiple service instances per host.

24.3 Service instance per VM

Within this approach we package the service as a virtual machine image and deploy each service instance as a separate VM.

The benefits of this approach include:

- Straightforward to scale the service by increasing the number of instances.
- VM encapsulates the details of the technology used to build the service.
- Each service instance is isolated.
- VM imposes limits on the CPU and memory consumed by the service instance.

The drawback of this approach include that building a VM image is slow and time consuming.

24.4 Service instance per container

Package the service as a container image ([Docker](#)) and deploy each service instance as a container. The benefits of this approach are:

- Straightforward to scale up and down a service by changing the number of container instances.
- Container encapsulates the details of the technology used to build the service. All services are, for example, started and stopped in exactly the same way.
- Each service instance is isolated.
- Container imposes limits on the CPU and memory consumed by a service instance.
- Containers are fast to build and start.

Drawback is that the infrastructure is not as rich as the infrastructure for deploying virtual machines.

TODO

25 Domain Driven Design

Introduction Domain Driven Design

The philosophy of Domain Driven Design (DDD), firstly described by Eric Evans, is about placing our attention at the heart of the application, focusing on the complexity that is intrinsic to the business domain itself.

Domain Driven Design consists of a set of [patterns](#) for building enterprise applications based on the domain model. Within this introduction we are going to run through some of concepts and terminology of DDD.

25.1 Domain Model

At the heart of DDD lies the concept of the domain model. This model is built by the team responsible for developing the system. The team consists of both domain experts from the business and software developers. The role of the domain model is to capture what is valuable or unique to the business. The domain model serves the following functions:

- It captures the relevant domain knowledge from the domain experts.
- It enables the team to determine the scope and verify the consistency of the knowledge.
- The model is expressed in code by the developers.
- It is constantly maintained to reflect evolutionary changes in the domain.

Domain models are typically composed of elements such as [entities](#), [value objects](#), [aggregates](#), and described using terms from a ubiquitous language.

25.2 Ubiquitous language

The ubiquitous language is very closely related to the domain model. One of the functions of the domain model is to create a common understanding between the domain experts and the developers. By having the domain experts and developers use the same terms of the domain model for objects and actions within the domain, the risk of confusion or misunderstanding is reduced.

25.3 Entities, value objects and services

DDD uses the following concepts to identify some of the building blocks that will make up the domain model:

- [Entities](#) are objects that are defined by their identity and that identity continuous through time.
- [Value objects](#) are objects which are not defined by their identity. Instead they are defined by the values of their attributes.
- [Services](#) is a collection of stateless methods that doesn't model properly to a entity or value object.

25.4 Aggregates and aggregate roots

DDD uses the term [aggregate](#) to define a cluster of related entities and value objects that form a consistency boundary within the system. That consistency boundary is usually based on transactional consistency. The aggregate root (also known as root entity) is the gatekeeper object for the aggregate. All access to the objects within the aggregate must occur through the aggregate root; external entities are only permitted to hold a reference to the aggregate root. In summary, aggregates are mechanism that DDD uses to manage the complex set of relationships between the many entities and value objects in a typical domain model.

25.5 Bounded context

For a large system, it may not be practical to maintain a single domain model; the size and complexity would make it difficult to keep it consistent. To address this, DDD introduces the concept of [bounded context](#) and multiple models. Within a system, you might choose to use multiple smaller models rather than a single large model, each focusing on a aspect or grouping of functionality within the overall system. A bounded context is the context for one particular domain model. Each bounded context has its own ubiquitous language, or at least its own dialect of the domain ubiquitous language.

25.5.1 Anti corruption layers

Different bounded contexts have different domains models. When your bounded contexts communicate with each other, you need to ensure that concepts specific to one domain do not leak into another domain model. An anti corruption layer functions as a gatekeeper between bounded contexts and helps you to keep the domain models clean.

25.5.2 Context maps

A large complex system can have multiple bounded contexts that interact with each other in various ways. A business entity, such as a customer, might exist in several bounded contexts. However, it may need to expose different facets or properties that are relevant to a particular bounded context. As a customer entity moves from one bounded context to another, you may need to translate it so that it exposes the relevant facets or properties for that particular context. A context map is the documentation that describes the relationships between these bounded contexts.

25.5.3 Bounded contexts and multiple architectures

Following the DDD approach, the bounded context will have its own domain model and its own ubiquitous language. Bounded contexts are also typically vertical slices through the system, meaning that the implementation of a bounded context will include everything from the data store, right up to the UI. One important consequence of this split is that you can use different implementation architectures in different bounded contexts. Due to this you can use an appropriate technical architecture for different parts of the system to address its specific characteristics.

25.5.4 Bounded contexts and multiple development teams

Clearly separating bounded contexts, and working with separate domain models and ubiquitous languages makes it possible to develop bounded contexts in parallel.

25.5.5 Maintaining multiple bounded contexts

It is unlikely that each bounded context will exist in isolation. Bounded contexts will need to exchange data with each other. The DDD approach offers a number of approaches for handling the interactions between multiple bounded contexts such as using anti-corruption layers or using [shared kernels](#).

26 Domain Model

Domain Model

Domain Driven Design is based on the idea of solving problems the organisations face through code. This is achieved by focusing on the domain of the business you are working with and the problems they want to solve. This will typically involve rules, processes and existing systems that need to be integrated as part of your solution; e.g. the domain is the ideas, entities and knowledge of the problem you are trying to solve. All of the knowledge around the company and how it operates will form the domain.

The Domain Model is your solution to the problem. The Domain Model usually represents an aspect of reality or something of interest. The Domain Model is often a simplification of the bigger picture, the important aspects of the solution are included while everything else is ignored.

The Domain Model should represent the vocabulary and key concepts of the problem domain and it should identify the relationships among all of the [entities](#) within the scope of the domain. The important this is that the Domain Model should be accessible and understandable by everyone who is involved with the project.

So the Domain is the world of the business and the Domain Model is the structured knowledge of the problem, but why is this important to Domain Driven Design? Well there are 3 reasons:

1. The Domain Model and implementation shape each other - The code that is written should be intimately linked to the Domain Model. Anyone on the team (business, developers, etc) should be able to look at your code and understand how it applies to the problem you are solving. Whenever a decision needs to be made during the course of the project everyone should refer to the Domain Model to look for the answer.
2. The Domain Model is the backbone of the language used by all team members - Every member of the team should use the [ubiquitous language](#). This ensures that technical and non-technical people have a common language so there is no loss of understanding between parties. The ubiquitous language should be directly derived from the Domain Model.
3. The Domain Model is distilled knowledge - The Domain Model is the outcome of an iterative discovering process where everyone on the team is involved to discuss the problem and how it should be solved. The Domain Model should capture how all think about the project and all of the distilled knowledge that has been derived from those collaboration sessions.

26.1 How do you create a Domain Model?

Creating a Domain Model is an iterative process that attempts to discover the real problem that is faced and the correct solution that is required. It's important to focus a Domain Model on one important problem. Trying to capture the entire scope of a business inside a single Domain Model will be far too overcomplicated and most likely work contradicting as concepts and ideas move around within the organisation, (see [bounded context](#) on how to cope with this).

The problem should be mapped out through talking to business experts to discover the problems they face. This should form the conceptual problem that you are looking to tackle. Business experts won't talk in terms of technical solutions, so during the process important aspects of the problem should be picket out from the language and given concrete definitions to form the ubiquitous language.

There should be a tight feedback loop where everyone on the project discuss the proposed Domain Model to get closer to the solution. This is an iterative approach that will likely encompass code and diagrams to really understand the problem and discover the correct solution.

27 Building blocks DDD

Building blocks Domain Driven Design

The diagram below is a navigational map. It shows the patterns that form the building blocks of Domain Driven Design and how they relate to each other.

By using these standard patterns we bring order in the design and make it easier for team members to understand each other's work. Using standard patterns also adds to the *ubiquitous language* which all team members can use to discuss model and design discussions.

28 Aggregates

Aggregates.

It is difficult to guarantee the consistency of changes to objects in a model with complex associations. Invariants need to be maintained that apply to closely related groups of objects, not just discrete objects. Yet cautious locking schemes cause multiple users to interfere pointlessly with each other and make a system unusable.

Therefore, cluster the `Entities` and `Value Objects` into `Aggregates` and define boundaries around each. Choose one `Entity` to be the root of each `Aggregate`, and control all access to the object inside the boundary through the root. Allow external objects to hold references to the root only. Transient references to internal members can be passed out for use within a single operation only. Because the root controls access, it cannot be blindsided by changes to the internals. This arrangement makes it practical to enforce all invariants for objects in the `Aggregate` and for the `Aggregate` as a whole in any state change.

29 Entities

Entities

Many object are not fundamentally defined by their attributes, but rather by a thread of continuity and identity.

Some objects are not defined primarily by their attributes. They represent a thread of identity that runs through time and often across distinct representations. Sometimes such an object must be matched with another object even though attributes differ. Mistaken identity can lead to data corruption.

Therefore, when an object is distinguished by its identity, rather than its attributes, make this primary to its definition in the model. Keep the class definition simple and focused on life cycle continuity and identity. Define a means of distinguishing each object regardless of its form or history. Be alert to requirements that call for matching objects by attributes. Define an operation that is guaranteed to produce a unique identity for each object. The model must define what it *means* to be the same thing.

30 Factories

Factories

When creation of an object, or an entire Aggregate, becomes complicated or reveals too much of the internal structure Factories provide encapsulation.

Creation of an object can be a major operation in itself, but complex assembly operations do not fit the responsibility of the created object. Combining such responsibilities can produce ungainly designs that are hard to understand. Making the client direct construction muddies the design of the client, breaches encapsulation of the assembled object or Aggregate, and overly couples the client to the implementation of the created object.

Therefore, shift the responsibility for creating instances of complex objects and Aggregates to a separate object, which may itself have no responsibility in the domain model but is still part of the domain design. Provide an interface that encapsulates all complex assembly and that does not require the client to reference the concrete classes of the object being instantiated. Create entire Aggregates as a piece, enforcing their invariants.

31 Layered Architecture

Layered Architecture

In an object oriented program UI, database, and other support code often gets written directly into the business objects. Additional business logic is embedded in the behavior of UI widgets and database scripts. This happens because it is the easiest way to make things work, in the short run. When the domain related code is diffused through such a large amount of other code, it becomes extremely difficult to see and reason about. Superficial changes to the UI can actually change business logic. To change a business rule may require meticulous tracing of UI code, database code, or other programming elements. Implementing coherent model driven objects impractical and automated testing is awkward. With all the technologies and logic involved in each activity, a program must be very simple or it becomes impossible to understand.

Therefore:

- Partition a complex program into 'layer's.
- Develop a design within each 'layer' that is cohesive and that depends only on the layers below.
- Follow standard architectural patterns to provide loose coupling to the layers above by means of a mechanism such as *Observer* or *Mediator*; there is never a direct reference from lower to higher.
- Concentrate all the code related to the domain model in one layer and isolate it from the user interface, application, and infrastructure code. The domain objects, free of the responsibility of displaying themselves, storing themselves, managing application tasks, and so forth, can be focused on expressing the domain model. This allows a model to evolve to be rich and clear enough to capture essential business knowledge and put it to work.

A typical enterprise application architecture consists of the following four conceptual layers:

- *Presentation Layer* - Responsible for presenting information to the user and interpreting user commands.
- *Application Layer* - Layer that coordinates the application activity. It doesn't contain any business logic. It does not hold the state of business objects, but it can hold the state of an application task's progress.
- *Domain Layer* - This layer contains information about the business domain. The state of business objects is held here. Persistence of the business objects, and possibly their state is delegated to the infrastructure layer.
- *Infrastructure Layer* - This layer acts as a supporting library for all the other layers. It implements persistence for business objects, contains supporting libraries, etc.

31.1 Application layer.

The application layer:

- Is responsible for the navigation between the UI screens in the bounded context as well as the interaction with application layers of other bounded contexts.

- Can perform the basic (non business related) validation on the user input data before transmitting it to the other (lower) layers of the application.
- Doesn't contain any business or domain related logic.
- Doesn't have any state reflecting a business use case but it can manage the state of the user session or the progress of a task.
- Contains [application services](#).

31.2 Domain layer.

The domain layer:

- Is responsible for the concepts of business domain, and the business rules. Entities encapsulate the state and behavior of the business domain.
- Manages the state of a business use case if the use case spans multiple user requests, (e.g. loan registration process which consists of multiple steps: user entering loan details, system returning products and rates, user selecting particular product, system locking the loan for selected rate).
- Contains [domain services](#).
- Is the heart of the bounded context and should be well isolated from the other layers. Also, it should not be dependent on the application frameworks used in the other layers, (Hibernate, Spring, etc).

31.3 CRUD operations

TODO: see [Domain Driven Design \(DDD\) architecture layer design for CRUD operation](#)

32 Modules

Modules

Everyone uses `Modules` but few treat them as full fledged part of the model. Code gets broken down into all sort of categories, from aspects of the technical architecture to developers work assignments. It is a truism that there should be low coupling between `Modules` and high cohesion within them. Explanations of coupling and cohesion tend to make them sound like technical metrics, to be judged mechanically based on the distributions of associations and interactions. Yet it isn't just code being divided into `Modules`, but concepts. There is a limit to how many things a person can think about at once (hence low coupling).

Therefore, choose `Modules` that tell the story of the system and contain a cohesive set of concepts. Seek low coupling in the sense of concepts that can be understood and reasoned about independently of each other. Refine the model until it partitions according to the high level domain concepts and the corresponding code is decoupled as well. Give the `Modules` names that become part of the ubiquitous language. `Modules` and their names should reflect insight into the domain.

33 Repositories

Repositories.

A client needs a practical means of acquiring references to preexisting domain objects. If the infrastructure makes it easy to do so, the developer of the client may add more traversable associations, muddling the model. On the other hand, they may use queries to pull the exact data they need from the database, or to pull a few objects rather than navigating from the ‘Aggregate’ roots. Domain logic moves into queries and client code, and the `Entities` and `Value Objects` become mere data containers. The sheer technical complexity of applying most database access infrastructure quickly swamps the client code, which leads developers to dumb down the domain layer, which makes the model irrelevant.

Therefore, for each type of object that needs global access, create an object that can provide the illusion of an in memory collection of all objects of that type. Setup access through a well known global interface. Provide methods to add and remove objects, which will encapsulate the actual insertion and removal of data in the store. Provide methods that select objects based on some criteria and return fully instantiated objects of object whose attribute values meet the criteria, thereby encapsulating the actual storage and query technology. Keep the client focused on the model, delegating all object storage and access to the `Repositories`.

34 Services

Services.

Some concepts from the domain aren't natural to model as objects. Forcing the required domain functionality to be the responsibility of an `Entity` or `Value Object` either distorts the definition of a model based object or adds meaningless artificial objects.

Therefore, when a significant process or transformation in the domain is not a natural responsibility of an `Entity` or `Value Object`, add an operation to the model as a standalone interface declared as a `Service`. Define the interface in terms of the language of the model and make sure the operation name is part of the ubiquitous language. Make the service stateless.

34.1 Service types.

Services exist in most layers of the DDD layered architecture:

- Application
- Domain
- Infrastructure

An infrastructure service would be something that communicates directly with external resources, such as file systems, databases, etc.

Domain services are the coordinators, allowing higher level functionality between many different smaller parts. Since domain services are first-class citizens of the domain model, their names and usage should be part of the ubiquitous language. Meanings and responsibilities should make sense to the stakeholders or domain experts.

The application service will be acting as façade and accepts any request from clients. Once the request comes, based on the operation, it may call a `Factory` to create domain object, or calls repository service to re-create existing domain objects. Conversion between, `DTO` (Data Transfer Object) to domain objects and Domain objects to `DTO` will be happening here. Application service can also call another application service to perform additional operations.

The differences between a domain service and an application service are subtle but critical:

- Domain services are very granular where as application services are a facade with as purpose providing an API.
- Domain services contain domain logic that can't naturally be placed in an entity or value object, whereas application services orchestrate the execution of domain logic and don't themselves implement any domain logic.
- Domain service methods can have other domain elements as operands and return values whereas application services operate upon trivial operands such as identity.
- Application services declare dependencies on infrastructural services required to execute domain logic.
- Command handlers are a flavor of application services which focus on handling a single command typically in a CQRS architecture.

35 Value Objects

Value Objects

Many objects have no conceptual identity. These objects describe some characteristic of a thing.

Tracking the identity of `Entities` is essential, but attaching identity to other objects can hurt system performance, add analytical work, and muddle the model by making all object look the same. Software design is a constant battle with complexity. We must make distinctions so that special handling is applied only where necessary. However if we think of this category of objects as just the absence of identity, we haven't added much to our toolbox or vocabulary. In fact, these objects have characteristics of their own, and their own significance to the model. These are the objects that describe things.

Therefore, when you care only about the attributes of an element of the model, classify it as a `Value Object`. Make it express the meaning of the attributes it conveys and give it related functionality. Treat the `Value Object` as immutable. Don't give it any identity and avoid the design complexities necessary to maintain `Entities`.

36 Context mapping

Context Mapping

TODO:

Include documentation regardign AcL Anti corruption Layer, and other approaches (hexa pattern?).

37 Enterprise Integration Pattern

Enterprise Integration Pattern

TODO

38 Correlation ID

Correlation ID

Microservices call each another and it will be difficult to figure out how one particular request got transformed and which services where called. Another usage of correlation ID is to track state and return address within asynchronous messaging. By using a correlation ID that is passed within the calls between services we enable the tracking of requests and its route, and we can return the response to the proper caller.

As far as the replier is concerned, the correlation ID is an opaque data type and all it has to do is copy the request ID from the request message to the correlation ID in the reply message.

Extra care has to be taken regarding the choice of request ID. The message ID seems a valid choice but in case of intermediaries, the intermediary consumes the original and sends a new message to the replier. During this step the new message might be assigned its own unique message ID. If the replier uses the message ID as the request ID and blindly copies it to the correlation ID, the requester will not be able to correlate the incoming reply message to the original request message.

To overcome this the message should not only contain a message identifier but also a conversation identifier. The requester who initiates the conversation picks a conversation ID while all intermediaries and repliers pass this ID along so that all messages belonging to the conversation carry a common conversation identifier.

TODO: determine whether we should use (custom) HTTP Header or URI parameter. Advantage of URI parameter is that it is easier log-able, however caching would become complexer.

38.1 See also

- [Implementing Correlation IDs in Spring Boot](#)

39 Exclusive Consumer

Exclusive Consumer

If there are multiple message consumer instances consuming from the same queue you will lose the guarantee of processing the messages in order since the messages will be processed concurrently in different threads. Sometimes it's important to guarantee the order in which messages are processed.

A common approach in this case is to pin one particular JVM in the cluster to have one consumer on the queue to avoid losing ordering. The problem with this is that if that particular JVM goes down, no one is processing the queue any more.

With an exclusive consumer we avoid pinning a particular JVM. Instead the message broker will pick a message consumer to get all the messages for that queue to ensure ordering. If that consumer fails, the broker will automatically failover and choose another consumer.

The effect is a heterogeneous cluster where each JVM has the same setup and configuration; the message broker is choosing one consumer to be the master and send all the messages to it until it dies; then you get immediately fail-over to another consumer.

39.1 Parallel exclusive consumer

By including the concept of message groups we can create a kind of parallel exclusive consumers. Rather than all messages going to a single consumer, a message group will ensure that all messages for the same message groups will be sent to the same consumer while that consumer stays alive. As soon as the consumer dies another will be chosen.

When a message is being dispatched to a consumer, the message group is checked. If there is a message group present the broker checks to see if there is a consumer that owns the message group. If no consumer is associated with a message group, a consumer is chosen. That message consumer will receive all further messages with the same message group until:

- The consumer closes.
- The message group is closed.

39.2 Resources

- [Exclusive Consumer within ApacheMQ](#)
- [Parallel exclusive consumers within ApacheMQ](#)

40 Message Router

Message Router

A message router consumes a [Message](#) from a [Message Channel](#) and republish it to a different message channel depending on a set of conditions.

A key property of the message router is that it does not modify the message content. It only concerns itself with the destination of the message.

41 Exception Handling

Exception Handling

TODO: Write nice intro.

41.1 Exception handling and logging

...

41.2 Best practices

- Use checked exceptions for recoverable errors and unchecked exceptions for programming errors
 - Checked exceptions ensure that you handle certain error conditions, but at the same time it also adds a lot of clutter in the code and might make it unreadable. Also it is only reasonable to catch exceptions if you have alternatives or recovery strategies.
- Avoid unnecessary exception handling - Exceptions are costly and can slow down your code. Don't just throw and catch exceptions, if you can use standard return values to indicate result of an operation. Also avoid unnecessary handling by fixing the root cause.
- Never swallow the exception in catch block.
- Declare specific checked exceptions that method can throw - Declare specific checked exceptions that can be thrown by your method. If there are too much checked exceptions, you should wrap them in your own exception and add information in the exception message.
- Never catch the `Exception` class - The problem with catching `Exception` is that if the method that you are calling adds a new checked exception to its method signature you will never know about it and the fact that your code now might be wrong and might break at any point in time.
- Never catch `Throwable` class - JVM errors are also subclass of `Throwable` and these errors are irreversible conditions that cannot be handled by the JVM itself. The JVM might even not invoke your catch clause on an error.
- Always wrap exceptions correctly - the original exception should always be included within the exception that wraps it.
- Either log the exception or throw it but never do both - Logging and throwing will result in multiple messages in the log for a single problem in the code.
- Never throw any exception from finally block - if we would throw an exception from the finally block we might hide the original first exception and correct reason would be lost forever.
- Only catch exception you can handle - Catch an exception only if you want to handle it or to provide additional contextual information in that exception.
- Don't print stack traces to the console.
- Use finally blocks instead of catch blocks if you are not going to handle exception. If a method throws some exception which you do not want to handle, but still want to perform some cleanup, then do this cleanup in the `finally` block. Do not use `catch` block.
- Prefer exceptions to return codes - it is preferable to throw understandable messages that are part of your API contract and guide the client programmer.
- Throw early catch late - This principle implicitly says that you will throw an exception in the low level methods and make the exception climb the stack trace until you reached a sufficient level of abstraction to be able to handle the problem.
- Always cleanup after handling exception - If you are using resources like databases make sure you clean them up. You can use the new java 7 auto-cleanup via `try-with-resources` statement.
- Exception names must be clear and meaningful - Use for specific exceptions clear names that states the cause of the exception. For example `AccountLockedException` instead of `AccountException`.

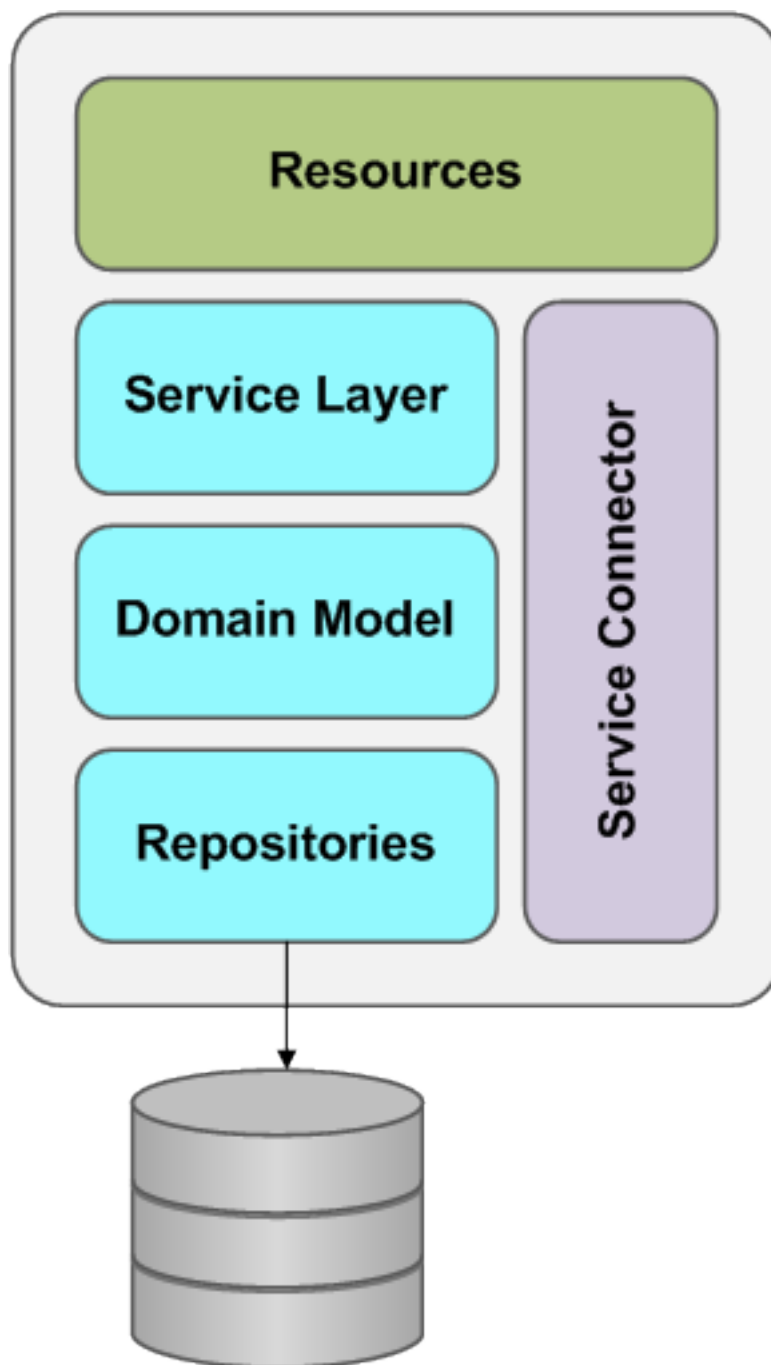
- Never use exceptions for flow control - it makes code hard to understand.
- Do not handle exceptions inside loops. Surround the loop with exception block instead.
- Always include all information about an exception in a single log message.
- Document all exceptions in your application in javadoc.

42 Microservice Architecture

Microservice architecture

A microservice architecture is the natural consequence of applying the single responsibility principle at architectural level. Within a microservice architecture the functionality is decomposed in a set of collaborating services and the [scale cube](#) is applied. Services communicate with each other either using synchronous protocols like HTTP/REST or asynchronous protocols such as AMQP/JMS. In this way, business domain concepts are modelled as resources with one or more of these managed by a microservice. Since a business request can span multiple microservices separated by network partitions, it is important to consider possible failures in the system. Techniques such as timeouts, [circuit breakers](#), [bulkheads](#) can help to maintain overall system uptime in spite of outage.

Often microservices display similar internal structure consisting of some or all of the layers as shown below.



- Resources act as a mapper between the application protocol as exposed by the micro service and messages to the [services/ entities](#) that are representing the domain. Typically, they are thin, responsible for sanity checking the request, and providing a protocol specific response according to the outcome of the request. If we for example expose the resources via REST and follow internally the [CQRS](#) approach, this layer would be responsible for converting the HTTP REST request into a `Command`.

- The microservice logic resides in the [domain model](#). [Services](#) coordinate across multiple domain activities, whilst [repositories](#) act on collections of domain entities and are often persistent. When a resource receives a request and has validated it, it either directly access the domain model, or, if many entities must be coordinated, it delegates the request to a service.
- If a service has another service as a collaborator, some logic is needed to communicate with the external service. A [service connector](#) encapsulates message passing with a remote service, marshalling requests and responses from and to other services. Service connectors should be resilient to outage of remote components.

Normally services are developed independently from another. Typically, a team will act as guardian to one or more microservices,

Each service has its own storage in order to decouple from other services. When necessary, consistency between services is maintained using application events.

Following this approach has a number of benefits:

- Each service is relatively small. Defining small is subjective but one of the rule of thumbs is that the microservice should be small enough to be owned by a small agile development team, re-writable within 1 or 2 sprint, and the complexity does not require refactoring or further division into another microservice. Advantage of such small microservices are:
 - Easier for developers to understand.
 - The web container starts fast, which makes developers more productive and speed up deployments.
- Each service can be deployed independently of other services; e.g. easier to deploy new versions of services frequently.
- Easier to scale development. It enables to organize the development effort around multiple teams. Each team can develop, deploy, and scale their service independently of all other teams.
- Improved fault isolation.
- Each service can be developed and deployed independently.
- Eliminates long term commitment to a technology stack. Each microservice has its own server, and network & hosting environment. Also the business logic, data model and the service interface(s) (API/UI) are all part of the entire system.

There are however also a number of drawbacks:

- Developers must deal with the additional complexity of creating distributed systems.
 - Testing is more difficult.
 - Developers must implement the inter-service communication mechanism
 - Implementing use cases that span multiple services without using distributed transactions is difficult (See [Try-Cancel/Confirm](#))

- Implementing uses cases that span multiple services requires careful coordination between the teams.
- Deployment complexity in production; there is additional operational complexity of deploying and managing a system comprised of many different services.
- Increased memory consumption. The microservice architecture replaces a monolithic application with N service instances. If each service runs in its own JVM (which is usually necessary to isolate instances) there is an overhead of N-1 JVM runtimes. Moreover, if each service runs on its own VM the overhead is even higher.

A challenge is deciding how to partition the system into microservices. One approach is to partition services by use case. We could for example have a micro service shipping within a partitioned e-commerce application that is responsible for shipping completed orders. Another approach is to partition by [entity](#). This kind of service is responsible for all operations that operate on entities of a given type.

Ideally, each service should have only a small set of responsibilities. The [Single Responsible Principle](#) states that every class should have responsibility over a single part of the functionality provided by the software, and that responsibility should be entirely encapsulated by the class. It makes sense to apply this principle to microservice design as well.

42.1 MicroService naming guidelines

There isn't a really straight forward approach regarding the naming of microservices but the following guidelines should help in creating a consistency in the naming of the microservices.

- *Use camel case for microservice name*
- *Don't reveal implementation details in the microservice name* - This not only has the potential to lead to confusion when you change the implementation of the microservice, but is also a security risk as it gives the microservice consumer an insight in how the microservice may be implemented which they may be able to exploit.
- *Don't include protocol information in the microservice name* - This is generally unnecessary as the service advertise itself at a particular endpoint which clearly defines the protocol to be used.
- *Don't include the word service in the microservice name*
- *Don't include a version in the microservice name*
- *Name microservice to entity at which it operates* - Microservices which operate on a specific entity should be named after the entity. For example if service that operates on customers may be simple named `Customer`.
- *Name microservice to functionality it performs* - Certain microservices are taking care of certain processes and therefore the use of verbs in the service name is common. For example a service that orchestrates the device registration could be called `RegisterDevice`.

Naming conventions may seem trivial at first, but as the number of microservices grow, so will the potential to reuse. In larger organizations, this means that more and more architects, analysts, and developers are discovering and then incorporating foreign services within their solution designs. The effort required to establish a consistent level of clarity across all microservices pays off quickly when interoperability and reuse opportunities are more easily recognized and seized.

42.2 Non-functional requirements

Non-functional requirements are important decision makers while designing micro services. The success of a system is largely dependent on its availability, scalability, performance, usability and flexibility.

42.2.1 Availability

The golden rule for availability says, anticipate failures and design accordingly so that the system will be available for 99.999%. It means that the system can only go down for 5.5 minutes a year. The cluster model is used to support such high availability (having group of services run in active-active mode or active-standby model).

So while designing microservices, it must be designed for appropriate clustering and high availability model. The basic properties of microservices such as stateless, independent and full stack will help to run multiple instances in parallel.

42.2.2 Scalability

Microservices must be scalable both horizontally and vertically. Being horizontally scalable means we can have multiple instances of the microservices to increase the performance of the system. The design of the microservices must support horizontal scaling (scale-out)

Also scaling vertically should be possible. If a microservice is hosted on a medium capacity and moved to a higher capacity the performance of the service should scale accordingly. Similarly downsizing the system capacity must be possible.

42.2.3 Performance

Performance is measured by throughput vs response time (TPS - Transactions Per Second). The performance requirements must be available in the beginning of the design phase itself. There are technologies and design choices that will affect the performance and to avoid rework at a later stage these should be known.

42.2.4 Usability

Usability aspects of the design focus on hiding the internal design, architecture, technology and other complexities to the end user of the system. Most of the time, microservices expose APIs to the end user, so the APIs must be designed in a normalized way so that it is easy to achieve the required functionality with a minimal number of API calls.

42.2.5 Flexibility

Flexibility measures the adaptability to change. In the microservices eco-system, where each microservice is owned (possibly) by different teams and developed in agile methodology changes will happen faster than any other system. The microservices may not inter-operate if they don't adapt or accommodate to changes in other systems. So there must be a proper mechanism in place to publish changes to API, functional changes, etc.

42.3 TODO

- Check whether OSGI could help with decreasing memory consumption.

42.4 Related Patterns

- [API Gateway](#) - defines how clients access the services in a microservice architecture.

- [Client side discovery](#) and [Service side discovery](#) - Patterns used to route requests for a client to an available microservice.
- [Command Query Response Separation](#) - Helps develop scalable, extensible and maintainable applications.
- [Domain Driven Design](#) - Set of patterns for building enterprise applications based on the domain model.
- [Monolithic Architecture](#) - alternative to the monolithic architecture.
- [Service Connector](#) - Provide high level interface that hides implementation details regarding communication, thereby making the use of the microservice easier.
- [Service Statelessness](#) - To have services scalable we should attempt to make them stateless.

42.5 See also

- [JBoss OSGI User guide](#) - Guide explaining how to deploy application as OSGI bundle on JBoss server.

43 Mocks

Mocks

Mocks are simulated objects/services that mimic the behavior of real objects/services in a controlled way. A mock is typically created to test the behavior of some other parts of the system.

43.1 See also

- [WireMock](#) - WireMock is a flexible library for stubbing and mocking web services.

44 Monitoring

Active Monitoring

Applications must expose runtime information that administrators and operators can use to manage and monitor the system.

The following patterns and guidances are related to monitoring applications:

- [Canary endpoint monitoring](#) - Implement functional checks within an application that external tools can access through exposed endpoints at regular intervals and verify whether the dependencies of the service are still operating as expected. This pattern can help to verify that applications and services are performing correctly.
- [Synthetic monitoring](#) - constantly checks the application for potential problems by mimicking a user.
- [Log aggregation](#) - Each service will create its own log and when an problem occurs you want to dig into them. Finding the problem in multiple logs can become a big pain. Solution for this is to aggregate the logs.
- [Service metering](#) - You may need to meter the use of services in order to plan future requirements, to gain knowledge on how they are used, to bill users, etc.

Dimensions of Service quality

See: * [reactive microservices monitoring](#) * [Problem 3 of problems micro services](#) * [Testing strategy in Micro services](#) * [Management and Monitoring Patterns and Guidance](#)

45 Canary endpoint monitoring

Canary endpoint monitoring

It is good practice - and often business requirement - to monitor web applications, middle-tier and shared services, to ensure that they are available and performing correctly. There are many factors that affect applications such as network latency, performance and availability of the underlying hardware and storage systems, and the network bandwidth between them. The service may fail entirely or partly due to any of these factors. Therefore, you must verify at regular intervals that the service is performing correctly to ensure the required level of availability.

To provide the functionality by which the application can be monitored one should implement an endpoint that performs the necessary checks and return an indication of its status. Such an endpoint, unlike other resources of a REST API, perform no business activity, but instead gathers status and latency of all dependencies of a service.

Checks that might be carried out by the monitoring code in the application include but not limited to:

- Checking storage or database for availability (and response time).
- Checking external distributed caches.
- Checking availability of Service bus
- Checking other resources or services used by the application.

The response code indicates the overall status of the application. If any of the dependencies has a non-success code, the returned status code will be non-success also. The payload contains the name of the dependencies it uses, and the status code that was returned by the canary endpoint of that particular service.

Implementation of the canary endpoint check is generally dependent on the underlying technology. For a cache service, it suffices to set a constant value and see it succeeding. For a SQL database a `select 1;` query is all that is needed. *Note that none of these are anywhere near a business activity, so that you could not, think that its success means your business is up and running.* A canary endpoint normally gets implemented as an HTTP GET call which returns a collection of connectivity check metrics.

Typical checks that should be performed by the monitoring tool include:

- Validating the response code. For example, an HTTP response of 200 (OK) indicates that the services dependencies are all ok.
- Checking the content of the response to detect errors, even when a 200 (OK) status code is returned.

45.1 Security of the canary endpoint.

Canary endpoints should be secured. A canary endpoint lists all internal dependencies, and potentially technologies of a system and this could be abused by hackers to target your system.

45.2 Performance impact.

Since canary endpoints do not trigger any business activity, its performance footprint should be minimal. However, since calling the canary endpoint generates a cascade of calls, calling all canary endpoints every few seconds might not be wise.

45.3 Related Patterns

- [Circuit breaker](#) - Based on the status code as exposed by the canary endpoint the circuit breaker can allow requests to the service or block.

45.4 See also

- [Spring Boot custom HealthIndicator](#)
- [AppDynamics](#) for monitoring tool.

46 Log aggregation

Log Aggregation

Logs are a critical part of any system, they give you insight into what a system is doing as well what happened. Most processes running on a system generates logs in some form or another and usually these logs are written to files on a local disk. When your system grows to multiple hosts, managing the logs and accessing them can get complicated. Searching for a particular error across hundreds of logs files on hundreds of servers is difficult without good tools. A common approach to this problem is to setup a centralized logging solution so that multiple logs can be aggregated in a central location.

46.1 Related to

- [Correlation ID](#) - To track requests there should be a common ID

46.2 See also

- [Log stash](#) - Allows to aggregate logs into 1 location.
- [Kibana](#) - To search (aggregated) logs.

47 Synthetic monitoring

Synthetic Monitoring

Synthetic (active) monitoring where an application that mimics a user constantly checks the application for potential problems.

47.1 Related Patterns

- [Try-Cancel/Confirm](#) - To verify business activity can be achieved by initiating 'business' functionality the Try Cancel/Confirm pattern can be (ab)used. After executing the business functionality, the state changes can be reverted by cancelling the action.

48 Monolithic Architecture

Monolithic Architecture

The application has either a layered or [hexagonal](#) architecture and consists of different types of components:

- Presentation components - responsible for handling HTTP requests and responding with either HTML or JSON/XML.
- Business logic- the application's business logic
- Database access logic
- Application integration logic - messaging layer

Problem is how we will deploy this application. Forces that influence the solution are:

- There is a team of developers working on the application.
- New team members must become quickly productive.
- The application must be easy to understand and modify.
- The application should be deployed using the continuous integration practice.
- Multiple instances of the application are running to satisfy requirements regarding scalability, and availability.

Most straightforward solution is to build an application with a monolithic architecture, (for example a single java WAR file). This approach has a number of benefits:

- Simple to develop
- Simple to deploy - you simply need to deploy the WAR file on the appropriate runtime.
- Simple to scale - You can scale the application by running multiple copies of the application behind a load balancer.

However once the application becomes large and teams grow in size, this approach has a number of drawbacks:

- The large monolithic code base intimidates developers and the application can be difficult to understand and modify. As a result, development typically slows down.
- The startup of the application will take longer. This will have an impact on developer productivity because of time wasted waiting for the web server to start. This is especially a problem for user interface developers since they usually need to iterate rapidly and redeploy frequently.
- A large monolithic application makes continuous deployment difficult because in case of an update the entire application has to be redeployed.
- A monolithic architecture can only scale in one dimension. This architecture can't scale with an increasing data volume. Each instance of the application will access all of the data, which makes caching less effective and increases memory consumption and I/O traffic. Also different application components have different resource requirements - one might be CPU intensive while another might be I/O intensive. With a monolithic architecture we cannot scale each component independently.

- A monolithic architecture is also an obstacle for scaling development. Once the application gets to a certain size it is useful to divide the development into several teams that focus on specific functional areas. The trouble with monolithic architecture is that it prevents teams from working independently. The teams must coordinate their development efforts and redeployments.
- A monolithic architecture requires a long term commitment to a technology stack.

48.1 Related Patterns

- [Microservice Architecture](#) is an alternative to the monolithic architecture.

49 Self registration

Self Registration

Service instances must be registered with the [service registry](#) on startup so that they can be discovered and unregistered on shutdown.

Registration of the service should happen with the service registry on startup and unregister on shutdown. If a service instance crashes the service must be unregistered from the service registry. The same applies for service instances that are running but are incapable of handling requests.

Within the self registration solution, a service instance is responsible for registering itself with the service registry. On startup the service instance register itself (location) with the service registry and makes itself available for discovery. The service instance must periodically renew it's registration so that the registry knows it is still alive.

[Eureka](#) is an example of a service registry which provides a service registry API and a client library that service instances use to (un)register themselves.

The benefit of the self registration pattern is that a service instance knows best it's own state so a state model that's more complex than UP/ DOWN should be possible.

There are also a couple of drawbacks:

- Couples the service to the service registry
- Service registration logic must be implemented for all languages that you are using within your environment/client.
- A service which is running but unable to process requests will often lack the self awareness to unregister from the service registry.

49.1 Related patterns

- [Active monitoring](#) - if a service crashes it might be that service is not properly unregistered. By actively monitoring the environment we might detect this and remove service instance from registry.
- [Service registry](#) - registry at which service instance registers itself, and un-registers when service is shutdown, crashes, or becomes unavailable.
- [3rd party registration](#) - alternative method to register services.

50 Server Side Discovery

Server Side Discovery

Services typically need to call one another. In a [monolithic](#), services invoke one another through direct calls. In a traditional distributed system deployment, services run at fixed, well known locations (host, port, etc) and can easily call one another using HTTP/REST or some RPC mechanism. A modern microservice based application typically runs in a virtualized or containerized environment where the number of instances of a service and their locations change dynamically. As a consequence of this, you must implement a mechanism that enables clients of services to make requests to a dynamically changing set of service instances.

Within a server side discovery approach a request to a service is going via a [message router](#) that runs at a well known location. The router queries a service registry, which might be built into the router, and forward the request to an available service instance.

The [AWS Elastic Load Balancing](#) is an example of a server side discovery solution. A client makes HTTP(S) requests to the ELB, which load balances the traffic amongst a set of EC2 instances. Within the Amazon solution ELB also functions as a [service registry](#). EC2 instances are either registered with the instance explicitly via an API call or automatically as part of an auto-scaling group.

Server side discovery has a number of benefits:

- Compared to [client side discovery](#), the client code is simpler since it does not have to deal with discovery. Instead a client simply makes a request to a router.

It however also has some drawbacks:

- Unless part of the cloud environment, the router is another system component that must be installed and configured. It will also need to be replicated for availability and capacity.
- More network hops are required compared to [client side discovery](#).

50.1 Related patterns

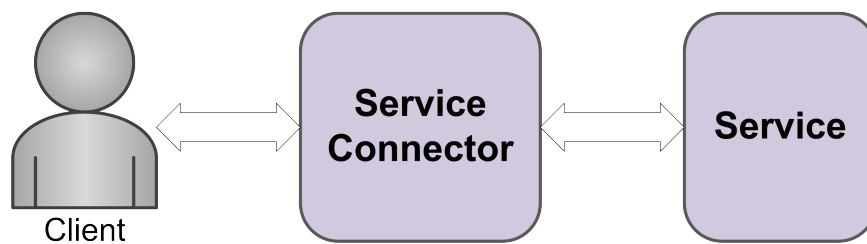
- [Client side discovery](#) is an alternative solution.
- [Message Router](#)
- [Service Registry](#)
- [System of record](#) - To assure that changes for a record are only happening on a single site.

51 Service Connector

Service Connector

Clients must know a lot about services in order to use them. REST APIs require clients to use specific media types, HTTP Headers, and specific request methods. Some web services require data to be encrypted or demand that clients submit authentication tokens. The client must also know the service's address, how to serialize requests, and parse responses. Each client could develop a custom solution to meet the specific requirements of the service but that would often result in duplicate code.

Instead we should create a service connector that encapsulates the logic a client can use to call a certain service.



Service connectors make services easier to use by hiding the communication specifics. The service connector encapsulate the generic communication logic required to use the service and also include the logic that is specific go the given service. Service connectors are typically responsible for service discovery and connection management, request dispatch, response handling, and some error handling

51.1 See also

- [Inter-Service communication](#) - There are multiple forms in which services communicate with each other and this describes the various forms.
- [Consumer-based testing](#) - To ensure that the service contract satisfies the consumer needs the Service Connector contains tests to confirm this.

52 Service Registry

Service Registry

Clients of a service user either [client side discovery](#) or [service side discovery](#) to determine the location of a service instance to send the requests. By providing a service registry, which is a database of service instances and their locations, the discovery mechanisms are able to find the available instances of a service.

The benefits of the service registry pattern is that the client of the service and/or [message router](#) can discover the location of service instances. There are also some drawbacks; Unless part of the cloud environment, the router is another system component that must be installed and configured. Moreover the service registry is a critical component. Although clients should [cache](#) data provided by the service registry, if the service registry fails that data will eventually become outdated. Consequently the service registry must be highly available.

Beside the service registry one also need to decide how service instances are registered with the service registry. There are two options:

1. [Self registration pattern](#) - service instances register themselves.
2. [3rd party registration pattern](#) - a 3rd party registers the service instances with the service registry.

The clients of the service registry need to know the location(s) of the service registry instances. This implicitly means that service registry instances must be deployed on fixed and well known locations. Clients are configured with those locations.

Examples of service registries include:

- [Eureka](#)
- [Zookeeper](#)
- [Consul](#)

Related patterns

- [Client side](#) - and [server side discovery](#) create the need for a service registry.
- [self registration](#) and [3rd-party-registration](#) are two different ways that service instances can be registered with the service registry.

53 Service Statelessness

Service Statelessness

The management of excessive state information can compromise the availability of a service and undermine its scalability principle. The services within a distributed application are deployed among multiple resources to benefit the scaling out functionality. The most significant factor complicating addition and removal of service instances is the internal state maintained by the service. In case of failure, this information might even be lost.

Services are therefore ideally designed to be stateless. Instead their state and configuration is stored externally in [storage offerings](#) or provided by the component with each request.

54 Single Responsibility Principle

Single Responsibility Principle

This principle states that, a subsystem, class or even function, should have only 1 reason to change. The classic example is a class that has methods that deal with business rules, reports, and persisting:

```
public interface Employee {  
    public Money calculatePay();  
    public int reportHours();  
    public void save();  
}
```

The problem with the interface shown above is that the functions change for entirely different reasons. The `calculatePay` function will change whenever the business rules for calculating pay change. The `reportHours` function will change whenever someone wants a different way of reporting hours. The `save` function will change whenever the database scheme is changed. These 3 reasons to change make `Employee` very volatile.

Applying the SRP principle means that we have to separate the interface into components that can be deployed independently. Independent deployment means that if we deploy 1 component we do not have to redeploy any of the others.

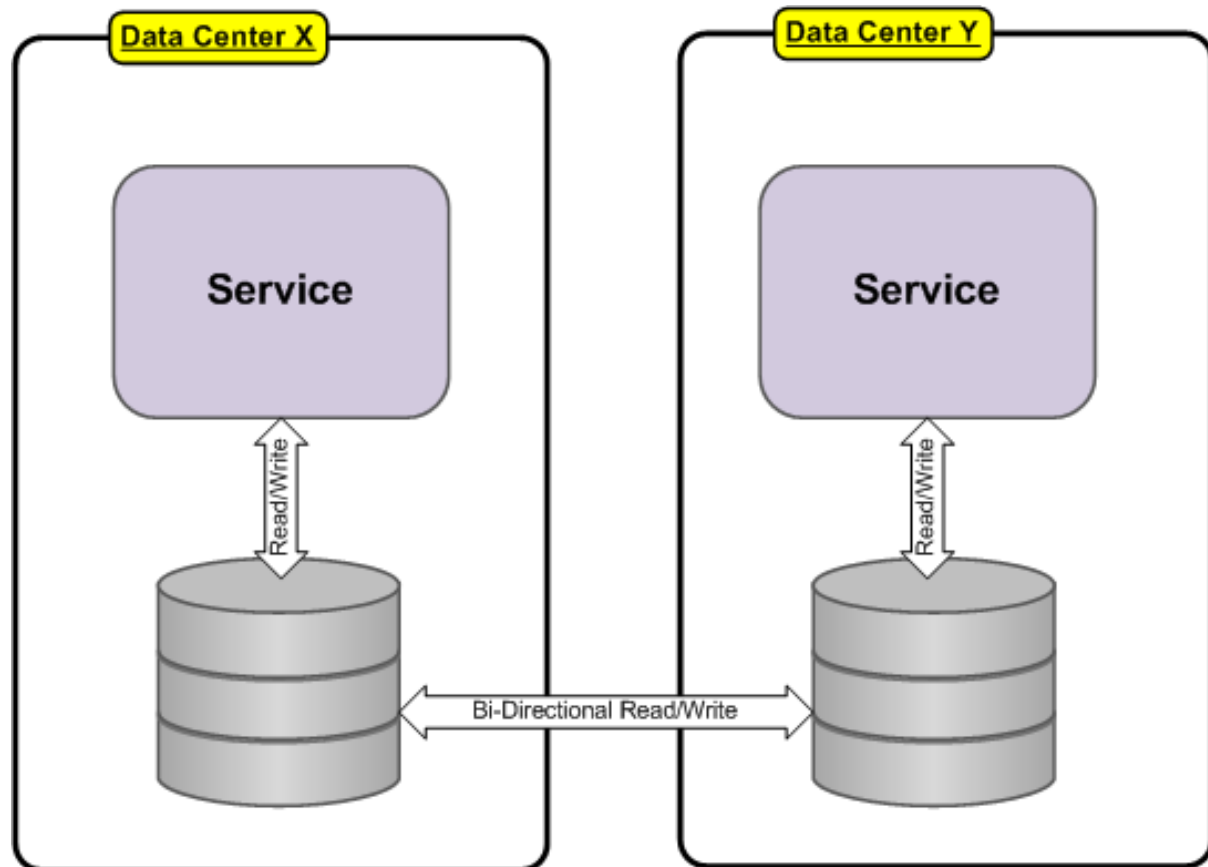
```
public interface Employee {  
    public Money calculatePay();  
}  
  
public interface EmployeeReporter {  
    public String reportHours(Employee e);  
}  
  
public interface EmployeeRepository {  
    public void save(Employee e);  
}
```

The simple partitioning shown above resolves the issue. We have to note that there is still a dependency from `Employee` to the other interfaces `EmployeeReporter` and `EmployeeRepository`. So if `Employee` is changed, the other classes will likely have to be recompiled and redeployed. We could prevent this through a careful use of the [Dependency Inversion Principle](#).

55 System of Record

System of Record

Multi master database systems that span sites have similar data replicated over multiple locations and the data is updated all the time. Goal of these kind of setups is little or no data loss on failure.



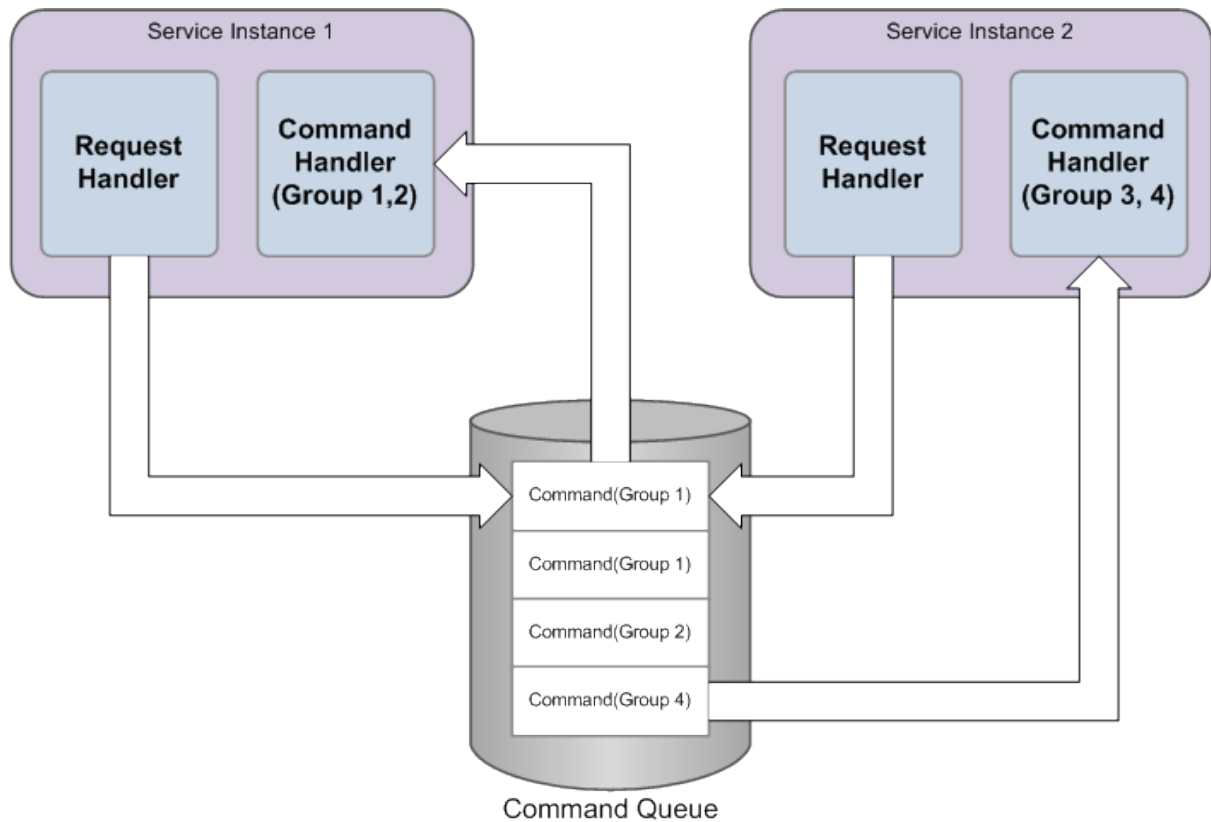
Problem with this solution is that it is impossible to build. The problem is the multi-master replication; updating the same record/table from two or more places on a LAN is already quite difficult, and the problem becomes unmanageable when you combine complex read/write operations, referential integrity, and high latency WAN connections.

Solution for this problem is system of record which states that individual records are updated in a single location only, but may have many copies elsewhere, (both locally as well as on other sites). When clients update particular information they do so on their 'own' master.

55.1 Related Patterns

- **Exclusive consumer** - To assure that only 1 instance is handling changes to certain entities we could use the exclusive consumer. In case of microservices this would mean that all instances of a particular service are listening to a particular queue and that the commands which are put on the queue are grouped by (for example) the ID of the **entity** at which the command applies. This would mean that the part that puts the command on the queue needs to know how to group the commands, (getting ID of entity at which the command applies from command). In the figure below we attempt to show how such a setup would look. Within this example there are 2 instances of the service whom both receive requests via their Request Handler. The Request Handler transforms the incoming requests into Commands, determine the group at which

they apply, and put them on the queue. The Command Handler will listen on the queue for commands within certain groups (standard functionality ActiveMQ) and process the commands within that group in the order that they were posted. Since there is only 1 command handler listening for a certain group we assure the system of record.



56 Timeouts

Timeouts

By defining timeouts we can specify how long to wait before a downstream system can be regarded as down. When the timeout is chosen too long you can slow down the whole system. Timeout too quickly, and you might have considered a call that might have worked as failed. Having no timeouts at all, and a downstream system being down could hang your whole system.

Put timeouts on all downstream calls and log when timeouts occur, and if required change them accordingly.

57 Try-Cancel/Confirm

Try-Cancel/Confirm.

The Try-Cancel/Confirm pattern focuses on transactions for microservices. A transaction is a set of related interactions (or operations) that may need to be cancelled after they were executed.

When a client initiates a state transition the service will return a handle by which the client can confirm or cancel the state change. If the service does not hear anything after some service specific timeout, it will cancel automatically. Once the workflow has completed successfully the set of returned handles is used to confirm the state transitions. If the service fails, the set of handles that has been collected until the failure is used to cancel the state transitions.

The timeout after which the service will automatically cancel the pending state transitions should be specified by the service.

57.1 Using RESTAT

You will need to deploy the coordinator as a war archive. The archive is contained in the `bin` folder of the [narayana download](#) (restat-web.war).

[RESTful transactions Quickstarts](#) [raw-xts-api-demo](#) [compensating transactions](#)

Related Patterns

- [API Gateway] - some use cases executed by/via an API gateway require multiple services. The API Gateway will orchestrate the sequence of service instantiations.

58 UUID

UUID

In an environment with multiple master relational databases, each master must be able to take updates, while also syncing correctly with others. One common concern is that unique identifiers must not conflict across master databases. Therefore server features such as table level auto-increment can not be used unless each table is only updated on one master database, or each server is given a distinct range of available auto-increment values.

One solution is for each database server or client to generate a UUID (Universal Unique Identifier) which is a string which is almost guaranteed to never conflict with another UUID generated on a different computer. Therefore each master database, or client connecting to any database, can generate a UUID and use it as the primary key on a table without much concern for replication conflicts.

Another advantage is that they can be freely exposed without disclosing sensitive information, and they are not predictable.

Downside to this choice is that every primary key is then a string which can be disadvantageous to performance on some systems and generating a massive amount of UUIDs can be expensive to process.

58.1 Related Patterns

- [Entities](#) - IDs of entities should be UUID.

59 3rd party registration

3rd Party Registration

Service instances must be registered with the [service registry](#) on startup so that they can be discovered and unregistered on shutdown.

Registration of the service should happen with the service registry on startup and unregister on shutdown. If a service instance crashes the service must be unregistered from the service registry. The same applies for service instances that are running but are incapable of handling requests.

Within the 3rd party registration a registrar is responsible for registering and un-registering a service instance with the service registry. When the service instance starts up, the registrar registers the service instance with the service registry and when the service shuts down the registrar un-registers the service from the service registry

The benefits of the 3rd party registration are:

- Service code is less complex than when using [self registration](#) since it is not responsible for registering itself.
- The registrar can perform health checks on a service instance and register/un-register the instance based on the health check.

Drawbacks of 3rd party registrations are:

- The registrar might only have a high level view of the service instance state (UP/DOWN) and so might not know whether it can handle requests.
- Unless the registrar is part of the infrastructure, it's another component that must be installed, configured and maintained. Also, since it is a critical system component it need to be highly available.

59.1 Related paterns

- [Service Registry](#)
- [Client](#) and [server side discovery](#).
- [Self registration](#) is an alternative solution.

60 Abbreviations

Abbreviations

- CQRS - Command-Query Responsibility Segregation
- DAS - Deployable Artifact Set
- DBC - Design By Contract
- DIP - Dependency Inversion Principle
- DTO - Data Transfer Object
- LOC - Lines Of Code
- LFU - Least Frequently Used
- LRU - Least Recently Used
- REST - REpresentational State Transfer
- SRP - Single Responsible Principle
- TCC - Try-Cancel/Confirm
- TPS - Transaction Per Second
- UUID - Unique Universal IDentifier
- VV - Version Vector