

# SPRING BOOT: THE HIDDEN MECHANICS

*Mental Models, Lifecycle Boundaries, and Failure Modes in  
Production*

DRAFT

**For experienced software engineers building Spring  
Boot applications in production**

Spring Boot is easy to use—until it isn't.

This book is written for experienced software engineers who build and operate Spring Boot systems in production. It focuses on how Spring Boot really behaves: auto-configuration mechanics, startup and bean lifecycles, proxy boundaries, concurrency, transactions, persistence, security, and failure handling at scale.

Rather than tutorials or API catalogs, the book presents mental models. It explains why certain problems appear, where behavior is decided, and how design choices interact with the framework under load and failure. Topics are concise, independent, and optimized for clarity over completeness.

You will not find setup instructions, beginner explanations, or exhaustive examples. You will find explanations that help you reason about Spring Boot systems during design reviews, incident analysis, and technical interviews.

AI-based writing tools were used selectively to support drafting and editing. All technical content, structure, and conclusions are the author's own.

### **This is not a guide to using Spring Boot.**

It is a guide to understanding what Spring Boot hides — and what happens when those abstractions fail.

**Spring and Spring Boot are trademarks of their respective owners; this book is an independent publication and is not affiliated with, endorsed by, or sponsored by VMware or the Spring team**

**Author: *Sujesh***

**© 2026 *Sujesh***

## Contents

Auto-Configuration: The Mental Model .....	3
Application Startup Lifecycle: From main() to Ready .....	6
@Component vs @Service vs @Repository vs @Bean .....	9
Circular Dependencies: Detection, Causes, and Design Fixes .....	12
Bean Scopes Explained with Real-World Use Cases .....	15
Bean Lifecycle Customization: Init and Destroy Done Right .....	18
Thread Safety and Singleton Beans .....	21
Designing a Scalable Spring Boot Microservice .....	23
Inter-Service Communication: REST vs Messaging vs gRPC .....	26
Failure Handling in Microservices .....	29
Distributed Tracing and Structured Logging .....	32
Diagnosing and Fixing Slow Spring Boot Applications .....	35
RestTemplate vs WebClient vs Reactive Programming .....	38
Concurrency and Async Processing in Spring Boot .....	41
Transaction Management Internals .....	44
Hibernate Performance Pitfalls .....	47
Safe Database Migrations in Production .....	50
Spring Security Internals .....	52
Designing JWT-Based Authentication .....	55
Testing Large Spring Boot Applications .....	58

## Auto-Configuration: The Mental Model

Spring Boot auto-configuration is not magic. It is a deterministic mechanism for registering beans based on a small set of rules. Once those rules are clear, auto-configuration becomes predictable.

### At a high level, Spring Boot auto-configuration works as follows:

- Discover candidate auto-configuration classes
- Import them into the ApplicationContext in a defined order
- Apply them conditionally based on classpath, existing beans, and properties
- Back off when the user provides their own beans

Everything else is an implementation detail.

### Discovery: where auto-configuration classes come from

Using `@SpringBootApplication` implicitly enables `@EnableAutoConfiguration`. This annotation instructs Spring Boot to load a list of auto-configuration class names from metadata provided by jars on the classpath.

### In Spring Boot 2.x, this metadata lives in:

META-INF/spring.factories

Key: `org.springframework.boot.autoconfigure.EnableAutoConfiguration`

### In Spring Boot 3.x, this moved to:

META-INF/spring/org.springframework.boot.autoconfigure.AutoConfiguration.imports

Starters such as `spring-boot-starter-web` contribute entries to this metadata. Adding a dependency is therefore not just adding code; it is signaling architectural intent. If a relevant class appears on the classpath, Spring Boot assumes you want the corresponding subsystem configured.

### Import phase: how auto-configurations enter the context

The discovered class names are imported into the ApplicationContext using an `ImportSelector` (historically `AutoConfigurationImportSelector`).

### Before importing, Spring Boot refines the list by:

- Removing duplicates
- Applying exclusions (exclude attribute or `spring.autoconfigure.exclude`)
- Applying ordering rules using annotations like `@AutoConfigureBefore` and `@AutoConfigureAfter`

Each auto-configuration class is a regular `@Configuration` class that defines one or more beans. There is no runtime scanning or late binding. The list is resolved deterministically during context creation.

## Conditional wiring: why not everything is created

Auto-configuration classes and their `@Bean` methods are guarded by conditional annotations. Common examples include:

- `@ConditionalOnClass` – only active if a class is present on the classpath
- `@ConditionalOnMissingBean` – only active if no user-defined bean exists
- `@ConditionalOnProperty` – only active if a property is set or matches a value
- `@ConditionalOnWebApplication` – only active in a web environment
- `@ConditionalOnBean` – only active if another bean already exists

Conditions are evaluated during the bean definition phase, not during bean instantiation. If a condition fails, the configuration silently backs off and is never reconsidered.

This is why adding or removing a dependency can change application behavior without modifying code. A condition flipped from false to true.

## The backoff rule (critical)

Spring Boot is designed to provide defaults, not override user intent.

For example, Spring Boot will create an `ObjectMapper` only if:

- Jackson is present on the classpath
- No `ObjectMapper` bean has already been defined

This backoff rule is the core contract of auto-configuration. Defaults exist until the user expresses intent. The framework optimizes for composability, not enforcement.

Many production issues stem from accidental intent: a bean defined indirectly or shared across modules disables an entire auto-configuration branch.

## Properties binding: how configuration becomes behavior

Auto-configuration relies heavily on `@ConfigurationProperties` to bind external configuration into typed objects.

For example, `spring.datasource.*` is bound into `DataSourceProperties`, which is then used to construct the `DataSource`. Configuration influences behavior indirectly through these bound objects, not through imperative logic.

This design keeps configuration declarative but also obscures cause-and-effect when defaults are overridden unintentionally.

## Debugging auto-configuration decisions

Most auto-configuration “failures” are silent backoffs. Debugging requires inspecting decisions, not stack traces.

## Practical tools include:

- Running with `--debug` or `debug=true` to print the Condition Evaluation Report

- The /actuator/conditions endpoint to see why configurations matched or did not
- The /actuator/beans endpoint to inspect what was actually created

The key debugging question is always the same: which condition caused this configuration to back off?.

DRAFT

## Application Startup Lifecycle: From main() to Ready

Spring Boot startup is often described as “slow” or “complex,” but most confusion comes from not knowing where time is spent and when decisions are made. The startup lifecycle is a sequence of well-defined phases. Understanding those phases is essential for debugging startup failures, tuning performance, and using lifecycle hooks correctly.

### At a high level, Spring Boot startup consists of:

- Launching SpringApplication
- Preparing the environment
- Creating and refreshing the ApplicationContext
- Registering bean definitions
- Instantiating beans
- Running startup callbacks
- Signaling readiness

Problems usually arise when code assumes a phase has happened earlier than it actually has.

### From main() to SpringApplication.run()

Application startup begins in main(), which delegates almost immediately to SpringApplication.run().

SpringApplication is not a thin wrapper. It orchestrates the entire bootstrap process, including environment loading, context creation, listener notification, and lifecycle callbacks. Most “Spring magic” happens here, before a single bean is instantiated.

At this point, no ApplicationContext exists. No beans exist. Only configuration metadata is being prepared.

Environment preparation happens first

Before the ApplicationContext is created, Spring Boot builds the Environment.

This includes:

- Loading property sources (application.yml, profiles, system properties, env vars)
- Resolving active profiles
- Applying property overrides and defaults

This ordering is intentional. Bean definitions, conditions, and auto-configuration decisions all depend on the Environment. If a property affects behavior, it must be available before bean registration begins.

A common mistake is assuming @Value or @ConfigurationProperties failures are “bean problems.” They are often environment problems that happened much earlier.

### ApplicationContext creation (but not population)

Once the Environment is ready, Spring Boot creates the ApplicationContext.

At this stage:

- The context object exists
- No beans are instantiated
- Bean definitions have not yet been fully loaded

The specific context type (Servlet, Reactive, or non-web) is chosen here. This choice affects which auto-configurations and conditions will match later.

Importantly, creating the context does not mean the application is usable. It is an empty container waiting to be filled.

### **Bean definition phase vs instantiation phase**

This distinction is critical and frequently misunderstood.

First, Spring registers bean definitions:

- From component scanning
- From `@Configuration` classes
- From auto-configuration

During this phase:

- `@Conditional` annotations are evaluated
- `@Bean` methods are analyzed
- Proxies are planned, not created

No constructors are called. No `@PostConstruct` methods run. Decisions made here are final for the lifetime of the context.

Only after all bean definitions are registered does Spring move to instantiation.

### **Bean instantiation and dependency wiring**

During instantiation:

- Singleton beans are created
- Dependencies are resolved and injected
- Proxies are created
- Lifecycle callbacks begin

This is where constructors execute and where most startup time is spent in real systems. Heavy constructors, eager initialization, and blocking I/O show up here.

`@PostConstruct` runs after dependency injection but before the application is considered ready. Code here blocks startup and should be treated as part of boot time, not runtime.

### **ApplicationContext refresh and startup hooks**

Once all singleton beans are created, the `ApplicationContext` is refreshed.

At this point:

- `ApplicationContextAware` callbacks have run

- BeanPostProcessors have completed
- The container is stable

Spring then invokes startup hooks such as:

- ApplicationRunner
- CommandLineRunner

These run after the context is fully initialized but before the application signals readiness. They are suitable for validation, warm-up, or verification logic, but not for long-running background work unless explicitly asynchronous.

### **Ready state and liveness implications**

After startup runners complete, Spring Boot signals that the application is ready.

In web applications, this typically means:

- The embedded server is accepting requests
- Health and readiness endpoints report UP

This boundary matters in containerized and orchestrated environments. Code that runs before readiness blocks traffic. Code that runs after readiness must be thread-safe and resilient.

Confusing startup work with runtime work is a common cause of slow deploys and unstable rollouts.

### **Where startup failures really come from**

Most startup failures fall into one of three categories:

- Environment misconfiguration (missing properties, wrong profiles)
- Bean definition conflicts (conditional mismatches, duplicate beans)
- Instantiation failures (constructor logic, missing dependencies)

Knowing which phase you are in drastically reduces debugging time. Stack traces alone are often misleading without lifecycle context.

## @Component vs @Service vs @Repository vs @Bean

These annotations look similar, behave similarly, and are often treated as interchangeable. They are not. The differences are not about syntax, but about semantic intent, proxy behavior, and where control lives.

At scale, misunderstanding these distinctions leads to fragile wiring, accidental behavior changes, and confusing proxy interactions.

### The core split is simple:

@Component, @Service, and @Repository are discovered by component scanning.  
@Bean is an explicit declaration inside a @Configuration class.

Everything else follows from that.

Component stereotypes: intent first, behavior second

@Component is the generic stereotype. @Service and @Repository are specializations of it. From the container's perspective, all three are component-scanned beans.

Their primary value is semantic, not functional.

@Service communicates business logic.

@Repository communicates data access.

@Component communicates “infrastructure or glue code.”

Spring itself uses these signals selectively, but humans use them constantly. In large codebases, semantic clarity matters more than the annotation's mechanics.

### @Repository and exception translation

@Repository is the only stereotype with consistent behavioral impact.

Classes annotated with @Repository are eligible for persistence exception translation. Low-level persistence exceptions are converted into Spring's DataAccessException hierarchy.

This matters because:

- It standardizes error handling across data access technologies
- It prevents leaking vendor-specific exceptions upward
- It only applies when the bean is detected as a repository

Annotating repositories as @Component technically works, but it opt-outs of a key Spring feature. That choice is rarely intentional and often forgotten.

### @Service is about boundaries, not mechanics

@Service does not change how the container creates the bean. There is no special lifecycle, proxying, or threading behavior attached to it.

Its value is architectural. It marks a boundary where:

- Transactions often start
- Business rules live
- Infrastructure concerns should not leak in

Senior engineers use `@Service` to communicate intent to both Spring and other developers. That signal becomes critical during refactoring and incident response.

## Component scanning vs explicit bean definitions

`@Component`-style annotations rely on classpath scanning. This is convenient, but implicit.

`@Bean` definitions are explicit. They live in `@Configuration` classes and declare exactly:

- What gets created
- With what dependencies
- Under what conditions

This distinction matters when behavior must be controlled precisely. Scanning optimizes for speed of development. Explicit beans optimize for correctness and predictability.

## @Bean and configuration control

`@Bean` methods are evaluated during the bean definition phase, just like scanned components. However, they give you tighter control over construction logic.

`@Bean` is the right choice when:

- You are wiring third-party classes
- Construction requires non-trivial logic
- You want configuration localized and reviewable
- You need conditional or profile-specific creation

In contrast, annotating a class with `@Component` pushes construction logic into the class itself, which often mixes configuration with behavior.

## Proxy behavior and method interception

Proxies are applied based on capabilities, not stereotypes.

`@Transactional`, `@Async`, and similar annotations create proxies regardless of whether the bean came from `@Component` or `@Bean`. The stereotype does not control proxying.

However, `@Bean` methods inside `@Configuration` classes are themselves proxied. This ensures:

- Singleton semantics for inter-bean method calls
- Correct lifecycle handling

Moving logic between `@Component` and `@Bean` can change proxy boundaries subtly, especially when self-invocation is involved.

## When explicit beats scanning

Component scanning works well until it doesn't.

Explicit `@Bean` definitions are usually superior when:

- Behavior must be obvious from configuration
- Multiple implementations exist
- Conditional creation matters
- Startup order is relevant
- You are debugging auto-configuration interactions

Senior teams often start with scanning and gradually move critical wiring into configuration as systems grow.

## Common failure patterns

Frequent production issues include:

- Missing exception translation due to incorrect repository annotation
- Accidental bean overrides caused by scanning collisions
- Configuration logic hidden inside component constructors
- Confusing proxy behavior after refactoring annotations

These failures are rarely caught by tests and usually surface during integration or production.

## A useful rule of thumb

Use component stereotypes to express role.

Use `@Bean` to express control.

If you need the container to “figure it out,” scanning is fine.

If you need to be sure, be explicit.

## Circular Dependencies: Detection, Causes, and Design Fixes

Circular dependencies are not a Spring quirk. They are a design signal. Spring detects them because it has to, not because it encourages them. When they appear, the container is telling you something about coupling and lifecycle assumptions.

A circular dependency exists when two or more beans depend on each other directly or indirectly during creation. The problem is not the cycle itself, but when the dependency is required.

### How Spring detects circular dependencies

Spring detects circular dependencies during bean instantiation, not during bean definition.

The container keeps track of beans currently in creation. If bean A requires bean B, and bean B requires bean A, Spring identifies the cycle and attempts resolution depending on injection style.

Constructor injection creates a hard requirement. Both beans must be fully constructed before either can exist. Spring cannot break this cycle and fails fast.

Setter and field injection are softer. Spring can create one bean, expose an early reference, and inject it later. This is why some cycles “work” while others fail.

### Early references and why they are dangerous

To resolve non-constructor cycles, Spring uses early references. A partially constructed bean is exposed before initialization completes.

This allows wiring to succeed, but it introduces risk:

- Methods may be called on uninitialized objects
- Proxies may not be fully applied
- Lifecycle callbacks may not have run

The system boots, but correctness becomes fragile. These issues surface later as intermittent bugs, not startup failures.

### Constructor injection vs setter injection

Constructor injection makes dependencies explicit and enforces immutability. It also forces Spring to detect cycles immediately.

Setter or field injection hides dependencies and allows Spring to defer wiring. This often masks architectural problems and delays failure until runtime.

This is why constructor injection is generally preferred. It trades convenience for correctness and pushes design issues to the surface early.

### The role of proxies in circular dependencies

Proxies complicate circular dependencies.

When A depends on B and B depends on A, and one or both are proxied (for transactions, async execution, etc.), the injected reference may be a proxy rather than the actual object.

This can result in:

- Self-invocation bypassing proxies
- Methods executing outside transactional boundaries
- Behavior changing after refactoring annotations

Cycles involving proxied beans are harder to reason about and more error-prone.

### **@Lazy as a mitigation, not a fix**

@Lazy delays bean instantiation until the dependency is actually needed.

This can break circular dependencies by replacing a direct reference with a proxy that resolves later. While useful in edge cases, it does not fix the underlying coupling.

Overusing @Lazy hides design issues and creates unpredictable runtime behavior. It should be treated as a tactical workaround, not a structural solution.

### **Why circular dependencies indicate design smells**

Circular dependencies usually indicate one of the following:

- Responsibilities are not clearly separated
- Business logic and orchestration are mixed
- Two components share too much knowledge

In healthy designs, dependencies tend to form a directed graph. Cycles imply mutual control, which makes change harder and testing more brittle.

### **Refactoring strategies that actually work**

Effective fixes usually involve:

- Introducing a third component to own shared behavior
- Inverting control via events or callbacks
- Splitting responsibilities along clearer boundaries
- Moving orchestration into a higher-level service

The goal is not to “satisfy Spring,” but to restore directional dependencies.

### **When cycles are sometimes acceptable**

Not all cycles are equally harmful.

Framework-level wiring, infrastructure components, or legacy boundaries may contain unavoidable cycles. In such cases, the goal is containment and documentation, not elimination.

However, cycles in core business services should be treated as defects, not features.

### **A practical rule of thumb**

If a circular dependency only works because of setter injection, early references, or @Lazy, the design is already compromised.

If breaking the cycle improves clarity, testability, and startup behavior, it is almost always the correct choice.

DRAFT

## Bean Scopes Explained with Real-World Use Cases

Bean scope defines how long a bean lives and who shares it. Most Spring applications rely almost entirely on singletons, often without realizing it. Scope mistakes rarely fail fast; they show up as race conditions, memory leaks, or subtle data corruption in production.

Understanding scope is less about memorizing options and more about knowing where state is allowed to exist.

### Singleton scope: the default and the trap

Singleton is the default scope in Spring. One instance exists per `ApplicationContext`.

This does not mean:

- One instance per request
- One instance per thread
- Thread-safe by default

Singleton beans are shared across all threads. Spring does not add synchronization or immutability guarantees. Any mutable state stored in a singleton must be explicitly protected.

Most production concurrency bugs in Spring applications come from accidental mutable state inside singleton beans.

### Prototype scope: new instance, new problems

Prototype scope creates a new bean instance every time it is requested from the container.

This sounds useful, but comes with important constraints:

- Spring manages creation, not destruction
- `@PreDestroy` is not called
- The container does not track lifecycle after injection

Injecting a prototype bean into a singleton does not create a new instance per use. The prototype is resolved once, at injection time. This surprises many developers and leads to false assumptions about isolation.

Prototype scope is only effective when the container is explicitly asked for new instances.

### Request and session scopes: web-only lifetimes

Request scope creates one instance per HTTP request.

Session scope creates one instance per HTTP session.

These scopes only exist in a web-aware `ApplicationContext`. Outside of a web request, they are meaningless.

Request-scoped beans are useful for:

- Per-request context
- Authentication or correlation data
- Request-local caching

Session-scoped beans introduce server-side state and should be used sparingly. They complicate horizontal scaling and increase memory pressure.

### Scoped proxies: bridging scope mismatches

Injecting a request-scoped bean into a singleton is not possible directly. The singleton lives longer than the request.

Spring solves this with scoped proxies. Instead of injecting the actual bean, Spring injects a proxy that resolves the real instance at runtime based on the active scope.

This allows scope mixing, but adds indirection and runtime resolution. Misuse can make execution flow harder to reason about and debug.

### Async execution and scope boundaries

Scopes interact poorly with async execution.

Request and session scopes are bound to the request thread. When execution moves to another thread via `@Async`, `CompletableFuture`, or reactive pipelines, those scopes are no longer active.

Accessing scoped beans asynchronously often results in:

- Scope not active exceptions
- Silent context loss
- Inconsistent behavior between sync and async paths

Passing data explicitly is usually safer than relying on scoped beans in async flows.

### Scope and memory behavior

Longer-lived scopes amplify memory issues.

Singletons retain references for the lifetime of the application. Session-scoped beans retain memory for the lifetime of the session. Leaks here accumulate slowly and are hard to detect.

Request-scoped beans are short-lived and GC-friendly, but only if they do not reference longer-lived objects.

Understanding scope lifetimes helps reason about memory retention, not just object creation.

### When to store state and when not to

A useful rule: treat singleton beans as stateless services.

If state is required:

- Prefer method-local variables
- Use request scope for request-bound state
- Externalize state to databases or caches

State inside long-lived beans should be immutable or carefully synchronized. Otherwise, behavior depends on timing and load.

## Common failure patterns

Frequent scope-related issues include:

- Mutable fields inside singleton services
- Assuming prototype beans are recreated per use
- Request-scoped beans leaking into async execution
- Session scope used as a convenience cache

These issues rarely appear in tests and usually surface under concurrency or load.

## A practical rule of thumb

If you are unsure which scope to use, singleton is usually correct—but only if the bean is stateless.

When scope becomes complicated, the design often is too.

DRAFT

## Bean Lifecycle Customization: Init and Destroy Done Right

Spring manages bean lifecycles, but it also allows customization at specific points. Misusing lifecycle hooks is a common cause of slow startup, unreliable shutdown, and hidden coupling. Understanding when these hooks run matters more than knowing how to declare them.

Lifecycle callbacks exist to integrate with the container, not to replace application logic.

The high-level lifecycle

### For a typical singleton bean, the lifecycle is:

- Bean definition registered
- Bean instantiated
- Dependencies injected
- Lifecycle callbacks executed
- Bean becomes usable
- Context shuts down
- Destruction callbacks executed

Customization hooks operate only at defined points in this sequence. Code that assumes otherwise usually behaves unpredictably.

Initialization callbacks: what runs and when

### Spring supports multiple initialization mechanisms:

`@PostConstruct`

`InitializingBean.afterPropertiesSet()`

Custom init methods defined on `@Bean`

All of these run after dependency injection but before the application is considered ready.

Ordering is fixed:

`@PostConstruct` runs first, then `afterPropertiesSet()`, then custom init methods.

Code here blocks startup. Any I/O, remote calls, or heavy computation increases boot time and delays readiness.

### Choosing the right initialization hook

`@PostConstruct` is declarative and concise. It is usually the right default.

`InitializingBean` is framework-coupled and rarely necessary outside infrastructure components.

Custom init methods on `@Bean` are useful when:

- You do not control the class
- Initialization must be conditional
- Configuration should own lifecycle behavior

Using multiple init mechanisms on the same bean is almost always a mistake.

## **Destruction callbacks and shutdown behavior**

Destruction callbacks run when the `ApplicationContext` shuts down gracefully.

Spring supports:

`@PreDestroy`

`DisposableBean.destroy()`

Custom destroy methods on `@Bean`

These callbacks are best-effort. They do not run on abrupt termination and should not be relied on for critical persistence.

Shutdown hooks should release resources, not perform business logic.

## **Shutdown timing and ordering**

Spring shuts down beans in reverse dependency order. This ensures that dependent beans are destroyed before the beans they rely on.

However, this ordering only applies to managed dependencies. External threads, executors, and resources must be handled explicitly.

Failing to shut down executors is a common cause of:

- Hanging shutdowns
- Forced termination
- Resource leaks

## **Lifecycle hooks vs application readiness**

Initialization does not mean readiness.

`@PostConstruct` and `init` methods run before:

- `ApplicationRunner`
- `CommandLineRunner`
- Readiness probes signaling UP

Work that must happen after the application is “live” should not run in lifecycle hooks. Mixing these phases causes slow deploys and failed rollouts.

## **Lifecycle hooks and proxies**

Lifecycle callbacks run on the target bean, not on the proxy.

This matters when:

- Initialization logic expects transactional behavior
- `@Async` or `@Transactional` annotations are present

Proxy-based features are not active during initialization callbacks. Assuming otherwise leads to subtle bugs.

## Common lifecycle anti-patterns

Frequent mistakes include:

- Remote calls inside `@PostConstruct`
- Long-running background threads started during init
- Business logic in destroy callbacks
- Multiple lifecycle mechanisms on the same bean

These patterns often work in development and fail under orchestration or load.

## A practical rule of thumb

Use lifecycle hooks to validate configuration and initialize local resources.

Do not use them for orchestration, remote coordination, or application logic.

If startup or shutdown behavior feels complex, the design usually is.

DRAFT

## Thread Safety and Singleton Beans

Spring's default bean scope is singleton, but Spring does not guarantee thread safety. This mismatch is one of the most common sources of production bugs in Spring Boot applications. The container manages lifecycle and wiring, not concurrent access.

Understanding this distinction is critical when building services that handle real load.

### What singleton actually means

In Spring, singleton means one instance per `ApplicationContext`.

It does not mean:

- One instance per request
- One instance per thread
- Serialized access

A singleton bean is shared across all threads handling requests. If it contains mutable state, that state is shared as well.

### Why Spring does not enforce thread safety

Spring intentionally does not make beans thread-safe.

Thread safety requirements vary by use case. Enforcing synchronization at the container level would impose unnecessary overhead and restrict design choices. Spring leaves concurrency control to the application.

This design favors flexibility but requires discipline.

### Common sources of shared mutable state

Typical sources include:

- Mutable fields used as temporary storage
- Caches implemented as instance variables
- `DateFormat`, `SimpleDateFormat`, or similar non-thread-safe utilities
- Reused buffers or collections

These issues often remain invisible until traffic increases or execution becomes parallel.

### Stateless services as the default

The safest singleton is a stateless singleton.

Business services should generally:

- Avoid storing request-specific data in fields
- Treat methods as pure functions over inputs
- Use method-local variables for transient state

Stateless design scales naturally and avoids most concurrency problems.

## When state is unavoidable

Sometimes state is necessary.

In those cases:

- Prefer immutable state initialized at startup
- Use thread-safe data structures where mutation is required
- Externalize state to databases, caches, or message queues

Synchronization should be explicit and intentional. Accidental synchronization is usually a performance bug.

## Synchronization strategies and trade-offs

Common approaches include:

- Synchronized blocks or methods
- Concurrent collections
- Atomic variables
- Explicit locks

Each has trade-offs in complexity, throughput, and contention. Over-synchronization often reduces scalability more than it improves correctness.

Thread safety should be designed, not patched.

## Interaction with async execution

Async execution multiplies thread-safety issues.

@Async, executor-based processing, and parallel streams increase concurrency beyond request handling. Code that “worked fine” synchronously often fails under async execution.

If a singleton is accessed asynchronously, it must be safe under concurrent access from multiple threads.

## Thread safety and proxies

Proxies do not make beans thread-safe.

@Transactional, @Async, and similar annotations add behavior around method calls, not around field access. They do not serialize execution or protect mutable state.

Assuming that “Spring handles concurrency” is a common and costly mistake.

## A practical rule of thumb

If a singleton bean has mutable fields, assume it is unsafe until proven otherwise.

If making it thread-safe feels complicated, the design likely needs to change.

## Designing a Scalable Spring Boot Microservice

Scalability is not a framework feature. Spring Boot can support scalable systems, but it does not make design decisions for you. Most scalability failures are architectural, not technical, and appear long before traffic becomes extreme.

A scalable microservice is predictable under load, tolerant of partial failure, and cheap to operate at higher throughput.

### Statelessness is the foundation

Horizontal scaling depends on stateless services.

A stateless service:

- Does not store request-specific data in memory
- Can handle any request on any instance
- Can be restarted or replaced without coordination

State that must persist beyond a request belongs in external systems: databases, caches, or message brokers. In-memory state ties requests to instances and limits scalability.

### Instance scaling vs workload scaling

Scaling instances increases capacity only if the workload scales linearly.

Common bottlenecks that prevent linear scaling include:

- Shared databases with limited write throughput
- Synchronized sections in application code
- Thread pool exhaustion
- Downstream services with lower capacity

Adding instances without addressing bottlenecks increases cost without improving throughput.

### Database design dominates scalability

For most services, the database is the real scaling limit.

Key considerations:

- Read/write ratios and contention hotspots
- Proper indexing and query shape
- Connection pool sizing
- Avoiding chatty transactional patterns

Caching often provides more scalability than adding application instances, but only when cache invalidation and consistency are understood.

### Idempotency under retries

Retries are inevitable in distributed systems. A scalable service must tolerate them.

Idempotent operations ensure that:

- Repeated requests do not cause duplicate side effects
- Clients can retry safely
- Failure handling does not corrupt state

This often requires idempotency keys, deduplication logic, or database constraints. Treat idempotency as a first-class design concern, not an afterthought.

### **Caching as a scaling lever**

Caching reduces load, but it also introduces complexity.

Effective caching requires:

- Clear ownership of cached data
- Defined expiration and invalidation rules
- Awareness of stale reads

Blind caching improves throughput but can undermine correctness. Cache what you can recompute or tolerate being slightly stale.

### **Threading and resource isolation**

Each microservice instance has finite resources.

Scalability depends on:

- Correct thread pool sizing
- Avoiding blocking calls in hot paths
- Isolating slow or unreliable dependencies

Unbounded queues and shared executors are common failure modes. They hide backpressure until the system collapses.

### **Backpressure and load shedding**

A scalable service must say “no” under load.

Mechanisms include:

- Bounded thread pools
- Timeouts on all remote calls
- Queue limits
- Explicit rejection strategies

Failing fast preserves overall system health. Accepting more work than you can handle does not increase throughput; it increases latency and failure rates.

### **Observability as a scaling requirement**

You cannot scale what you cannot observe.

At minimum, a scalable service exposes:

- Latency percentiles

- Error rates
- Saturation indicators (threads, pools, queues)

Scaling decisions based on averages are misleading. Tail latency drives user experience and system instability.

### **A practical rule of thumb**

If a service cannot be safely restarted at any time, it is not scalable.

If scaling increases operational complexity faster than capacity, the design is wrong.

DRAFT

## Inter-Service Communication: REST vs Messaging vs gRPC

Inter-service communication is where microservice designs usually succeed or fail. The choice is not about preference or trend; it is about coupling, latency, failure modes, and evolution over time.

Spring Boot supports multiple communication styles, but it does not remove the trade-offs between them.

### Synchronous vs asynchronous as a design choice

The first decision is not REST vs messaging vs gRPC. It is synchronous vs asynchronous.

Synchronous communication couples:

- Availability
- Latency
- Deployment timing

Asynchronous communication trades immediacy for decoupling and resilience.

Most systems use both. Problems arise when synchronous calls are used where asynchronous behavior is required.

### REST: simple, visible, tightly coupled

REST over HTTP is the default choice.

Strengths:

- Simple mental model
- Human-readable
- Easy to debug
- Wide tooling support

Weaknesses:

- Tight runtime coupling
- Latency amplification across call chains
- Fragile under retries and partial failures

REST works best for request-response interactions where low latency and immediate consistency are required.

### Feign vs WebClient

Feign provides declarative HTTP clients. It optimizes for readability and consistency but hides execution details.

WebClient is explicit and non-blocking. It provides fine-grained control over timeouts, retries, and threading, but requires more discipline.

Key distinction:

Feign favors developer ergonomics.

WebClient favors control and scalability.

In high-throughput or latency-sensitive services, implicit blocking becomes a scalability bottleneck.

### **Messaging: decoupling through time**

Messaging systems (Kafka, queues, event streams) decouple producers from consumers.

Benefits:

- Loose coupling
- Natural buffering and backpressure
- Better failure isolation

Costs:

- Eventual consistency
- Operational complexity
- Harder debugging

Messaging is ideal when services should not depend on each other's availability or response time.

### **At-least-once delivery and idempotency**

Most messaging systems guarantee at-least-once delivery.

This means:

- Duplicate messages are normal
- Consumers must be idempotent
- Side effects must be repeatable

Designing consumers without idempotency is a correctness bug, not a messaging problem.

### **gRPC: efficient but constraining**

gRPC provides:

- Binary serialization
- Strongly typed contracts
- Low latency and high throughput

It excels in internal service-to-service communication where performance matters.

Trade-offs include:

- Reduced debuggability
- Tighter schema coupling
- More rigid evolution

gRPC works best when teams align on contracts and versioning discipline.

## Coupling analysis over protocol choice

Protocol choice matters less than coupling shape.

Key questions:

- Does this call require immediate consistency?
- What happens if the dependency is slow or down?
- Can this interaction tolerate retries or reordering?
- Who owns the data and the contract?

Many failures blamed on “the protocol” are actually coupling failures.

## Backpressure and failure propagation

Synchronous calls propagate pressure upstream.

Without:

- Timeouts
- Circuit breakers
- Bounded thread pools

a slow downstream service can exhaust resources across the system.

Asynchronous systems absorb pressure but require explicit capacity planning and monitoring.

## A practical rule of thumb

Use REST for simple, synchronous queries with clear ownership.

Use messaging for workflows, fan-out, and decoupled processing.

Use gRPC when performance dominates and contracts are stable.

If a communication choice feels “convenient,” it is probably hiding a trade-off.

## Failure Handling in Microservices

In distributed systems, failure is normal. What distinguishes reliable systems is not the absence of failure, but how failures are isolated, absorbed, and recovered from. Spring Boot does not solve this for you; it only provides integration points.

Most production outages are caused by unbounded failure propagation, not by a single component failing.

### Timeouts are mandatory, not optional

Every remote call must have a timeout.

Without timeouts:

- Threads block indefinitely
- Resource pools are exhausted
- Failures cascade upstream

Timeouts define failure boundaries. They should be set based on business requirements, not defaults. A fast failure is usually preferable to a slow one.

### Retries: useful but dangerous

Retries increase resilience but also amplify load.

Safe retries require:

- Idempotent operations
- Bounded retry counts
- Backoff strategies

Blind retries turn transient failures into sustained outages. Retrying without understanding downstream capacity is a common failure pattern.

### Circuit breakers and failure isolation

Circuit breakers stop repeated calls to unhealthy dependencies.

They:

- Prevent resource exhaustion
- Allow systems to degrade reminders gracefully
- Create space for recovery

Circuit breakers are not about hiding failures. They are about preventing failures from spreading.

### Bulkheads and resource partitioning

Bulkheads isolate failures by partitioning resources.

Examples include:

- Separate thread pools per dependency

- Connection pool limits
- Queue boundaries

Without bulkheads, one slow dependency can starve unrelated traffic. Isolation preserves partial functionality under stress.

### **Cascading failures and amplification**

Cascading failures occur when:

- A slow service causes retries
- Retries increase load
- Increased load causes more failures

This feedback loop can bring down healthy services. Breaking the loop requires early rejection and load shedding.

### **Graceful degradation**

Not all functionality is equally critical.

Graceful degradation means:

- Serving partial responses
- Falling back to cached or stale data
- Disabling non-essential features

Designing fallbacks after an outage is too late. They must be part of the original design.

### **Resilience4j and integration points**

Spring integrates with libraries like Resilience4j for:

- Timeouts
- Retries
- Circuit breakers
- Rate limiting

These tools are powerful, but configuration matters more than the library choice. Poorly tuned resilience mechanisms are worse than none.

### **Observability of failure handling**

Failure handling without observability is guesswork.

You must be able to see:

- Timeout rates
- Retry counts
- Circuit breaker states
- Saturation indicators

Without visibility, resilience mechanisms mask problems instead of revealing them.

### **A practical rule of thumb**

If failure handling increases complexity more than it reduces risk, the design is wrong.

Systems should fail fast, isolate damage, and recover predictably.

DRAFT

## Distributed Tracing and Structured Logging

In distributed systems, failures rarely occur in isolation. A single user request may traverse multiple services, threads, and hosts. Without correlation, logs become noise and metrics lack context. Observability is not about collecting more data, but about making causality visible.

Tracing and logging serve different purposes. Confusing them leads to blind spots.

### Why logs alone stop scaling

Traditional logging assumes a single process boundary.

In microservices:

- Requests span services
- Execution is concurrent
- Failures propagate indirectly

Plain logs answer “what happened here,” but not “why did this request fail.” As systems grow, log volume increases faster than insight.

### Correlation IDs as the foundation

Correlation starts with a request identifier.

A correlation ID:

- Is generated at the edge
- Is propagated across service boundaries
- Appears in every log line related to the request

Without consistent propagation, tracing collapses. With it, even basic logs become navigable.

### Distributed tracing: seeing the full path

Distributed tracing tracks a request across services and time.

Key concepts:

- Trace: the full request lifecycle
- Span: a unit of work within a trace
- Context propagation: passing trace metadata downstream

Tracing reveals latency sources, fan-out patterns, and unexpected dependencies that are invisible in isolated logs.

### Trace propagation and thread boundaries

Trace context must cross:

- HTTP boundaries
- Messaging systems
- Async execution
- Thread pools

Losing context breaks traces. Async code, executor misconfiguration, and custom threading are common places where propagation fails silently.

If traces look incomplete, context loss is usually the cause.

## OpenTelemetry as the standard

OpenTelemetry defines a vendor-neutral model for:

- Traces
- Metrics
- Logs

It separates instrumentation from backend choice. This avoids lock-in and allows teams to change observability backends without rewriting code.

The key value is standardization, not tooling.

## Zipkin-style tracing and limitations

Tracing systems such as Zipkin visualize request flows and timings.

They excel at:

- Identifying slow dependencies
- Understanding call graphs
- Debugging latency spikes

They do not explain business intent or correctness. Tracing shows where time went, not whether behavior was right.

## Structured logging over free text

Structured logs are machine-readable.

Instead of free text, logs emit:

- Consistent fields
- Explicit identifiers
- Typed values

This enables filtering, aggregation, and correlation at scale. Free-text logs optimize for humans reading files, not systems analyzing streams.

## Log levels and signal discipline

More logs do not mean more observability.

Effective logging requires:

- Clear distinction between info, warn, and error
- Avoiding debug logs in hot paths
- Logging decisions, not noise

Logs should explain why something happened, not restate what the system already knows.

### **A practical rule of thumb**

If you cannot answer “what happened to this request” in minutes, observability is insufficient.

If tracing exists but is rarely used, it is probably too noisy or incomplete.

DRAFT

## Diagnosing and Fixing Slow Spring Boot Applications

Performance problems rarely have a single cause. They emerge from interactions between the JVM, the framework, the database, and the workload. Treating slowness as a tuning problem instead of a diagnostic problem usually makes things worse.

The goal is not to make the system fast in isolation, but predictable under load.

### Start with symptoms, not solutions

Performance debugging begins with observation.

Before tuning anything, answer:

- Is latency high, or throughput low?
- Is the problem constant or load-dependent?
- Does it affect all requests or specific paths?

Optimizing without isolating the symptom often hides the real bottleneck.

JVM behavior matters more than most code

In many Spring Boot services, JVM behavior dominates performance.

Common JVM-related issues include:

- Excessive garbage collection
- Memory pressure from object churn
- Thread contention

Heap sizing, GC choice, and allocation patterns usually matter more than micro-optimizations in application code.

### Startup performance vs runtime performance

Startup slowness and runtime slowness are different problems.

Startup issues typically come from:

- Eager bean initialization
- Heavy `@PostConstruct` logic
- Classpath scanning and auto-configuration

Runtime issues usually involve:

- Blocking calls
- Database latency
- Thread pool exhaustion

Optimizing one does not fix the other.

Database access is the usual bottleneck

Most “slow” applications are waiting on the database.

Common causes:

- Missing or ineffective indexes
- N+1 query patterns
- Chatty transactional workflows
- Connection pool misconfiguration

Database latency propagates upward and consumes application threads while waiting.

### **Caching as a performance multiplier**

Caching can dramatically improve performance, but only when used deliberately.

Effective caching requires:

- Identifying read-heavy paths
- Choosing appropriate cache scope and TTL
- Understanding consistency trade-offs

Caching the wrong layer or data often increases complexity without improving latency.

### **Thread pools and blocking behavior**

Thread pools define concurrency limits.

Problems arise when:

- Blocking I/O occurs on limited pools
- Pools are undersized for peak load
- Work is submitted faster than it can be processed

Symptoms include high latency, request timeouts, and uneven throughput.

### **Async execution is not free**

Async processing increases throughput only if it reduces blocking.

Using `@Async` or reactive APIs while still performing blocking I/O shifts the bottleneck instead of removing it. This often makes debugging harder without improving performance.

Concurrency increases coordination costs. It should be applied intentionally.

### **Profiling mindset over tools**

Tools matter less than approach.

Effective profiling focuses on:

- Where threads are waiting
- What resources are saturated
- Which code paths dominate time

Sampling profilers, thread dumps, and slow request traces usually provide more insight than fine-grained metrics alone.

### **A practical rule of thumb**

If performance improves after increasing resources, the system was under-provisioned.

If it does not, the design is the bottleneck.

DRAFT

## RestTemplate vs WebClient vs Reactive Programming

The choice between RestTemplate, WebClient, and reactive programming is often framed as “old vs new.” That framing is misleading. The real question is whether your system benefits from non-blocking execution and whether you are prepared to pay the complexity cost.

Using reactive APIs without a reactive design rarely improves performance and often makes systems harder to reason about.

Blocking vs non-blocking is the core distinction

Blocking I/O ties up a thread while waiting for a response.

Non-blocking I/O frees the thread and resumes work when data is available.

Blocking scales by adding threads.

Non-blocking scales by multiplexing work over fewer threads.

Neither model is universally better. Each optimizes for different constraints.

### RestTemplate: simple and thread-bound

RestTemplate is synchronous and blocking.

Strengths:

- Simple mental model
- Easy debugging
- Predictable execution

Costs:

- One thread per in-flight request
- Limited scalability under high concurrency
- Thread pool pressure during downstream slowness

RestTemplate works well for low to moderate concurrency and simple call chains. Its limitations appear when downstream latency increases or fan-out grows.

### WebClient: API choice, not architecture

WebClient is non-blocking by default, but using it does not automatically make an application reactive.

WebClient can be used in:

- Blocking mode (calling `block()`)
- Non-blocking mode (reactive chains)

Calling `block()` negates most scalability benefits while retaining complexity. This hybrid approach is common and usually accidental.

## Reactive programming: end-to-end commitment

Reactive programming changes how control flow, error handling, and backpressure are expressed.

Benefits:

- Better resource utilization under high concurrency
- Natural backpressure
- Improved resilience for I/O-bound workloads

Costs:

- Steeper learning curve
- Harder debugging
- Less intuitive stack traces
- Limited benefit for CPU-bound work

Reactive systems must be reactive end-to-end. Mixing blocking dependencies into reactive pipelines creates hidden bottlenecks.

### When reactive helps

Reactive programming is most effective when:

- Workloads are I/O-bound
- Concurrency is high
- Latency variance is significant
- Backpressure is required

Typical examples include API gateways, aggregation services, and fan-out-heavy workloads.

### When reactive hurts

Reactive programming often hurts when:

- Business logic is CPU-heavy
- Most dependencies are blocking
- Team familiarity is low
- Operational debugging matters more than raw throughput

In these cases, reactive complexity outweighs its benefits.

### Threading models and resource usage

Blocking applications scale by increasing thread pools, which increases memory usage and context switching.

Reactive applications use fewer threads but require careful control of execution contexts. Misconfigured schedulers or blocking calls can collapse performance.

Non-blocking does not mean “no limits.” It shifts where limits must be enforced.

## **Error handling and observability trade-offs**

Reactive error handling is explicit and compositional, but unfamiliar.

Exceptions become signals.

Stack traces become fragmented.

Logs lose linear flow.

Observability must be designed intentionally. Without tracing and correlation, reactive systems are harder to debug than blocking ones.

## **A practical rule of thumb**

Choose the simplest model that meets your scalability requirements.

If you do not have a clear concurrency or backpressure problem, blocking I/O is often the better choice.

Adopting reactive programming without a concrete scaling goal is a liability.

DRAFT

## Concurrency and Async Processing in Spring Boot

Concurrency increases throughput, but it also increases complexity. Spring provides multiple ways to execute work asynchronously, but none of them remove the need to reason about threads, resources, and failure modes. Most concurrency bugs come from assuming the framework “handles it.”

Async execution should be a deliberate design choice, not an optimization reflex.

### Threads are the real currency

Every asynchronous mechanism ultimately consumes threads.

Concurrency is limited by:

- Thread pool size
- Blocking behavior
- Downstream latency

Increasing concurrency without understanding thread usage leads to saturation, not scalability.

### @Async: simple but constrained

@Async allows methods to run on a separate thread pool.

Strengths:

- Minimal code changes
- Clear intent
- Works well for fire-and-forget tasks

Limitations:

- Requires proxy invocation
- No implicit backpressure
- Easy to misuse for long-running work

@Async is best suited for short, isolated tasks that can fail independently.

### Task executors and pool sizing

Spring uses TaskExecutor abstractions backed by thread pools.

Key considerations:

- Core and max pool sizes
- Queue capacity
- Rejection policies

Unbounded queues hide overload by accumulating work until memory or latency collapses. Bounded queues force explicit backpressure.

### CompletableFuture and composition

CompletableFuture enables asynchronous composition without full reactive adoption.

Benefits:

- Explicit async boundaries
- Better control than `@Async`
- Familiar imperative style

Costs:

- Error handling complexity
- Manual executor management
- Risk of blocking inside async stages

`CompletableFuture` works well for moderate concurrency and structured parallelism.

### **Async does not mean non-blocking**

Async execution often moves blocking work to a different thread rather than eliminating blocking.

Common mistakes include:

- Performing blocking I/O inside async tasks
- Using async to hide slow dependencies
- Assuming parallelism increases throughput

If all threads are blocked, async provides no benefit.

### **Backpressure and overload handling**

Async systems need explicit overload control.

Without backpressure:

- Queues grow unbounded
- Latency increases
- Failures propagate late

Backpressure mechanisms include bounded queues, timeouts, and rejection strategies. Failing fast is usually safer than processing everything slowly.

### **Context propagation pitfalls**

Async execution breaks implicit context.

Security context, MDC logging context, and tracing information do not automatically propagate across threads. Losing context leads to:

- Incomplete logs
- Broken traces
- Security bugs

Context propagation must be handled explicitly.

### **Testing and debugging async code**

Async code is harder to test and debug.

Common issues:

- Race conditions
- Timing-dependent failures
- Flaky tests

Deterministic execution, explicit synchronization, and time-bound assertions are essential for reliable tests.

### **A practical rule of thumb**

If async logic makes the system harder to reason about than the problem it solves, it is the wrong tool.

Concurrency should reduce waiting, not clarity.

DRAFT

## Transaction Management Internals

Spring's transaction management is simple to use and easy to misunderstand. `@Transactional` looks declarative, but its behavior is driven by proxies, method boundaries, and propagation rules. Most transaction bugs come from assuming it behaves like a language feature rather than a framework mechanism.

Transactions are not attached to code blocks. They are attached to method invocations.

### What `@Transactional` actually does

`@Transactional` does not start a transaction by itself. It marks a method as transactional, which causes Spring to wrap the bean in a proxy.

When a transactional method is invoked through the proxy:

- A transaction is started or joined
- The method executes
- The transaction is committed or rolled back

If the method is not invoked through the proxy, none of this happens.

### Proxy boundaries and the self-invocation trap

Transactional behavior applies only when calls cross the proxy boundary.

When a method within a class calls another `@Transactional` method on the same class, the call bypasses the proxy. No transaction is created, changed, or propagated.

This is the self-invocation trap. It is one of the most common causes of “`@Transactional` not working” bugs.

The fix is architectural, not annotational: move transactional boundaries to external services or ensure calls go through the proxy.

### Transaction propagation: joining vs creating

Propagation defines how a method behaves when a transaction already exists.

The most common modes are:

- `REQUIRED` – join an existing transaction or create one
- `REQUIRES_NEW` – suspend the existing transaction and create a new one
- `SUPPORTS` – run with or without a transaction

Propagation is about composition. Misusing it leads to unexpected commits, rollbacks, or lock contention.

### Isolation levels and what they really protect

Isolation controls how concurrent transactions interact.

Higher isolation levels reduce anomalies but increase locking and reduce throughput. Lower levels increase concurrency but allow inconsistent reads.

The important point: isolation is enforced by the database, not Spring. Spring only requests a level; the database decides what is actually supported.

Tuning isolation without understanding database behavior is ineffective.

### **Rollback rules and exception semantics**

By default, Spring rolls back transactions only on unchecked (runtime) exceptions.

Checked exceptions do not trigger rollback unless explicitly configured. This design assumes checked exceptions represent recoverable conditions.

Misaligned exception hierarchies lead to partial commits and data inconsistency. Rollback behavior must match business semantics, not Java exception taxonomy.

### **Transaction boundaries vs business boundaries**

Transactions should align with business invariants.

Common mistakes include:

- Long-running transactions spanning remote calls
- Transactions used as workflow coordination
- Transactions wrapping async execution

Long transactions increase lock duration and reduce system throughput. Transactions are for consistency, not orchestration.

### **Transactions and proxies with other concerns**

Transactional proxies often coexist with other proxies such as security or async execution.

Order matters:

- `@Transactional` with `@Async` rarely behaves as expected
- Async execution runs outside the original transaction
- Security context may or may not propagate

Mixing concerns without understanding proxy order produces subtle bugs.

### **Testing transactional behavior**

Transactional tests can hide real behavior.

Tests often:

- Run in a single transaction
- Roll back automatically
- Mask commit-time failures

What passes in tests may fail in production. Critical transactional paths should be tested with real commits.

### **A practical rule of thumb**

If you cannot explain where a transaction starts, where it ends, and what happens on failure, the design is unsafe.

Transactions should be small, explicit in intent, and aligned with business consistency.

DRAFT

## Hibernate Performance Pitfalls

Hibernate is often blamed for performance problems that are actually caused by misunderstanding its behavior. The ORM trades explicit control for convenience. When that trade-off is not managed deliberately, performance degrades in non-obvious ways.

Most Hibernate performance issues are invisible until the system is under real load.

### The N+1 query problem

The N+1 problem occurs when Hibernate loads a collection or association lazily and then issues one query per element.

This usually happens when:

- Lazy associations are accessed inside loops
- Serialization triggers entity traversal
- Business logic assumes data is already loaded

N+1 problems scale with data size and are rarely noticeable in small test datasets.

### Fetch strategies and their trade-offs

Hibernate supports eager and lazy fetching.

Eager fetching:

- Reduces query count
- Increases data volume
- Can create Cartesian product explosions

Lazy fetching:

- Reduces upfront cost
- Defers loading
- Risks N+1 queries

There is no universally correct strategy. Fetching must be shaped around access patterns, not entity relationships.

### Join fetching and query shape

Join fetching allows related entities to be loaded in a single query.

It is effective for:

- Read-heavy queries
- Known access paths
- Bounded result sets

However, join fetching large collections multiplies rows and increases memory pressure. It must be used selectively and tested with realistic data volumes.

## First-level cache behavior

Hibernate's first-level cache is the persistence context.

Within a transaction:

- Each entity is loaded once
- Repeated lookups return the same instance
- Changes are tracked automatically

This cache improves consistency but can increase memory usage and dirty checking costs in long-running transactions.

## Second-level cache realities

The second-level cache is optional and global.

Benefits:

- Reduces database load
- Speeds up repeat reads

Costs:

- Cache invalidation complexity
- Consistency trade-offs
- Operational overhead

Second-level caching only helps when access patterns are stable and read-heavy. It is not a universal performance solution.

## Batching and write performance

By default, Hibernate writes entities one at a time.

Batching:

- Reduces round-trips
- Improves throughput
- Requires careful configuration

Batching works best for bulk operations. Mixing batching with cascades and entity graphs can reduce its effectiveness.

## Dirty checking and flush behavior

Hibernate tracks changes to managed entities and flushes them to the database.

Unexpected flushes can occur:

- Before query execution
- At transaction boundaries
- When auto-flush is triggered

This can cause performance spikes and lock contention. Understanding flush timing is essential for predictable performance.

## Pagination pitfalls

Pagination with joins is tricky.

Common mistakes include:

- Paginating after join expansion
- Fetching large result sets and slicing in memory
- Using OFFSET on large tables

Efficient pagination often requires query redesign, keyset pagination, or denormalization.

## A practical rule of thumb

If you cannot predict how many queries a request executes, you do not control performance.

Measure query count, not just execution time.

DRAFT

## Safe Database Migrations in Production

Database migrations are one of the highest-risk activities in production systems. Unlike application code, database changes are shared across versions, services, and deployments. A safe migration strategy assumes partial failure, mixed versions, and the need for rollback.

Most migration incidents come from treating schema changes as atomic and reversible. They are neither.

### Schema changes are not code deploys

Application code can be rolled back quickly. Database schema changes often cannot.

Once a column is dropped or data is transformed, rollback becomes expensive or impossible. Production migrations must be designed to be forward-compatible and survivable under failure.

Flyway vs Liquibase: control vs abstraction

Flyway emphasizes explicit, versioned SQL migrations.

Liquibase provides higher-level abstractions and supports rollback metadata.

Trade-offs:

- Flyway favors predictability and simplicity
- Liquibase favors flexibility and database portability

The tool choice matters less than discipline. Unsafe migrations are unsafe regardless of tooling.

### Backward-compatible migrations

Safe migrations are backward compatible.

Common patterns include:

- Add new columns before using them
- Avoid dropping or renaming fields immediately
- Deploy code that tolerates both old and new schemas

Backward compatibility allows old and new application versions to coexist during rolling deployments.

### Zero-downtime migration patterns

Zero-downtime migrations require separating schema change from behavior change.

Typical sequence:

- Add new schema elements
- Deploy code that writes to both old and new structures
- Backfill existing data
- Switch reads to new structures
- Remove old schema later

This process is slower but dramatically reduces risk.

### **Data backfills and long-running migrations**

Large data migrations should not run inside deployment transactions.

Risks include:

- Long locks
- Transaction log growth
- Replication lag

Backfills should be incremental, resumable, and observable. Treat them as production jobs, not schema changes.

### **Locks, indexes, and traffic impact**

Schema changes can block writes and reads.

Index creation, column type changes, and constraint validation may lock tables. Migration impact must be tested under realistic data volume and traffic.

“Fast in staging” often means “dangerous in production.”

### **Rollback strategy in practice**

True rollback often means forward fixes.

Rather than undoing schema changes, production rollbacks typically:

- Revert application behavior
- Leave schema changes in place
- Apply corrective migrations later

Design migrations with the assumption that rollback means coexistence, not reversal.

### **Migration ownership and safety checks**

Migrations should be:

- Reviewed like production code
- Tested against production-like data
- Owned by the service team

Automated checks help, but human review is essential for high-risk changes.

### **A practical rule of thumb**

If a migration cannot safely run while multiple application versions are live, it is unsafe.

If rollback requires dropping data, it is not a rollback.

## Spring Security Internals

Spring Security is powerful, flexible, and frequently misunderstood. Most confusion comes from treating it as a set of annotations instead of a request-processing pipeline. Security behavior emerges from filter ordering, context propagation, and decision points—not from individual annotations.

Understanding the internals turns security from trial-and-error into a controlled system.

### Security is a filter chain, not an annotation

Spring Security is built around a servlet filter chain.

Every incoming request passes through a sequence of security filters before reaching application code. These filters handle authentication, context setup, authorization checks, and cleanup.

Annotations like `@PreAuthorize` do not secure endpoints by themselves. They are evaluated only after the request has already passed through the filter chain.

### Authentication vs authorization

Authentication answers: “Who is the caller?”

Authorization answers: “What is the caller allowed to do?”

Spring Security separates these concerns deliberately.

Authentication typically occurs early in the filter chain and results in an `Authentication` object. Authorization happens later, using that authentication to make access decisions.

Mixing these concepts leads to brittle and confusing security rules.

### The SecurityContext lifecycle

The `SecurityContext` holds the current `Authentication`.

During request processing:

- The context is created or loaded
- Authentication is stored in it
- Authorization checks read from it

By default, the `SecurityContext` is thread-bound. Once execution leaves the request thread—via async execution or custom threading—the context is no longer available unless explicitly propagated.

This is a common source of “works in sync, fails in async” security bugs.

### Filter ordering and why it matters

Filter order defines behavior.

For example:

- Authentication must happen before authorization

- Exception translation must wrap protected resources
- Context cleanup must happen last

Custom filters placed incorrectly can:

- Bypass authentication
- Break exception handling
- Leak security context across requests

Security issues caused by filter misordering are subtle and hard to detect in tests.

## Authorization decisions and voters

Authorization is evaluated using access decision logic.

Historically, Spring Security used voters and decision managers. Modern configurations abstract this, but the model remains the same: multiple checks contribute to a final allow or deny decision.

Understanding this helps when:

- Multiple authorization rules apply
- Behavior changes unexpectedly after refactoring
- Method-level and URL-level security interact

Authorization failures are often rule conflicts, not missing permissions.

Method security and proxies

## Method-level security uses proxies.

This implies:

- Calls must cross the proxy boundary
- Self-invocation bypasses security checks
- Proxy order matters when combined with transactions or async

Security annotations on private methods or internal calls often provide a false sense of protection.

## Exception handling and error translation

Spring Security translates security exceptions into HTTP responses.

Authentication failures, access denials, and internal errors are handled by different components.

Misconfiguration here leads to:

- Incorrect status codes
- Information leakage
- Inconsistent client behavior

Security is not only about access control, but also about how failures are exposed.

## Common internal failure patterns

Frequent issues include:

- Assuming annotations secure endpoints without filters

- Losing SecurityContext in async execution
- Misordered custom filters
- Self-invocation bypassing method security

These failures rarely surface as obvious security holes. They appear as inconsistent behavior.

### **A practical rule of thumb**

If you cannot trace a request through the security filter chain, you do not understand its security.

Security should be explainable as a sequence of steps, not as a collection of annotations.

DRAFT

## Designing JWT-Based Authentication

JWT-based authentication is often adopted to achieve statelessness, but stateless does not mean simple. JWTs move complexity from the server to the system design. Most problems arise from treating tokens as sessions rather than as signed claims.

JWTs are a tool, not an authentication strategy.

### What a JWT actually represents

A JWT is a signed set of claims.

It represents:

- Who the caller is
- What the caller is allowed to do
- Under what constraints (time, issuer, audience)

It does not represent server-side state. Once issued, a JWT is valid until it expires unless additional controls exist.

### Statelessness and its real cost

Stateless authentication removes server-side session storage.

Benefits:

- Horizontal scalability
- Easier load balancing
- Reduced shared state

Costs:

- Harder revocation
- Token leakage risk
- Longer blast radius for compromised credentials

Statelessness simplifies infrastructure but complicates security guarantees.

### Token lifetime and expiration strategy

Token lifetime is a core design decision.

Short-lived tokens:

- Reduce damage from compromise
- Increase refresh frequency

Long-lived tokens:

- Reduce auth overhead
- Increase security risk

Most production systems use short-lived access tokens combined with longer-lived refresh tokens to balance usability and risk.

## Refresh tokens and rotation

Refresh tokens allow access tokens to remain short-lived.

Safe refresh token designs include:

- Server-side storage or tracking
- Rotation on use
- Revocation on suspicious behavior

Refresh tokens are credentials. Treating them as harmless metadata is a security bug.

## Revocation is not optional

JWTs are not revocable by default.

Common revocation strategies:

- Short expiration times
- Token blacklists
- Key rotation
- Audience or version checks

Each approach trades simplicity for control. Systems that require immediate revocation must reintroduce server-side state.

## Claims design and overloading

JWT claims should be minimal and stable.

Common mistakes:

- Encoding business logic into claims
- Including large or frequently changing data
- Treating claims as a database

Claims are authorization hints, not authoritative state. Overloaded tokens create tight coupling and frequent reissuance.

## Signature validation and trust boundaries

JWTs are only as secure as their validation.

Critical checks include:

- Signature verification
- Issuer validation
- Audience validation
- Expiration and clock skew

Skipping or relaxing these checks is equivalent to disabling authentication.

## JWTs in microservices

In microservice architectures, JWTs are often propagated downstream.

This enables:

- End-to-end identity propagation
- Decentralized authorization

It also spreads trust. Every service that accepts the token becomes part of the security boundary. Compromise of one service can affect others.

### **A practical rule of thumb**

If you need immediate revocation, fine-grained control, or server-enforced sessions, pure JWT-based authentication is the wrong tool.

JWTs work best when identity is stable, lifetimes are short, and trust boundaries are clear.

DRAFT

## Testing Large Spring Boot Applications

Testing is not about coverage. In large Spring Boot systems, tests exist to preserve behavior while the system evolves. The challenge is not writing tests, but choosing the right boundaries so tests remain fast, meaningful, and maintainable.

Most testing pain comes from mixing concerns and testing too much at once.

### Test types and their roles

Different tests answer different questions.

Unit tests:

- Validate local logic
- Run fast
- Do not require Spring

Integration tests:

- Validate wiring and configuration
- Exercise framework behavior
- Catch misconfiguration

End-to-end tests:

- Validate critical flows
- Are slow and brittle
- Should be few and targeted

Confusing these layers leads to slow builds and fragile test suites.

### Unit tests: isolate aggressively

Unit tests should not start Spring.

They should:

- Test business logic in isolation
- Mock external dependencies
- Avoid framework concerns

If a unit test requires `ApplicationContext` startup, it is not a unit test.

### Integration tests: test the wiring, not the world

Integration tests validate that components work together under Spring.

Good integration tests:

- Load a limited context
- Focus on one slice of functionality
- Assert behavior, not implementation

Spring's test slicing helps limit scope, but slices must be chosen intentionally.

## Testcontainers and realism

Mocks hide infrastructure problems.

Testcontainers allows tests to run against real dependencies such as databases or message brokers. This increases confidence at the cost of speed.

Use Testcontainers when:

- Behavior depends on database semantics
- Integration bugs are costly
- Production parity matters

Not every test needs production realism.

## Mocking boundaries and ownership

Mocking is about ownership.

Mock:

- External services you do not control
- Unreliable or slow dependencies

Do not mock:

- Core domain logic
- Code you own and want to refactor safely

Over-mocking creates tests that validate assumptions instead of behavior.

## Transactional tests and false confidence

Transactional tests often roll back by default.

This:

- Keeps tests isolated
- Hides commit-time failures
- Masks database constraints and triggers

Critical persistence logic should be tested with real commits to catch production-only failures.

## Managing test performance

Slow tests kill feedback loops.

Common causes include:

- Starting full application contexts
- Reusing expensive test fixtures
- Overusing end-to-end tests

Optimizing test performance is a design activity, not a tooling tweak.

## **CI-friendly testing strategies**

Tests must run reliably in CI.

This requires:

- Deterministic execution
- Explicit timeouts
- Minimal environmental dependencies

Flaky tests erode trust and are worse than missing tests.

## **A practical rule of thumb**

Tests should make refactoring safer, not scarier.

If changing code breaks many unrelated tests, the test boundaries are wrong.

DRAFT