# Android
## In The Real World

Fabio Chiarani

# 3

# So you think you know Android?

*Every time you call an API, you are not executing code, you are awakening a system.*

Your request travels through layers you will never see, negotiating with services, runtimes, drivers, and silicon. From your perspective it is a function call; from the system's perspective, it is a disturbance in the whole machine.

## 3.8 A Journey Through the Android Stack: When Your App Starts a Bluetooth Scan

It begins with a single line of Kotlin. A quiet, harmless instruction, sitting somewhere in your code: maybe inside a ViewModel, maybe inside a Composable callback, maybe behind a UI button that the user just tapped:

```
bluetoothAdapter.startDiscovery()
```

To an application developer, it feels small. It is just a function call. An intention. A polite request to *please look for nearby devices.*

But Android is not a polite system. As we have seen in (Figure 3.1), it is a layered organism, and every request you make must pass through

each layer, each one interpreting, validating, translating, and enforcing its own rules.

This example Bluetooth scan is not a local action. It is a **journey**:

## The App Layer, where Intent begins

The moment your app executes `startDiscovery()`, nothing actually happens in the physical world. No radio waves move, and no frequency bands are scanned. Not yet.

At this point, you are simply handing a high level intention to Android: "*I would like to scan for Bluetooth devices*".

Your code is declarative, abstract, human. But Android is not. Inside the `BluetoothAdapter`, your intent is wrapped, packaged, and prepared for going down to the layers of the system.

The real work is about to begin.

Your request enters the Framework, a cathedral of rules and constraints carved across years of OS evolution.

The Framework never trusts your app. It checks:

- you have permission to scan?

- are you scanning too frequently?

- is the system in a restricted mode?

- are you respecting background limits?

Only once you pass the tests does the Framework forward your intention to a lower, quieter part of the system: the system services, where the bureaucracy of Android lives: `BluetoothManagerService`, `AdapterService`, and others.

They are the administrators of your request. Your function call becomes an entry in their ledger. And then it moves further down.

## Binder: the Crossing Border Between Worlds

Your process is not the one that controls Bluetooth. The system does. To submit your request, Android uses Binder, its custom inter process communication system. Imagine your request as a stamped envelope:
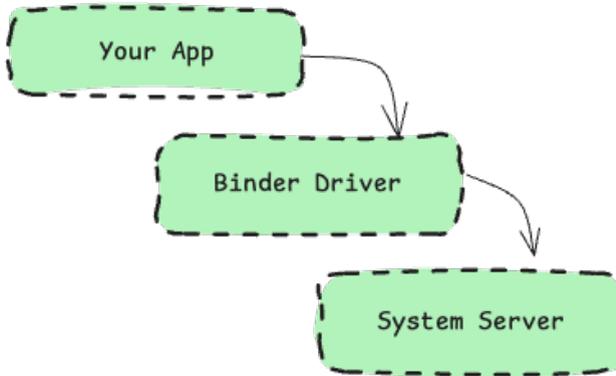
Figure 3.3: ART

So Binder serializes your intention, transports it across a boundary your code cannot cross alone, enforces identity and permissions, wakes the receiving service, and delivers the message. Now your request is no longer yours.

It belongs to the system.

## The Native Bluetooth Stack, Where Software Meets Reality

Deep inside the system sits a second engine: the native Bluetooth stack, written in C/C++.

This layer has no Kotlin, no coroutines, no ViewModels. It is procedural, stateful, mechanical.

It decides which mode of scanning to use, how long to scan, which filters to apply, how to manage hardware power levels, and how to parse incoming HCI events. For the first time, your request becomes something more concrete:

A command that will eventually reach hardware.

But it still is not there yet. The native stack prepares low level instructions and pushes them downward, toward another boundary: the HAL.

## The HAL: the translator of worlds

The Hardware Abstraction Layer exists because no two devices are the same. Samsung's Bluetooth chip is not Google's. Qualcomm firmware is not Broadcom firmware. Each device hides its own scars, its own quirks, its own silicon personality.

The Android Bluetooth HAL translates Android's universal commands into vendor-specific instructions.

It speaks two languages.

Your request changes shape here. It becomes raw, technical, stripped of abstractions. The HAL passes it to the entity that finally understands it: the vendor driver.

## The Vendor Driver: where electricity begins to move

Now we are inside the kernel, at the edge of hardware.

Vendor drivers receive the command, wrap it into HCI (Host Controller Interface) packets, and deliver it to the firmware running inside the Bluetooth chip.

Here, the request becomes real: frequencies shift, antennas activate, the radio listens, BLE advertising channels are scanned, and packets are captured out of the air.

This is where the physical world responds to your code. The chip listens to the invisible chatter around it: beacons, wearables, headsets, IoT devices, and then sends discoveries upward, back into the layers of Android and up to your app.

Okay, we have talked enough about how Android is structured beneath your application. I know you are eager to discover all the tips and tricks hidden in this book, but let us take it step by step.

[...]

# 7

# Everything Scales Until it Doesn't

*You wake up one morning, and you have two million users.*

You are on the sofa on a Sunday morning, coffee in hand, scrolling through notifications, when your phone buzzes nonstop. Two million users. Overnight, your Android app, the one you built with an MVP architecture on a shoestring schedule, has exploded. Downloads spike, reviews pop up with five stars and one star alike, and the analytics dashboard looks like a firehose. What do you do?

Your architecture, which was "good enough" for hundreds of users, is suddenly stress tested in ways you never imagined. Users are experiencing crashes on specific devices and OS versions you never saw before. You need to analyze metrics, diagnose crashes, test hypotheses, and deploy patches before the bad reviews cascade and retention plummets. On top of that, you want to roll out A/B tests for new conversion flows, analyze backend performance, and support a backend that is now carrying millions of requests per day. This is not theoretical; this is real scalability in Android development.

Real scalability problems involve multiple dimensions: handling large user bases, maintaining code quality across a growing codebase, managing team coordination as more engineers join, and ensuring responsiveness across diverse hardware. In this section we address the first dimension:

what it means for an Android app to scale beyond initial expectations, and the immediate tactical and architectural questions you face when the numbers go from hundreds to millions.

## 7.1  What We Mean When We Say "It Scales"

In Android development, saying an app "scales" is shorthand for a set of capabilities: the application continues to function correctly and responsively across:

- a rapidly increasing number of users;

- a rapidly increasing number of crashes and ANR;

- a growing, more complex codebase;

- a wider range of hardware and OS versions;

- and a development team that expands over time.

A scalable system must handle not only *more users* but also *more variation.* Android's ecosystem spans from low-end devices to high-end flagships, and users may be on different API levels and hardware configurations. This diversity means that issues will surface that never appeared in your early QA. For example, crash reporting at scale becomes essential: integrating a robust crash analytics system such as Google Crashlytics allows you to collect, analyze, and prioritize crash data across millions of sessions. Crashlytics is designed to gather crash reports from a global user base, providing insights into the most common failure paths and the devices most affected.

When you wake up to millions of users, the app must remain stable, responsive, and informative. Performance regressions that were invisible at small scale become glaring: UI jank, main thread contention, slow network handling, database lock contention, and even memory leaks become user-facing issues. The definition of scalability in this context is not just *"can the server handle it?"*, but *"can the entire system, client, server, and pipeline, continue to deliver a predictable experience"*?

[...]

These metrics should be sliced by device tier, Android version, and release cohort. A 500ms regression on low-end devices may be invisible in the overall median but lethal for retention.

**Technical checklist:**

- Do you monitor performance by device class rather than global averages?

  - *Yes: metrics are sliced by device tier (low, mid, high), RAM/CPU class, Android version, and OEM cohorts, and you review distributions (P50/P95/P99) instead of just means. This makes "small cohort, high pain" visible.*

  - *No: you track global averages, which are dominated by modern devices. Regressions on low-end hardware and poor networks remain invisible until ratings, churn, or support tickets spike.*

- Are startup regressions blocked by CI or release gates?

  - *Yes: you have thresholds and regression budgets, and you block rollout when cold start time or time-to-first-frame drifts beyond acceptable deltas. The gate is automated, and exceptions require explicit sign-off.*

  - *No: startup performance is monitored passively. Regressions are discovered after release, when rollback is more expensive and the feedback loop is dominated by external signals rather than engineering data.*

- Can you map jank spikes to specific screens and UI interactions?

  - *Yes: frame metrics are labeled by screen and interaction context, so you can pinpoint where jank happens, reproduce it with targeted traces, and validate fixes with before/after comparisons.*

  - *No: you see a global jank number without attribution. You know the app is slower, but you cannot identify which UI path is responsible, so performance work becomes guesswork.*

## 9.4  Tracing Critical Flows End-to-End

Logs and metrics show symptoms. Traces show causality. End-to-end tracing connects a user action to its downstream effects: network calls, database writes, and UI state transitions. This is how you diagnose invisible failures like slow syncs or inconsistent caches.

Tracing does not need to be heavy. Even lightweight tracing that measures the latency of key flows (login, search, purchase, sync) will tell you where product reliability is leaking. The key is to trace at the product boundary, not just at the technical boundary.

**Technical checklist:**

- Do you trace the time from user intent to visible result?

  - *Yes: critical flows define clear start and end markers (tap, navigation, "content visible") and record end-to-end latency. This makes reliability measurable as perceived by users, not just as internal timings.*

  - *No: you only measure internal segments (network time, DB time) without connecting them to user-perceived completion, so slow paths and coordination delays remain unaccounted for.*

- Are traces linked to release versions and experiment buckets?

  - *Yes: every trace carries version/build metadata and experiment assignment, so you can answer "did variant B increase login latency on Android 12 devices?" without reconstructing context manually.*

  - *No: traces exist but are not attributable to configuration, which makes them poor evidence for rollout decisions and slows regression isolation.*

- Can you identify which downstream dependency caused the delay?

  - *Yes: traces include spans for network, database, disk, and major subsystems, plus reason codes for retries, cache misses, and timeouts. When latency spikes, you can assign blame to a dependency rather than the entire app.*

[...]

## 5.7 Frame Time Anatomy: A 16ms Budget Autopsy

The sixteen millisecond number is often repeated without context, but the real constraint is more brutal: you do not own sixteen milliseconds. On Android, the display refresh interval (e.g., 16.67 ms at 60 Hz) is shared with the OS, SurfaceFlinger, the scheduler, drivers, background work, and thermals. Your code executes in what is left, and what is left varies by device, refresh rate, and current load.

Figure 5.3: Frame on time

A single frame is not one operation. It is a pipeline of dependent stages that must complete early enough for SurfaceFlinger to latch your buffer, compose, and present it at the next vsync. Missing any stage means the new buffer cannot be used for that refresh. The user does not see a partially updated frame; they see the previous frame repeated, which is perceived as jank.

Figure 5.3 illustrates the healthy case: input is dispatched, the UI thread performs its work (animations, measure/layout when needed, and draw recording), RenderThread submits rendering commands, and the GPU finishes execution soon enough that the frame is ready for presentation at VSync #(N+1). In other words, the pipeline for *Frame K* completes before the vsync-driven deadline, so the display advances smoothly.

[...]

# 4

# The Compose Illusion: Everything Works, Until It Doesn't

*The real art in Compose isn't making everything work, it's discovering why sometimes it doesn't.*

Oh yes, that is the chapter you want to read. We will dive into the core and highly technical aspects of Jetpack Compose. At first glance, Compose can feel almost magical in its simplicity: you define functions, describe what should be on the screen, and everything seems to work effortlessly.

But behind this declarative facade lies a sophisticated mechanism of state handling, composition, and recomposition, capable of producing subtle and unexpected behavior when you least expect it. We'll explore in detail how Compose updates the UI in response to state changes, why this process is efficient yet sometimes deceptive, and what key concepts (such as recomposition, stability, and state management) you need to master in order to truly understand what's happening under the hood.

## 4.1  Thinking in compose:  from declaration to pixel

Jetpack Compose is a declarative UI toolkit: you describe what should appear on screen as a function of state, and Compose turns that state into UI. In the classic imperative model, UI is a tree of objects (`Views`) with internal mutable state that you manually update via setters (`setText()`, `addView()`, etc.). The real problem with that world is not just verbosity, it's brittleness. The more nodes you must update, the easier it is to forget one, update at the wrong time, or update a view that no longer exists.

```
           Imperative (layout.xml)                              Declarative (formScreen.kt)

<?xml version="1.0" encoding="utf-8"?>              @Composable
<androidx.constraintlayout.widget.ConstraintLayout   fun SimpleFormScreen() {
    xmlns:android="http://schemas.android.com/apk/res/android"
    xmlns:app="http://schemas.android.com/apk/res-auto"    var text by rememberSaveable { mutableStateOf("") }
    android:layout_width="match_parent"                val btnNextEnabled = text.isNotBlank()
    android:layout_height="match_parent"
    android:padding="24dp">                            Column(
    <LinearLayout                                          modifier = Modifier
        android:id="@+id/content"                              .fillMaxSize()
        android:layout_width="0dp"                             .padding(24.dp),
        android:layout_height="wrap_content"               verticalArrangement = Arrangement.Center
        android:orientation="vertical"                 ) {
        app:layout_constraintTop_toTopOf="parent"
        app:layout_constraintBottom_toBottomOf="parent"    OutlinedTextField(
        app:layout_constraintStart_toStartOf="parent"          value = text,
        app:layout_constraintEnd_toEndOf="parent">             onValueChange = { text = it },
                                                               modifier = Modifier.fillMaxWidth(),
        <EditText                                              label = { Text("Your Name") },
            android:id="@+id/nameInput"                    )
            android:layout_width="match_parent"
            android:layout_height="wrap_content"           Spacer(Modifier.height(16.dp))
            android:hint="Your Name"
            android:inputType="text"                       Button(
            android:padding="16dp" />                          onClick = {  /* dispatch Action */  },
        <Button                                                enabled = btnNextEnabled,
            android:id="@+id/btnNext"                          modifier = Modifier.fillMaxWidth()
            android:layout_width="match_parent"            ) {
            android:layout_height="wrap_content"               Text("Next")
            android:text="Next"                            }
            android:layout_marginTop="16dp"
            android:enabled="false" />                 }
    </LinearLayout>                                  }

</androidx.constraintlayout.widget.ConstraintLayout>
```

Figure 4.1: Imperative vs Declarative

Compose shifts the axis of the system: you do not "mutate UI", you conceptually regenerate the UI description when state changes, and Compose applies only the required differences. That is the promise: fewer invalid UI states, fewer sync bugs, and UI as a pure reflection of data.

```kotlin
package chiarani.it

import realworld.production.BugsThatOnlyHappenOnFriday
import caffeine.levels.StillNotEnough
import systems.thinking.PatternsYouOnlyLearnInProduction



@Suppress("ThisBookMayContainExperience")
@Composable
public fun BackCover() {

    val message = """
```

**Most Android books teach APIs.**

**This is not a book about how to build Android applications.**

**It is a book about what happens after you already know how to build them.**

```kotlin
    """

    Text(
        modifier = Modifier.padding(24.dp),
        text = message,
        style = MaterialTheme.typography.bodyLarge,
    )
}
```