

O'REILLY®

# Mastering API Architecture

Defining, Connecting, and Securing Distributed  
Systems and Microservices



Early  
Release

RAW &  
UNEDITED

James Gough,  
Daniel Bryant &  
Matthew Auburn

# **Mastering API Architecture**

Defining, Connecting, and Securing Distributed Systems and Microservices

With Early Release ebooks, you get books in their earliest form—the authors' raw and unedited content as they write—so you can take advantage of these technologies long before the official release of these titles.

**James Gough, Daniel Bryant, and Matthew Auburn**

# **Mastering API Architecture**

by James Gough, Daniel Bryant, and Matthew Auburn

Copyright © 2021 James Gough Ltd, Big Picture Tech Ltd, and Matthew Auburn Ltd. All rights reserved.

Printed in the United States of America.

Published by O'Reilly Media, Inc., 1005 Gravenstein Highway North, Sebastopol, CA 95472.

O'Reilly books may be purchased for educational, business, or sales promotional use. Online editions are also available for most titles (<http://oreilly.com>). For more information, contact our corporate/institutional sales department: 800-998-9938 or [corporate@oreilly.com](mailto:corporate@oreilly.com).

Acquisitions Editor: Melissa Duffield

Development Editor: Nicole Tache

Production Editor: Kate Galloway

Interior Designer: David Futato

Cover Designer: Karen Montgomery

Illustrator: Kate Dullea

November 2022: First Edition

## **Revision History for the Early Release**

- 2021-03-18: First Release
- 2021-10-29: Second Release
- 2022-01-12: Third Release

See <http://oreilly.com/catalog/errata.csp?isbn=9781492090632> for release details.

The O'Reilly logo is a registered trademark of O'Reilly Media, Inc.  
*Mastering API Architecture*, the cover image, and related trade dress are trademarks of O'Reilly Media, Inc.

While the publisher and the authors have used good faith efforts to ensure that the information and instructions contained in this work are accurate, the publisher and the authors disclaim all responsibility for errors or omissions, including without limitation responsibility for damages resulting from the use of or reliance on this work. Use of the information and instructions contained in this work is at your own risk. If any code samples or other technology this work contains or describes is subject to open source licenses or the intellectual property rights of others, it is your responsibility to ensure that your use thereof complies with such licenses and/or rights.

978-1-492-09056-4

[LSI]

# Preface

---

## A NOTE FOR EARLY RELEASE READERS

With Early Release ebooks, you get books in their earliest form—the authors' raw and unedited content as they write—so you can take advantage of these technologies long before the official release of these titles.

If you have comments about how we might improve the content and/or examples in this book, or if you notice missing material within this chapter, please reach out to the author at [feedback@mastering-api.com](mailto:feedback@mastering-api.com).

## Why Did We Write This book?

In early 2020 we attended O'Reilly Software Architecture in New York, where Jim and Matt gave a workshop on APIs and a presentation on API Gateways. Jim and Daniel know each other from the London Java Community, and like at many architecture events, had got together to talk about our thoughts and understanding around API architectures. As we were talking on the hallway track, several conference delegates came up to us and chatted about their experiences with APIs. People were asking for our thoughts and guidance on their API journey. It was at this point that we thought writing a book on the topic of APIs would help share our discussions from conferences with other architects.

## Why Should You Read This Book?

This book has been designed to provide a complete picture on designing, building, operating and evolving an API Architecture. We will share

through both writing and case study key focus areas for consideration for getting the best architectural results building an API architecture.

We also believe in allowing you to form your own decisions, to support this we will:

- Be clear when we have a strong recommendation or guidance
- Highlight areas of caution and problems that you may encounter
- Where the answer is “it depends” we will supply an ADR Guideline, to help inform the best possible decision given the circumstances of your architecture and give guidance on what to consider
- Highlight references and useful articles where you can find out more in-depth content

The book is not just a greenfield technology book. We felt that covering existing architectures with an evolutionary approach towards more suitable API architectures would have the most use for you. We have also tried to balance this with looking forward to newer technologies and developments in the API architecture domain.

## **The Core Personas For This Book**

### **Developer**

You have most likely been coding professionally for several years and have a good understanding of common software development challenges, patterns, and best practices. You are increasingly realizing that the software industry’s march towards building service-oriented architecture and adopting cloud services means that building and operating APIs is fast becoming a core skill. You are keen to learn more about designing effective APIs and testing them. You are wanting to explore the various implementation choices (e.g. synchronous versus asynchronous

communication) and technologies (e.g. REST, gRPC, messaging), and learn how to ask the right questions and evaluate which approach is best for a given context.

## **Accidental Architect**

You have most likely been developing software for many years, and have often operated as a team lead or resident software architect (even if you don't have the official titles). You understand core architectural concepts, such as designing for high cohesion and loose coupling, and apply these to all aspects of software development, including design, testing, and operating systems. You are realizing that your role is increasingly focused on combining systems to meet customer requirements, both internally built applications and also third-party COTS and SaaS-type offerings, and the APIs of these systems play a big role in the successful integration. You want to learn more about the supporting technologies (e.g. API gateway, service mesh etc) and also understand how to operate and secure API-based systems.

## **Solutions/Enterprise Architect**

You have been designing and building enterprise software systems for several years and most likely have the word "architect" in your job title or role description. You are responsible for the big picture of software delivery, and typically work within the context of a large organization or a series of large inter-connected organizations. You recognize the changes that the latest iteration of service-based architectural styles are having on the design, integration, and governance of software, and you see APIs are pivotal to the success of your organization's software strategy. You are keen to learn more about evolutionary patterns and understand how the choice of API design and implementation will impact this. You also want to focus on the cross-functional "ilities" and understand how to build API-based systems that exhibit such properties as usability, maintainability, scalability, availability and security.

## What This Book is Not

We realise that APIs encompass a vast market space and we want to be clear what this book will not cover. It doesn't mean to say that we believe these topics are not important, however if we tried to cover everything we wouldn't be able to share our knowledge effectively with you.

We will cover application patterns for migration and modernization that will include taking advantage of cloud, but the book is not wholly focused on cloud. Many of you will have hybrid architectures or even have all of your systems hosted in data centres. We want to ensure that we cover the design and operational factors of API architectures that support both approaches.

The book is not tied to a specific language, but will use some lightweight examples to demonstrate approaches to building/designing APIs and their corresponding infrastructure. The book will focus more on the approach and significant code examples will be available in the accompanying GitHub repository.

The book does not favour one style of architecture over another, however we will discuss situations where architectural approaches may cause limitations to the API offering presented.

## Conventions Used in This Book

The following typographical conventions are used in this book:

### *Italic*

Indicates new terms, URLs, email addresses, filenames, and file extensions.

### *Constant width*

Used for program listings, as well as within paragraphs to refer to program elements such as variable or function names, databases, data types, environment variables, statements, and keywords.

### ***Constant width bold***

Shows commands or other text that should be typed literally by the user.

### ***Constant width italic***

Shows text that should be replaced with user-supplied values or by values determined by context.

#### **TIP**

This element signifies a tip or suggestion.

#### **NOTE**

This element signifies a general note.

#### **WARNING**

This element indicates a warning or caution.

## **Using Code Examples**

Supplemental material (code examples, exercises, etc.) will be available for download at a later date.

If you have a technical question or a problem using the code examples, please send email to [bookquestions@oreilly.com](mailto:bookquestions@oreilly.com).

This book is here to help you get your job done. In general, if example code is offered with this book, you may use it in your programs and documentation. You do not need to contact us for permission unless you're reproducing a significant portion of the code. For example, writing a program that uses several chunks of code from this book does not require

permission. Selling or distributing examples from O'Reilly books does require permission. Answering a question by citing this book and quoting example code does not require permission. Incorporating a significant amount of example code from this book into your product's documentation does require permission.

We appreciate, but generally do not require, attribution. An attribution usually includes the title, author, publisher, and ISBN. For example: "*Book Title* by Some Author (O'Reilly). Copyright 2012 Some Copyright Holder, 978-0-596-xxxx-x."

If you feel your use of code examples falls outside fair use or the permission given above, feel free to contact us at [permissions@oreilly.com](mailto:permissions@oreilly.com).

## O'Reilly Online Learning

### NOTE

For more than 40 years, *O'Reilly Media* has provided technology and business training, knowledge, and insight to help companies succeed.

Our unique network of experts and innovators share their knowledge and expertise through books, articles, conferences, and our online learning platform. O'Reilly's online learning platform gives you on-demand access to live training courses, in-depth learning paths, interactive coding environments, and a vast collection of text and video from O'Reilly and 200+ other publishers. For more information, please visit <http://oreilly.com>.

## How to Contact Us

Please address comments and questions concerning this book to the publisher:

O'Reilly Media, Inc.

1005 Gravenstein Highway North

Sebastopol, CA 95472

800-998-9938 (in the United States or Canada)

707-829-0515 (international or local)

707-829-0104 (fax)

We have a web page for this book, where we list errata, examples, and any additional information. You can access this page at  
<https://learning.oreilly.com/library/view/mastering-api-architecture/9781492090625>.

Email [bookquestions@oreilly.com](mailto:bookquestions@oreilly.com) to comment or ask technical questions about this book.

For more information about our books, courses, conferences, and news, see our website at <http://www.oreilly.com>.

Find us on Facebook: <http://facebook.com/oreilly>

Follow us on Twitter: <http://twitter.com/oreillymedia>

Watch us on YouTube: <http://www.youtube.com/oreillymedia>

# Part I. API Fundamentals

---

## A NOTE FOR EARLY RELEASE READERS

With Early Release ebooks, you get books in their earliest form—the authors' raw and unedited content as they write—so you can take advantage of these technologies long before the official release of these titles.

If you have comments about how we might improve the content and/or examples in this book, or if you notice missing material within this chapter, please reach out to the author at [feedback@mastering-api.com](mailto:feedback@mastering-api.com).

This section will introduce the architectural foundation for Mastering API Architecture. Architecture can be quite abstract and potentially difficult to set the context, to avoid this the conference system case study is introduced from the outset. We will document the current architecture using C4 diagrams and explain the plans for the conference system. Using the case study you will learn about API interactions, capturing decisions and taking our first evolutionary step towards an API based architecture.

This section wraps up by discussing the challenges of moving to an API based architecture and sets the agenda for the rest of the book.

**Chapter 1**

# Chapter 1. API Architecture Primer

---

## A NOTE FOR EARLY RELEASE READERS

With Early Release ebooks, you get books in their earliest form—the authors' raw and unedited content as they write—so you can take advantage of these technologies long before the official release of these titles.

This will be Chapter 1 of the final book.

If you have comments about how we might improve the content and/or examples in this book, or if you notice missing material within this chapter, please reach out to the author at [feedback@mastering-api.com](mailto:feedback@mastering-api.com).

Anyone who has taken a long journey will no doubt have encountered the question (and possibly persistently) “*are we there yet?*”. For the first few inquiries you look at the GPS or a route planner and provide an estimate - hoping that you don’t encounter any delays along the way. Similarly, the journey to building API based architectures can be complex for developers and architects to navigate, even if there was an *Architecture GPS* what would your destination be? Architecture is a journey without a destination, and you cannot predict how technologies and architectural approaches will change. For example, you may not have been able to predict Service Mesh technology would become so widely used, but once you learn about its capabilities it may cause you to think about evolving your existing architecture. The culminating effect of delivering incremental value combined with new emerging technologies leads to the concept of Evolutionary Architecture. Along the way, we ask you to keep the following advice in approaching API architecture in mind:

Though architects like to be able to strategically plan for the future, the constantly changing software development ecosystem makes that difficult. Since we can't avoid change, we need to exploit it.

—Building Evolutionary Architectures - Neal Ford,  
Rebecca Parsons and Patrick Kua

In many projects APIs themselves may be evolutionary, at the time of connecting software services or systems specific functionality was more important than the definition of the API. The authors have built services focussing on a single function and not considering the broader API reuse. API-First design is an approach where developers and architects consider the functionality of their service and design an API suitable for any consumer. The API consumer could be a mobile application, another service or even an external customer. In [Chapter 2](#) we will analyze the API-First approach and discover how we build APIs that are durable to change and delivery value to a broad consumer base.

The good news is that you can start an API-First architecture journey at any point, and as authors we have experienced many points along the journey. If you are responsible for pre-existing technical inventory, we will show you techniques to evolve your architecture to promote the use of APIs in your platform. On the other hand, if you are lucky and have a blank canvas to work with, we will share with you the benefit of adopting API architectures based on our years of experience, while also highlighting key factors in decision making. All the authors actively work in the API multiverse and our journey is still underway, we hope to share our perspective to help you with Mastering API Architecture.

## API Architecture Case Study

Throughout the book, to demonstrate approaches to building API based architectures, we will focus on a case study based on a conference system. [Figure 1-1](#) visualizes the conference system at a high level. The system is used by an external customer to create their attendee account, review the conference sessions available and book their attendance. The intention of

this diagram is to set context for both a technical and non technical audience. Many architecture conversations dive straight into the low level detail and miss setting the context of the high-level interactions. Consider the implications of getting a system context diagram wrong, the benefit of summarizing the approach may save months of work to correct a misunderstanding.

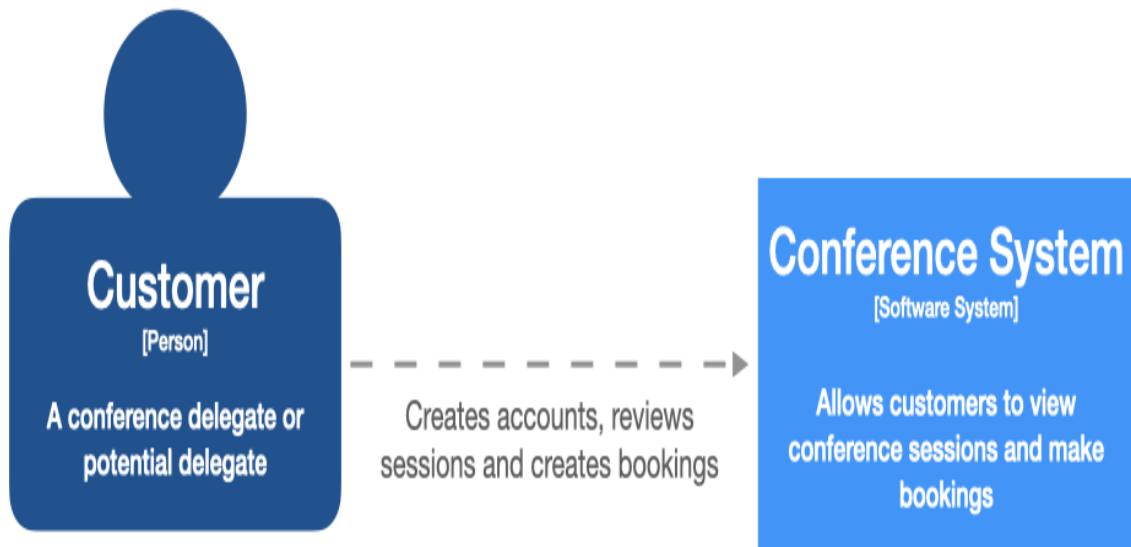
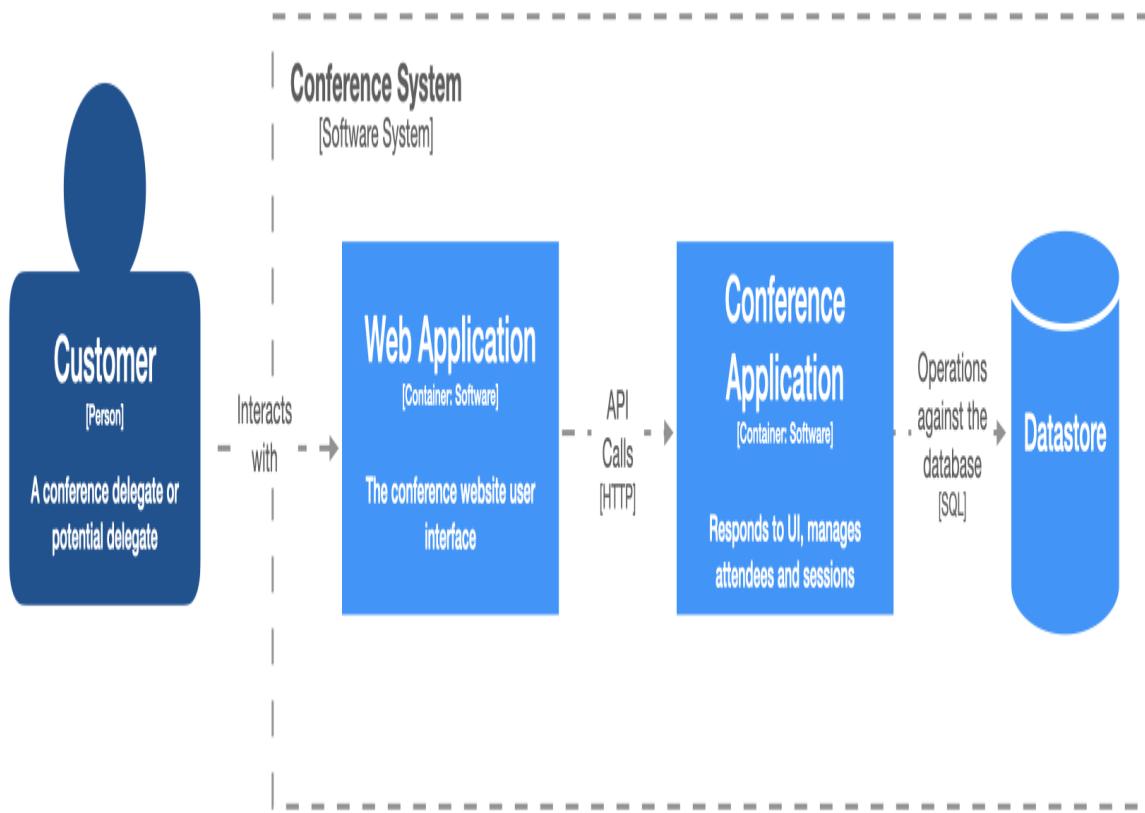


Figure 1-1. C4 Conference System Context Diagram

We have chosen C4 diagrams because we believe C4 is the best documentation standard for communicating architecture, context and interactions to a diverse set of stakeholders. You may be wondering what about UML? The Unified Modeling Language (UML) provides an extensive dialect for communicating software architectures. A major challenge is that the majority of what UML provides is not committed to memory of architects and developers, people quickly revert to boxes/circles/diamonds. It becomes a real challenge to understand the structure of diagrams before getting onto the technical content of the discussion. Many diagrams are only committed to a project history if someone accidentally uses a permanent marker instead of dry wipe marker by mistake. The C4 Model provides a simplified set of diagrams that act as a guide to your project architecture at various levels of detail.

While [Figure 1-1](#) provides the big picture of the conference system, a container diagram helps describe the technical breakout of the major participants in the architecture. A container in C4 is defined as “*something that needs to be running in order for the overall system to work*”, for example the conference database. Container diagrams are technical in nature and build on the higher level system context diagram. [Figure 1-2](#), a container diagram, documents the detail of a customer interacting with the conference system. The user interacts with the Web Application, which invokes APIs on the conference application. The conference application uses SQL to query the backing database.



*Figure 1-2. C4 Conference System Container Diagram*

## NOTE

The conference application container in [Figure 1-2](#) is documented as simply *software*. Normally a C4 container would provide more detail into the type of container e.g. *Java Spring Application*. However in this book we will be avoiding technology specifics, unless it helps to demonstrate a specific solution. The advantage of APIs and indeed modern applications is there is a significant amount of flexibility in the solution space.

The current architectural approach has worked for the conference system for many years, however the conference owner has asked for two improvements, which are driving architectural change:

- The conference organizers would like to build a mobile application.
- The conference organizers plan to go global with their system, running ten conferences instead of one per year. In order to facilitate this expansion they would like to integrate with an external Call for Papers (CFP) system for managing speakers and their application to present sessions at the conference.

Our goal is to migrate the conference system to be able to support the new requirements, without impacting the existing production system or rewriting everything in one go. In this chapter we will introduce the concept of traffic patterns in API Architecture and the opportunities to use this concept to evolve our architecture. As the case study evolves, we will share techniques for documenting architectures and recording decisions.

## A Brief Introduction to APIs

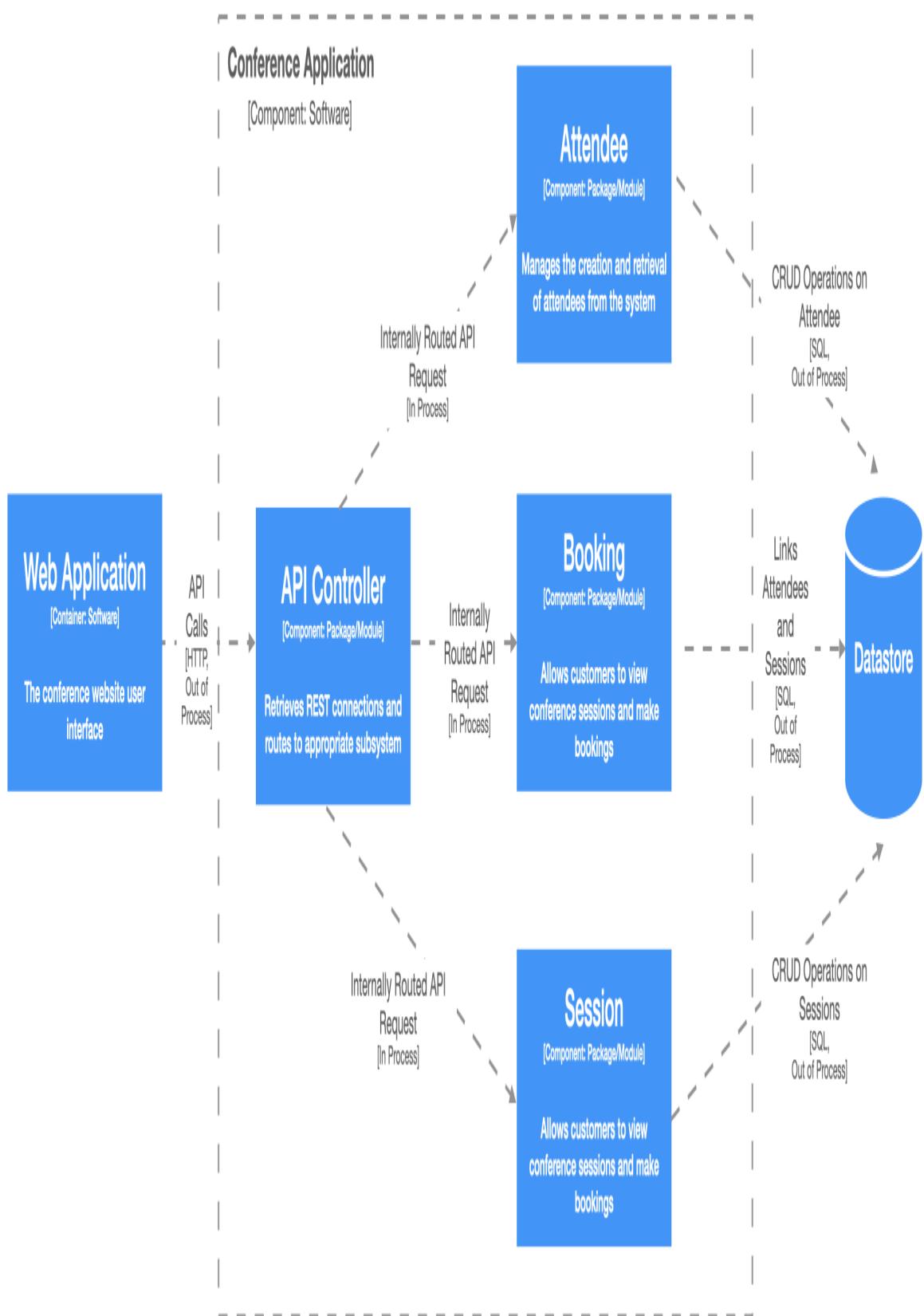
In the field of software architecture there are a handful of terms that are incredibly difficult to define. The term API (Application Programming Interface) falls into this categorization, as the concept first surfaced as many as 80 years ago. Terms that have been around for a significant amount of

time end up being overused and having multiple meanings in different problem spaces. In this book we will refer to an API to mean the following:

- An API represents an abstraction of the underlying implementation.
- An API is represented by a specification that introduces types. Developers can understand the specifications and use tooling to generate code in multiple languages to implement an API
- Consumer (software that consumes an API).
- An API has defined semantics or behavior to effectively model the exchange of information.
- Effective API design enables APIs to be extended to customers or third parties for a business integration.

Broadly speaking APIs can be broken into two general categories depending on whether the API invocation is *in process* or *out of process*. In order to explore this concept further with our C4 diagrams we need to zoom into the Conference Application Container box in [Figure 1-2](#).

The component diagram in [Figure 1-3](#) helps to define the roles and responsibilities within each container, along with the internal interactions. This diagram is useful if the detail of a container is queried, it is also provides a very useful map to the codebase. Think about the first time starting work on a new project, browsing a *self documenting* codebase is one approach - but it can be difficult to piece everything together. A component diagram reveals the detail of the language/stack you are using to build your software. In order to remain technology agnostic we have used the term *package/module*.



*Figure 1-3. C4 Conference System Component Diagram*

The *web application* to *API Controller* arrow is an out of process call, whereas the *API Controller* to *Attendee Component* arrow is an example of an in process call. All interactions within the Conference Application boundary are examples of in process calls. The in process invocation is well defined and restricted by the programming language used to implement the Conference Application. The invocation is compile time safe, the conditions under which the exchange mechanism are enforced at the time of writing code.

An *out of process* exchange is significantly more complicated, the external process can be written in any language and potentially located anywhere. There needs to be an agreement on the protocol, expected behavior and compatibility. Out of process API modelling will be the focus of [Chapter 2](#).

[Figure 1-3](#) describes 4 major components and the database involved in the current system. The **API Controller** faces all incoming traffic from the UI and make a decision about where to route the request in the system. This component would also be responsible for marshalling from the *on the wire* network level representation to an object or representation in code. The API Controller component is intriguing from the perspective of *in process* routing and acting as a junction point or *front controller* pattern. For API requests and processing this is an important pattern, all requests pass through the controller which makes a decision on where the request is directed. In [Chapter 4](#) we will look at the potential for taking the controller out of process.

The **Attendee**, **Booking** and **Session** packages are involved in translating the requests into queries and execute SQL against the database out of process. In the existing architecture the database is an important component, potentially enforcing relationships for example constraints between bookings and sessions.

## **From Tiered Architecture to Modelling APIs**

The starting point of the case study is a typical 3-tier architecture, composed of a UI, a server-side processing tier and a datastore. To begin to discuss an evolutionary architecture we need a model to think about the way API requests are processed by the components. We need a model/abstraction that will work for both public cloud, virtual machines in a data center and a hybrid approach.

The abstraction of traffic will allow us to consider *out of process* interactions between an API consumer and an API service, sometimes referred to as the API producer. With architectural approaches like Service Oriented Architecture (SOA) and Microservices based Architecture the importance of modelling API interactions is critical. Learning about API traffic and the style of communication between components will be the difference between realizing the advantages of increased decoupling or creating a maintenance nightmare.

### WARNING

Traffic patterns are used by data center engineers to describe network exchanges within data centers and between low level applications. At the API level we are using traffic patterns to describe flows between groups of applications. For the purposes of this book, we are referring to application and API level traffic patterns.

## An Evolutionary Step

In order to start to consider traffic pattern types it will be useful to take a small evolutionary step in our case study architecture. In [Figure 1-4](#) a step has been taken to refactor the **Attendee** component into an independent service, as opposed to a package or module within the legacy conference system. The conference system now has two traffic flows, the interaction between the customer and the legacy conference system and the interaction between the legacy system and the attendee system.

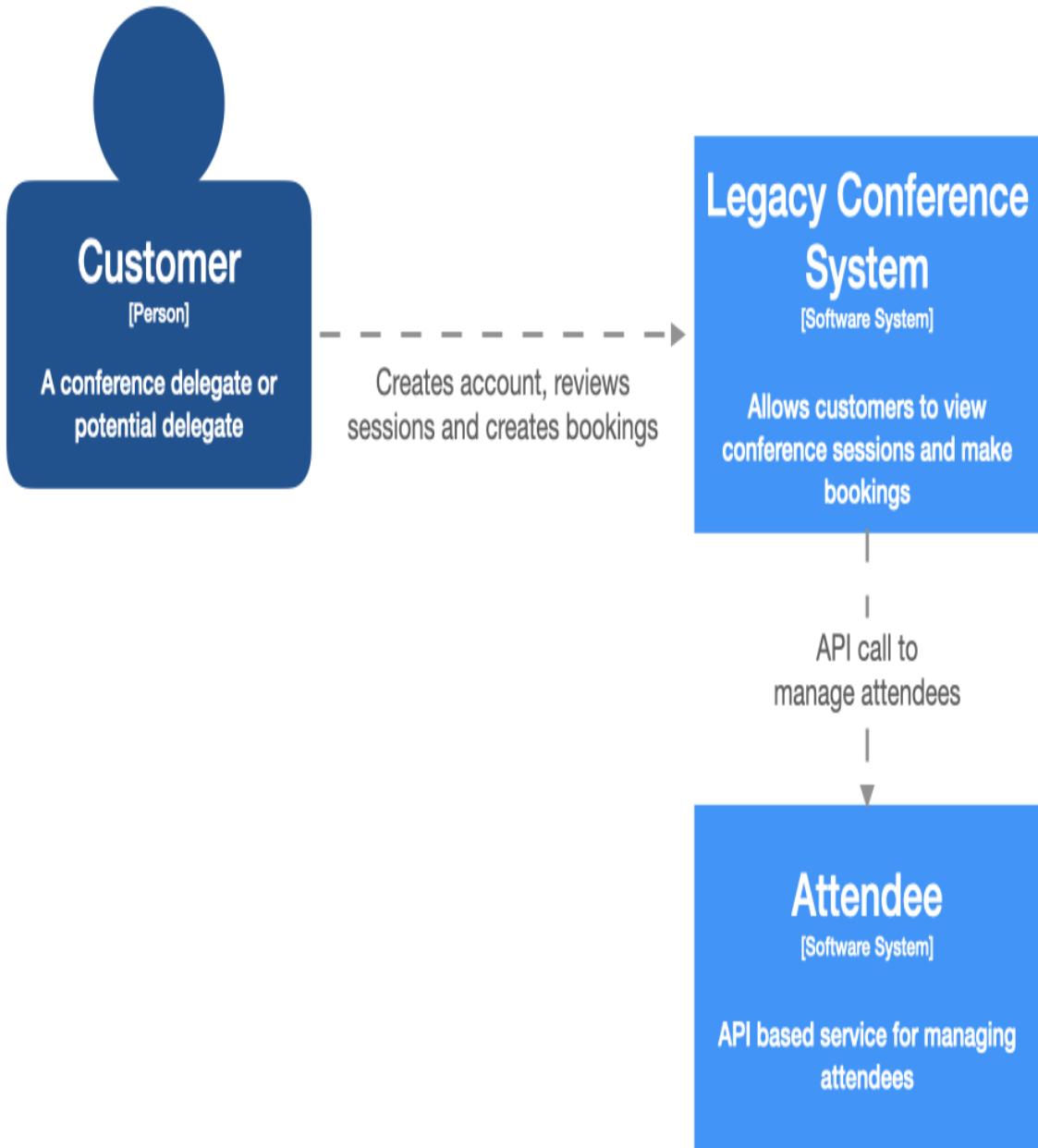


Figure 1-4. C4 Conference System Context - Evolutionary Step

## North - South Traffic

The interaction between the customer and the legacy conference system is referred to a north → south traffic, and represents an ingress flow. The customer is using the UI which is sending requests to the legacy conference system over the internet. This represents a point in our network that is exposed publicly, and will be accessed by the UI<sup>1</sup>. This means that any

component handling north → south traffic must make concrete checks about client identity and also include appropriate challenges before allowing traffic to progress through the system. *Chapter 7* will go into detail about securing north → south API traffic.

## East - West Traffic

The new interaction between the legacy conference system and the **Attendee** service introduces an east → west traffic flow to our system. East → west traffic can be thought of as *service-to-service* style of communication within a group of applications. Most east → west traffic, particularly if the origin is within your wider infrastructure, can be trusted to some degree. Although we can trust the source of the traffic, it is still necessary to consider securing east → west traffic.

## API Infrastructure and Traffic Patterns

There are two key infrastructure components present in API based architectures, which are key to controlling traffic. Controlling and coordinating traffic is often described as Traffic Management. Generally north → south traffic will be controlled by API Gateways, the key subject for [Chapter 4](#).

East → west traffic will often be handled by infrastructure components like Kubernetes or Service Mesh, the key subject for *Chapter 5*. Infrastructure components like Kubernetes and Service Mesh use network abstractions to route to services, requiring services to run inside a managed environment. In some systems east → west traffic is managed by the application itself and service discovery techniques are implemented to locate other systems.

## Using Architecture Decision Records

As developers, architects and indeed humans we have all been in the position where we ask the question “*what were they thinking??*”. If you have ever driven on the M62 between Leeds and Manchester in the United

Kingdom you may have been baffled by the construction of the motorway. As you climb the hill on the 3 lane highway it starts to deviate away from the traffic contraflow, until eventually Scott Hall Farm emerges surrounded by around 15 acres of farming land nestled between the carriages. Local legend of what happened described the owner of the land as stubborn and refused to move or handover his land, so the engineers simply built around him <sup>2</sup> 50 years later a documentary surfaced revealing that the real reason for this was a geological fault beneath the land, meaning the motorway had to be built that way. When people guess why something was done in particular way expect rumour, humour and criticism to emerge.

In software architecture there will be many constraints that we have to build around, so it is important to ensure our decisions are recorded and transparent. Architecture Decision Records (ADRs) help make decisions clear in software architecture.

An Architecture Decision Record (ADR) is a document that captures a decision, including the context of how the decision was made and the consequences of adopting the decision.

—Engineering at Spotify Blog

An ADR is created in a proposed state and based on discussion will usually be either accepted or rejected, it is also possible that later the decision may be superseded by a new ADR. The context helps to set the scene and describe the problem or the bounds in which the decision will be made. The context is not intended to be a blog post or detailed description, often creating a blog post ahead of the ADR and then linking from the ADR helps the community to follow your working. The decision clearly sets out what you plan to do and how you plan to do it. All decisions carry consequences or trade offs, in architecture these can sometimes be incredibly costly to get wrong.

When reviewing an ADR it is important to see if you agree with the decision in the ADR or if there is an alternative approach. An alternative approach that has not been considered may cause the ADR to be rejected. There is a lot of value in a rejected ADR and most teams choose to keep

ADRs immutable to capture the change in perspective. ADRs work best when they are presented in a location where key participants can view the ADR, comment and help move the ADR to accepted.

### TIP

A question we often get asked is at what point should the team create an ADR? It is useful to ensure that there has been discussion ahead of the ADR and the record is a result of collective thinking in the team.

## Attendees Evolution ADR

In [Figure 1-4](#) we made the decision to take an evolutionary step in the conference system architecture. This is a major change and would warrant an ADR. Below is an example ADR that might have been proposed by the engineering team owning the conference system.

*T*

*a*

*b*

*l*

*e*

*l*

-  
*l*

.

*A*

*D*

*R*

*0*

*0*

*l*

*S*

*e*

*p*

*a*

*r*

*a*

*t*

*i*

*n*

*g*

*A*

*t*

*t*

*e*

*n*

*d*

*e*

*e*

*s*

*f  
r  
o  
m*

*t  
h  
e  
L  
e  
g  
a  
c  
y  
C  
o  
n  
f  
e  
r  
e  
n  
c  
e  
S  
y  
s  
t  
e  
m*

Status	Proposed
Context	The conference owners have requested two new major features to the current conference system, that need to be implemented without disrupting the current system.

	<p>The conference system will need to be evolved to support a mobile application and an integration with an external CFP system.</p> <p>Both the mobile application and the external CFP system need to be able to access attendees to log in users to the third party service.</p>
Decision	<p>We will take an evolutionary step as documented in <a href="#">Figure 1-4</a> to split out the <b>Attendee</b> component into a standalone service.</p> <p>This will allow API-First development against the <b>Attendee</b> service and allow the API to be invoked from the legacy conference service.</p> <p>This will also support the ability to design for direct access to the <b>Attendee</b> service to provide user information to the external CFP system.</p>
Consequences	<p>The call to the <b>Attendee</b> service will not be <i>out of process</i> and may introduce a latency that will need to be tested.</p> <p>The <b>Attendee</b> service could become a single point of failure in the architecture and we may need to take steps to mitigate the potential impact of running a single <b>Attendee</b> service.</p>

Some of the consequences in the ADR are fairly major, and definitely require further discussion. We are going to defer some of the consequences to later chapters.

## Mastering API - ADR Guidelines

Within Mastering API Architecture we will be supplying *ADR Guidelines* to help collect important questions to ask when decisions on the topic covered. Making decisions about an API based Architecture can be really tough, and in a lot of situations the answer is *it depends*. Rather than say it depends without context, the ADR Guidelines will help describe what it depends on and help inform your decisions. The ADR Guidelines can be used as a reference point to come back to or to read ahead to if you're facing a specific challenge.

*T  
a  
b  
l  
e  
I  
-  
2  
.A  
D  
R*

*G  
u  
i  
d  
e  
li  
n  
e  
s  
-  
F  
o  
r  
m  
a  
t*

Decision      Describes a decision that you might need to make when considering an aspect of this book.

Discussion Points      This section helps to identify the key discussions that you should be having

when making a decision about your API Architecture.

In this section we will reveal some of the authors experiences that may have influenced the decision. We will help you to identify the key information to inform your decision making process.

Recommendations	We will make specific recommendations that you should consider when creating your ADR, explaining the rationale behind why we are making a specific recommendation.
-----------------	---

## Roadmap for our Journey

In “An Evolutionary Step” we started our journey to evolve the conference system as an API based architecture, but how should we design our API? What are our different options around modelling APIs and what are the benefits and drawbacks of each? Designing an API is an important undertaking, particularly if it will be used by an external customer. In [Chapter 2](#) you will discover how to design your API, the different styles of API and the considerations you must make with each.

One of the benefits of an API based architecture, which we revealed in this chapter is the abstraction of implementation detail, which leads to being able to rapidly change your application. However APIs also introduce the possibility of impacting vast numbers of consumers if an error or compatibility issue occurs. Out of process API interactions are a lot more susceptible to runtime errors than in process exchanges within a compiled application. In order to ensure that an API responds in the expected manner it is essential that appropriate testing is applied on the API. In [Chapter 3](#) you will discover the different approaches to testing APIs within the context of the test pyramid.

Once we have more APIs in our system we can look to at opening these up to the outside world. In this chapter we introduced the model of north → south traffic, which usually involves entering the system from the internet. How do you go about creating that entry point and setting up a pattern that will work for routing to multiple APIs? In [Chapter 4](#) we will introduce API

Gateways and demonstrate the importance of getting started with an API based architecture.

In this chapter we also introduced the model of east → west traffic, or services communicating laterally within a group of services. The pattern can introduce challenges with locating services, scaling up services and even deployments if configured manually. Service Mesh represents a best in class solution to this type of problem and is the focus of *Chapter 5*.

Once we have patterns and approaches for building APIs and managing traffic the next step is to consider how we operationally manage API based architectures in production. In *Chapter 6* you will learn about some of the approaches to deploying and releasing APIs and how we can potentially use traffic management concepts to reduce the risk of releasing into production.

The security conscious may now be in wild panic about services that we can rapidly release as APIs to the internet! An API based architecture with an internet facing component must consider security from the outset in the design and architecture of the system. In *Chapter 7* we will demonstrate approaches that help to interrogate API Architectures and ensure that the overall architecture has security considerations built in from the outset.

*Chapter 8* takes security a step further, introducing the worst case scenario - the architecture is under attack. How do you find out that the system is under attack and what mitigation can be put in place to prevent the entire system falling over? We will also explore how you can attack and test your own system to ensure that the architecture holds up to expected volumes.

In the final section we start to explore evolutionary architectures with APIs. In *Chapter 9* we will cover patterns and approaches to redesigning monolithic application to use an API based series of services. We will discuss how you structure those services and consider a service oriented approach and a microservices based approach.

*Chapter 10* will build on *Chapter 9* and discuss how the content of the book can be applied to help migrate to a hybrid-cloud or even a pure cloud deployment for an API ecosystem. We will talk about some of the

approaches to cloud migration and how this does not need to be an all or nothing approach.

## Summary

In this chapter we have introduced a variety of key concepts that are an important foundation to building API based architectures that will be revisited throughout the book:

- Architecture is an endless journey and the conference based case study will demonstrate this throughout the book.
- C4 Architecture diagrams allows us to explore the conference system at different levels.
- Briefly introduced APIs and considered the approach of modelling API based systems using traffic patterns.
- Introduced an evolutionary step to the conference system architecture and documented this in C4 with a sample ADR. We also introduced the concept of ADR Guidelines that will guide decisions that are context dependant.
- Looked at the challenges of API Architecture and how the book is structured to help you work through the challenges.

---

<sup>1</sup> The intention is it will be the UI accessing the ingress point. However it is open for potential exploit.

<sup>2</sup> Local stubborn traits fueled this likely explanation.

# Part II. Designing, Building and Testing APIs

---

## A NOTE FOR EARLY RELEASE READERS

With Early Release ebooks, you get books in their earliest form—the authors' raw and unedited content as they write—so you can take advantage of these technologies long before the official release of these titles.

If you have comments about how we might improve the content and/or examples in this book, or if you notice missing material within this chapter, please reach out to the author at [feedback@mastering-api.com](mailto:feedback@mastering-api.com).

This section provides the foundational building blocks for API Ecosystems. In general, an API Ecosystem will have API services and infrastructure-based routing, this section will focus on design and testing of APIs before we consider how architectural patterns are applied.

In the first chapter you will learn the foundation of REST APIs and why without following API guidelines things can go wrong quite quickly for your customers. You will learn the practical application for OpenAPI Specifications and why they have rapidly become an industry standard. You will learn about why REST is a low barrier to entry for an API ecosystem and we will review the best way to revise and evolve REST APIs. Not all API exchanges will be best modelled with REST and additional constraints might require alternative considerations. You will learn about gRPC and the pros and cons vs a REST approach. You will also look at GraphQL and consider where this might fit into an API ecosystem. Finally, we will

consider if there is a way to mix and match our approach to a service modelling different styles of API.

The first chapter presents a variety of choices for modelling APIs, and we need to review the practical applications to start to break out parts of our conference monolith. You will look at the attendees service and consider the implications of modelling API exchanges in the conference ecosystem. The second chapter leads straight into learning practices to test the APIs that we have seen in the first chapter. To understand how to get the most value out of testing APIs you will learn what the testing quadrant and the test pyramid is. From these two testing strategies you will be able to devise how they can apply this to your own environment and understand why they provide value. A practical introduction to Unit tests, Service tests and UI tests will be given next. By learning about the fundamentals of testing you will be shown the value that testing provides and how it need not be difficult to start and introduce into your architecture. The Conference System will be used to demonstrate this.

One of the toughest parts of testing is to validate that two services can communicate with one another before actually being deployed. With this in mind, we go into detail explaining different strategies to solve this problem, however, the authors think that one of the best ways to solve this is using Contract testing. You will come away armed with knowledge about what tests should be applied and at what level.

[Chapter 2](#)  
[Chapter 3](#)

# Chapter 2. Defining Rest APIs and Alternative API Approaches

---

## A NOTE FOR EARLY RELEASE READERS

With Early Release ebooks, you get books in their earliest form—the authors' raw and unedited content as they write—so you can take advantage of these technologies long before the official release of these titles.

This will be Chapter 2 of the final book.

If you have comments about how we might improve the content and/or examples in this book, or if you notice missing material within this chapter, please reach out to the author at [feedback@mastering-api.com](mailto:feedback@mastering-api.com).

## WARNING

Dear Reader,

Thank you for taking the time to read and review our early access content. Over the past few months we have received some excellent feedback on the structure of the book. As a result, we will be making a few changes to this chapter, including updating the content and adding a case study. Once the update is complete, we will remove this warning.—  
James, Daniel and Matthew

Microservices based architectures and Service Oriented Architectures promote an increased number of independent services. Often services are running multi-process, on multiple machines in data centers around the globe. This has led to an explosion in mechanisms to communicate between processes and deal with the challenges of distributed communication between services. The software development community has responded by

building a variety of creative API protocols and exchanges including REST, gRPC and GraphQL learning from vintage protocols. With a range of options available an API architect needs an overall understanding of the exchange styles available and how to pick the right technologies for their domain.

In this chapter we will explore how REST addresses the challenges of communicating across services, techniques for specifying REST APIs using OpenAPI and the practical applications of OpenAPI specifications. With modern tooling building REST APIs is easy, however building a REST API that is intuitive and follows recommended practices requires careful design and consideration. In order to be successful with microservices, a focus on continuous integration, deployment speed and safety is essential. REST APIs need to complement DevOps and rapid deployments rather than become a bottleneck or deployment constraint. In order to keep our system free to evolve and not disrupt consumers, it is also necessary to consider versioning in the release process of APIs.

Due to the simplicity of REST and wide support it is usually one of the best approaches to create an API. However there will be situations and constraints where we must consider alternative approaches for our API implementation. This may come down to performance or the requirement to create a query layer across multiple APIs.

The chapter concludes by looking at whether we can combine multiple API formats in a single service. We will also discuss whether the complexity of multi specifications takes us further away from the fundamentals of REST and other API approaches.

## Introduction to REST

Roy Fielding's dissertation [Architectural Styles and the Design of Network-based Software Architectures](#) provides an excellent definition of the architectural applications of REST. REpresentational State Transfer (REST) is one of the primary mechanisms for exchanging data between client and server and is extremely popular within the API ecosystem and beyond. If

we adjust “*distributed hypermedia system*” in Roy Fielding’s definition below to “*an API*” it describes the benefit that we get from using REST APIs.

The Representational State Transfer (REST) style is an abstraction of the architectural elements within a distributed hypermedia system. REST ignores the details of component implementation and protocol syntax in order to focus on the roles of components, the constraints upon their interaction with other components, and their interpretation of significant data elements. It encompasses the fundamental constraints upon components, connectors, and data that define the basis of the Web architecture, and thus the essence of its behavior as a network-based application.

—Roy Thomas Fielding

To explore REST in more practical terms [Figure 2-1](#) describes a typical REST exchange and the key components involved in the exchange over HTTP. In this example the consumer requests information about the **attendees resource**, which is known by the **resource identifier** <http://mastering-api.com/attendees>. Modelling an exchange in REST involves a request method (also known as the request verb), in this case a GET to retrieve data, and the resource identifier to describe the target of our operation. REST defines the notion of a **representation** in the body and allows for **representation metadata** to be defined in the headers. In this example we are informing the server we are expecting *application/json* by stating this in the **Accept** header.

The response includes the status code and message from the server, which enables the consumer to interrogate the result of the operation on the server-side resource. In the response body a JSON representation containing the conference attendees is returned. In [“REST API Standards and Structure”](#) we will explore approaches to modelling the JSON body for API compatibility.

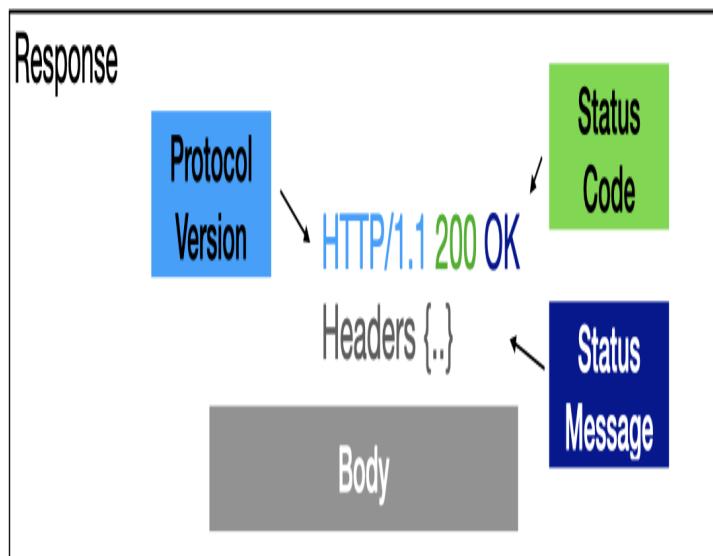
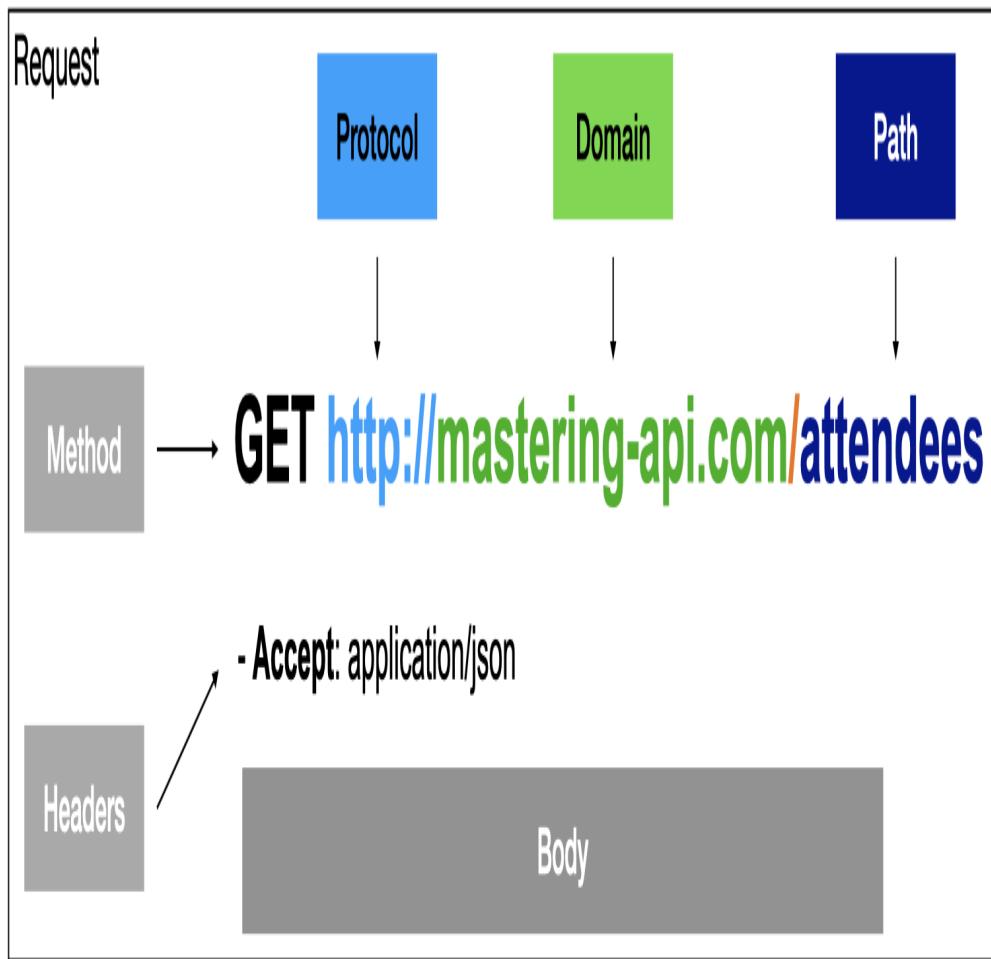


Figure 2-1. Anatomy of a RESTful Request and Response over HTTP

## NOTE

Although there is nothing in the REST specification that states HTTP is required, the data elements and their usage were designed with HTTP in mind.

## The Richardson Maturity Model

The REST specification does not enforce a high degree of restriction in terms of how developers and architects build and design REST APIs. Speaking at [QCon in 2008](#) Leonard Richardson presented his experiences of reviewing many REST APIs. Martin Fowler also covered Leonard Richardson's maturity model on [his blog](#). The model presents levels of adoption that teams apply to building APIs.

*T*  
*a*  
*b*  
*l*  
*e*

*2*  
-  
*l*  
. *R*  
*i*  
*c*  
*h*  
*a*  
*r*  
*d*  
*s*  
*o*  
*n*

*M*  
*a*  
*t*  
*u*  
*r*  
*i*  
*t*  
*y*

*M*  
*o*  
*d*  
*e*

*l  
L  
e  
v  
e  
l  
s*

Level 0 - HTTP/RPC	<p>Establishes that the API is built using HTTP and has the notion of one URI and one HTTP method.</p> <p>Taking our example above of <code>GET /attendees</code>, if we stopped there the example would match Level 0.</p> <p>Essentially this represents an RPC (Remote Procedure Call) implementation, a consideration we will compare and contrast throughout this chapter.</p>
Level 1 - Resources	<p>Establishes the use of resources, and starts to bring in the idea of modelling resources in the context of the URI.</p> <p>In our example if we added <code>GET /attendees/1</code> returning a specific attendee it would start look like a level 1 API.</p> <p>Martin Fowler draws an analogy to the classic object oriented world of introducing identity.</p>
Level 2 - Verbs (Methods)	<p>Starts to introduce the correct modelling of multiple resources URIs accessed by different request methods (also know as HTTP Verbs) based on the effect of the resources on the server.</p> <p>An API at level 2 can make guarantees around <code>GET</code> methods not impacting server state and presenting multiple operations on the same resource URI.</p> <p>In our example adding <code>DELETE /attendees/1</code>, <code>PUT /attendees/1</code> would start to add the notion of a level 2 compliant API.</p>
Level 3 - Hypermedia Controls	<p>This is the epitome of REST design and involves navigable APIs by the use of <b>HATEOAS (Hypertext As The Engine Of Application State)</b>.</p> <p>In our example when we call <code>GET /attendees/1</code> the response would contain the actions that are possible on the object returned from the server.</p> <p>This would include the option to be able to update the attendee or delete the attendee and what the client is required to invoke in order to do so.</p>

Most APIs reach API maturity level 2, but with this wide range of options there is a lot for the developer to consider. The practical applications of

HATEOAS are limited in modern RESTful APIs and this is where the theoretical ideal REST concept meets the battle testing of production code. In order to be consistent we need to establish baseline expectations of the APIs implemented using REST and provide that consistency across all the APIs offered.

## REST API Standards and Structure

REST has some very basic rules, but for the most part the implementation and design of an API is left as an exercise for the developer. It is useful to have a more practical definition around the APIs to provide a uniformity and expectations across different APIs. This is where REST API Standards or Guidelines can help, however there are a variety of sources to choose from. For the purposes of discussing REST API design we will use the [Microsoft REST API Guidelines](#), which represent a series of internal guidelines that have been OpenSourced to create a dialog in the community. The guidelines use RFC-2119 which defines terminology for standards such as MUST, SHOULD, SHOULD NOT, MUST NOT etc allowing the developer to determine whether requirements are optional or mandatory.

### TIP

As REST API Standards are evolving, an open list of API Standards are available on the [book's Github](#). Please contribute via Pull Request any open standards you think would be useful for other readers to consider.

Lets evolve our *attendees* API using the Microsoft REST API Guidelines and introduce an endpoint to create a new `attendee`. If you are familiar with REST the thought will immediately be to use `POST`, but the recommended response header might not be so obvious.

```
POST http://mastering-api.com/attendees
{
  "displayName": "Jim",
  "givenName": "James",
```

```
        "surname": "Gough",
        "email": "jim@mastering-api.com"
    }
---
201 CREATED
Location: http://mastering-api.com/attendees/1
```

The *Location* header reveals the location of the new resource created on the server, and in this API we are modelling a unique ID for the user. It is possible to use the *email* field as a unique ID, however the Microsoft REST API guidelines recommend in [section 7.9](#) that PII should **not** be part of the URL.

### WARNING

The reason for removing sensitive data from the URL is paths or query parameters might be inadvertently cached in the network, for example in server logs or elsewhere.

Another aspect of APIs that can be difficult is naming, as we will discuss in [“API Versioning”](#) something as simple as changing a name can break compatibility. There is a short list of standard names that should be used in the Microsoft Rest API Guidelines, however teams should expand this to have a common domain data dictionary to supplement the standards. Lets now take a look at patterns for retrieving data.

## Collections and Pagination

It seems reasonable to model the GET */attendees* request as a response containing a raw array. The source snippet below shows an example of what that might look like as a response body.

```
GET http://mastering-api.com/attendees
---
200 OK
[
    {
        "displayName": "Jim",
```

```

        "givenName": "James",
        "surname": "Gough",
        "email": "jim@mastering-api.com",
        "id": 1,
    },
{
    "displayName": "Matt",
    "givenName": "Matthew",
    "surname": "Auburn",
    "email": "matt@mastering-api.com",
    "id": 2,
}
]

```

Let's consider at an alternative model to the `GET /attendees` that nests the array of attendees inside an object. It may seem strange that an array response is returned in an object, however the reason for this is that allows for us to model bigger collections and pagination. This is reaping the benefits of hindsight, adding pagination later and converting from an array to a object in order to add a `@nextLink` (as recommended by the standards) would break compatibility.

```

GET http://mastering-api.com/attendees
---
200 OK
{
    "value": [
        {
            "displayName": "Jim",
            "givenName": "James",
            "surname": "Gough",
            "email": "jim@mastering-api.com",
            "id": 1,
        }
    ],
    "@nextLink": "{opaqueUrl}"
}

```

## Filtering Collections

Our conference is looking a little lonely with only two attendees, however when collections grow in size we may need to add filtering in addition to

pagination. The **filtering** standard provides an expression language within REST to standardize how filter queries should behave, based upon the **OData Standard**. For example we could find all attendees with the `displayName` Jim using

```
GET http://mastering-api.com/attendees?$filter=displayName eq 'Jim'
```

It is not necessary to complete all filtering and searching features from the start. Designing an API in line with the standards will allow the developer to support an evolving API architecture without breaking compatibility for consumers.

## Updating Data

When designing an API the developer would need to make an active decision on whether to use **PUT** or **PATCH** to update an attendees details. A **PUT** is used to replace the resource entirely with the content of the request, whereas a **PATCH** would only update the attributes specified by the request. Depending on the number of fields and expected parallel updates the API should be designed accordingly. For example two **PUT** operations would likely involve a lost update whereas two **PATCH** requests may be successful independently.

## Error Handling

An important consideration when extending APIs to consumers is defining what should happen in various error scenarios. **Error standards** are useful to define upfront and share with API producers to provide consistency. It is important Errors describe to the API Consumer exactly what has gone wrong with the request, this will avoid increased support of the API.

## **WARNING**

Ensure that the error messages sent back to the consumer does not contain stack traces and other sensitive information. This information can help an hacker aiming to compromise the system.

We have just scratched the surface on building REST APIs, but clearly there are important decisions to be made at the beginning of the project to build an API. If we combine the desire to present intuitive APIs that are consistent and allow for an evolving compatible API, it is worth adopting an API Standard early.

## **Checklist: Choosing an API Standard**

*T*  
*a*  
*b*  
*l*  
*e*

*2*  
-  
*2*  
. *A*  
*P*  
*I*

*S*  
*t*  
*a*  
*n*  
*d*  
*a*  
*r*  
*d*  
*s*

*C*  
*h*  
*e*  
*c*  
*k*  
*l*  
*i*  
*s*  
*t*

Decision	Which API standard should we adopt?
Discussion Points	<p>Does the organization already have other standards within the company? Can we extend those standards to external consumers?</p> <p>Are we using any third party APIs that we will need to expose to a consumer (e.g. Identity Services) that already have a standard?</p> <p>What does the impact of not having a standard look like for our consumers?</p>
Recommendations	<p>Pick an API standard that best matches the culture of the Organization and formats of APIs you may already have in the inventory.</p> <p>Be prepared to evolve and add to a standard any domain/industry specific amendments.</p> <p>Start with something early to avoid having to break compatibility later for consistency.</p> <p>Be critical of existing APIs, are they in a format that consumers would understand or is more effort required to offer the content?</p>

## Specifying REST APIs

As we have seen the design of an API is fundamental to the success of an API platform. The next challenge to overcome is sharing the API with developers consuming our APIs.

API marketplaces provide a public or private listing of APIs available to a consumer. The developer can browse documentation and quickly trial an API in the browser to explore the API behavior and functionality. Public and private API marketplaces have placed REST APIs prominently into the consumer space. Architecturally REST APIs are increasingly important in support of both microservices based architectures and Service Oriented Architectures. The success of REST APIs has been driven by both the technical landscape and the low barrier to entry for both the client and server.

Prior to 2015 there was no standard for specifying REST APIs, which may be quite surprising given that API specifications are not a new concept. XML had XML Schema Definitions (or XSD), which were a core

mechanism in ensuring compatibility of services. However, it is important to remember that REST was designed for the web, rather than specifically for APIs. As the number of APIs grew it quickly became necessary to have a mechanism to share the *shape* and structure of APIs with consumers. This is why the [OpenAPI Initiative](#) was formed by API industry leaders to construct the OpenAPI Specification (OAS). The OAS was formerly known as Swagger and documentation and implementation use OAS and Swagger interchangeably.

## OpenAPI Specification Structure

Let's explore an example OpenAPI Specification for the attendees API.

```
"openapi": "3.0.3",
"info": {
    "title": "Attendees Mastering API",
    "description": "Example accompanying the Mastering API Book",
    "version": "1.0"
},
"servers": [
    {
        "url": "http://mastering-api.com",
        "description": "Demo Server"
    }
],
```

The specification begins by defining the OAS version, information about the API and the servers the API is hosted on. The `info` attribute is often used for top level documentation and for specifying the version of the API. The version of the API is an important attribute, which we will discuss in more detail in [“API Versioning”](#). The `servers` array is one of the new additions in OpenAPI Specification 3, prior to this only a single host could be represented. The `openapi` object key is named `swagger` in older versions of the specification.

## NOTE

As well as defining the shape of an API the OpenAPI Specification often conveys full documentation about the API.

```
"paths": {
  "/attendees": {
    "get": {
      "tags": [
        "attendees-controller"
      ],
      "summary": "Retrieve a list of all attendees stored within the
system",
      "operationId": "getAttendees",
      "responses": {
        "200": {
          "description": "OK",
          "content": {
            "/*": {
              "schema": {
                "$ref":
                  "#/components/schemas/AttendeeResponse"
              }
            }
          }
        }
      }
    }
  },
},
```

The `paths` tag conveys the possible operations on the RESTful API and the expected request and response object formats. In this example on a 200 status response the consumer can expect to receive an object `AttendeeResponse`. The `components` object will describe the response has a key `value` containing the `Attendee` array. The `$ref` tag indicates that this will be represented elsewhere in the specification.

```
"components": {
  "schemas": {
    "Attendee": {
      "title": "Attendee",
```

```
        "required": [
            "email",
            "givenName",
            "surname",
            "displayName"
        ],
        "type": "object",
        "properties": {
            "displayName": {
                "type": "string"
            },
            "email": {
                "maxLength": 254,
                "minLength": 0,
                "type": "string"
            },
            "givenName": {
                "maxLength": 35,
                "minLength": 0,
                "type": "string"
            },
            "id": {
                "type": "integer",
                "format": "int32"
            },
            "surname": {
                "maxLength": 35,
                "minLength": 0,
                "type": "string"
            }
        }
    },
    "AttendeeResponse": {
        "title": "AttendeeResponse",
        "type": "object",
        "properties": {
            "value": {
                "type": "array",
                "items": {
                    "$ref": "#/components/schemas/Attendee"
                }
            }
        }
    }
}
```

## NOTE

In addition to documentation the author can also supply example responses to demonstrate how the API should be used.

**Components** holds the various object schemas <sup>1</sup> for the specification and in this case defines our `Attendee` and `AttendeeResponse` object.

OpenAPI specifications can also include a wide range of additional metadata and useful features for developers and API consumers. In the `Attendees` example the required fields of an `Attendee` are all fields except for `id`, at the time of generating an `Attendee` the consumer does not know the id. The specification also sets out maximum lengths for some of the strings, it is also possible to set a regular expression to pattern match `email`. Full details of the [OpenAPI Specification](#) are hosted on the book's GitHub.

Example requests and responses, like the one we've shown here, demonstrate a typical data exchange supported by the API. The OAS also documents the OAuth2 flows that are supported by an API, which we will explore further in *Chapter 9*. Over the course of the chapter it should become clear how important the OpenAPI Specification is to offering any type of REST API platform.

## Visualizing OpenAPI Specifications

It's quite difficult to read a specification in JSON or in YAML (which is also supported by OAS), especially as APIs grow beyond a handful of operations. The example specification above includes no user documentation. When documentation is added specifications can rapidly become thousands of lines long, which makes the specification difficult to read without a tool. One of the big success stories of OpenAPI Specifications has been the number of tools available in many different languages. There are tools that enable the developer to generate OpenAPI Specifications directly from their code or use the Swagger Editor in *Figure 2* (TK).

# Practical Application of OpenAPI Specifications

Once an OpenAPI Specification is shared the power of the specification starts to become apparent. [OpenAPI Tools](#) documents a full range of open and closed source tools available. In this section we will explore some of the practical applications of tools based on their interaction with the OpenAPI Specification.

## Code Generation

Perhaps one of the most useful features of an OpenAPI specification is allowing the generation of client side code to consume the API. As discussed in “[Specifying REST APIs](#)” we can include the full details of the server, security and of course the API structure itself. With all this information we can generate a series of model and service objects that represent and invoke the API. The [OpenAPI Generator](#) project supports a wide range of languages and tool chains. For example, in Java you can choose to use Spring or JAX-RS and in Typescript you can choose a combination of Typescript with your favorite framework. It is also possible to generate the API implementation stubs from the OpenAPI Specification.

This raises an important question about what should come first the specification or the server side code? In the next chapter we are going to discuss “[Contract testing](#)” which presents a behavior driven approach to testing and building APIs. The challenge with OpenAPI Specifications is that alone they only convey the shape of the API. OpenAPI specifications do not fully model the semantics (or expected behavior of the API) under different conditions. If you are going to present an API to external users it is important that the range of behaviors is modelled and tested to help avoid having to drastically change the API later.

A common challenge with API modelling, as discussed in “[The Richardson Maturity Model](#)”, is determining whether you need a RESTful API or whether you need RPC. We will explore this idea further in “[Alternative](#)

**API Formats**”. It is important that this is an active decision, as delivering RPC over REST can result in a modelling mistake and a challenge for consumers. APIs should be designed from the perspective of the consumer and abstract away from the underlying representation behind the scenes. It is important to be able to freely refactor components behind the scenes without breaking API compatibility, otherwise the API abstraction loses value.

## OpenAPI Validation

OpenAPI Specifications are useful for validating the content of an exchange to ensure the request and response match the expectations of the specification. At first it might not seem apparent where this would be useful, if code is generated surely the exchange will always be right? One practical application of OpenAPI validation is in securing APIs and API infrastructure. In many organizations a zonal architecture is common, with a notion of a DMZ (Demilitarized Zone) used to shield a network from inbound traffic. A useful feature is to interrogate messages in the DMZ and terminate the traffic if the specification does not match.

Atlassian, for example, Open Sourced a tool called the [swagger-request-validator](#), which is capable of validating JSON REST content. The project also has adapters that integrate with various mocking and testing frameworks to help ensure that API Specifications are conformed to as part of testing. The tool has an `OpenApiInteractionValidator` which is used to create a `ValidationReport` on an exchange.

```
//Using the location of the specification create an interaction validator
//The base path override is useful if the validator will be used behind a
gateway/proxy
final OpenApiInteractionValidator validator = OpenApiInteractionValidator
    .createForSpecificationUrl(specUrl)
    .withBasePathOverride(basePathOverride)
    .build;

//Requests and Response objects can be converted or created using a builder
final ValidationReport report = validator.validate(request, response);
```

```
if (report.hasErrors()) {  
    // Capture or process error information  
}
```

## Examples and Mocking

The OpenAPI Specification can provide example responses for the paths in the specification. Examples, as we've discussed, are useful for documentation to help developers understand the API usage. Some products have started to use examples to allow the user to query the API and return example responses from a mock service. This can be really useful in features such as a Developer Portal, which allows developers to explore documentation and invoke APIs.

Examples can potentially introduce an interesting problem, which is that this part of the specification is essentially a string (in order to model XML/JSON etc). `openapi-examples-validator`<sup>2</sup> validates that an example matches the OpenAPI Specification for the corresponding request/response component of the API.

## Detecting Changes

OpenAPI Specifications can also be useful in detecting changes in an API. This can be incredibly useful as part of a DevOps pipeline. Detecting changes for backwards compatibility is incredibly important, but first it is useful to understand versioning of APIs in more detail.

## API Versioning

We have explored the advantages of sharing an OpenAPI specification with a consumer, including the speed of integration. Consider the case where multiple consumers start to operate against the API. What happens when there is a change to the API or one of the consumers requests the addition of new features to the API?

Let's take a step back and think about if this was a code library built into our application at compile time. Any changes to the library would be packaged as a new version and until the code is recompiled and tested against the new version, there would be no impact to production applications. As APIs are running services, we have a couple of upgrade options that are immediately available to us when changes are requested:

- **Release a new version and deploy in a new location.** Older applications continue to operate against the older version of the APIs. This is fine from a consumer perspective, as the consumer only upgrades to the new location and API if they need the new features. However, the owner of the API needs to maintain and manage multiple versions of the API, including any patching and bug fixing that might be necessary.
- **Release a new version of the API that is backwards compatible with the previous version of the API.** This allows additive changes without impacting existing users of the API. There are no changes required by the consumer, but we may need to consider downtime or availability of both old and new versions during the upgrade. If there is a small bug fix that changes something as small as an incorrect fieldname, this would break compatibility.
- **Break compatibility with the previous API and all consumers must upgrade code to use the API.** This seems like an awful idea at first, as that would result in things breaking unexpectedly in production.<sup>3</sup> However a situation may present itself where we cannot avoid breaking compatibility with older versions. One example where APIs have had to break compatibility for a legal reasons was the introduction of GDPR (General Data Protection Regulation) in Europe.

The challenge is that each of these different upgrade options offer advantages, but also drawbacks either to the consumer or the API owner. The reality is that we want to be able to support a combination of all three

options. In order to do this we need to introduce rules around versioning and how versions are exposed to the consumer.

## Semantic Versioning

**Semantic Versioning** offers an approach that we can apply to REST APIs to give us a combination of the above. Semantic versioning defines a numerical representation attributed to an API release. That number is based on the change in behavior in comparison to the previous version, using the following rules.

- A **Major** version introduces non-compatible changes with previous versions of the API. In an API platform upgrading to a new major version is an active decision by the consumer. There is likely going to be a migration guide and tracking as consumers upgrade to the new API.
- A **Minor** version introduces a backwards compatible change with the previous version of the API. In an API service platform it is acceptable for consumer to receive minor versions without making an active change on the client side.
- A **Patch** version does not change or introduce new functionality, but is used for bug fixes on an existing Major.Minor version of functionality.

Formatting for semantic versioning can be represented as **Major.Minor.Patch**. For example 1.5.1 would represent major version 1, minor version 5 with patch upgrade of 1. Whilst reading the above the reader may have noticed that with APIs running as services there is another important aspect to the story. Versioning alone is not enough, an element of deployment and what is exposed to the consumer at what time is part of the challenge. This is where the API Lifecycle is important to consider, n terms of versioning.

## API Lifecycle

The API space is moving quickly, but one of the clearest representations of version lifecycle comes from the now archived [PayPal API Standards](#). The lifecycle is defined as follows:

*T*  
*a*  
*b*  
*l*  
*e*

*2*  
-  
*3*  
. *A*  
*P*  
*I*

*L*  
*i*  
*f*  
*e*  
*c*  
*y*  
*c*  
*l*  
*e*

(  
*a*  
*d*  
*a*  
*p*  
*t*  
*e*  
*d*

*f*

*r*  
*o*  
*m*

*P*  
*a*  
*y*  
*P*  
*a*  
*l*  
*A*  
*P*  
*I*

*S*  
*t*  
*a*  
*n*  
*d*  
*a*  
*r*  
*d*  
*s*  
)

Planned	<p>Exposing an API from a technology perspective is quite straightforward, however once it is exposed and consumed we have multiple parties that need to be managed.</p> <p>The planning stage is about advertising that you are building an API to the rest of the API program.</p> <p>This allows a discussion to be had about the API and the scope of what it should cover.</p>
Beta	<p>Involves releasing a version of our API for users to start to integrate with, however this is generally for the purpose of feedback and improving the API. At this stage the producer reserves the right to break compatibility, it is not a versioned API.</p>

	<p>This helps to get rapid feedback from consumers about the design of the API before settling on a structure.</p> <p>A round of feedback and changes enables the producer to avoid having many major versions at the start of the APIs lifetime.</p>
Live	<p>The API is now versioned and live in production.</p> <p>Any changes from this point onward would be versioned changes.</p> <p>There should only ever be one live API, which marks the most recent major/minor version combination.</p> <p>Whenever a new version is released the current live API moves to deprecated.</p>
Deprecated	<p>When an API is deprecated it is still available for use, but significant new development should not be carried out against it.</p> <p>When a minor version of a new API is released an API will only be deprecated for a short time, until validation of the new API in production is complete.</p> <p>After the new version is successfully validated a minor version moves to retired, as the new version is backwards compatible and can handle the same features as the previous API.</p> <p>When a major version of the API is released the older version becomes deprecated.</p> <p>It is likely that will be for weeks or months, as an opportunity must be given to consumers to migrate to the new version.</p> <p>There is likely going to be communication with the consumers, a migration guide and tracking of metrics and usage of the deprecated API.</p>
Retired	The API is retired from production and is no longer accessible.

The lifecycle helps the consumer fully understand what to expect from the API. The main question is what does the consumer see with respect to the versioning and lifecycle? With Semantic Versioning combined with the API Lifecycle the consumer only needs to be aware of the major version of the API. Minor and patch versions will be received without updates required on the consumers side and won't break compatibility.

One often controversial question is how should the major version be exposed to the user. One way is to expose the major version in the URL i.e. <http://mastering-api.com/v1/attendees>. From a purely RESTful perspective however the version is not part of the resource. Having the major version as

part of the URL makes it clear to the consumer what they are consuming. A more RESTful way is to have the major version as part of the header, e.g. VERSION: 1. Having the version in a header may be slightly hidden from the consumer. A decision would need to be made to be consistent across APIs.

You may be wondering how APIs with multiple versions can be presented side-by-side during deployments and route to specific API services. We will explore this further in [Chapter 4](#) and in [Chapter 5](#).

# OpenAPI Specification and Versioning

Now that we have explored versioning we can look at examples of breaking changes and non breaking changes using the attendees API specification. There are several tools to choose from to compare specifications, in this example we will use [openapi-diff](#) from OpenAPITools.

We will start with a breaking change. We will change `givenName` to be a field called `firstName`. We can run the diff tool from a docker container using the following command:

```
$docker run --rm -t \
-v $(pwd):/specs:ro \
openapitools/openapi-diff:latest /specs/original.json /specs/first-
name.json
=====
==                               API CHANGE LOG                         ==
=====
=====Attendees Mastering API=====
--                               What's Changed
-- 
- GET      /attendees
  Return Type:
    - Changed 200 OK
      Media types:
        - Changed */*
          Schema: Broken compatibility
          Missing property: [n].givenName (string)
  -----
  --                               Result
  --
```

---

---

API changes broke backward compatibility

---

---

**TIP**

The `-v $(pwd):/specs:ro` adds the present working directory to the container under the `/specs` mount as read only.

We can try to add a new attribute to the `/attendees` return type to add an additional field `age`. Adding new fields does not break existing behavior and therefore does not break compatibility.

```
$ docker run --rm -t \
-v $(pwd):/specs:ro \
openapitools/openapi-diff:latest --info /specs/original.json /specs/age.json
=====
==                         API CHANGE LOG                         ==
=====
                        Attendees Mastering API
-----
--                         What's Changed -- 
-----
- GET      /attendees
  Return Type:
    - Changed 200 OK
    Media types:
      - Changed */*
    Schema: Backward compatible
-----
--                         Result -- 
-----
                        API changes are backward compatible
-----
```

It is worth trying this out to see what would be compatible changes and what would not. Introducing this type of tooling as part of the API pipeline is going to help avoid unexpected non compatible changes for consumers. OpenAPI specifications are an important part of an API program, and when combined with tooling, versioning and lifecycle they are invaluable.

## Alternative API Formats

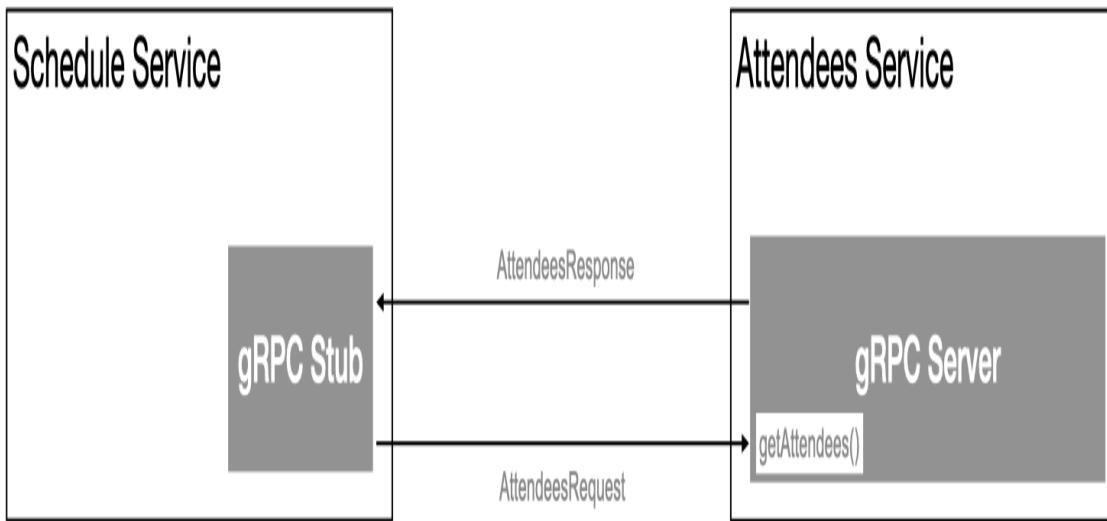
REST APIs work incredibly well for extending services to external consumers. From the consumer perspective the API is clearly defined, won't break unexpectedly and all major languages are supported. But, is using a REST API for every exchange in a microservices based architecture the right approach? For the remainder of the chapter we will discuss the various API formats available to us and factors that will help determine the best solution to our problem.

### Remote Procedure Call (RPC)

Remote Procedure Calls (RPC) are definitely not a new concept. RPC involves executing code or a function of another process. It is an API, but unlike a REST API it generally exposes the underlying system or function internals. With RPC the model tends to convey the exact functionality at a method level that is required from a secondary service.

RPC is different from REST as REST focuses on building a model of the domain and extending an abstraction to the consumer. REST hides the system details from the user, RPC exposes it. RPC involves exposing a method from one process and allows it to be called directly from another.

gRPC is a modern open source high performance Remote Procedure Call (RPC). gRPC is under stewardship of the Linux Foundation and is the defacto standard for RPC across most platforms. [Figure 2-2](#) describes an RPC call in gRPC, which involves the Schedule Service invoking the remote method on the Attendees Service. The gRPC Attendees Service creates a server, allowing methods to be invoked remotely. On the client side, the Schedule Service, a stub is used to abstract the complexity of making the remote call into the library.



*Figure 2-2. Example Attendees with RPC using gRPC*

Another key difference between REST and RPC is state, REST is by definition stateless - with RPC state depends on the implementation. The authors have seen huge systems built around SOAP, which was the XML-RPC successor of the noughties! Many systems have been built around content based routing of messages to specific services with the right cached state to handle the request. In state based RPC systems the developer must have a detailed understanding of each exchange and expected behavior. In order to scale, systems start to leak business logic into routing, which if not carefully managed can lead to increased complexity.

## Implementing RPC with gRPC

The Attendees service could model either a North→South or East→West API. In addition to modelling a REST API we are going to evolve the Attendees service to support **gRPC**.

## TIP

East→West such as Attendees tend to be higher traffic, and can be implemented as microservices used across the architecture. gRPC may be a more suitable tool than REST for East→West services, owing to the smaller data transmission and speed within the ecosystem. Any performance decisions should always be measured in order to be informed.

Let's explore using a [Spring Boot Starter](#) to rapidly create a gRPC server. The Java code below demonstrates a simple structure for implementing the behavior on the generated gRPC server classes.

```
@GrpcService
public class AttendeesServiceImpl extends
AttendeesServiceGrpc.AttendeesServiceImplBase {

    @Override
    public void getAttendees(AttendeesRequest request,
        StreamObserver<AttendeeResponse> responseObserver) {
        AttendeeResponse.Builder responseBuilder =
            AttendeeResponse.newBuilder();

        //populate response
        responseObserver.onNext(responseBuilder.build());
        responseObserver.onCompleted();
    }
}
```

The following .proto file defines an empty request and returns a repeated **Attendee** response. In protocols used for binary representations it is important to note that position and order of fields is important, as they govern the layout of the message. Adding a new service or new method is backward compatible as is adding a field to a message, but care is required. Removing a field or renaming a field will break compatibility, as will changing the datatype of a field. Changing the field number is also an issue as field numbers are used to identify fields on the wire. The restrictions of encoding with gRPC mean the definition must be very specific. REST and OpenAPI are quite forgiving as the specification is only a guide <sup>4</sup>. Extra

fields and ordering do not matter in OpenAPI, versioning and compatibility is therefore even more important when it comes to gRPC.

The following .proto file models the same `attendee` object that we explored in our OpenAPI Specification example.

```
syntax = "proto3";
option java_multiple_files = true;
package com.masteringapi.attendees.grpc.server;

message AttendeesRequest {
}

message Attendee {
    int32 id = 1;
    string givenName = 2;
    string surname = 3;
    string email = 4;
}

message AttendeeResponse {
    repeated Attendee attendees = 1;
}

service AttendeesService {
    rpc getAttendees(AttendeesRequest) returns (AttendeeResponse);
}
```

The Java service modelling this example can be found on the [Book GitHub page](#). gRPC cannot be queried directly from a browser without additional libraries, however you can install [gRPC UI](#) to use the browser for testing. `grpcurl` also provides a command line tool:

```
$ grpcurl -plaintext localhost:9090 \
    com.masteringapi.attendees.grpc.server.AttendeesService/getAttendees
{
  "attendees": [
    {
      "id": 1,
      "givenName": "Jim",
      "surname": "Gough",
      "email": "gough@mail.com"
    }
  ]
}
```

```
    ]  
}
```

gRPC gives us another option for querying our service and defines a specification for the consumer to generate code. gRPC has a more strict specification than OpenAPI and requires methods/internals to be understood by the consumer.

## GraphQL

RPC offers access to a series of individual functions provided by a producer, but does not usually extend a model or abstraction to the consumer. REST extends a resource model for a single API provided by the producer. It is possible to offer multiple APIs on the same base URL by combining REST APIs together using API Gateways. We will explore this notion further in [Chapter 4](#). If we offer multiple APIs in this way the consumer will need to query the APIs sequentially to build up state on the client side. This approach is also wasteful if the client is only interested in a subset of fields on the API response. Consider a user interface that models a dashboard of data on our conference system using visual tiles. Each individual tile would need to invoke each API to populate the UI to display the tile content the user is interested in. [Figure 2-3](#) shows the number of invocations required.

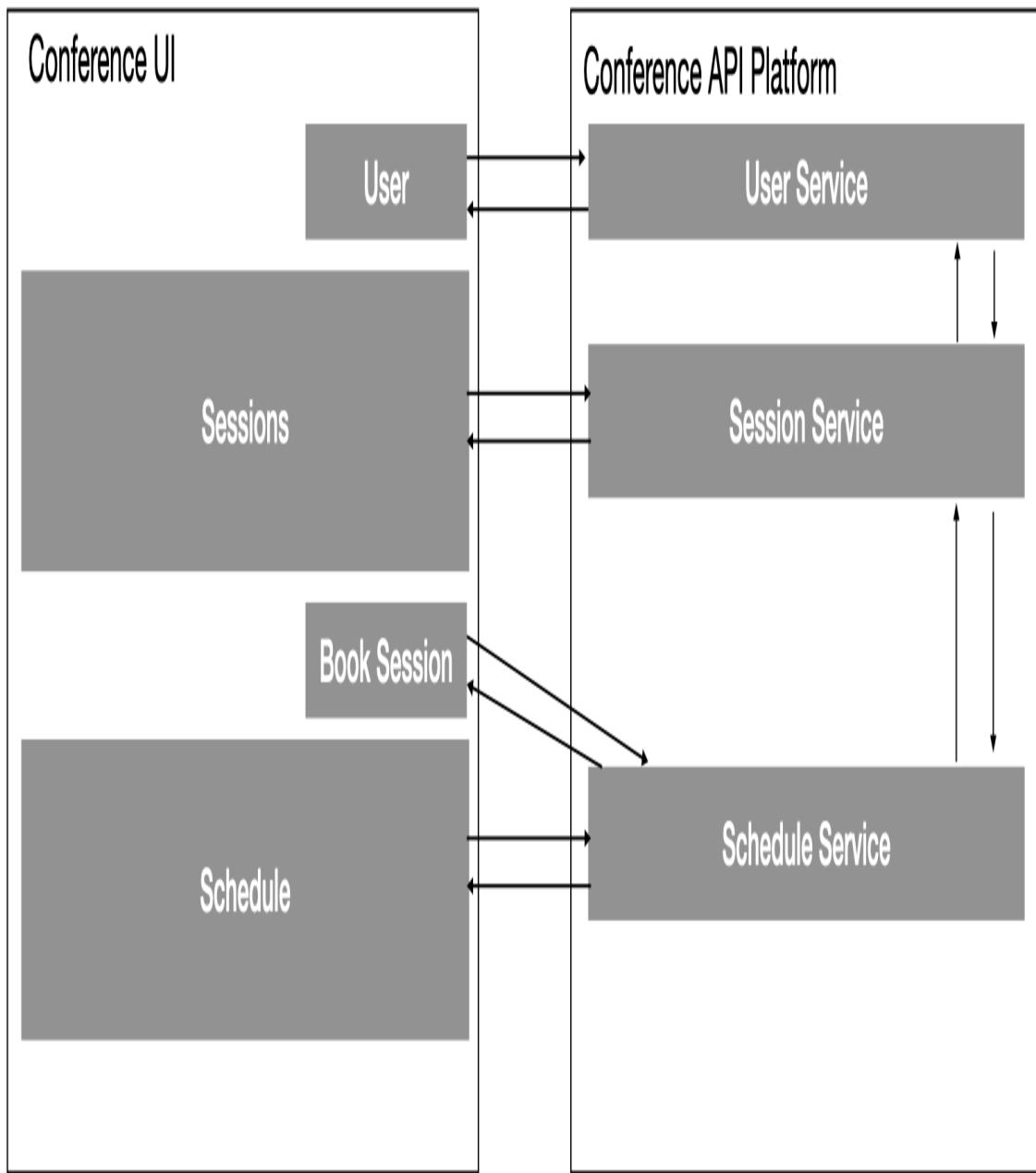


Figure 2-3. Example User Interface model

GraphQL introduces a technology layer over existing services, data stores and APIs that provides a query language to query across multiple sources. The query language allows the client to ask for exactly the fields required, including fields that span across multiple APIs.

GraphQL introduces the [GraphQL schema language](#), which is used to specify the types in individual APIs and how APIs combine. One major advantage of introducing a GraphQL schema is the ability to provide a single version across all APIs, removing the need for potentially complex version management on the consumer side.

After defining the schema the next consideration is the implementation of behavior and defining how data is retrieved and if necessary converted to match the schema. A [Resolver](#) is a function that a GraphQL server implementor creates to source the data for data elements in the GraphQL schema. Mark Stuart has written an excellent [blog on GraphQL resolver best practices for PayPal Engineering](#).

### TIP

One mistake that API developers and architects often make is assuming that GraphQL is only a technology used with User Interfaces. In systems where vast amounts of data is stored across different subsystems GraphQL can provide an ideal solution to abstracting away internal system complexity.

Lets implement a very trivial GraphQL schema for the `/attendees` call to look at what GraphQL looks like to the consumer.

```
var schema = buildSchema(`  
  type Attendee {  
    givenName: String  
    surname: String  
    displayName: String  
    email: String  
    id: Int  
  }  
  
  type Query {  
    attendees: [Attendee]  
  }  
`);  
  
// Logic to resolve and fetch content
```

```
var app = express();
app.use('/graphql', graphqlHTTP({
  schema: schema,
  rootValue: root,
  graphiql: true,
}));
app.listen(4000);
```

GraphQL has a single POST `/graphql` endpoint (not too dissimilar from RPC), however unlike RPC it is not single a method with set behavior that is invoked but a declarative request for specific elements of the schema.

```
curl -X POST -H "Content-Type: application/json" \
-d '{"query": "{ attendees { email } }"}' http://localhost:4000/graphql

{"data": {"attendees": [{"email": "jgough@gmail.com"}]}}
```

[Figure 2-4](#) shows how GraphiQL (a UI tool for GraphQL) provides the consumer with a graphical mechanism for building up queries, along with the ability to explore the schema data types.

GraphQL ➔ Prettify Merge Copy History

⟨ Query Attendee X

```
1 {
2   attendees {
3     id
4     email
5     givenName
6     surname
7   }
8 }
9 }
```

```
{
  "data": {
    "attendees": [
      {
        "id": 1,
        "email": "jgough@gmail.com",
        "givenName": "Jim",
        "surname": "Gough"
      }
    ]
  }
}
```

Q Search Attendee...

No Description

FIELDS

givenName: String

surname: String

displayName: String

email: String

id: Int

*Figure 2-4. Executing a Query from the GraphiQL UI*

GraphQL is a fascinating technology and offers a complement to REST and RPC and in some cases will be a better choice. [Learning GraphQL by Eve Porcello and Alex Banks](#) offers an in-depth exploration for the curious architect. GraphQL works very well when the data and services that a team or company present are from a specific business domain. In the case where disparate APIs are presented externally GraphQL could introduce a complicated overhead to maintain if it tried to do too much. Whilst you can use GraphQL to normalize access to a domain, maintenance may be reduced if the services behind the scenes have already been normalized.

## Exchanges and Choosing an API Format

In [Chapter 1](#) we discussed the concept of traffic patterns, and the difference between requests originating from outside the ecosystem and requests within the ecosystem. Traffic patterns are an important factor in determining the appropriate format of API for the problem at hand. When we have full control over the services and exchanges within our microservices based architecture, we can start to make compromises that we would not be able to make with external consumers.

It is important to recognize that the performance characteristics of an East→West service are likely to be more applicable than a North→South service. In a North→South exchange traffic originating from outside the API producer's environment will generally involve the exchange using the internet. The internet introduces a high degree of latency, and an API architect should always consider the compounding effects of each service. In a microservices based architecture it is likely that one North→South request will involve multiple East→West exchanges. High traffic East→West exchange need to be efficient to avoid cascading slow-downs propagating back to the consumer.

## High Traffic Services

In our example Attendees is a central service. In a microservices based architecture components will keep track of an `attendeeId`. APIs offered to consumers will potentially retrieve data stored in the Attendees service, and at scale it will be a high traffic component. If the exchange frequency is high between services, the cost of network transfer due to payload size and limitations of one protocol vs another will be more profound as usage increases. The cost can present itself in either monetary costs of each transfer or the total time taken for the message to reach the destination.

## Large Exchange Payloads

Large payload sizes may also become a challenge in API exchanges and are susceptible to increasing transfer performance across the wire. JSON over REST is human readable, and will often be more verbose than a fixed or binary representation.

### TIP

A common misconception is that “human readability” is quoted as a primary reason to use JSON in data transfers. The number of times a developer will need to read a message vs the performance consideration is not a strong case with modern tracing tools. It is also rare that large JSON files will be read from beginning to end. Better logging and error handling can mitigate the human readable argument.

Another factor in large payload exchanges is the time taken by components to parse the message content into language level domain objects. Performance time of parsing data formats varies vastly depending on the language a service is implemented in. Many traditional server side languages can struggle with JSON compared to a binary representation for example. It is worth exploring the impact of parsing and include that consideration when choosing an exchange format.

## HTTP/2 Performance Benefits

Using HTTP/2 based services can help to improve performance of exchanges by supporting binary compression and framing. The **binary framing layer** is transparent to the developer, but behind the scenes will split and compress the message into smaller chunks. The advantage of binary framing is it allows for a full request and response multiplexing over a single connection. Consider processing a list in another service and the requirement is to retrieve 20 different attendees, if we retrieved these as individual HTTP/1 requests it would require the overhead of creating 20 new TCP connections. Multiplexing allows us to perform 20 individual requests over a single HTTP/2 connection.

gRPC uses HTTP/2 by default and reduces the size of exchange by using a binary protocol. If bandwidth is a concern or cost gRPC will provide an advantage, in particular as content payloads increase significantly in size. gRPC may be beneficial compared to REST if payload bandwidth is a cumulative concern or the service exchanges large volumes of data. If large volumes of data exchanges are frequent it is also worth considering some of the asynchronous capabilities of gRPC, which we will cover in *Chapter 10*.

## Vintage Formats

Not all services in an architecture will be based on a modern design. In *Chapter 6* we will look at how to isolate and evolve vintage components, however as part of an evolving architecture older components will be an active consideration. Many older services will use formats such as SOAP/XML over HTTP/TCP. It is important that an API architect understands the overall performance impact of introducing vintage components.

## Performance Testing Exchanges

Recognizing the performance characteristics of exchanges is a useful skill for an API Architect to develop. Often it is not the network alone that needs to be considered - the test should include the network, parsing, responding to the query and returning a response. Smaller benchmarks do not capture

the full picture, so it is important to look at the performance in the context of your system. Let's explore the approach of a simple end-to-end performance test in the gRPC ecosystem.

Performance is at the heart of every build of the gRPC libraries, and a **Performance Dashboard** monitors each build and the impact of changes on performance. Buying into the gRPC ecosystem will provide a full stack of complementing libraries that work together in the target language for building services.

If performance is a primary concern for a service it is important to build a benchmark that can be used to test changes to code and libraries over time. We can use a gRPC benchmarking tool `ghz` to get an idea of the performance of the attendees service.

```
brew install ghz
ghz --insecure --format=html --total=10000 \
  --proto ./attendees.proto \
  --call
com.masteringapi.attendees.grpc.server.AttendeesService.getAttendees \
  -d {} localhost:9090 > results.html
```

**Figure 2-5** shows a graphical representation of the the performance of 10,000 requests. The average response time was 2.49 milliseconds, the slowest response time was 14.22 milliseconds and the fastest was 0.16ms.

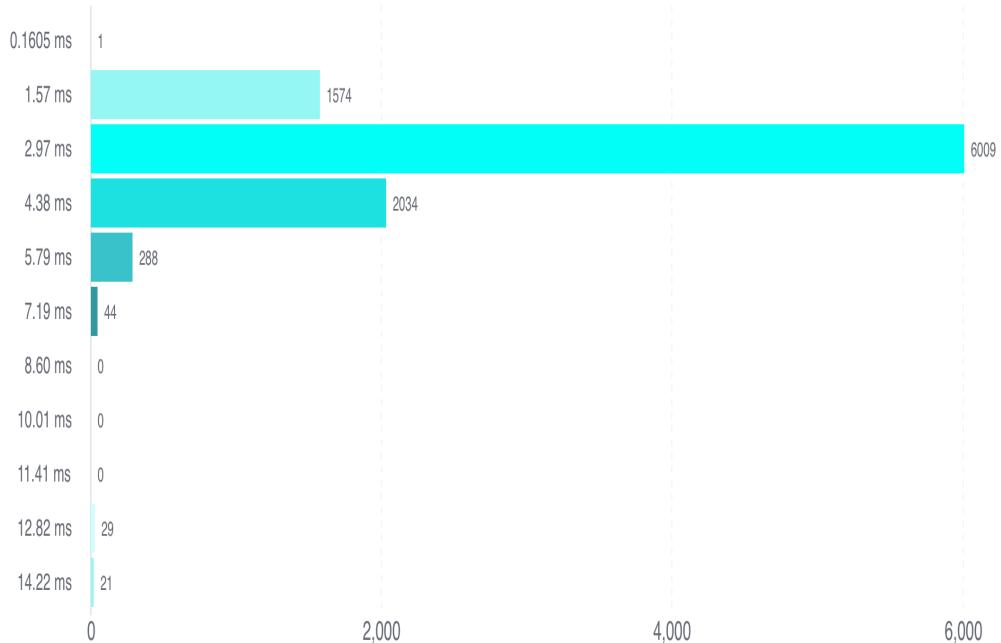
## Summary

Count	10000
Total	519.74 ms
Slowest	14.22 ms
Fastest	0.16 ms
Average	2.49 ms
Requests/sec	19240.53

## Options

```
{  
  "host": "localhost:9090",  
  "proto": "./attendees.proto",  
  "import-paths": [  
    ".."  
  ],  
  "call": "com.masteringapi.attendees.grpc.server.AttendeesService.getAttendees",  
  "insecure": true,  
  "total": 10000,  
  "concurrency": 50,  
  "connections": 1,  
  "timeout": 2000000000,  
  "dial-timeout": 1000000000,  
  "data": {},  
  "binary": false,  
  "CPUs": 8  
}
```

## Histogram



*Figure 2-5. GHZ gRPC Benchmark Tool - <https://github.com/bojand/ghz>*

**TIP**

We want to avoid premature optimization and benchmarks without analysis can lead to a confirmation bias. gRPC will provide performance benefits, but it is important to consider the consumer and their expectations of the API.

gRPC also supports asynchronous and streaming APIs, we will spend *Chapter 10* discussing asynchronous approaches to APIs. If services are constantly exchanging information an open asynchronous pipe would offer advantages over an individual request/response model. In *Chapter 5* we will explore alternative approaches to testing and monitoring the behavior of applications and exchanges in production.

## **Checklist: Modelling Exchanges**

*T*  
*a*  
*b*  
*l*  
*e*

*2*  
-  
*4*  
. *M*  
*o*  
*d*  
*e*  
*l*  
*l*  
*i*  
*n*  
*g*

*E*  
*x*  
*c*  
*h*  
*a*  
*n*  
*g*  
*e*  
*s*

*C*  
*h*  
*e*  
*c*

*k*  
*l*  
*i*  
*s*  
*t*

Decision	What format should we use to model the API for our service?
Discussion Points	<p>Is the exchange a North→South or East→West exchange? Are we in control of the consumer code?</p> <p>Is there a strong business domain across multiple services or do we want to allow consumers to construct their own queries?</p> <p>What versioning considerations do we need to have?</p> <p>What is the deployment/change frequency of the underlying data model.</p> <p>Is this a high traffic service where bandwidth or performance concerns have been raised?</p>
Recommendations	<p>If the API is consumed by external users REST is the lowest barrier to entry and provides a strong domain model.</p> <p>If the APIs offered connect well together and users are likely to use the API to query across APIs frequently consider using GraphQL.</p> <p>If the API is interacting between two services under close control of the producer or the service is proven to be high traffic consider gRPC.</p>

## Multiple Specifications

In this chapter we have explored a variety of API formats for an API Architect to consider and perhaps the final question is “*Can we provide all formats?*”. The answer is yes we can support an API that has a RESTful presentation, a gRPC service and connections into a GraphQL schema. However, it is not going to be easy and may not be the right thing to do. In this final section we will explore some of the options available for a multi-format API and the challenges it can present.

## The Golden Specification

The `.proto` file for attendees and the OpenAPI specification do not look too dissimilar, they contain the same fields and both have data types. Is it possible to generate a `.proto` file from an OpenAPI specification using the `openapi2proto` tool? Running `openapi2proto --spec spec-v2.json` will output the `.proto` file with fields ordered alphabetically by default. This is fine until we add a new field to the OpenAPI specification that is backwards compatible and suddenly the ID of all fields changes, breaking backwards compatibility.

The sample `.proto` file below shows that adding `a_new_field` would be alphabetically added to the beginning, changing the binary format and breaking existing services.

```
message Attendee {  
    string a_new_field = 1;  
    string email = 2;  
    string givenName = 3;  
    int32 id = 4;  
    string surname = 5;  
}
```

OpenAPI specifications support the idea of extensions, and by using the `openapi2proto` specific OpenAPI extensions it is possible to generate the compatibility between the two mechanisms of specification.

### NOTE

There are other tools available to solve the specification conversion problem, however it is worth noting that some tools only support OpenAPI Specification version 2. The time taken to move between version 2 and 3 in some of the tools built around OpenAPI has led to many products needing to support both versions of the OAS.

An alternative option is `grpc-gateway`, which generates a reverse-proxy providing a REST facade in front of the gRPC service. The reverse proxy is generated at build time against the `.proto` file and will produce a best

effort mapping to REST, similar to `openapi2proto`. You can also supply extensions within the `.proto` file to map the RPC methods to a nice representation in the OpenAPI specification.

```
import "google/api/annotations.proto";
//...
service AttendeesService {
    rpc getAttendees(AttendeesRequest) returns (AttendeeResponse) {
        option(google.api.http) = {
            get: "/attendees"
        };
    }
}
```

Using `grpc-gateway` gives us another option for presenting both a REST and gRPC service. However, `grpc-gateway` involves several commands and setup that would only be familiar to developers who work with the go language or build environment.

## Challenges of Combined Specifications

It's important to take a step back here and consider what we are trying to do. When converting from OpenAPI we are effectively trying to convert our RESTful representation into a gRPC series of calls. We are trying to convert an extended hypermedia domain model into a lower level function to function call. This is a potential conflation of the difference between RPC and APIs and is likely going to result in wrestling with compatibility.

With converting gRPC to OpenAPI we have a similar issue, the objective is trying to take gRPC and make it look like a REST API. This is likely going to create a difficult series of issues when evolving the service.

Once specifications are combined or generated from one another, versioning becomes a challenge. It is important to be mindful of how both the gRPC and OpenAPI specifications maintain their individual compatibility requirements. An active decision should be made as to whether coupling the REST domain to an RPC domain makes sense and adds overall value. Rather than generate RPC for East→West from North→South, what makes more sense is to carefully design the

microservices based architecture (RPC) communication independently from the REST representation, allowing both APIs to evolve freely.

GraphQL offers a mechanism that is version-less from the consumers perspective, they interact with only the data that they wish to retrieve. This is at the cost to the producer in maintaining a GraphQL Schema and logic that is used to fetch and resolve data from the underlying services. It is possible to offer REST APIs to external users and then use the separate GraphQL server to aggregate together APIs that have combined domains. It is also possible to use GraphQL to present RPC based services in a normalized schema to clients.

## Summary

In this chapter we have scratched the surface of a variety of topics that an API Architect, API developer or individuals involved in an API program perspective must appreciate.

- The barrier to building a REST API is really low <sup>5</sup> in most technologies.
- REST is a fairly loose standard and for building APIs, conforming to an agreed API Standards ensures our APIs are consistent and have the expected behavior for our consumers.
- OpenAPI specifications are a useful way of sharing API structure and automating many coding related activities. Teams should actively select OpenAPI features within their platform and choose what tooling or generation features will be applied to projects.
- Versioning is an important topic that adds complexity for the API producer but is necessary to ease API usage for the API consumer. Not planning for versioning in APIs exposed to consumers is dangerous. Versioning should be an active decision in the product feature set and a mechanism to convey versioning to consumers should be part of the discussion. Versioning alone is usually not

enough and ensuring we have an API Lifecycle to help govern versioning will lead to a successful API offering.

- REST is great, but is it not always the best option especially when traffic patterns and performance concerns are factored in. It is important to consider how we approach and model exchanges of information in our microservices based architecture. gRPC and GraphQL provide options that need to be considered when we design our exchange modelling.
- Modelling multiple specifications starts to become quite tricky, especially when generating from one type of specification to another. Versioning complicates matters further but is an important factor to avoid breaking changes. Teams should think carefully before combining RPC representations with RESTful API representations, as there are fundamental differences in terms of usage and control over the consumer code.

The challenge for an API architect is to drive the requirements from a consumer business perspective, create a great developer experience around APIs, and avoid unexpected compatibility issues.

---

1 The schema object is an extended subset of the [JSON Schema Specification Wright Draft 00](#).

2 <https://github.com/codekie/openapi-examples-validator>

3 The authors have been in this situation many times, usually first thing on a Monday!

4 Validation of OpenAPI specifications at runtime helps enforce a greater strictness of OpenAPI Specifications.

5 Anecdotally, whilst researching this chapter some developers claim it is 10 minutes.

# Chapter 3. Testing APIs and the Test Pyramid

---

## A NOTE FOR EARLY RELEASE READERS

With Early Release ebooks, you get books in their earliest form—the authors' raw and unedited content as they write—so you can take advantage of these technologies long before the official release of these titles.

This will be Chapter 3 of the final book.

If you have comments about how we might improve the content and/or examples in this book, or if you notice missing material within this chapter, please reach out to the author at [feedback@mastering-api.com](mailto:feedback@mastering-api.com).

In [Chapter 2](#) we covered the different types of APIs and the value that they provide to your architecture. This chapter closes out the Designing, Building and Testing APIs section of this book, we finish by reviewing approaches to testing APIs. In [Chapter 1](#) we explained that we were taking the evolutionary step to separate the Attendee functionality into its own separate API, and we want to make this a testable solution. The authors believe that, testing is core to building APIs and that testing helps deliver a quality product to the consumer of your API.

If we compare building APIs to any sort of product, such as a mouthguard <sup>1</sup>, the only way to verify that the product works as expected is to test it. In the case of the mouthguard this can mean stretching, hitting, pushing and pulling the product, or even running simulations. Similarly, testing an API service gives us confidence that it operates correctly under varying conditions as expected.

As discussed in Chapter 2 “[Specifying REST APIs](#)”, an API should not return anything unexpected from its documented results. It is also infuriating when an API introduces breaking changes or causes network timeouts due to the large duration of time to retrieve a result. These types of issues drive customers away and are entirely preventable by creating quality tests around the API service. Any API built should be able to handle a variety of scenarios, including sending useful feedback to users who provide a bad input, being secure and returning results within a specified SLO (Service Level Objectives) based on our SLIs (Service Level Indicator). <sup>2</sup> that are agreed.

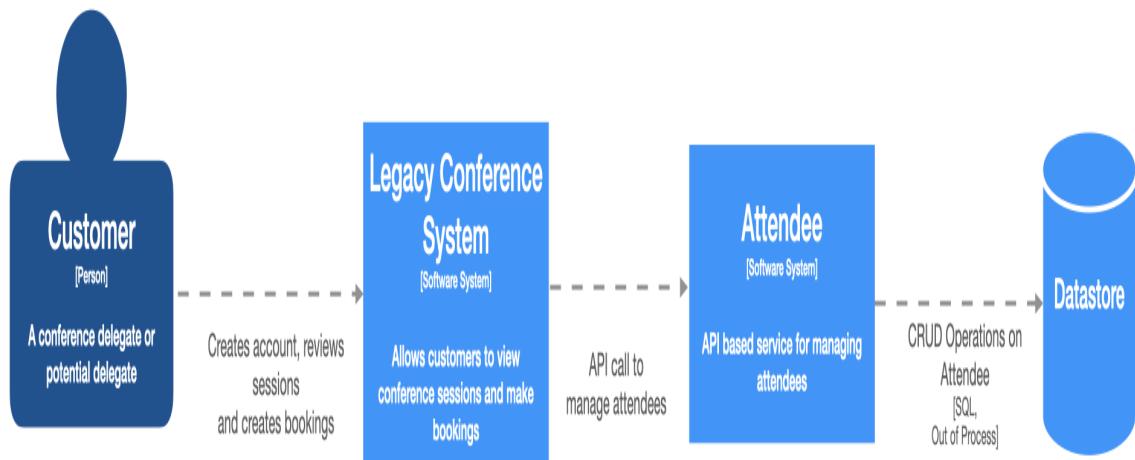
This chapter will introduce the different types of testing that can be applied to your API to help avoid these unforeseen issues. We will highlight the positives and the negatives of each type of testing to understand where the most time should be invested. Throughout the chapter we will reference some recommended resources for those readers seeking to gain a significantly more in depth and specialist knowledge about a subject area.

## Scenario for this chapter

In Chapter 1, “[Attendees Evolution ADR](#)”, we cover the reasons that we took the evolutionary approach to separate the Attendee API service from the rest of the Conference System. This separation introduces new interactions between services, the new Attendee API service will still be used by the Legacy Conference System, as well as the new use case, in which it will be used by the external CFP system. We can see how the Legacy Conference System will work with the new Attendee API service here [Figure 3-1](#). We will spend this Chapter covering the testing that will be needed for this new service and how we can use testing to verify the interactions between the Legacy Conference System and the Attendee API. As a collective we have seen enough APIs that become inconsistent or produce accidental breaking changes as new releases are made, this is all due to a lack of testing. We want to make sure that the new Attendee API is not like this and that as the Attendee API will be functionally correct. It is

important to ensure that the Attendee API returns the correct results, a service that is not correct is useless, so we need to get this right.

# Chapter 3 - Testing Introduction



*Figure 3-1. Scenario for the chapter*

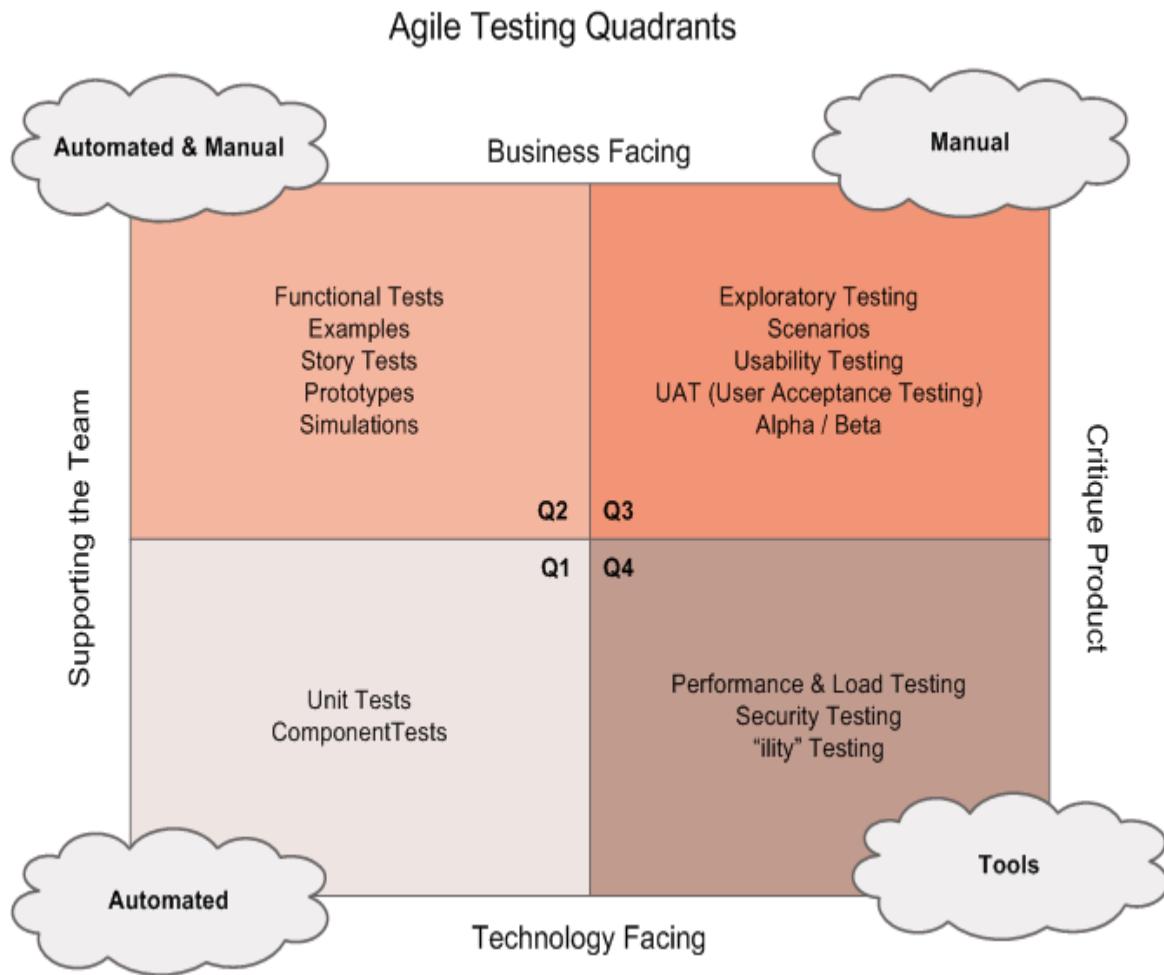
Testing can be applied at different levels to an API, this starts with the individual building blocks that make up the API itself to verifying that it works as part of the entire ecosystem. Before implementing some tests and researching what testing tools and frameworks are available, it is important to understand the strategies that can be used when testing. This is where we believe the Testing Quadrant is important to discuss.

## Testing Quadrant

The Attendee API service is being separated out from the Legacy Conference System and as such when we do start testing the new setup we want to know that we are performing the correct tests. We use the Testing Quadrant to help guide our testing efforts. If we start testing the wrong thing then we waste time or produce incorrect results.

The testing quadrant was first introduced by Brian Marick in his blog series on [agile testing](#). This became popularized in the book *Agile Testing* by Lisa Crispin and Janet Gregory (Addison-Wesley). The testing quadrant brings together Technology and the Business. Technology cares that the service has been built correctly, that its pieces (e.g. functions or endpoints) respond as expected, and that it is resilient and continues to behave under abnormal circumstances. The Business cares that the right service is being developed i.e. in the case of the Attendee API service do the right attendees get returned when querying the API? To clarify by the term *The Business* we mean that this someone that has a clear understanding of what should be developed to take the product in the right direction. For our case we want to ensure that we are building the right API.

There is a distinction between the priorities from each perspective and the Testing Quadrant brings these together to create a nice collaborative approach to develop testing. The popular image of the test quadrant is shown in [Figure 3-2](#).



*Figure 3-2. Agile Testing Quadrants from Agile Testing (Addison-Wesley) by Lisa Crispin and Janet Gregory*

The testing Quadrant does **not** depict any order, they are labeled for convenience and this is a common source of confusion that Lisa describes in one of her [blog posts](#). The four quadrants can be generally described as follows:

- Q1 Unit and Component tests for technology, what has been created works as it should, automated testing.
- Q2 Tests with the business to ensure what is being built is serving its purpose, combination of manual and automated.
- Q3 Testing for the business, meeting functional requirements, exploratory testing and expectations being fulfilled.

- Q4 The system will work as expected from a technical standpoint, including aspects such as; security, SLA integrity, the service scales effectively, whether horizontally or vertically. These are all non-functional requirements. We will cover more of the Q4 concerns in Chapter 6.

The left side of quadrant (Q1, Q2) is all about supporting the product, it helps guide the product and prevent defects. The right side (Q3, Q4) is about critiquing the product and finding defects in the product. The top of the pyramid (Q2,Q3) is the external quality of your product, what the business finds important, in our case the users of our APIs. The bottom of the pyramid (Q1, Q4) are the Technology facing tests so that we maintain the internal quality of our application.<sup>3</sup>

If our conference system was looking at selling tickets, which is a system that must handle large traffic spikes, it may be best to start with Q4. When building APIs you will look at each of the quadrants, though in this chapter the primary focus will be on automating the testing. We want to create the right tests for our API service and automated testing allows for faster feedback.

## Test Pyramid

In addition to the testing quadrants we can use the Test pyramid as part of our strategy for test automation. The Testing pyramid was first introduced in the book **Succeeding with Agile** by **Mike Cohn**. This pyramid gives the notion of how much time should be spent on a testing area, the difficulty and the returned value that it brings. If you search online for images of a testing pyramid thousands will appear, they all consist of different blocks, colours, arrows drawn on them and some even have clouds at the top. However, the Testing Pyramid at its core has remained unchanged. The Testing Pyramid has Unit Tests as its foundation, Service Tests as the middle block and UI tests at the peak of the pyramid.

Figure 3-3 shows a Testing pyramid which comes from Martin Fowler's online post <https://martinfowler.com/bliki/TestPyramid.html>

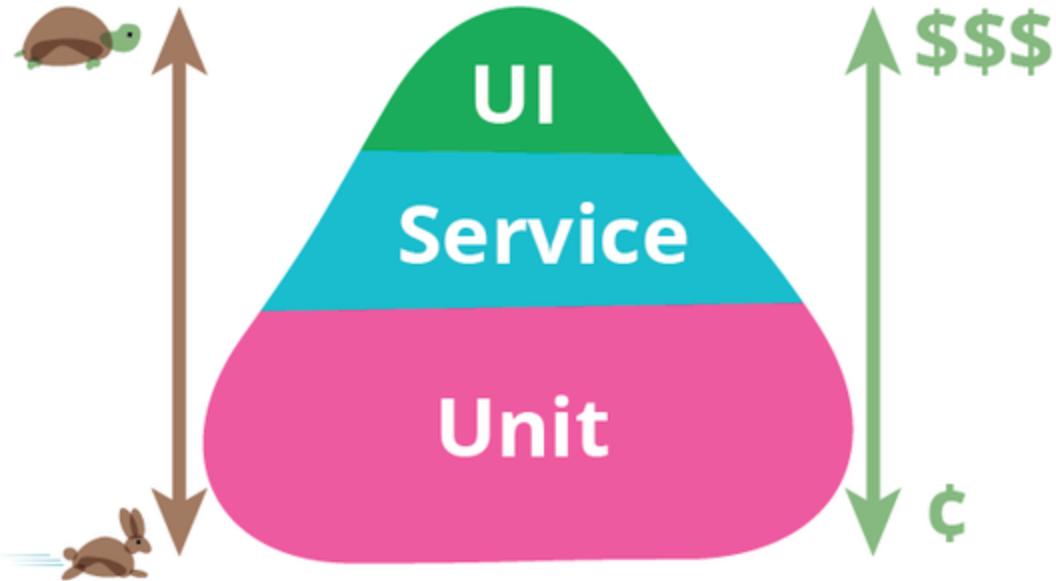


Figure 3-3. *Martin Fowler Testing Pyramid*

An important point to note is that when we look at this image we see intuitive icons that allow us to understand that the UI tests at the top of the pyramid are slower and more costly than the faster and cheaper unit tests at the bottom of the pyramid.

When testing APIs start at the bottom of the pyramid with the Unit tests, which form the foundation of your automated testing. These fit into Q1 of the testing quadrant, they are used to provide quality to the internals of your API. Unit tests are testing small, isolated units of your code to ensure that your defined unit is running as expected.

Service tests are next in the pyramid with a higher development and maintenance cost than unit tests and a slower run time than unit tests. The reason for the increased maintenance cost in a service test is that they are more complex, service tests. They consist of multiple units tested together to verify they integrate correctly. For APIs this is about ensuring that your API can integrate with other components, maybe even other APIs, which comes under Q1 of the testing quadrant. Service tests also check that your API is responding to requests and responses as expected, ensuring that it is returning the payloads as expected, this fits into Q2 of the testing quadrant. This can cross over into Q4 of the quadrant if we want to verify that the API is working as expected with using automated performance tests. As

these tests are more complex this will mean that they will run more slowly than unit tests.

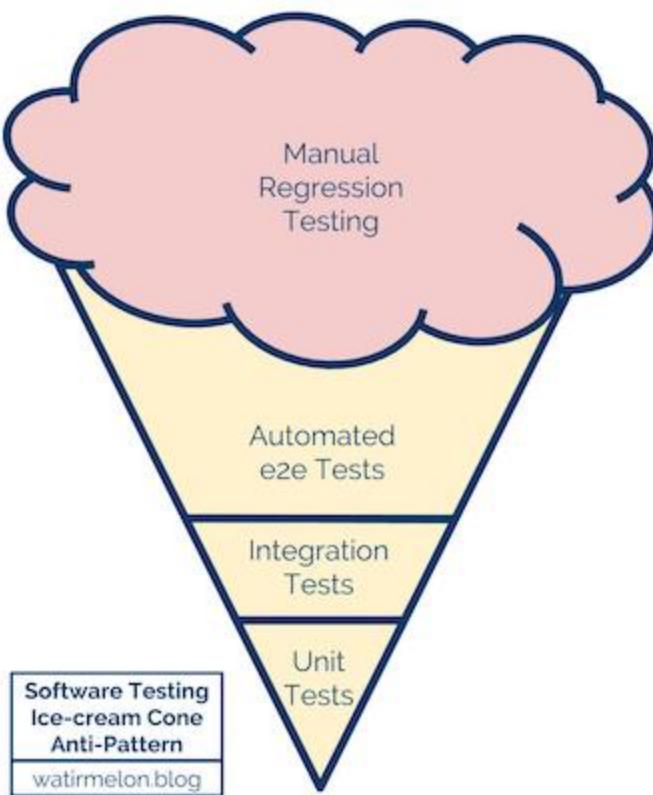
Finally, we have UI tests which are the peak of the pyramid. We want to state now that we will not be using the term UI tests and instead be using the term end-to-end tests. In the olden days of software development where everything was a **LAMP stack**, the web UI was the only way to test your application all the way through. In the world of APIs we do not care what the callee is , it can be a UI, a process or anything for that matter. Therefore, the term End-to-End covers the same ground, of a request flowing from a start to an end point, but does not enforce that it must be a UI. These are the most complex tests, so they will have the most cost in terms of creation and maintenance and will also run the most slowly. The End-to-end tests will test that an entire module is working together with all its integrations. This high cost and low speed demonstrates why UI tests are the peak of the pyramid as the cost to benefit diminishes.

This does not mean that one type of test is better than another. When building an API it is important to perform all of these types of test. We do not want to end up with only unit tests or end-to-end tests.

The testing pyramid is a guide to the proportions of each type of testing that should be done. Most of your time/expense should be on Unit tests, and less on Service/end-to-end tests

The reason that we need a testing pyramid is that from a high level, as an Architect or a project owner standpoint, UI tests are the most tangible. Having something tangible makes these tests feel the safest and by giving someone a list of step by step tests to follow this will catch all the bugs. This gives the false sense of security that these higher level tests are of higher quality/value than unit tests, which are not in an architects control. This fallacy gives rise to the ice cream cone representation of testing, which is the opposite of a testing pyramid. For a robust argument on this please read Steve Smiths' blog post **End-To-End Testing considered harmful**

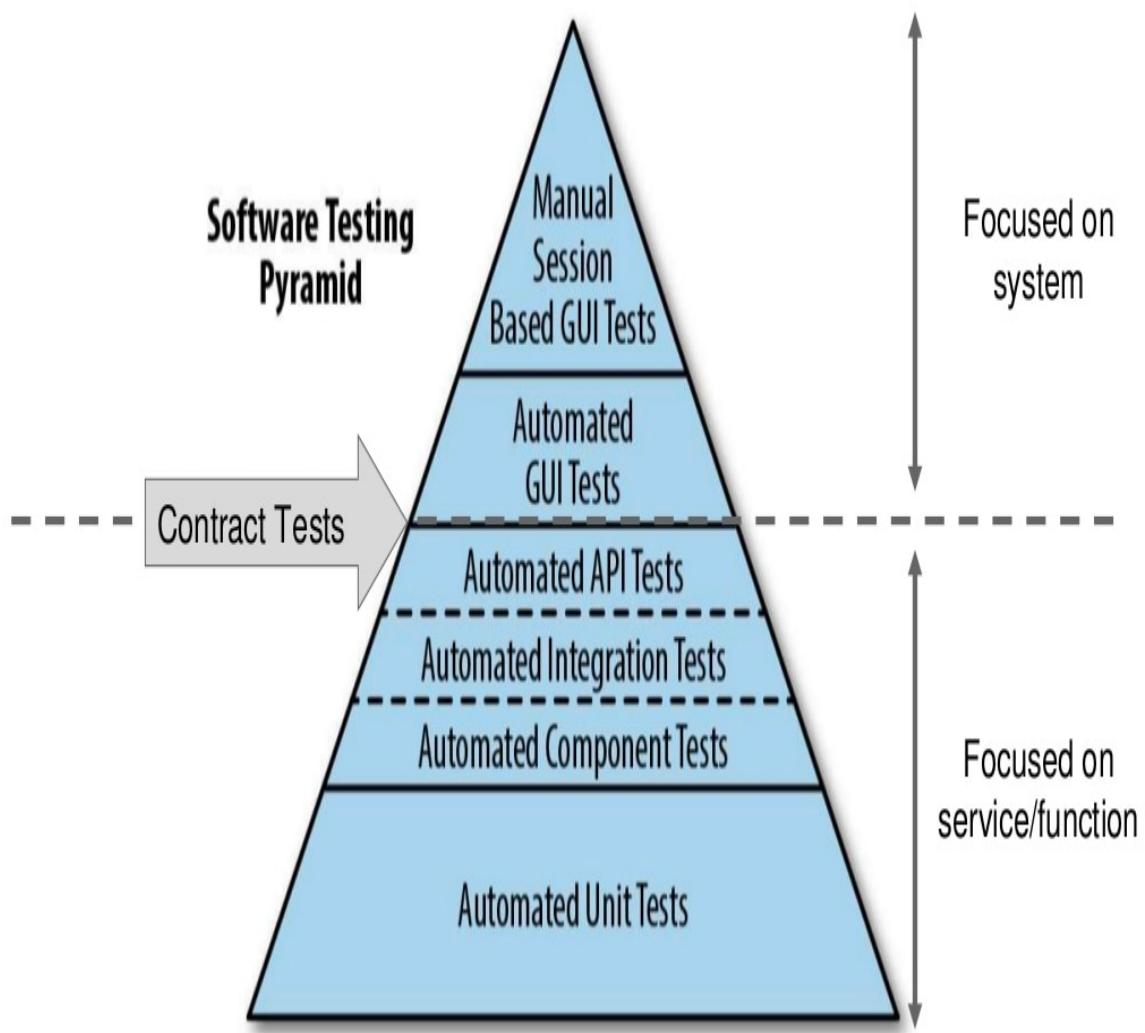
Here we have an image of the anti-pattern testing ice cream cone  
<https://alisterbscott.com/kb/testing-pyramids/>.



*Figure 3-4. Software Testing Ice-Cream Cone by Alister B Scott*

As can be seen in [Figure 3-4](#) the emphasis here is on the high level manual and end-to-end testing over the lower level and more isolated unit tests. This is exactly what you do not want to do, to end up with lots of automated tests that are slow and have a high cost to maintain. Martin Fowler wrote an updated piece on testing shapes and covered why he feels that testing that is guided by any shape other than the pyramid is incorrect., <https://martinfowler.com/articles/2021-test-shapes.html>.

For this Chapter we would like to show the following test pyramid [Figure 3-5](#) that we feel helps breakdown further the test pyramid and show the parts that we are more interested in. It should help guide you in each area that we will cover and give you a clearer picture of where these tests fit in.



*Figure 3-5. Broken down test pyramid*

## Unit Testing

As we've mentioned, the foundation of the testing pyramid is unit testing, we have outlined some characteristics that unit tests should have, they should be fast numerous and cheap to maintain. However, what is the unit? The typical example of a Unit in Object-Oriented (OO) languages is a Class, though some make unit tests per public method of a Class. For a declarative language, such as Haskell, a unit may be defined as a function or some Logic.<sup>4</sup>

At this point you may be wondering why are we spending time looking at unit tests in this discussion of creating our new Attendee APIs? Unit tests are for testing logic at small isolated parts. This gives a developer, who is implementing the unit, the confidence that this piece has been built “correctly” and doing what it is supposed to do. When building an API it must be reliable for the customer and so it is essential that all these foundational units are working as expected. We want our new Attendees API service to return correct results.

Let us take an example where we have a class that handles verifying the input data of a request made to our Attendee API. In this case we are not concerned about how the input was deserialized, there is another part of the service that handles that, we want to validate the data that has been passed in.

```
boolean valid = attendeeValidator.validate(attendeeData);
```

We can write unit tests to ensure that the data we care about exists, in this case we care that the attendee data includes the event attended and the time that they arrived and left the event. Other checks include the time that the attendee left the event after the attended time. This is a test to ensure the internal quality of your API works as expected. The nice thing about this is that it is very self-contained, however not all cases will be like that.

This is where we can use stubs and/or mocks.<sup>5</sup> A stub is an object used for a test that can return hard coded responses. Mocks are pre-programmed objects that work on behavior verification. In the following example an attendee record should be added to the database if the data is valid. For the test we have the option to stub or mock the database and the validator.

```
Response createAttendeeRecord(AttendeeData attendee){
    if(!attendeeValidator.validate(attendee)){
        return new ClientRequestError400();
    }
    var newAttendeeRecord = attendeeDB.addRecord(attendee);
    return new CreatedRecord201(newAttendeeRecord);
}
```

This is really valuable for our use case as we are testing core logic of our API, however, we are not escaping this unit. With these unit tests we can perform lots of different verification such as checking:

- When non-valid attendee data is presented
- When the database fails to save
- When valid data is saved to the database

Stubs and mocks are known as types of doubles. To learn a lot more on this subject of mocks and stubs please read this excellent article by Martin Fowler titled [Mocks Aren't Stubs](#)

## Test Driven Development

Test Driven Development also known as TDD is a style of development used to write tests before writing any code. This simple iterative process was discovered/re-discovered by Kent Beck, for an in depth learning about TDD refer to the works of Kent Beck's *Test Driven Development: By Example*(Addison-Wesley).

The TDD process (also known as Red, Green, Refactor) is defined as:

- Write Test - Write the test that you wish your unit to fulfill
- Run all the Tests - Your code should compile and run, however, your tests should not pass.
- Write the code to make test pass
- Run all the Tests - Validates that your code is working and nothing is broken
- Refactor code - Tidy up the code
- Repeat

So why is TDD important to building APIs? We find this a useful technique for development to deliver the correct outcome and having tests help guide

development. Knowing the outcomes upfront helps the developer to avoid making assumptions and encourages them to ask questions early if they are unsure of what the actual end result should be. Writing tests also ensures that your code keeps functioning as expected, if it is modified in the future and breaks the implementation then a test should catch the error. There are different ways to perform TDD. However, for the Attendee API service we would use an Outside-In approach. This approach entails asking what is required, i.e what is the functionality that should be implemented, and focusing on writing tests that verify that behaviour is implemented. We like this approach as it requires the developer to ask questions about what is actually required. This also means that you do not write functionality that is not required which saves time.

A consequence of TDD is that as you write all of your tests first, you end up with lots of unit tests, which is essential to the base of the test pyramid. This is a benefit — if you do not use TDD then you will need to have another process to add unit tests to your API. We do recommend TDD as a general development approach, but it is not necessary.

## Service Testing

As part of developing your API you want to ensure that it can be called by a consumer and that the correct data is being produced. We also want to validate that if our API is calling other services that it can successfully integrate with these services as well. For APIs these service level tests are important, in order to confirm responses that match what is expected.

Here we see the Attendees API service in our C4 diagram, for our service tests, the scope of testing (the center area within the dotted lines) is just the Attendees API.

# Chapter 3 - Service Testing



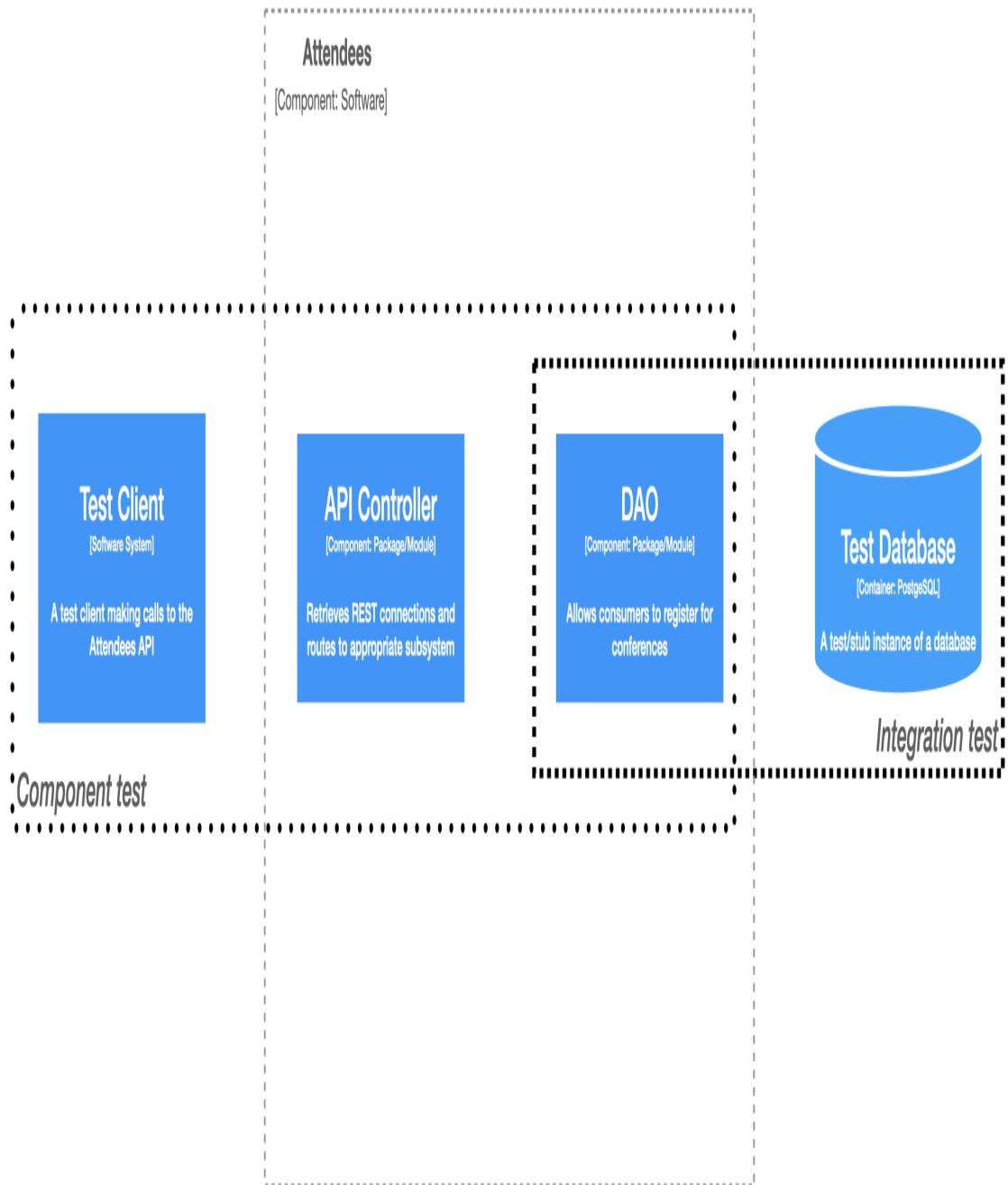
*Figure 3-6. Scope of service tests*

Since we want tests to validate what our API produces and how our API service integrates with other services, we will have two definitions for our service level tests.

- Component tests, which are testing our API and do not include external services (if they are needed then doubles are used).
- Integration tests, which are tests that verify egress communication across a module boundary (e.g. DB Service, External service, other Module)

The following image clarifies the scope of each of these

# Chapter 3 - Service Testing - Component and integration



We start by demonstrating how we can validate the Attendees API is working as expected when requests are made to it. Then we will want to show an integration test, however, instead of showing the integration between the DAO (**Data Access Object**) and the Test DB shown in the image we first want to show how the Legacy Conference System can write tests to integrate with the new Attendee API. After which we will discuss how we can improve testing between services by using contract tests and then move onto showing how to test the DAO and test database using Testcontainers.

## Component tests

Component testing to see how our module works with multiple units is really valuable. It validates that the modules work as expected and our API can return the correct response from incoming consumer requests.

This is where using a library or testing framework that wraps a request client can be really useful, as these libraries usually have a DSL and make it easy to analyse responses from the API. An example in Java is to use RestAssured which is a testing library for REST services or using the `httpptest` package that comes with Go. Depending on the language or framework that you use there should be something available, otherwise, creating a small wrapper around a standard client can make things considerably easier to integrate responses when writing tests.

The type of tests that we want to trigger in this scope varies based on the business case, however, we do want to validate cases such as:

- Is the correct status code returned when I make a request?
- Does the response contain the correct data?
- When I send a request where the accepted content type is XML will the data return the expected format?

- If a request is made by a user who does not have the correct entitlements, what will the response be?
- What will happen if an empty dataset is returned, is this a 404 or is it an empty array?
- When creating a resource does the location header point to the new asset created?

From this selection of tests we see how these bleed into two areas of the testing quadrant. This includes Q1 where we are checking that we built the API correctly as we want the components of our API to all work together, Q2 where we have automated tests to verify that we built the right API (as we want to check the response) and that it is functioning correctly (as it is rejecting unauthorized requests).

For component tests we do not include external dependencies, any that are required should be test doubles. In the Attendees component tests we will want to use a test double for the DAO, there is no need to have a real instance that communicates with a database as that is crossing a boundary. Also, this is useful if you want to force a module to fail and trigger issues such as client errors and internal service errors which can help verify that responses are generated correctly.

Let's look at an example of a case for our Attendees API. We have an endpoint `/conference/{conference-id}/attendees` which returns a list of the attendees at a conference event. The things we want to test for this endpoint are:

- Requests that are successful have response of 200 (OK)
- Users without the right level of access will return a status of 403 (Forbidden)
- When a conference has no attendees an empty array will be returned

Here we use **REST-Assured** to show an example of calling out to the Attendees API to verify these cases

```
@Test
void response_for_attendees_should_be_200() {
    given()
    .when()
        .get("/conference/conf-1/attendees")
    .then()
        .statusCode(HttpStatus.OK.value());
}

@Test
void response_for_attendees_should_be_403() {
    given()
    .when()
        .get("/conference/conf-1/attendees")
    .then()
        .statusCode(HttpStatus.FORBIDDEN.value());
}
```

We see that running this type of test gives us confidence that our API is behaving correctly. This is what we are checking for primarily, the behaviour of your components is appropriate.

If you are not going to do contract testing, which we will see later on in this chapter “**Contract testing**”, then it is important that you check the response body of your endpoints. This is going to give you a faster feedback loop to validate the response than if done as part of end-to-end testing. You will need to ensure that the response fits with the API specification that you have laid out. If your tests do not cover this and your response is different from the API Specification then users of your API will quickly become very frustrated.

## Integration tests

Integration tests in our definition are tests across boundaries between the module being developed and any external dependencies.

The Attendees API has a single other external services which is the PostgreSQL Database. For communicating with the PostgreSQL database

the Attendees API has a DAO component, which will be demonstrated at the end of this section (“[Testcontainers](#)”). We first need to examine how the Legacy Conference System will interact with the Attendees API, this will demonstrate the tests that can be used to interact with an external API.

It is important to clarify that when performing an integration test we only want to test our communication communicate with the external dependency, as normally it is not feasible to launch the whole external dependency. This can be seen in the case of the Legacy Conference System that wishes to have an integration test with the Attendees API. The Attendees Service has a dependency on a Database as well, so bringing up the whole system, (Legacy Conference System, Attendees API service and the database ) would not be an integration test but becomes an end-to-end test. All that the Legacy Conference System cares about when running integration tests with the Attendees API is if the requests it makes are correct and what the responses will be.

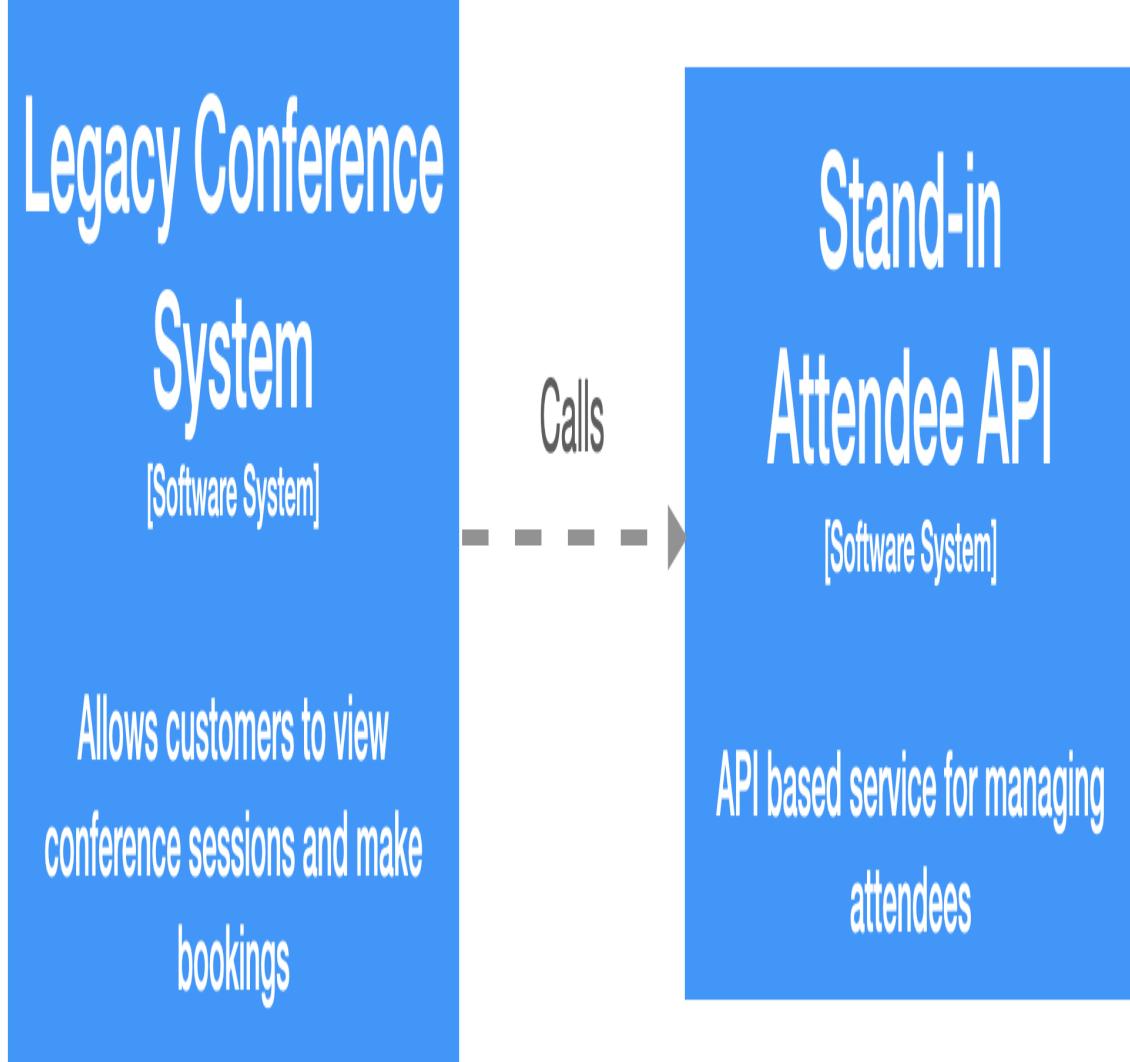


Figure 3-7. Legacy Conference System wanting to test the Attendees API

## External Service integration

When testing across an external boundary the validation that should be performed is to confirm that the communication across the boundary is correct. The types of things we want to verify that we have successfully integrated with an external service are the following:

- Ensuring that an interaction is being made correctly e.g. for a RESTful service this may be specifying the correct URL or the payload body is correct
- Is the data being sent in the correct format?
- Can the unit that is interacting with external service handle the responses that are being returned?

For the Legacy Conference System, it will need to verify it can communicate with the new the Attendees API.

So, how can we validate that the two services can talk to each other?

There are many options available to us, the simplest one though is to hand roll a stub service that mimics the requests and responses of the service we interact with (hence, you'll see why we label a *stand-in* API in the [Figure 3-7](#)). This is certainly a viable option as in your chosen language and framework it is usually very easy for developer to create a stub service with canned responses that integrate with tests. The key considerations when hand rolling a stub service is to make sure that the stub is accurate. It can be very easy to make mistakes, such as inaccurately portraying the URL or making mistakes in the response property names and values, can you see the errors in this hand typed response? <sup>6</sup>

```
{  
  "values": [  
    {  
      "id": 123456,  
      "givenName": "James",  
      "familyName": "Gough"  
    },  
    {  
    }
```

```
        "id": 123457,
        "givenName": "Matthew",
        "familyName": "Auburn"
    },
    {
        "id": 123456,
        "givenName": "Daniel",
        "familyName": "Bryant"
    }
]
```

This should still not put off the reader as this is still a good solution and one of the authors has hand rolled a stub server for an IDP to replicate logging into a service and this solution worked extremely well. We just want to highlight that care must be taken with hand-rolling stubs.

A way to avoid these inaccuracies and to ensure that requests to URLs are accurately captured along with the responses is use a recorder. It is possible to use a tool, that will record the requests and responses to an endpoint and generate files that can be used for stubbing. The following image shows how this works.

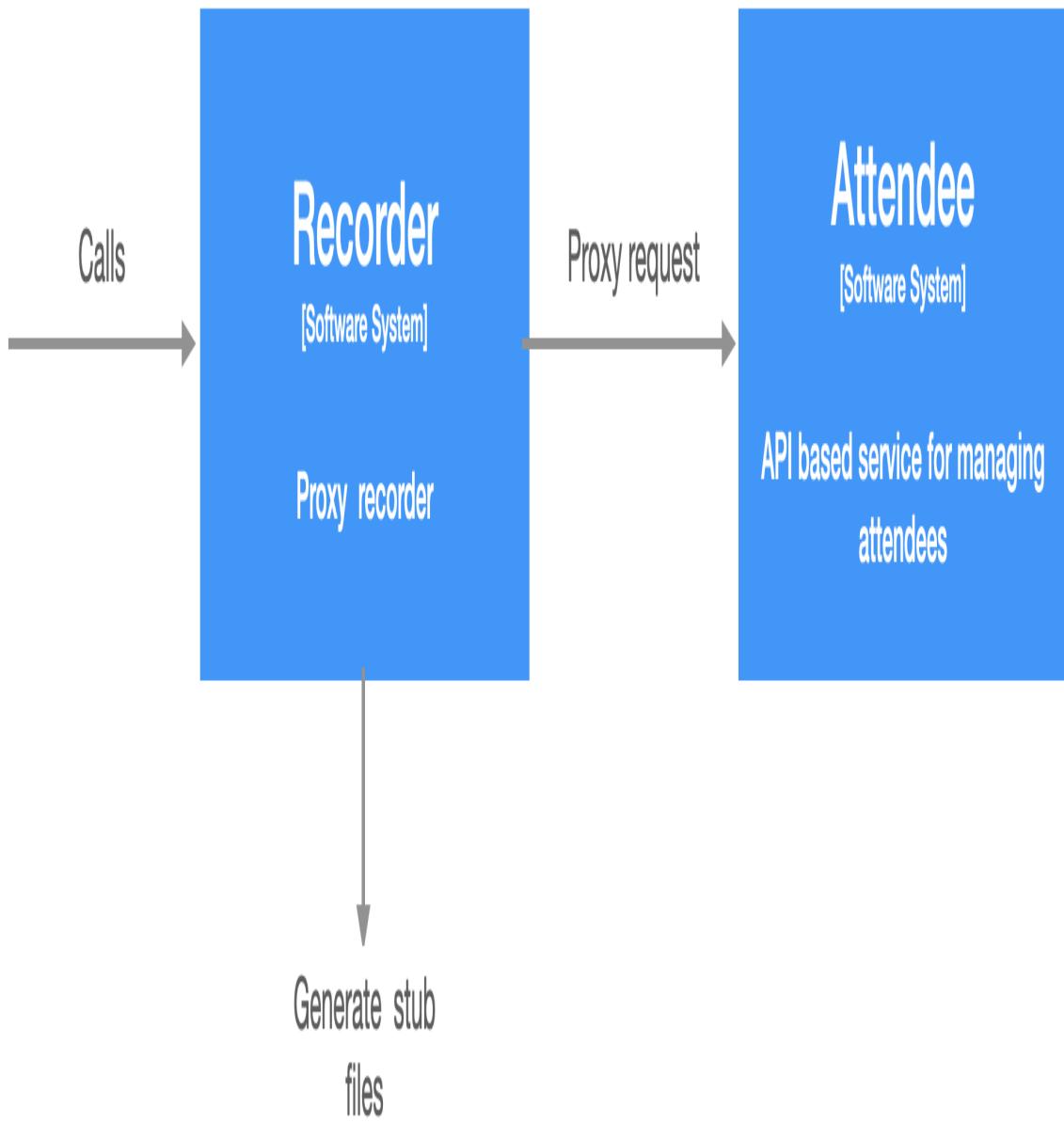


Figure 3-8. How a recorder works

These generated files are mappings that can then be used for tests to accurately portray requests and responses and as they are not hand rolled we have guarantees that they are accurate. To use these generated files, when a stub server is launched it reads the mappings files. When a request it

made to the stub server it looks at the requests and checks to see if it matched any of the mappings, if it matches then the mapped response will be returned. **Wiremock** is a tool that can do this and it can be used as a standalone service so is language agnostic, though as it is written in Java there are specific Java integrations that can be taken advantage of. There are many other tools available that have a similar capability in other languages such as **camouflage** which is written in Typescript. Recording calls to APIs for use of stubbing, in our opinion, this is usually a better way to go than hand rolling a stub as it is going to be more accurate than writing them.

What you really want is a golden source of defined interactions that can be used, and this is where contract testing can help. We'll discuss contract testing next, but first want to summarize the ADR guidelines for integration testing.

*T*

*a*

*b*

*l*

*e*

*3*

-

*l*

.

*A*

*D*

*R*

*G*

*u*

*i*

*d*

*e*

*li*

*n*

*e*

*s*

-

*I*

*n*

*t*

*e*

*g*

*r*

*a*

*ti*

*o*

*n*

*t*

*e  
s  
ti  
n  
g*

Decision      Should we add integration testing to our API Testing.

Discussion Points

If your API is integrating with any other service what level of integration test should you use?

- Do you feel confident that you can just mock responses and do not need to perform integration tests?
- For creating a stub service to test against, are we able to accurately craft the request and responses or should they be recorded?
- Will you be able to keep stub services up to date and recognise if an interaction is incorrect?

If your stubs are incorrect or become out of date, this means it is possible to have tests that pass against your stub service, however, when you deploy your service to interact with another API it can fail due to the service changing.

Recommendations      We do recommend having integration testing using recordings of interactions. Having integration tests that can be ran locally and give confidence in an

---

integration, especially when refactoring, helps to ensure that changes have not broken an integration.

---

## Contract testing

Contract testing is a written definition of an interaction between two entities a Consumer and a Producer and is about providing guarantees about how two parties can interact. The beauty of defining an interaction with a contract is that it is possible to generate tests and stubs from this.

Automated tests and stub services are what give us the ability to perform local testing without reaching out to the actual service, though still allow us to verify an interaction. This is really useful for APIs as that is a key requirement when publishing an API, you want it to keep working and not change by accident for user.

If we take the Attendee API service we want to ensure that the endpoint we mentioned earlier `/conference/{conference-id}/attendees`, continues to return the correct value. Contract testing is not the same as saying that an API conforms to a schema. A system is either compatible with a schema (like OpenAPI Spec,) or it is not, while a contract is about a defined interaction between parties and provide examples. A Matt Fellows has an excellent piece on this titled [Schema-based contract testing with JSON schemas and Open API \(Part 1\)](#).

There are only two roles for contract testing.

- A Consumer requests data from an API. e.g. web client, terminal shell.
- A Producer responds to the API requests, it is producing data. e.g. a RESTful Web Service.

Producers are also known as Providers. In this book the word Producer will be used, though the reader should be aware that they are interchangeable and provider is the term of choice in some contract frameworks. In the case

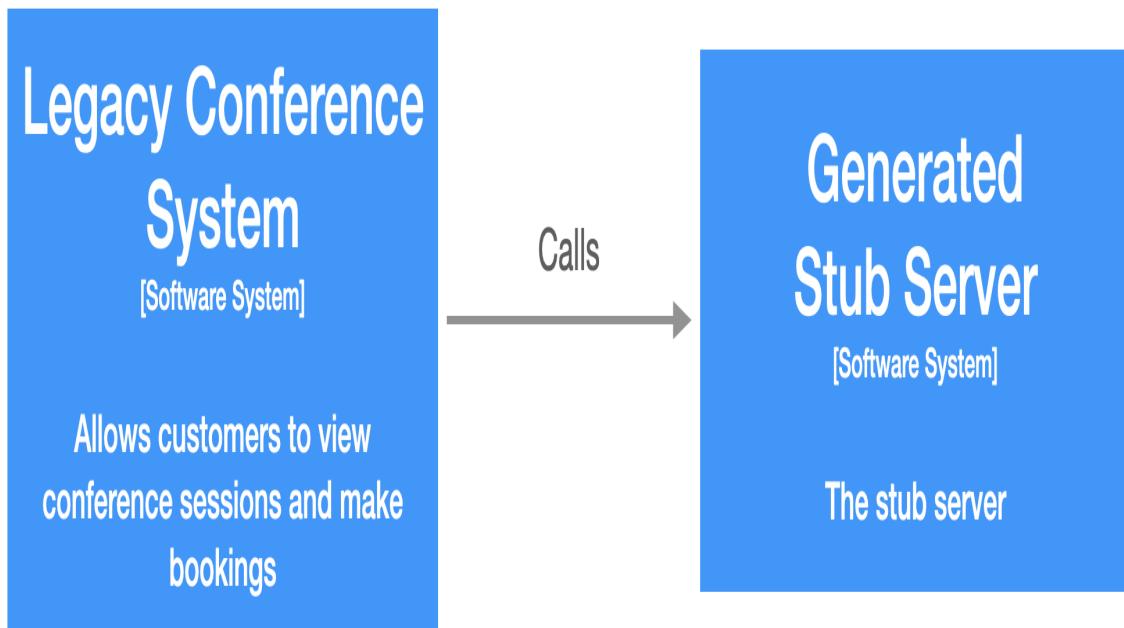
of the Attendee API service we say that it is a producer and the Legacy Conference System that is calling a consumer.

So how does a contract work? As was mentioned a contract is a shared definition of how two entities interact. Once a definition is defined of how two parties should interact then we can generate tests and stubs from this. If a contract is made for the following GET request to the url endpoint `/conference/{conference-id}/attendees` we state that we expect a response that has a property “value” that contains an array of values about the attendees. This is a sample definition of what an interaction looks like.

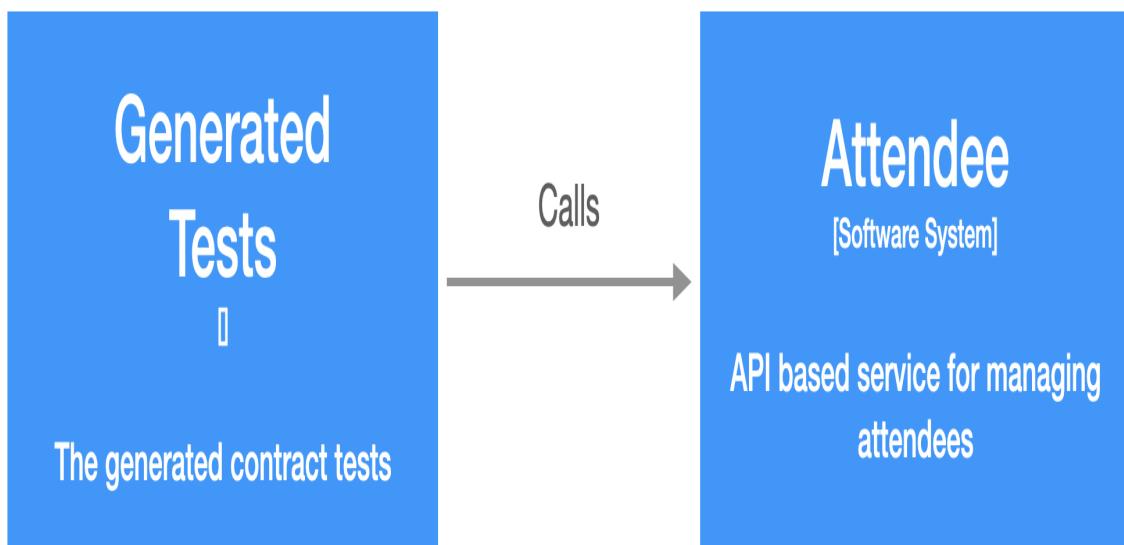
```
Contract.make {
    request {
        description("""Get a list of all the attendees at a conference""")
        method GET()
        url '/conference/1234/attendees'
        headers {
            contentType('application/json')
        }
    }
    response {
        status OK()
        headers {
            contentType('application/json')
        }
        body(
            value: [
                ${
                    id: 123456,
                    givenName: "James",
                    familyName: "Gough"
                },
                ${
                    id: 123457,
                    givenName: "Matthew",
                    familyName: "Auburn"
                }
            ]
        )
    }
}
```

From this we can generate a test for the producer. This is possible as we know when a request is made to the producer for at this endpoint we expect the response to match the response in the body and the content type. We can generate stubs as well, as we can have a stub server respond when it receives a request that matches the input and the response returned will be the response defined in the body. For a consumer this test allows the consumer to validate that it is making the right request to the producer and can handle the response. For the producer this test confirms that when the request is made that it can generate the correct response. We see this in the following image

# Generated stub for the consumer



# Test for the producer



*Figure 3-9. Generated stubs and tests*

## WARNING

It is tempting to use contracts for scenario tests. e.g. first add an attendee to a conference, second, use the get attendees to check the behavior is correct in the auto generated tests. Frameworks do support this though also discourage it. A producer should verify this type of behavior in component and unit tests and not in the Contracts, Contracts are for testing interactions.

A key benefit of using contracts is that once it is confirmed the contract will be fulfilled then this decouples the dependency of building. The consumer has a stub service to develop against and the producer has tests to ensure that they are building the right interaction. The contract test process saves time as when both the consumer and producer are deployed they should integrate seamlessly.

## TIP

We have used generated stub servers to run demos for stakeholders. This was useful as the producer for a new API was still being developed but had agreed contracts.

In order to understand how we generate our contracts, let's look at the two main contract methodologies.

### Producer Contracts

Producer contract testing is when a producer defines its own contracts. If you are just starting with contract testing or wanting to introduce them into your APIs then this is a great place to start. New contracts can be created for the producer to ensure that that service fulfils the interaction criteria and will continue to fulfil them. If consumers are complaining that a producers API is breaking on new releases then introducing contracts can help with this issue.

The other common reason for using producer contract testing is when there is a large audience for your API, such as when an API is being used outside your immediate organization and by unknown users, i.e. external third parties. When developing an API that has a wide public audience it will need to maintain its integrity and though it is something that is updated and improved, immediate feedback and individual feedback will not be applied.

A concrete example of such an API in the real world is the Microsoft Graph API, it can be used to look at the users registered in an Active Directory Tenant. A consumer may find it preferable to have an additional property for adding in a preferred name on the **Users** endpoint of the API. This is something that can be messaged to Microsoft as a suggestion, however, this is not likely to be changed and if the suggestion was seen as a good idea it would unlikely happen quickly. This is something that would have to be weighed up and considered. Is this something that will be useful for others, is it backwards compatible change, how does this change the interactions across the rest of the service?

If we take the same approach with our Attendee API and say that we want to make it available for public consumption, we do not want consumers suggested changes to the API that only benefits them. What we care more about in this case is using the contract to ensure that our interactions and the data returned are consistent. Breaking changes involves versioning which we read about in “[API Versioning](#)”

## Consumer Driven Contracts

Consumer Driven Contracts (CDC) is when a consumer drives the functionality that they wish to see. Consumers submit contracts to the producer for new API functionality and the producer will choose to accept or reject the contract.

CDC is very much an interactive and social process. The owners of the applications that are consumers and producers should be within reach, e.g. in the same organization as one another. When a consumer would like a new interaction (e.g., API call,) or have an interaction updated (e.g., a new property added,) then they submit a request for that feature. In our case, this

may mean that a pull request is submitted from the Legacy Conference System to the new Attendee API service. The request for the new interaction is then reviewed and discussion takes place about this new functionality to ensure that this is something that the Attendees service should and will fulfil and that the contract is correct. For example, let's say a contract is suggested for a PUT request, a discussion can take place if this should be in fact a PATCH request. This is where a good part of the value of contracts comes from, this discussion for both parties about what the problem is, and using a contract to assert that this is what the two parties accept and agree to. Once the contract is agreed, the producer (Attendees service) accepts the contract as part of the project and can start fulfilling it.

## **Contracts methodology overview**

These methodologies should hopefully give an overview of how to use contracts as part of the development process. This should not be taken as gospel as variations do exist on the exact steps. For example, one process can ask that the consumer when writing the contract should also create a basic implementation of the producer code to fulfil the contract. In another example, the consumer should TDD the functionality they require and then create the contract before submitting the pull request. The exact process that is put in place may vary by team. Once you understand the core concepts and patterns of CDC the exact process that is put in place is just an implementation detail.

If starting out on a journey to add contracts it should be noted that there is a cost to it. This cost is the setup time to incorporate contracts into a project and also the cost of writing the contracts. It is worth looking at tooling that can create contracts for you based on an OpenAPI Specification.<sup>7</sup>

*T*

*a*

*b*

*l*

*e*

*3*

-

*2*

.

*A*

*D*

*R*

*G*

*u*

*i*

*d*

*e*

*li*

*n*

*e*

*s*

-

*C*

*o*

*n*

*t*

*r*

*a*

*c*

*t*

*T*

*e*

*s*

*ti  
n  
g*

Decision      When building an API should you use contract testings and if so should you use Consumer Driven Contracts or Producer Contracts?

Discussion Points

Determine whether you are ready to include contract testing as part of your API testing.

•

Do you want to add an extra layer of testing to your API that developers will be required to learn about?

If contracts have not been used before then it requires time to decide how you will use them.

•

Should contracts be centralized or in a project?

•

Do we need to bring in additional tools and training to help people with contracts?

If deciding to use contracts, then which methodology should be used? CDC or Producer Contracts.

- Do you know who will use this API?
- Will this API be used just within our organization?
- Does the API have consumers that are willing to engage with us to help drive our functionality?

Recommendations	<p>We recommend using contract testing when building an API. Even if there is a developer learning curve and deciding how you are going to set up your contracts for the first time, we believe it is worth the effort. Defined interactions that are tested save so much time when integrating services together.</p> <p>If you are exposing your API to a large external audience then it is important to use Producer contracts. Again having defined interactions that help ensure that your API does not break compatibility is crucial.</p> <p>If you're building an internal API then the ideal is to work towards CDC, even if you have to start with producer contracts and evolve over to CDC.</p>
-----------------	--

## Contracts testing frameworks

The likelihood is that when it comes to contract testing frameworks for HTTP you are best to start looking at **PACT**. PACT has evolved into the default contract testing framework due to the ecosystem that has been built around it and the sheer number of languages it supports. There are other

contract testing frameworks available, and they can be opinionated. PACT is opinionated that you should perform CDC and is specifically designed for that. A test is written by a consumer and that test generates a contract, this is an intermediate representation of what an interaction should be. It is language agnostic, which is why PACT has such wide language usage. Contract generation certainly is a good way to go and means that when using PACT you have no real need to ever see a contract as it is a language agnostic representation of the interaction. There are other frameworks that have differing opinions, an example is Spring Cloud Contracts that does not have a strong opinion of CDC or provider contracts, either can be achieved. This is possible as with Spring Cloud Contracts you write the contracts by hand opposed to having it generated. Though Spring Cloud Contracts is language agnostic by using a containerized version of the product, to get the most out of it you need to really be using the Spring and JVM ecosystem.<sup>8</sup>

There are options for contract testing for other protocols, contract testing is not exclusively for HTTP communications.

## API Contracts Storage and Publishing

Having seen how contracts work and methodologies of incorporating them to the development process the next question becomes where are contracts stored and how should they be published.

There are a few options for storing and publishing contracts and these again depend on the setup that is available to you and your organization.

Contracts can be stored alongside the Producer code in version control (e.g. git). They can also be published alongside your build into an artifact repository such as [Artifactory](#).

Ultimately the contracts should be easily available for the producer to auto generate tests and the consumer to generate stubs and be able to submit new contracts for the project. The Producer should have control over which contracts are accepted in the project and can ensure that undesired changes aren't made or additional contracts are added. The downside to this

approach is that in a large organization it can be difficult to find all the API services that use contracts.

Another option is to store all the contracts in a centralized location to enable visibility into other API interactions that are available. This central location could be a git repository, though if well organized could also be a folder structure. The downside to this approach is that unless organized and setup correctly it is possible and likely that contracts get pushed into a module that the producer has no intention on fulfilling.

Yet another option for storing contracts is to use a broker. The PACT contract framework has a **broker product** that can be used as a central location to host contracts. A broker can show all contracts that have been validated by the producer as the producer will publish those contracts that have been fulfilled. A broker can also see who is using a contract to produce a network diagram, integrate with CI/CD pipelines and provide even more valuable information. This is the most comprehensive solution available and if you use a framework that is compatible with the PACT Broker then it is recommended.

## Testing during the building of your application

As we learned in Chapter 2 “[OpenAPI Specification and Versioning](#)”, having an OpenAPI specification informs a user of your API how it will work. We update our OpenAPI spec with changes made to our API to ensure that we do not break backwards compatibility. Using tools such as OpenAPI-Diff as part of your build offers the opportunity to validate whether your API is still backwards compatible. This check can be integrated into your build pipeline to diff the current version and the deployed version of your API. OpenAPI diff tools will give a pass or fail to declare if a schema is backwards compatible or not. If a comparison fails then we know that we have broken the API Specification and can fix this early in our development. This is much better than releasing a change and a customer finding that we have broken the API spec.

At this point we may ask if we need to diff our OpenAPI specs as we have contract tests. This is a valid point, if your contract tests are well-defined then this might certainly be enough, however, you are putting more responsibility on your contract tests to catch breaking changes. Adding a diff tool to the build pipeline provides a nice clear-cut yes or no if your API spec is backwards compatible.

## Testcontainers

We have seen how we can communicate with other APIs and use contracts to help our development. Another useful practice to aid verifying communications across a boundary is to perform testing with a deployable solution.

The Attendees API service will provide gRPC interface as well as the RESTful interface we have been developing. The gRPC interface is to be developed after the RESTful interface, but we want to provide the developers with a gRPC interface to test against. The decision is made to provide a smart stub service, this will be a stub that provides a few pre-canned responses and stores information in memory. To achieve this goal a barebones' application is made that fulfills this objective.<sup>9</sup> This gRPC stub is then packaged up, containerized and published. This stub can be now used by developers for testing across a boundary, i.e. they can make real calls to this service in their tests, and this containerized stub can run locally on their machine. The difficult thing with this is that the developer will need to orchestrate the lifecycle of this container for their test, which is non-trivial.

This is where **Testcontainers** can help. Testcontainers is a library that integrates with your testing framework to orchestrate containers. Testcontainers will start and stop and generally organise the lifecycle of containers you use with your tests, in our case the stub gRPC server. It can be beneficial to use a containerized solution as it provides a simple and common method that others can use to deploy it.

So why is this important to the development of APIs? Having libraries and tools that run services locally are extremely valuable when testing across boundaries for your API. Being able to package stubs, or whole services, and distribute them over a common mechanism, like Docker, makes it easier to perform testing.

## Other containerized services

We have seen that the Attendee API service has a connection to an external database, and so it is important to perform an integration test. The options for testing integration boundaries for a database would be mocking out the Database, use an inMemory Database e.g. H2 or run a local version of the database using Testcontainers. Using a real instance of the database in your test provides a lot of value as with mocks you can mock the wrong return value or make an incorrect assumption. With an inmemory DB you are assuming that the implementation matches the real DB and when you use H2 enough you find that it normally does not. Using a real instance of the product and the same version that you run in production means that you get reliable testing across a boundary to be sure your integration will work when going live.

We see in [Figure 3-10](#) how we would structure the test to confirm proper integration across a boundary with a database.

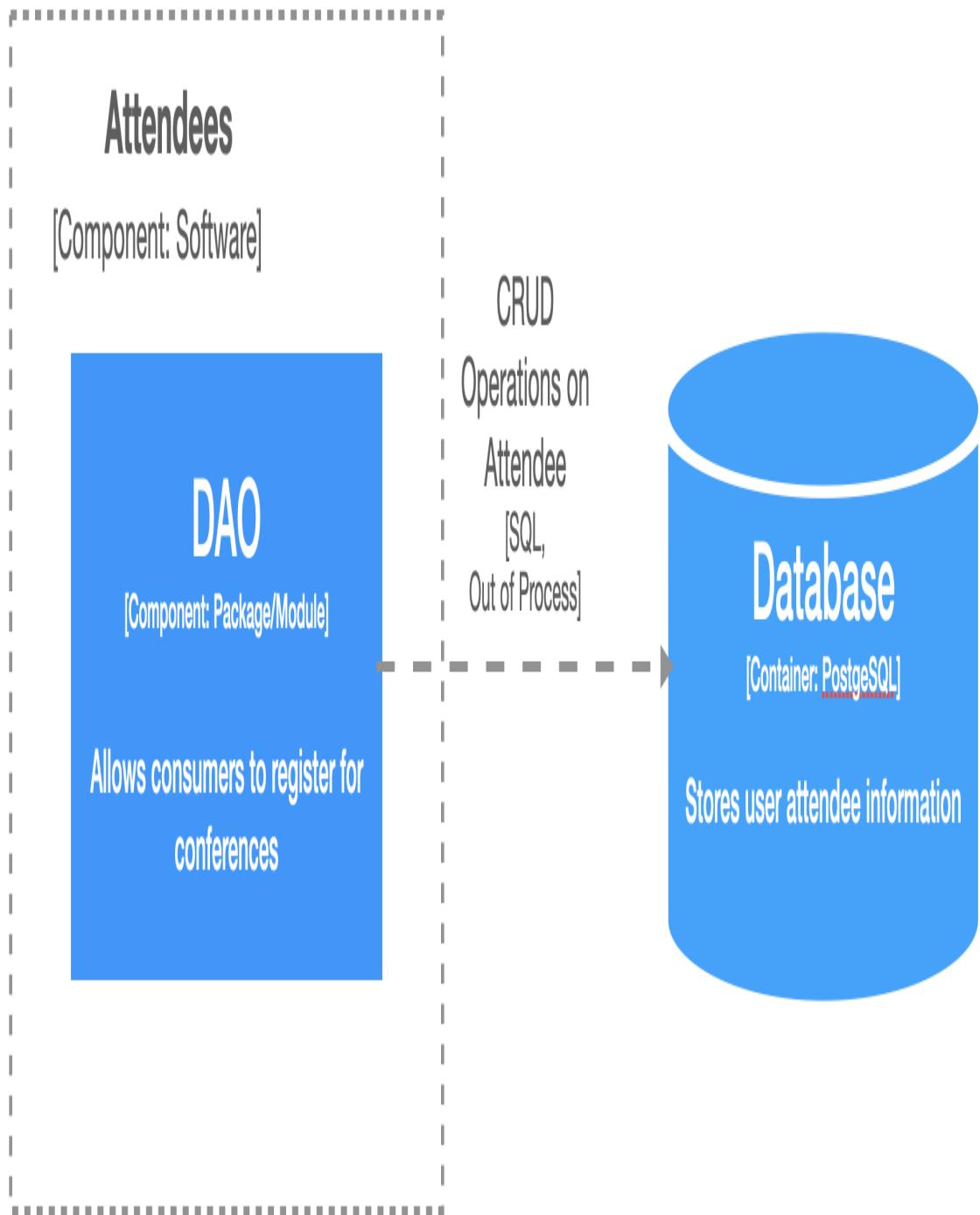


Figure 3-10. Testcontainers DAO test

Testcontainers is a powerful tool and should be considered when testing boundaries between any external services. Other common external services that benefit from using Testcontainers include Kafka, Redis, NGINX.

You may be asking if this type of testing is considered end-to-end testing as we have a real instance of another service that we are testing against. This is a valid question, for integration tests we are concerned about the interactions across the boundaries and not what the external service is doing. We want to call the service and get a response that is the piece we want to check, we are not performing tests to confirm that the services are working together to validate any behaviour within both services. That is why the boundary of what you are testing matters, so in this case we are not end-to-end testing, though it is possible and shall cover in the next section “End-to-end Testing”

## End-to-end Testing

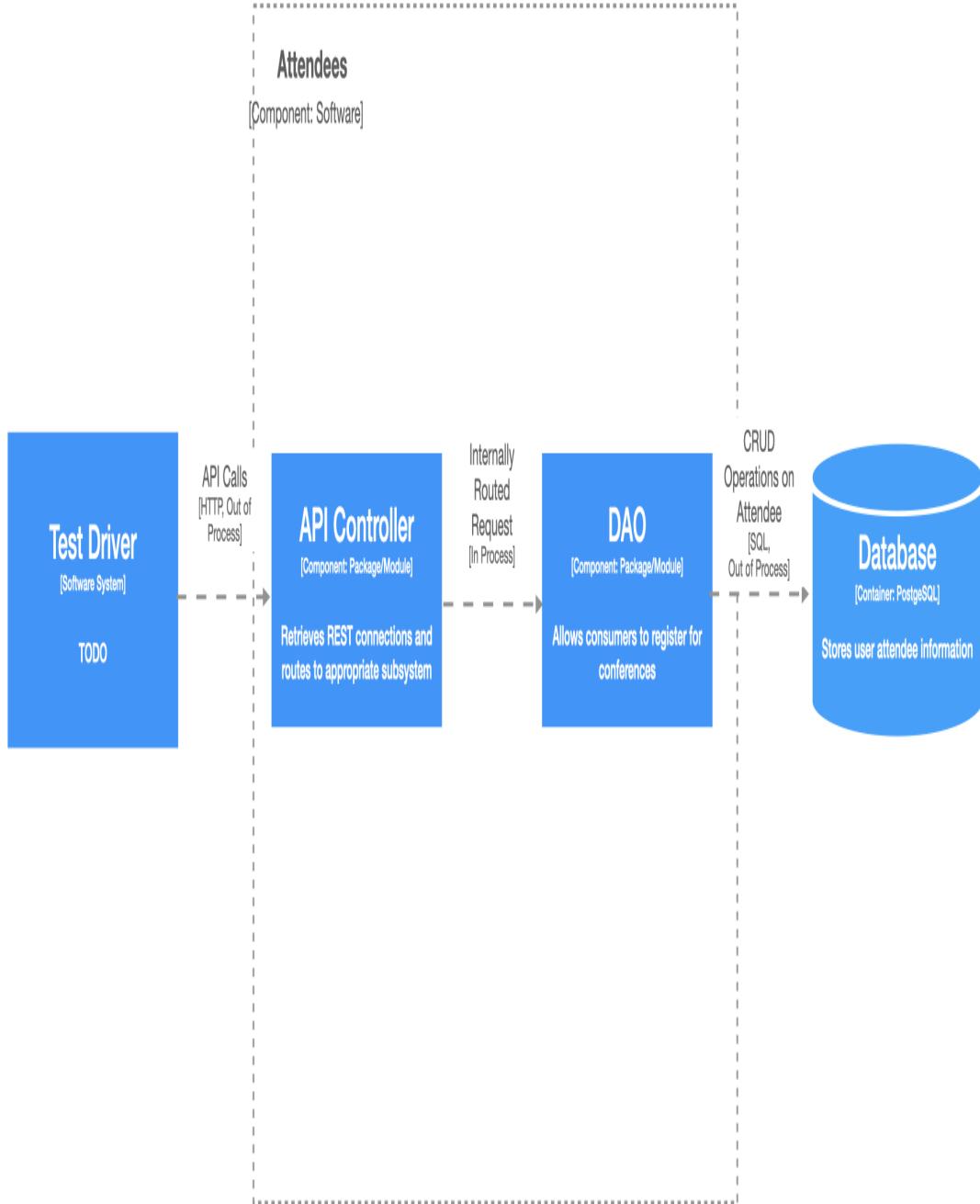
The essence of end-to-end testing is to test services and their dependencies together to verify it works as expected. It is important to validate that when a request is made, and it reaches the front door (i.e. the request reaches your infrastructure), the request flows all the way through and the correct response is given. This validation gives confidence that all these systems work together as expected. We see these tests as being part of Q2 of the testing quadrant, these are automated tests that test core cases and validate business functionality. As we are building and exposing APIs we care only about the API, we want to validate this as our public interface. We do not care about what is calling the API, we should not need to validate the Mobile Apps or Web-UIs’ that are communicating with your API as they are not part of our domain. <sup>10</sup>.

## WARNING

If you are building an external facing API and you have multiple third parties that are consuming it, don't try and copy the third party UI to replicate how their UI works. Doing so will mean that huge amounts of time will be spent trying to replicate something out of your domain.

We do want to clarify that the ideal is to have real versions of *your* services running and interacting together, however, sometimes this is not always feasible. Therefore, it is okay to stub out some entities of a system that are outside your organizations domain and are provided by an external party. Such as AWS S3. You really do not want to run your tests and have to rely on an external entity where you are connecting over the internet and the network can go down. This is why for our end-to-end tests we need to ensure that we define our test boundary, we also do not want to launch services that are not needed. For an end-to-end test of our Attendee service, we need to launch the database and the attendee service, we do not need to launch the whole conference system. We have defined a boundary here, we could have an end-to-end test that includes the Legacy Conference System, this would also be valid for an end-to-end.

# Chapter 3 - End to End Testing



*Figure 3-11. End-to-end test scope*

Managing and co-ordinating multiple systems together is not easy to automate and end-to-end tests can be brittle. Running end-to-end tests locally is becoming easier, as we have just seen in “**Testcontainers**” with containerization which allow you to spin up multiple systems locally. However, the process to launch is slow and these tests are still difficult to maintain, hence they are at the top of the testing pyramid.

One of the other parts of end-to-end testing that provides value is using realistic payloads. This is extremely valuable, we have seen end-to-end tests where payloads have been created that are small and concise, then when investigating why APIs are breaking it is found that the Clients are regularly sending very large payloads, larger than the buffers support. Though we did say that we should not care what is interacting with our API, you still need to know who the consumers are. If a WebUI is using your API are you needing to support requests for **CORS** or **CSRF**?

As we look to start writing our end-to-end tests, a key point is to ensure that you are only testing core user journeys and not testing edge cases.

This is where **Behavior Driven Development (BDD)** can be used as a nice way to write your user stories as part of your business facing tests. This may mean testing multiple endpoints to show a core user journey.

## **Can you test too much?**

While it is recommended more time be spent on writing tests over business logic there is a point where too much testing can be done. If too much time is spent working on tests then the module will never be delivered. A balance is always needed to have a good test coverage and provide confidence. Being smart about what should be tested and what is irrelevant, such as creating tests that duplicate scenarios, is a waste of resources.

An Architect must be able to recognize where the boundary is for excessive testing for an API/module/application, as its value for customers and the

business is only realized when running in production. What we are implying is that not all of this testing needs to be implemented to be able to release an API, it may not be feasible due to time constraints and business demands. Though what we do want the reader to take away is that ideally you would be able to incorporate a lot of this testing at all these levels to give yourself the best API you can.

## Summary

In this chapter we have covered the core types of testing for APIs, including what should be tested and where time should be dedicated. Key takeaways are as follows:

- Stick to the fundamentals of the Unit Testing and TDD for the core of your API.
- Perform service tests on your component and isolate the integrations to validate incoming and outgoing traffic.
- Contract testing can help you develop a consistent API and test with other APIs.
- Using end-to-end tests replicate core user journeys to help validate your APIs all integrate correctly.
- Use the ADR Guidelines as a way to work out if you should add different tests to your API.

While we've given you lots of information, ideas and techniques for building a quality API, this is by no means an exhaustive list of tools available. We encourage you to do some research on testing frameworks and libraries that you may want to use, to ensure you are making an informed decision.

However, no matter how much testing is done upfront, nothing is as good as seeing how an application actually runs in production. To learn more about testing in production refer next to Chapter 6

- 
- 1 The author's friend owns a mouthguard company. One of the authors was on the receiving end of hearing about the arduous process for testing the integrity of the product. No one wants a mouthguard where the only testing takes place during the match!
  - 2 SLOs and SLIs will be discussed in more detail in Chapter 6
  - 3 To learn more on agile testing check out the books [Agile Testing](#), [More Agile Testing](#) or the video series [Agile Testing Essentials](#)
  - 4 A Unit need not be defined exactly as a function or a class, however, you will need to define a unit for your case.
  - 5 Historical fact: Stubbing is known as the Classicist approach or the Chicago School and Mocking is known as the Mockist approach or also the London School.
  - 6 replicated id and familyNane
  - 7 At the time of writing there are a few projects that are available, though none are actively maintained, so it is difficult to recommend any.
  - 8 PACT does a good job of comparing itself to other contract frameworks  
[https://docs.pact.io/getting\\_started/comparisons](https://docs.pact.io/getting_started/comparisons)
  - 9 Please refer to the book [github](#) for an implementation example
  - 10 If these are part of the domain, and you are responsible for them then these should be considered for your end-to-end testing. Martin Fowler discusses testing below the UI, these are known as Subcutaneous Tests are are useful “[when you want to test end-to-end behavior, but it's difficult to test through the UI itself](#)”

# Part III. Traffic Patterns and Management

---

## A NOTE FOR EARLY RELEASE READERS

With Early Release ebooks, you get books in their earliest form—the authors’ raw and unedited content as they write—so you can take advantage of these technologies long before the official release of these titles.

If you have comments about how we might improve the content and/or examples in this book, or if you notice missing material within this chapter, please reach out to the author at [feedback@mastering-api.com](mailto:feedback@mastering-api.com).

This section explores how API traffic is managed. This includes both traffic originating externally from end users that is entering (ingressing) into the system and traffic originating internally from services that is travelling across (service-to-service) the system.

The first chapter in this section, and where most API architects should start their traffic management journey, examines API gateway technology for managing ingress, or “north-south”, traffic. You will learn about core motivations, patterns, and functionality provided by an API gateway, alongside developing an understanding of the larger context of where this technology sits within a typical edge stack, which often comprises of CDNs, firewalls, and load balancers. You will explore how architectural choices (monolith, microservices etc) impact the role of an API gateway and learn key questions to ask when selecting your technology choice.

Although service mesh technology is relatively new, the functionality provided for managing service-to-service, or east west, traffic is compelling for many API architects to explore. The next chapter in this section examines the service mesh pattern, and you will explore the fundamental differences between this and API gateways (and older technology, such as an ESB). You will learn about the core traffic routing, observability, and security functionality, and explore the emerging service mesh interface (SMI) specification. You will wrap up the chapter with a speculative look into the future as to where this technology is heading. It is the job of an API architect to not only choose technologies based on current requirements, but also to understand the longer-term impact and opportunities.

[Chapter 4](#)

[Chapter 5](#)

# Chapter 4. API Gateways

---

## A NOTE FOR EARLY RELEASE READERS

With Early Release ebooks, you get books in their earliest form—the authors’ raw and unedited content as they write—so you can take advantage of these technologies long before the official release of these titles.

This will be Chapter 4 of the final book.

If you have comments about how we might improve the content and/or examples in this book, or if you notice missing material within this chapter, please reach out to the author at [feedback@mastering-api.com](mailto:feedback@mastering-api.com).

Now that you have a good understanding of the life of an API, the protocols involved, and how to begin testing, we can turn our attention to platforms and tooling that are responsible for delivering APIs to end users in production. An API gateway is a critical part of any modern technology stack, sitting at the “edge” of systems and acting as a management tool that mediates between a client and a collection of backend services.

In this chapter you will learn about the “why?”, “what, and “where” of API gateways, and also explore the history of the API gateway and other edge technologies. In order to build upon the success of previous generations of engineers and developers, it is important that you understand the evolution of this technology and examine how this relates to important architectural concepts like coupling and cohesion.

You will also explore the taxonomy of API gateways, and learn how these fit into the bigger picture of system architecture and deployment models. You will revisit your earlier exploration of ingress (“north-south”) traffic and service-to-service (“east-west”) traffic, and explore the technologies that can be used to manage each traffic type. You will also explore the

challenges and potential pitfalls when mixing traffic patterns and technologies.

An API gateway is also a very useful tool for migrating and evolving systems, as it can act as a selective router between old and new systems and provide a level of abstraction or encapsulation between clients and backend systems. You'll explore these concepts further by evolving the Conference System case study.

Building on all of the topics above, you will conclude the chapter by learning how to select an appropriate API gateway based on your requirements, constraints, and use cases.

## Why Use an API Gateway?

A big part of the modern software architect's role is asking the hard questions about design and implementation. This is no different when dealing with APIs and traffic management and related technologies. You need to balance both short term implementation and long term maintainability. There are many API-related cross-cutting concerns that you might have, including: maintainability, extensibility, security, observability, product lifecycle management, and monetization. An API gateway can help with all of these!

This section of the chapter will provide you with an overview of the key problems that an API gateway can address, such as:

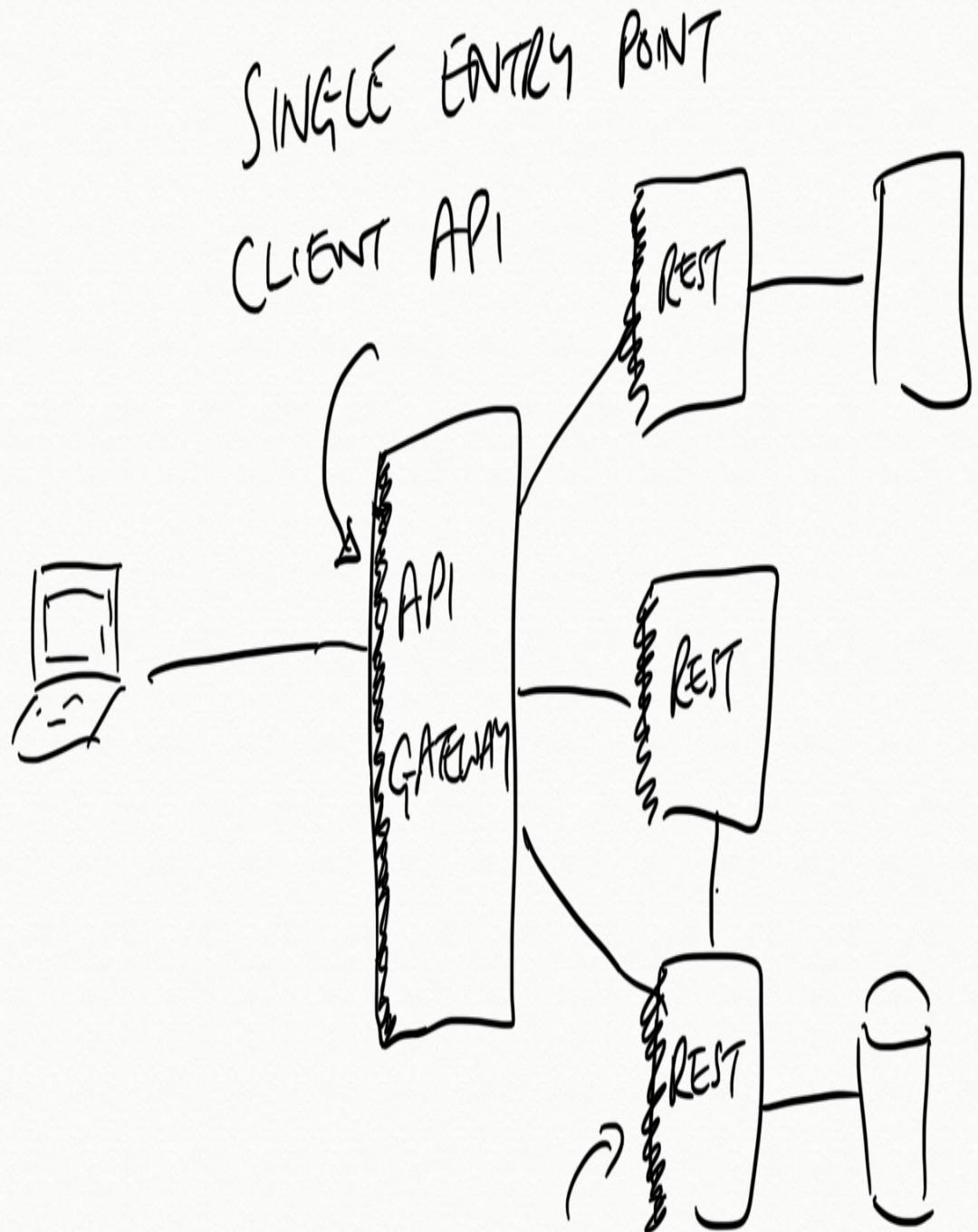
- Reduce Coupling: Adapter / Facade Between Front Ends and Back Ends
- Simplify Consumption: Aggregating / Translating Back End Services
- Protect APIs from Overuse and Abuse: Threat Detection and Mitigation
- Understand How APIs Are Being Consumed: Observability

- Manage APIs as Products: API Lifecycle Management
- Monetize APIs: Account Management, Billing, and Payment

## **Reduce Coupling: Adapter / Facade Between Front Ends and Back Ends**

Two fundamental concepts that every software architect learns about early in their career are coupling and cohesion. You are taught that systems that are designed to exhibit loose coupling and high cohesion will be easier to understand, maintain and modify. Loose coupling allows different implementations to be swapped in easily, and internals to be modified without experiencing unintended cascading effects on surrounding modules or systems. High cohesion promotes understandability — i.e all code in a module or system supports a central purpose — and reliability and reusability. In our experience, APIs are often the locations in a system in which the architectural theory meets the reality; an API is quite literally and figuratively an interface that other engineers integrate with.

An API gateway can act as a single entry point and a **facade** or an **adapter**, and hence promote loose coupling and cohesion. A facade defines a new simpler interface for a system, whereas an adapter re-uses an old interface with the goals of supporting interoperability between two existing interfaces. Clients integrate with the API exposed at the gateway, which, providing the agreed upon contract is maintained, allows components at the backend to change location, architecture, and implementation (language, framework etc) with minimal impact.



BACKEND APIs CAN CHANGE  
LOCATION, ARCHITECTURE, IMPL

*Figure 4-1. An API gateway providing a facade between front ends and back ends*

## Simplify Consumption: Aggregating / Translating Back End Services

Building on the discussion of coupling in the previous section, it is often the case that the API you want to expose to the front end systems is different than the current interface provided by a back end or composition of backend systems. For example, you may want to aggregate the APIs of several back end services that are owned by multiple owners into a single client-facing API in order to simplify the mental model for front end engineers, streamline data management, or hide the back end architecture. [GraphQL](#) is often used for exactly these reasons.

### ORCHESTRATING CONCURRENT API CALLS

A popular simplification approach implemented in API gateways is orchestrating concurrent backend API calls. This is where the gateway orchestrates and coordinates the concurrent calling of multiple independent backend APIs. Typically you want to call multiple independent and non-coupled APIs in parallel rather than sequentially in order to save time when gathering results for the client. Providing this in the gateway removes the need to independently implement this functionality in each of the clients.

It is also common within an enterprise context that some protocol translation will be required. For example, you may have several “heritage” (money making) systems that provide SOAP-based APIs, but you only want to expose REST-like APIs to clients. Or your legacy systems may only support HTTP/1.1, but clients require HTTP/2 connections. Some organizations may implement all internal service APIs via gRPC and Protobuf, but want to expose external APIs using HTTP and JSON. The list goes on; the key point here is that some level of aggregation and translation is often required to meet externals requirement or provide further loose coupling between systems.

An API gateway can provide this aggregation and translation functionality. Implementations vary and can be as simple as exposing a single route and composing together (“mashing”) the responses from the associated multiple internal system, or providing a protocol upgrade from HTTP/1.1 to HTTP/2, all the way through to mapping individual elements from an internal API to a completely new external format and protocol.

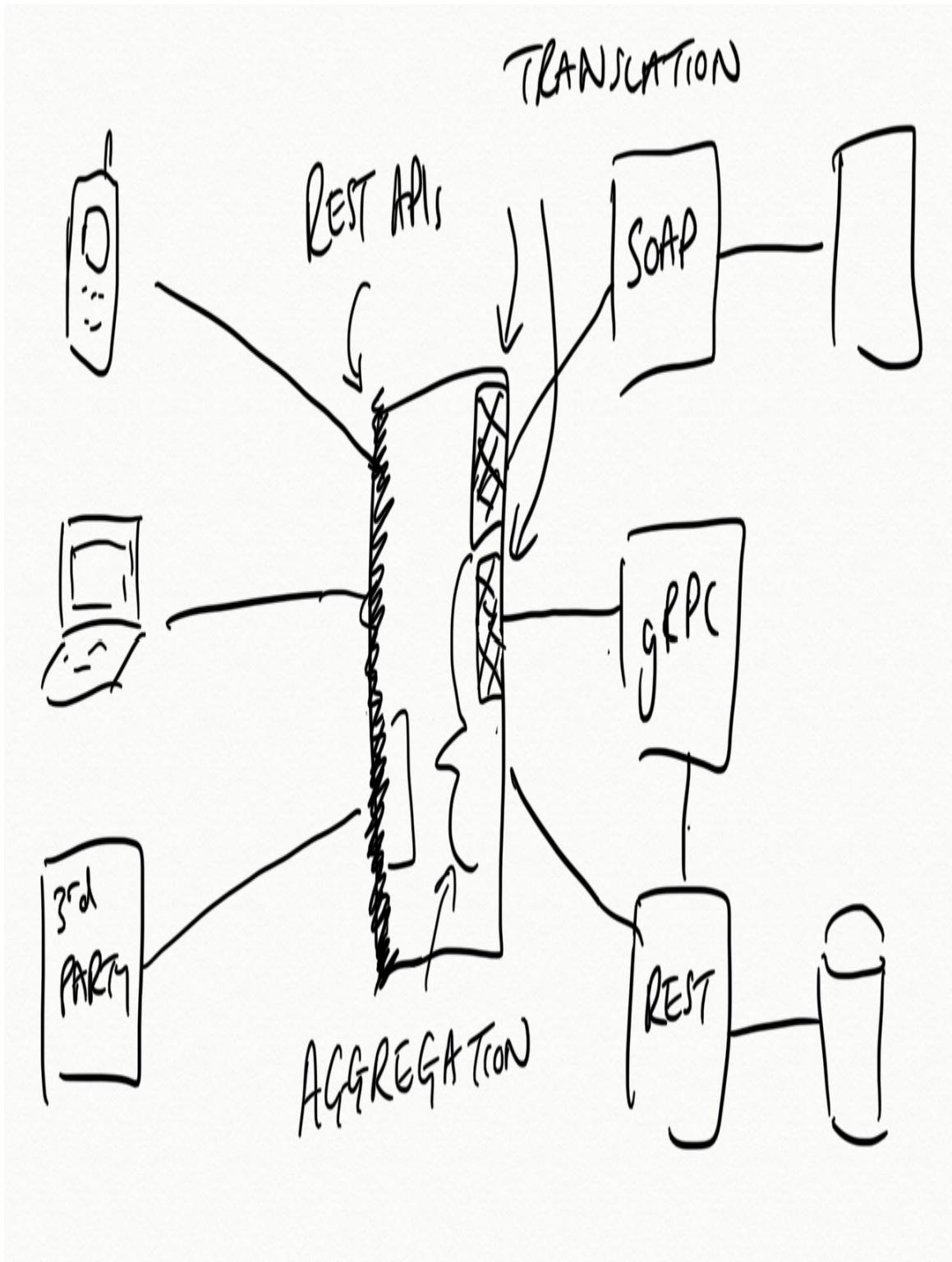
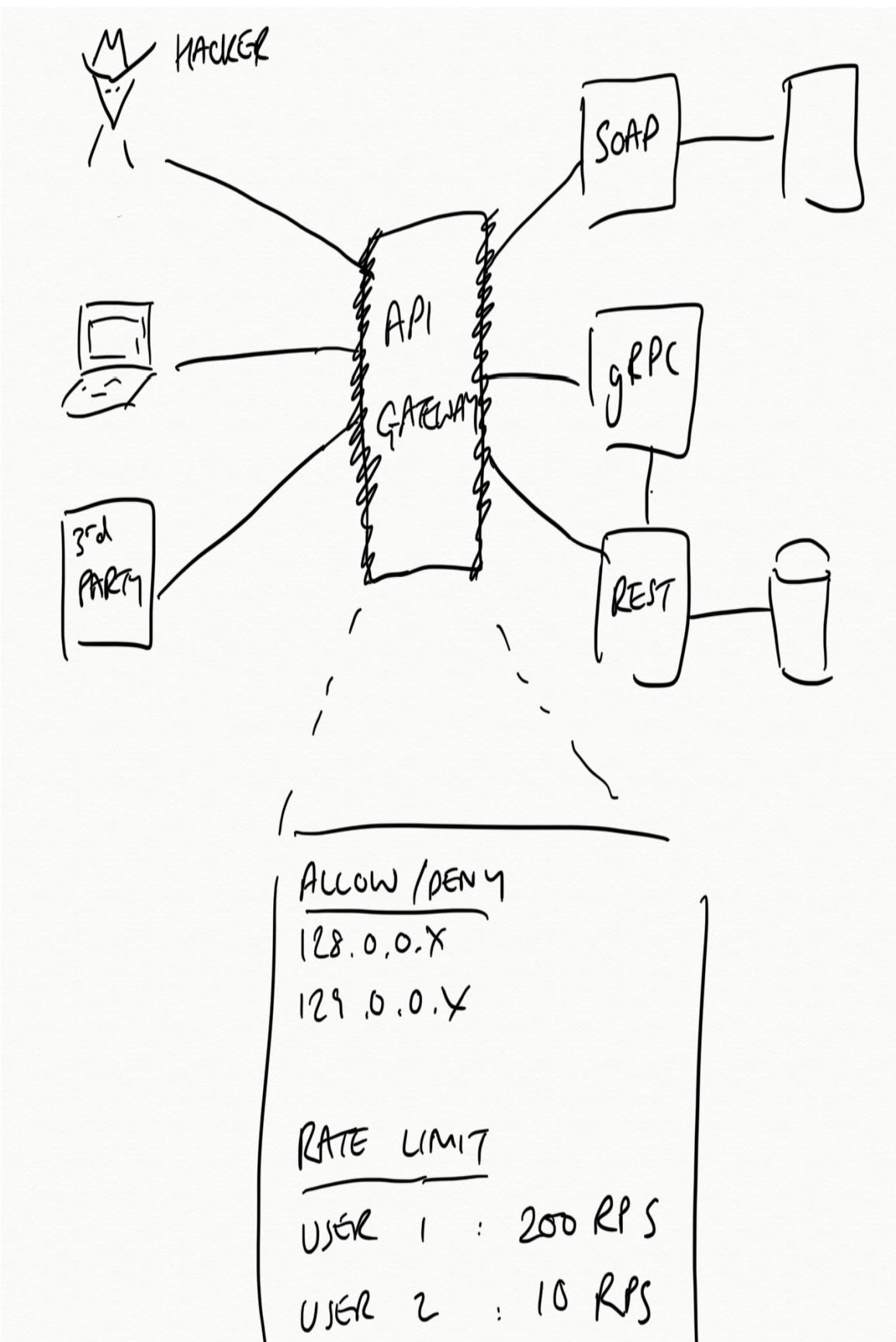


Figure 4-2. An API gateway providing aggregation and translation

## **Protect APIs from Overuse and Abuse: Threat Detection and Mitigation**

The edge of a system is where your users first interact with your applications. It is also often the point where bad actors and hackers first encounter your systems. Although the vast majority of enterprise organizations will have multiple security-focused layers to their edge stack, such as a content delivery network (CDN) and web application firewall (WAF), and even a perimeter network and dedicated demilitarised zone (DMZ), for many smaller organizations the API gateway can be the first line of defense. For this reason many API gateways include security-focused functionality, such as TLS termination, authentication/authorization, IP allow/deny lists, WAFs (either inbuilt or via external integration), and rate limiting and load shedding.



*Figure 4-3. API gateway overuse and abuse*

A big part of this functionality is the capability to detect API abuse, either accidental or deliberate, and for this you will need to implement a comprehensive observability strategy.

## **Understand How APIs Are Being Consumed: Observability**

Understanding how systems and applications are performing is vitally important for ensuring business goals are being met and that customer requirements are being satisfied. It is increasingly common to measure business objectives via key performance indicators (KPIs), such as customer conversion, revenue per hour, stream starts per second etc.

Infrastructure and platforms are typically observed through the lens of **service level indicators (SLIs)**, such as latency, errors, queue depth etc. As the vast majority (if not all) of user requests flow through the edge of a system, this is a vital point for observability. It is an ideal location to capture top-line metrics, such as the number of errors, throughput, and latency, and it is also a key location for identifying and annotating requests (potentially with application-specific metadata) that flow throughout the system further upstream. Correlation identifiers (such as **OpenZipkin b3** headers) are typically injected into a request via the API gateway and are then propagated by each upstream service. These identifiers can then be used to correlate log entries and request traces across services and systems.

Although the emitting and collecting of observability data is important at the system-level, you will also need to think carefully how to process, analyse, and interpret this data into actionable information that can then be used to drive decision making. Creating dashboards for visual display and manipulation, and also defining alerts are vital for a successful observability strategy.

## ADDITIONAL READING: OBSERVABILITY

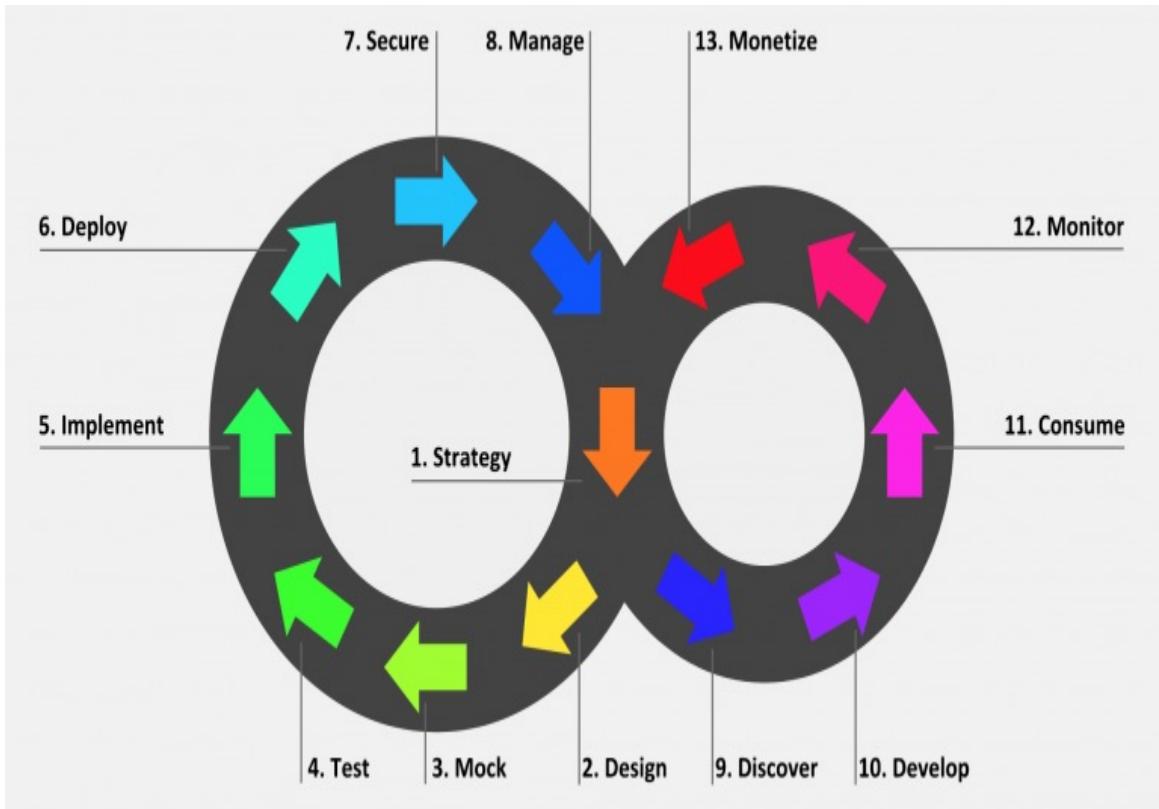
Cindy Sridharan's O'Reilly book [Distributed Systems Observability](#) is a great primer for learning more about the topic of observability.

## Manage APIs as Products: API Lifecycle Management

Modern APIs are often designed, built, and run as products that are consumed by both internal and third-parties, and must be managed as such. Many large organizations see APIs as a critical and strategic component within their business, and as such will create a API program strategy, and set clear goals, constraints, and resources. With a strategy set, the day-to-day tactical approach is often focused on [application lifecycle management](#). Full lifecycle API management spans the entire lifespan of an API that begins at the planning stage and ends when an API is retired. Engineers interact with an associated API gateway, either directly or indirectly, within many of these stages, and all user traffic flows through the gateway. For these reasons, choosing an appropriate API gateway is a critical decision.

There are multiple definitions for key API lifecycle stages. The [Swagger](#) and [SmartBear](#) communities define the five key steps as: planning and designing the API, developing the API, testing the API, deploying the API, and retiring the API.

The [3Scale](#) and [Red Hat](#) teams define thirteen steps:



*Figure 4-4. The 3Scale and Red Hat teams approach to full API lifecycle management*

The Axway team strike a good balance with 3 key components — create, control, and consume — and 10 top stages of an API lifecycle:

### *Building*

Designing and building your API.

### *Testing*

Verifying functionality, performance, and security expectations.

### *Publishing*

Exposing your APIs to developers.

### *Securing*

Mitigating security risks and concerns.

### *Managing*

Maintaining and managing APIs to ensure they are functional, up to date, and meeting business requirements.

### *Onboarding*

Enabling developers to quickly learn how to consume the APIs exposed. For example, offering OpenAPI or ASyncAPI documentation, and providing a portal and sandbox.

### *Analyzing*

Enabling observability, and analyzing monitoring data to understand usage and detect issues.

### *Promoting*

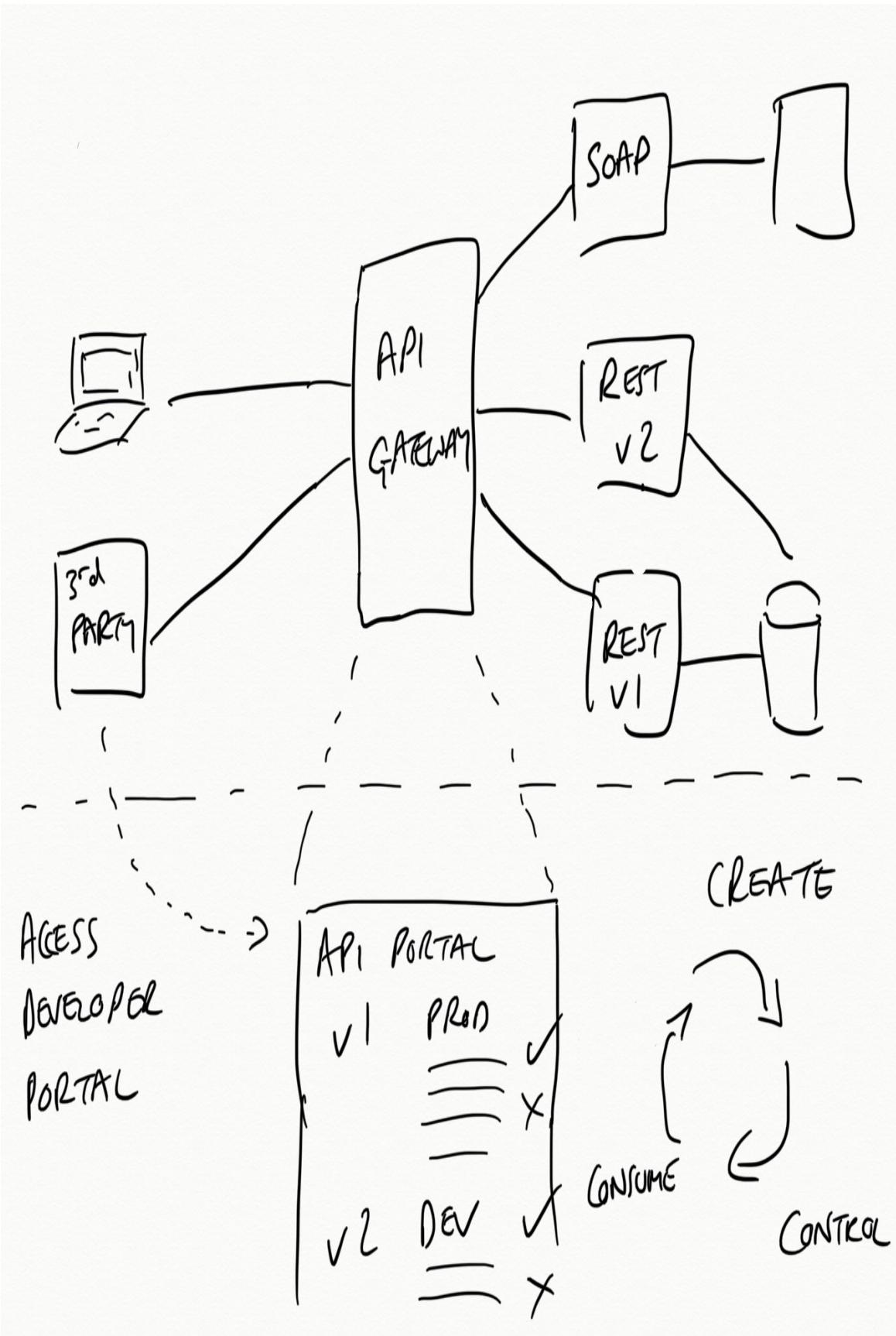
Advertising APIs to developers, for example, listing in an API Marketplace.

### *Monetizing*

Enabling the charging for and collection of revenue for use of an API. We cover this aspect of API lifecycle management as a separate stages in the next section.

### *Retirement*

Supporting the deprecation and removal of APIs, which happens for a variety of reasons including, business priority shifts, technology changes, and security concerns.



*Figure 4-5. API gateway lifecycle management*

## Monetize APIs: Account Management, Billing, and Payment

The topic of billing monetized APIs is closely related to API lifecycle management. The APIs being exposed to customers typically have to be designed as a product, and offered via a developer portal that also includes account management and payment options. Many of the enterprise API gateways include monetization, such as [Apigee Edge](#) and [3Scale](#). These payment portals often integrate with payment solutions, such as PayPal or Stripe, and enable the configuration of developer plans, rate limits, and other API consumption options.

## What is an API Gateway?

With the “why” of an API gateway defined, it’s now important to ask the “what”. In a nutshell, an API gateway is a management tool that sits at the edge of a system between a client and a collection of backend services and acts as a single point of entry for a defined group of APIs. The client can be an end-user application or device, such as a single page web application or a mobile app, or another internal system or third-party application or system.

An API gateway is implemented with two high-level fundamental components, a control plane and data plane. These components can typically be packaged together or deployed separately. The control plane is where operators interact with the gateway and define routes, policies, and required telemetry. The data plane is the location where all of the work specified in the control plane occurs; where the network packets are routed, the policies enforced, and telemetry emitted.

## What Functionality Does an API Gateway Provide?

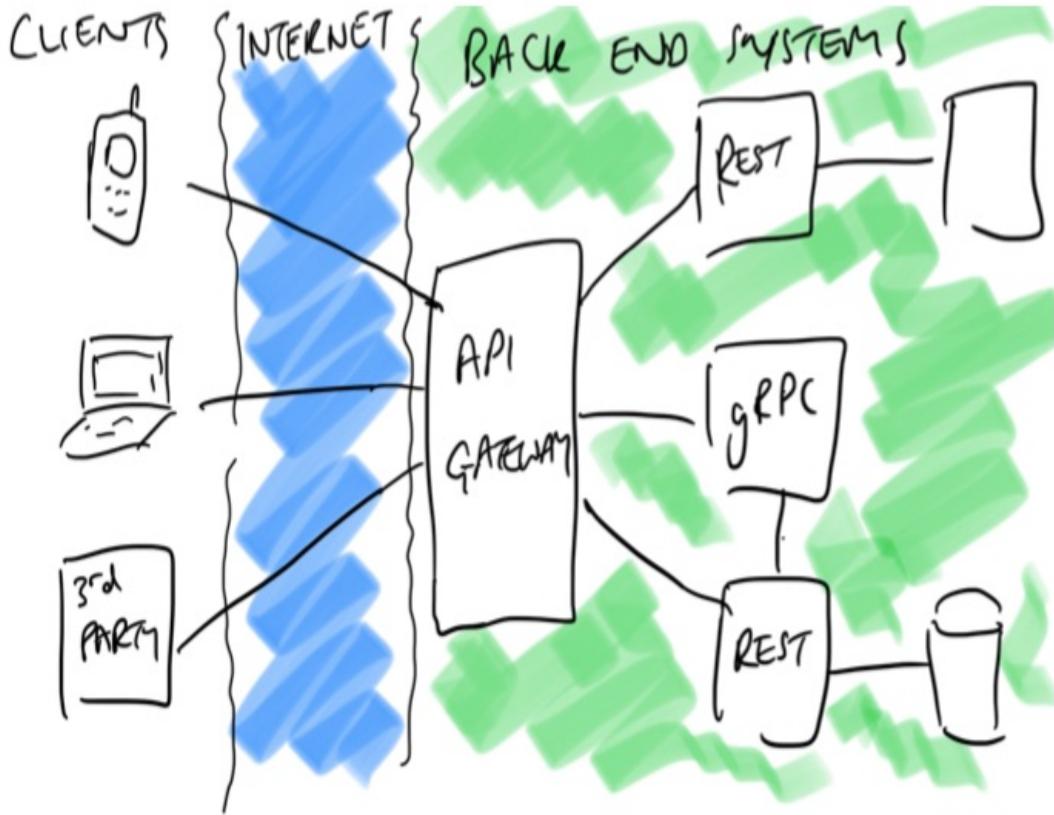
At a network level an API gateway typically acts as a reverse proxy to accept all of the API requests from a client, calls and aggregates the various application-level backend services (and potentially external services) required to fulfill them, and returns the appropriate result. An API gateway provides cross-cutting requirements such as user authentication, request rate limiting, and timeouts/retries, and can provide metrics, logs, and trace data in order to support the implementation of observability within the system. Many API gateways provide additional features that enable developers to manage the lifecycle of an API, assist with the onboarding and management of developers using the APIs (such as providing a developer portal and related account administration and access control), and provide enterprise governance.

### **WHAT IS A PROXY? WHAT IS A REVERSE PROXY?**

A proxy server is an intermediary server that forwards requests for content from multiple clients to different servers across the Internet. For instance, a business may have a proxy that routes and filters employee traffic to the public Internet. A reverse proxy server, on the other hand, is a type of proxy server that typically sits behind the firewall in a private network and routes client requests to the appropriate backend server.

## **Where is an API Gateway Deployed?**

An API gateway is typically deployed at the edge of a system, but the definition of “system” in this case can be quite flexible. For startups and many small-medium businesses (SMBs) an API gateway will often be deployed at the edge of the data center or cloud. In these situations there may only be a single API gateway (deployed and running via multiple instances for high availability) that acts as the front door for the entire back end estate, and the API gateway will provide all of the edge functionality discussed in this chapter via this single component.



*Figure 4-6. A typical startup/SMB API gateway deployment*

For large organizations and enterprises an API gateway will typically be deployed in multiple locations, often as part of the initial edge stack at the perimeter of a data center, and additional gateways may be deployed as part of each product, line of business, or organizational department. In this context these gateways would more typically be separate implementations, and may offer differing functionality depending on geographical location (e.g. required governance) or infrastructure capabilities (e.g. running on low-powered edge compute resources).

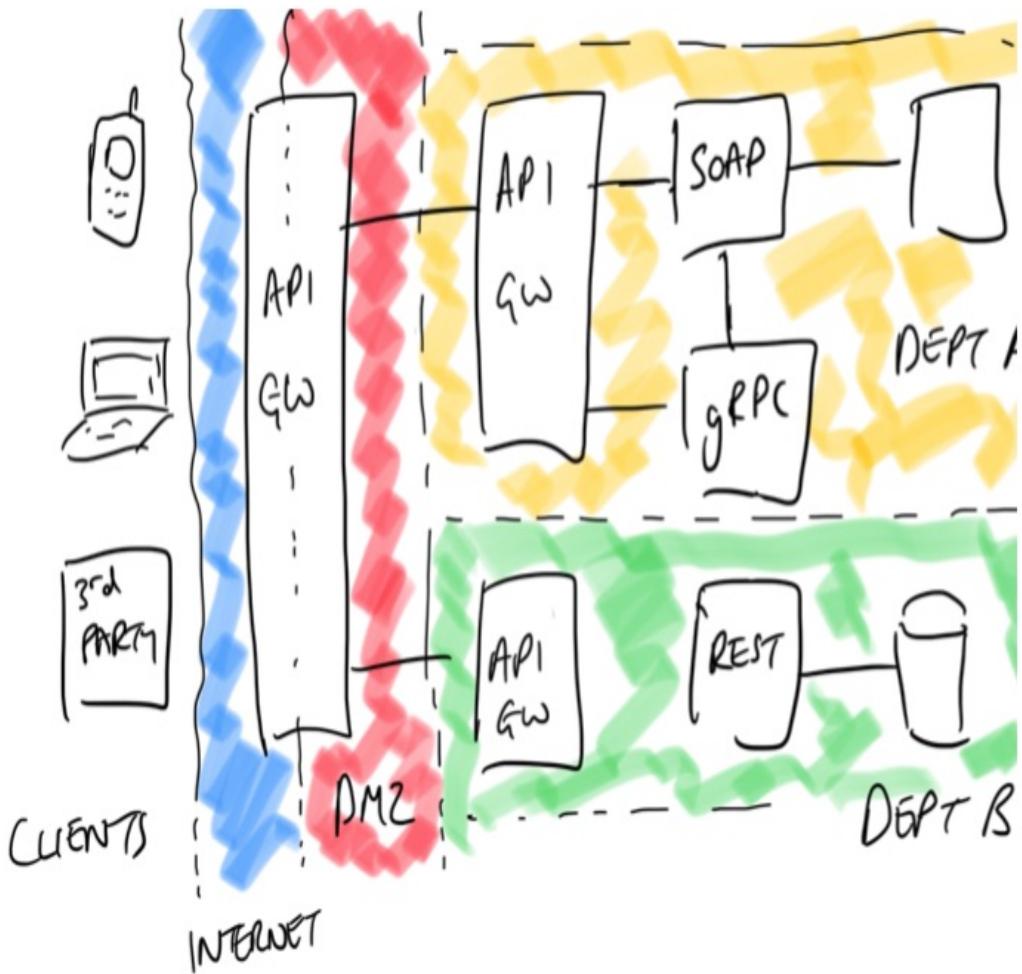


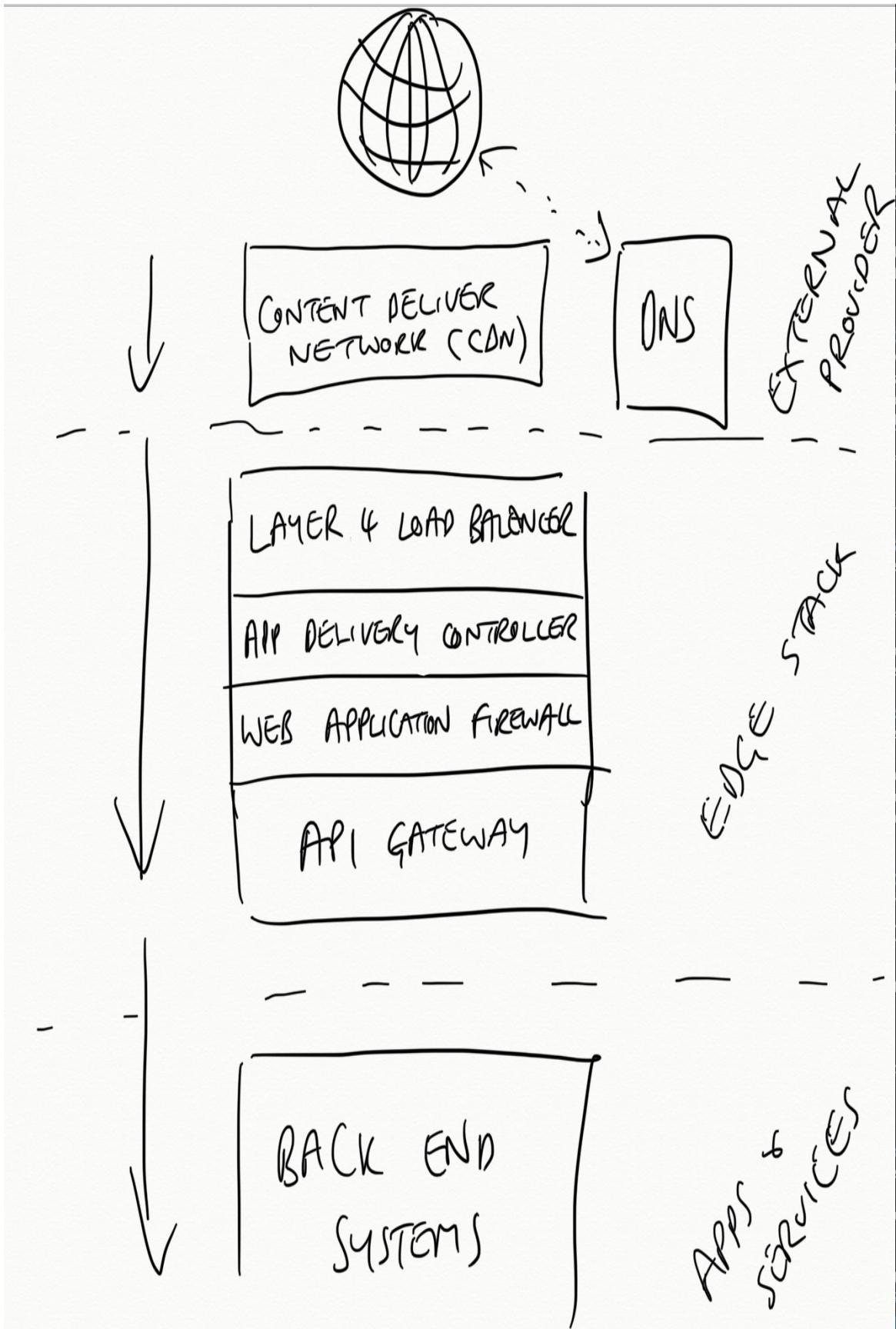
Figure 4-7. A typical large/enterprise API gateway deployment

As you will learn later in this chapter, the definition and exact functionality offered within an API gateway isn't always consistent across implementations, and so the diagrams above should be thought of as more conceptual rather than an exact implementation.

## How Does an API Gateway Integrate with Other Technologies at the Edge?

There is typically many components deployed at the edge of an API-based system. This is where the clients and users first interact with the backend, and hence many cross-cutting concerns are best addressed here. Therefore, a modern edge technology stack or “edge stack” provides a range of

functionality that meets essential cross functional requirements for API-based applications. In some edge stacks each piece of functionality is provided by a separately deployed and operated component, and in others the functionality and/or components are combined. You will learn more about the individual requirements in the next section of the chapter, but for the moment the diagram below should highlight the key layers of a modern edge stack.



*Figure 4-8. A modern edge stack*

Now that you have a good idea about the “what” and “where” of an API gateway, let’s now look at why an organization would use an API gateway.

## A Modern History of API Gateways

Now that you have a good understanding of the “what”, “where”, and “why” of API gateways, it is time to take a glance backwards through history before looking forward to current API gateway technology. As Mark Twain was alleged to have said, “history doesn’t repeat itself, but it often rhymes”, and anyone who has worked in technology for more than a few year will definitely appreciate the relevance this quote has to the general approach seen in the industry. Architecture style and patterns repeat in various “cycles” throughout the history of software development, as do operational approaches. There is typically progress made between these cycles, but it is also easy to miss the teachings that history has to offer.

This is why it is important that you understand the historical context of API gateways and traffic management at the edge of systems. By looking backwards we can build on firm foundations, understand fundamental requirements, and also try to avoid repeating the same mistakes.

## Late 1990s Onwards: Hardware Load Balancers

The concept of World Wide Web (www) was proposed by Tim Berners-Lee in the late 1980s, but this didn’t enter the consciousness of the general public until the mid 1990s, where the initial hype culminated in the dotcom boom and bust of the late 90s. This “Web 1.0” period drove the evolution of web browsers (Netscape Navigator was launched late 1994), the web server (Apache Web Server was released in 1995), and hardware load balancers (F5 was founded in 1996). The Web 1.0 experience consisted of users visiting web sites via making HTTP requests using their browser, and the entire HTML document for each target page being returned in the response. Dynamic aspects of a website were implemented via Common Gateway

Interface (CGI) in combination with scripts written in languages like Perl or C. This was arguably the first incantation of what we would call “Function as a Service (FaaS)” today.

As an increasing number of users accessed each website this strained the underlying web servers. This added the requirement to design systems that supported spreading the increased load and also provide fault tolerance. Hardware load balancers were deployed at the edge of the data center, with the goal of allowing infrastructure engineers, networking specialists, and sysadmins to spread user requests over a number of web server instances. These early load balancer implementations typically supported basic health checks, and if a web server failed or began responding with increased latency then user requests could be routed elsewhere accordingly. Hardware load balancers are still very much in use today. The technology may have improved alongside transistor technology and chip architecture, but the core functionality remains the same.

## **Early 2000s Onwards: Software Load Balancers**

As the Web overcame the early business stumbles from the dotcom bust, the demand for supporting a range of activities, such as users sharing content, ecommerce and online shopping, and businesses collaborating and integrating systems, continued to increase. In reaction, web-based software architectures began to take a number of forms. Smaller organizations were building on their early work with CGI and were also creating monolithic applications in the emerging web-friendly languages such as Java and .NET. Larger enterprises began to embrace modularization (taking their cues from David Parnas’ work in the 1970s), and Service Oriented Architecture (SOA) and the associated “Web Service” specifications (WS-\*) enjoyed a brief moment in the sun.

The requirements for high availability and scalability of web sites were increasing, and the expense and inflexibility of early hardware load balancers was beginning to become a constraining factor. Enter software load balancers, with HAProxy being launched in 2001 and NGINX in 2002. The target users were still operations teams, but the skills required meant

that sysadmins comfortable with configuring software-based web servers were increasingly happy to take responsibility for what used to be a hardware concern.

### **SOFTWARE LOAD BALANCERS: STILL A POPULAR CHOICE TODAY**

Although they have both evolved from initial launches, NGINX and HAProxy are still widely in use, and they are still very useful for small organizations and simple API gateway use cases (both also offer commercial variants more suitable for enterprise deployment). The rise of cloud (and virtualization) cemented the role of software load balancers, and we recommend learning about the basics of this technology.

This time frame also saw the rise of other edge technologies that still required specialized hardware implementation. Content Delivery Networks (CDNs), primarily driven by the Telco organizations, began to be increasingly adopted in order to offload requests from origin web servers. Web Application Firewalls (WAFs) also began to see increasing adoption, first implemented using specialized hardware, and later via software. The open source [ModSecurity](#) project, and the integration of this with the Apache Web Server, drove mass adoption of WAFs.

## **Mid 2000s: Application Delivery Controllers (ADCs)**

The mid 2000s continued to see the increasing pervasiveness of the web in everyday life. The emergence of Internet-capable phones only accelerated this, with BlackBerry initially leading the field, and everything kicking into a higher gear with the launch of the first iPhone in 2007. The PC-based web browser was still the de facto method of accessing the www, and the mid 2000s saw the emergence of “Web 2.0”, triggered primarily by the widespread adoption in browsers of the XMLHttpRequest API and the corresponding technique named “asynchronous JavaScript and XML (AJAX)”. At the time this technology was revolutionary. The asynchronous nature of the API meant that no longer did an entire HTML page have to be returned, parsed, and the display completed refreshed with each request. By

decoupling the data interchange layer from the presentation layer, AJAX allowed web pages to change content dynamically without the need to reload the entire page.

All of these changes placed new demands on web servers and load balancers, for yet again handling more load, but also supporting more secure (SSL) traffic, increasingly large (media rich) data payloads, and different priority requests. This led to the emergence of a new technology named Application Delivery Controllers (ADCs). As these were initially implemented using specialized hardware this led to the existing networking players like F5 Networks, Citrix, Cisco dominating the market. ADCs provided support for compression, caching, connection multiplexing, traffic shaping, and SSL offload, combined with load balancing. The target users were once again infrastructure engineers, networking specialists, and sysadmins.

### THE BENEFITS, AND COSTS, OF SPECIALIZATION

By the mid 2000s nearly all of the components of a modern traffic management edge stack were widely adopted across the industry. The benefits of the separation of concerns were becoming clear (e.g. each edge technology had a clear and focused purpose), but this was increasingly the siloing between teams. If a developer wanted to expose a new application within a large organization this typically meant many separate meetings with the CDN vendors, the load balancing teams, the InfoSec and WAF teams, and the web/application server team. Movements like DevOps emerged, partly driven by a motivation to remove the friction imposed by these silos. If you still have a large number of layers in your edge stack and are migrating to the cloud or a new platform, now is the time to potentially think about the tradeoffs with multiple layers and specialist teams.

## Early 2010s: First Generation API Gateways

The late 2000s and early 2010s saw the emergence of the API economy and associated technologies. Organizations like Twilio were disrupting telecommunications, with their founder, Jeff Lawson, [pitching](#) that “We have taken the entire messy and complex world of telephony and reduced it to five API calls.” The Google Maps API was enabling innovative user

experiences, and Stripe was enabling organizations to easily charge for access to services. Founded in late 2007, Mashape was one of the early pioneers in attempting to create an API marketplace for developers.

Although this exact vision didn't pan out (arguably it was ahead of its time, looking now to the rise of "no code"/"low code" solutions), a byproduct of the Mashape business model was the creation of the Kong API Gateway, built upon **OpenResty** and the open source NGINX implementation. Other implementations included WSO2 with **Cloud Services Gateway**, Sono Systems with Apigee, and Red Hat with 3Scale Connect.

These were the first edge technologies that were targeted at developers in addition to platform teams and sysadmins. A big focus was on managing the software development lifecycle (SDLC) of an API and providing system integration functionality, such as endpoints and protocol connectors, and translation modules. Due to the range of functionality offered, the vast majority of first generation API gateways were implemented in software. Developer portals sprang up in many products, which allowed engineers to document and share their APIs in a structured way. These portals also provided access controls and user/developer account management, and publishing controls and analytics. The theory was that this would enable the easy monetization of APIs, and the management of "APIs as a product".

During this evolution of developer interaction at the edge there was increasing focus on the application layer (layer 7) of the **OSI Networking model**. The previous generations of edge technologies often focused on IP addresses and ports, which primarily operate at the transport layer (layer 4) of the OSI model. Allowing developers to make routing decisions in an API gateway based on HTTP metadata such as path-based routing or header-based routing provided the opportunity for richer functionality.

There was also an emerging trend towards creating smaller service-based architectures, and some organizations were extracting single-purpose standalone applications from their existing monolithic code bases. Some of these monoliths acted as an API gateway, or provided API gateway-like functionality, such as routing and authentication. With the first generation of API gateways it was often the case that both functional and cross-functional

concerns, such as routing, security, and resilience, were performed both at the edge and also within the applications and services.

### **WATCH FOR RELEASE COUPLING WHEN MIGRATING AWAY FROM A MONOLITH**

Extracting standalone services from a monolithic application and having the monolith acting as a simple gateway can be a useful migration pattern towards adopting a microservices architecture. However, beware of the costs related to coupling between the application gateway and services that are introduced with this pattern. For example, although the newly extracted services can be deployed on demand, their release is often dependent on the release cadence of the monolith, as all traffic must be routed through this component. And if the monolith has an outage, then so do all of your services operating behind it.

## **2015 Onwards: Second Generation API Gateways**

The mid 2010s saw the rise of the next generation of modular and service-oriented architectures, with the concept of “microservices” firmly entering the zeitgeist by 2015. This was largely thanks to “unicorn” organizations like Netflix, AWS, and Spotify sharing their experiences of working with these architectural patterns. In addition to back end systems being decomposed into more numerous and smaller services, developers were also adopting container technologies based on Linux LXC. Docker was released in March of 2013, and Kubernetes followed hot on its heels with a v1.0 release in July of 2015. This shift in architectural style and changing runtimes drove new requirements at the edge. Netflix released their bespoke JVM-based API gateway, [Zuul](#), in mid 2013. Zuul supported service discovery for dynamic back end services and also allowed Groovy scripts to be injected at runtime in order to dynamically modify behaviour. This gateway also consolidated many cross cutting concerns into a single edge component, such as authentication, testing (canary releases), rate limiting and load shedding, and observability. Zuul was a revolutionary API gateway in the microservices space, and it has since evolved into a second version, and Spring Cloud Gateway has been built on top of this.

## **WATCH FOR COUPLING OF BUSINESS LOGIC IN THE GATEWAY**

When using Zuul it was all too easy to accidentally highly couple the gateway with a service by including related business logic in both a Groovy script (or scripts) and the service. This meant that the deployment of a new version of a service often required modifications to the code running in the gateway. The pinnacle of this bad coupling can be seen when a microservice team decides to reuse an existing deployed Groovy script for their service, and then at some arbitrary time in the future the script is modified by the original owning team in an incompatible way. This can quickly lead to confusion as to why things are broken, and also to whack-a-mole type fixes. This danger wasn't unique to Zuul, and nearly all proxies and many gateways allowed the injection of plugins or dynamic behaviour, but this was often only accessible to operations teams or written in obscure (or unpopular) languages. The use of Groovy in Zuul made this very easy for application developers to implement.

With the increasing adoption of Kubernetes and the open source release of the Envoy Proxy in 2016 by the Lyft team, many API gateways were created around this technology, including Ambassador, Contour, and Gloo. This drove further innovation across the API gateway space, with Kong mirroring functionality offered by the next generation of gateways, and other gateways being launched, such as Traefik, Tyk, and others. Increasingly, many of the Kubernetes Ingress Controllers called themselves "API gateways", regardless of the functionality they offered, and this led to some confusion for end users in this space.

## CONFUSION IN THE CLOUD: API GATEWAYS, EDGE PROXIES, AND INGRESS CONTROLLERS

As Christian Posta noted in his blog post [API Gateways Are Going Through an Identity Crisis](#), there is some confusion around what an API gateway is in relation to proxy technologies being adopted within the cloud computing domain. Generally speaking, in this context an API gateway enables some form of management of APIs, ranging from simple adaptor-style functionality operating at the application layer (OSI layer 7) that provides fundamental cross-cutting concerns, all the way to full lifecycle API management. Edge proxies are more general purpose traffic proxies or reverse proxies that operate at the network and transport layers (OSI layers 3 and 4 respectively) and provide basic cross-cutting concerns, and tend not to offer API-specific functionality. [Ingress controllers](#) are a Kubernetes-specific technology that controls what traffic enters a cluster, and how this traffic is handled.

The target users for the second generation of API gateways was largely the same as the cohort for the first generation, but with a clearer separation of concerns and a stronger focus on developer self-service. The move from first to second generation of API gateways saw increased consolidation of both functional and cross-functional requirements being implemented in the gateway. Although it became widely accepted that microservices should be built around the idea espoused by James Lewis and Martin Fowler of “smart endpoints and dumb pipes”, the uptake of polyglot language stacks mean that “microservice gateways” emerged (more detail in the next section) that offered cross-cutting functionality in a language-agnostic way.

## Coda (2017 Onwards): Service Mesh and/or API Gateway?

In early 2017 there was increasing buzz about the use of “service meshes”, a new communication framework for microservice-based systems. William Morgan, CEO of Buoyant and ex-Twitter engineer, is [largely credited](#) with coining the term. The early service meshes exclusively targeted the east-west, service-to-service communication requirements. As Phil Calçado writes in [Pattern: Service Mesh](#), service meshes evolved from early in-process microservice communication and resilience libraries, such as

Netflix’s Ribbon, Eureka, and Hystrix. Buoyant’s Linkerd was the first service mesh to be hosted by the CNCF, but Google’s Istio captured the attention of many engineers (largely thanks to Google’s impressive marketing machine).

Although there is a lot of functional overlap between service meshes and API gateways, and much of underlying proxy (data plane) technology is identical, the use cases and control plane implementations are at the moment quite different. Several service meshes offer ingress and gateway-like functionality, for example, [Istio Gateway](#) and [Consul Ingress Gateways](#). However, this functionality is typically not as feature-rich or as mature as that currently offered by existing API gateways. This is why the majority of modern API gateways integrate effectively with at least one service mesh.

### START AT THE EDGE AND WORK INWARDS

Every production system that exposes an API will require an API gateway, but not necessarily a service mesh. Many existing systems, and also simple greenfield applications, are implemented as a shallow stack of services, and therefore the added operational complexity of deploying and running a service mesh does not provide a high enough level of benefit. For this reason it is typically the case that we recommend engineers “start at the edge and work inwards”, i.e. select and deploy an API gateway, and only when the number of services (and their interactions) grows consider selecting a service mesh.

## Current API Gateway Taxonomy

As can be the case with terminology in the software development industry, there often isn’t an exact agreement on what defines or classifies an API gateway. There is broad agreement in regards to the functionality this technology should provide, but different segments of the industry have different requirements, and hence different views, for an API gateway. This has led to several sub-types of API gateway emerging and being discussed. In this section of the chapter you will explore the emerging taxonomy of

API gateways, and learn about their respective use cases, strengths, and weaknesses.

## **Traditional Enterprise Gateways**

The traditional enterprise API gateway is typically aimed at the use case of exposing and managing business-focused APIs. This gateway is also often integrated with a full API lifecycle management solution, as this is an essential requirement when releasing, operating, and monetizing APIs at scale. The majority of gateways in this space may offer an open source edition, but there is typically a strong usage bias towards the open core/commercial version of the gateway.

These gateways typically require the deployment and operation of dependent services, such as data stores. These external dependencies have to be run with high availability to maintain the correct operation of the gateway, and this must be factored into running costs and DR/BC plans.

## **Micro/Microservices Gateways**

The primary use case of microservices API gateway, or micro API gateway, is to route ingress traffic to backend APIs and services. There is typically not much provided for the management of an API's lifecycle. These types of gateway are often available fully-featured as open source or are offered as a lightweight version of a traditional enterprise gateway.

They tend to be deployed and operated as standalone components, and often make use of the underlying platform (e.g. Kubernetes) for the management of any state. As microservices gateways are typically built using modern proxy technology like Envoy, the integration capabilities with service meshes (especially those built using the same proxy technology) is typically good.

## **Service Mesh Gateways**

The ingress or API gateway included with a service mesh is typically designed to provide only the core functionality of routing external traffic into the mesh. For this reason they often lack some of the typical enterprise features, such as comprehensive integration with authentication and identity provider solutions, and also integration with other security features, such as a WAF.

The service gateway typically manages state using the capabilities of the mesh itself or underlying deployment fabric (e.g. Kubernetes). This type of gateway is also implicitly coupled with the associated service mesh (and operational requirements), and so if you are not yet planning to deploy a service mesh, then this is most likely not a good first choice of API gateway.

## **Comparing API Gateway Types**

The table below highlights the difference between the two most widely deployed API gateway types across six important criteria.

*T*

*a*

*b*

*l*

*e*

*4*

-

*l*

.

*C*

*o*

*m*

*p*

*a*

*r*

*i*

*s*

*o*

*n*

*o*

*f*

*E*

*n*

*t*

*e*

*r*

*p*

*r*

*i*

*s*

*e*

,

*M*

*i*

*c*  
*r*  
*o*  
*s*  
*e*  
*r*  
*v*  
*i*  
*c*  
*e*  
*s*

,

*a*  
*n*  
*d*  
*S*  
*e*  
*r*  
*v*  
*i*  
*c*  
*e*

*M*  
*e*  
*s*  
*h*  
*A*  
*P*  
*I*  
*g*  
*a*  
*t*  
*e*  
*w*  
*a*  
*y*

Use case	Traditional Enterprise API gateway
Microservices API gateway	Service Mesh Gateway
Primary Purpose	Expose, compose, and manage internal business APIs and associated services.
Expose, compose, and manage internal business services.	Expose internal services within the mesh.
Publishing Functionality	API management team or service team registers / updates gateway via admin API.
Service team registers / updates gateway via declarative code as part of the deployment process.	Service team registers / updates mesh and gateway via declarative code as part of the deployment process.
Monitoring	Admin and operations focused e.g. meter API calls per consumer, report errors (e.g. internal 5XX).
Developer focused e.g. latency, traffic, errors, saturation.	Platform focused e.g. utilization, saturation, errors.
Handling and Debugging Issues	L7 error-handling (e.g. custom error page). For troubleshooting, run gateway/API with additional logging and debug issue in staging environment.
L7 error-handling (e.g. custom error page, failover, or payload). For debugging issues configure more detailed monitoring, and enable traffic shadowing and / or canarying to recreate the problem.	L7 error-handling (e.g. custom error page or payload). For troubleshooting, configure more detailed monitoring and/or utilise traffic “tapping” to view and debug specific service-to-service communication.
Testing	Operate multiple environments for QA, Staging, and Production. Automated integration testing, and gated API deployment. Use client-driven API versioning for compatibility and stability (e.g. semver)
Enables canary routing and dark launching	Facilitate canary routing for dynamic testing.

for dynamic testing.  
Use contract testing for upgrade management.

#### Local Development

Deploy gateway locally (via installation script, Vagrant or Docker), and attempt to mitigate infrastructure differences with production. Use language-specific gateway mocking and stubbing frameworks.

Deploy gateway locally via service orchestration platform (e.g. container, or Kubernetes)

Deploy service mesh locally via service orchestration platform (e.g. Kubernetes)

## Evolving the Conference System Using an API Gateway

In this section of the chapter you will learn how to install and configure an API gateway to route traffic directly to the Attendee service that has been extracted from the monolithic Conference System. This will demonstrate how you can use the popular “strangler fig” pattern to evolve your system from a monolith to a microservices-based architecture over time by gradually extracting pieces of an existing system into independently deployable and runnable services.

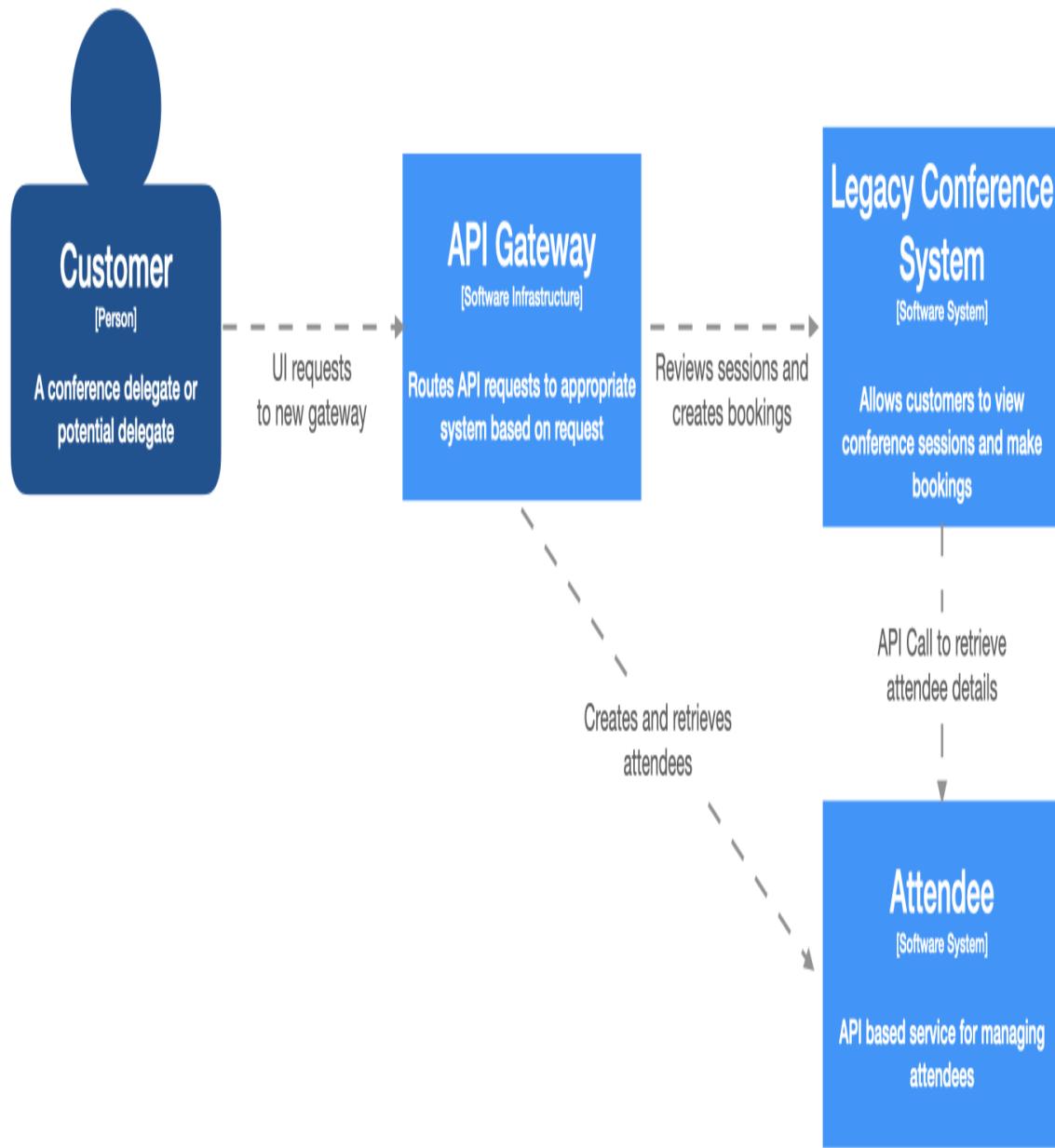


Figure 4-9. Using an API gateway to route to a new Attendees service running independently from the monolith

Many organisations often start such a migration by extracting services but have the monolithic application perform the routing and other cross-cutting concerns for the externally running services. This is often the easy choice, as the monolith already has to provide this functionality for internal functions. However, this leads to tight coupling between the monolith and services, with all traffic flowing through the monolithic application and the

configuration cadence determined by the frequency of deployment of the monolith. From a traffic management perspective, both the increased load on the monolithic application and increased blast radius if this does fail mean the operational cost can be high. And being limited in updating routing information or cross-cutting configuration due to a slow release train or a failed deployment can prevent you from iterating at speed. Because of this we generally do not recommend using the monolith to route traffic in this fashion; particularly if you plan to extract many services within a relatively short time scale.

As long as the gateway is deployed to be highly available and developers have access to self-serve with routing and configuration, extracting and centralising application routing and cross-cutting concerns to an API gateway provides both safety and speed. Let's now walk through a practical example of deploying an API gateway within the Conference System and using this to route to the new Attenee service.

## Installing Ambassador Edge Stack in Kubernetes

As you are deploying the Conference System into a Kubernetes cluster, you can easily install an API gateway using the standard Kubernetes-native approaches, such as applying YAML config or using Helm, in addition to using command line utilities. For example, the [Ambassador Edge Stack API gateway](#) can be installed using the `edgectl` tool, which will also create a sample “edgestack.me” sub-domain and configure a TLS certificate using LetsEncrypt.

Make sure you have your local Kubernetes context configured correctly so that you can execute commands using the `kubectl` tool, such as `kubectl get svc -A` successfully. Visit <https://www.getambassador.io/docs/edge-stack/latest/tutorials/getting-started/> and download the `edgectl` binary that is suitable for your operating system. Now run `edgectl install`. You'll be prompted to enter your email address for the purpose of LetsEncrypt being able to notify you when your TLS certificate is about to expire.

```
$ edgetl install  
Installing the Ambassador Edge Stack  
  
Please enter an email address for us to notify you before your TLS certificate  
and domain name expire. In order to acquire the TLS certificate, we share this  
email with Let's Encrypt.  
Email address [daniel.bryant@redacted.co.uk]:  
=====  
Beginning Ambassador Edge Stack Installation  
-> Finding repositories and chart versions  
-> Installing Ambassador Edge Stack 1.14.1  
-> Checking the AES pod deployment  
-> Provisioning a cloud load balancer  
-> Your Ambassador Edge Stack installation\'s address is 34.70.47.47  
-> Checking that Ambassador Edge Stack is responding to ACME challenge  
-> Automatically configuring TLS  
-> Acquiring DNS name practical-archimedes-7849.edgestack.me  
-> Obtaining a TLS certificate from Let\'s Encrypt  
-> TLS configured successfully  
  
Ambassador Edge Stack Installation Complete!  
=====  
Congratulations! You\'ve successfully installed the Ambassador Edge Stack in  
your Kubernetes cluster. You can find it at your custom URL:  
https://practical-archimedes-7849.edgestack.me/
```

With the API gateway up and running and providing an HTTPS connection the Conference System application no longer needs to be concerned with terminating TLS connections or listening to multiple ports. Similarly, authentication and rate limiting can also be easily configured without having to reconfigure or deploy your application.

## Configuring Mappings from URL Paths to Backend Services

You can now use an Ambassador Edge Stack Mapping Custom Resource to map the root of your domain to the “conferencesystem” service listening on port 8080 and running in the “legacy” namespace within the Kubernetes cluster. This Mapping should be familiar to anyone that has configured a

web application or reverse proxy to listen for user requests. The `metadata` provides a name for the Mapping, and the `prefix` determines the path (the “`/`” root in this case) that is mapped to the target `service` (with the format `service-name.namespace:port`). An example is shown below:

```
---
apiVersion: getambassador.io/v2
kind: Mapping
metadata:
  name: legacy-conference
spec:
  prefix: /
  rewrite: /
  service: conferencesystem.legacy:8080
```

Another Mapping can be added to route any traffic sent to the “`/attendees`” path to the new (“nextgen”) attendees microservice that has been extracted from the monolith. The information included in the Mapping should look familiar from the previous example. Here a `rewrite` is specified which “rewrites” the matching `prefix` path in the URL metadata before making the call to the target attendees service. This makes it appear to the attendees service that the request originated with the “`/`” path, effectively stripping out the “`/attendees`” part of the path.

```
---
apiVersion: getambassador.io/v2
kind: Mapping
metadata:
  name: legacy-conference
spec:
  prefix: /attendees
  rewrite: /
  service: attendees.nextgen:8080
```

This pattern of creating additional Mappings as each new microservice is extracted from the legacy application can continue. Eventually the legacy application becomes a small shell with only a handful of functions, and all of the other functionality is handled by microservices, each with their own Mapping.

## Configuring Mappings Using Host or Path-based Routing

Matching prefixes can be nested e.g. `/attendees/affiliation` or use regular expressions e.g. `/attendees/^a-z].*`. Most API gateways will also let you perform host-based routing e.g. `host: attendees.conferencesystem.com`. An example of this using Ambassador Edge Stack Mappings can be seen below:

```
---
apiVersion: getambassador.io/v2
kind: Mapping
metadata:
  name: attendees-host
spec:
  prefix: /
  host: attendees.conferencesystem.com
  service: attendees.nextgen:8080
```

Query parameter-based routing is often supported, too: e.g. `/attendees?specialcase=true`. An example of this using Ambassador Edge Stack Mappings can be seen below:

```
---
apiVersion: getambassador.io/v2
kind: Mapping
metadata:
  name: attendees-param
spec:
  prefix: /attendees
  query_parameters:
    specialcase: true
  service: attendees.nextgen:8080
```

Some API gateways will enable routing based on the payload or body of a request, but this should generally be avoided for two reasons: first, this often leaks highly-coupled domain-specific information into the API gateway config (e.g. a payload often conforms to a schema/contract that may change in the application, which the gateway will now need to be synchronized with); and second, it can be computationally expense to

deserialize and parse a large payload in order to extract the required information for routing.

## **Deploying API Gateways: Understanding and Managing Failure**

Regardless of the deployment pattern and number of gateways involved within a system, an API gateway is typically on the critical path of many, if not all, user requests entering into your system. An outage of a gateway deployed at the edge typically results in the unavailability of the entire system. And an outage of a gateway deployed further upstream typically results in the unavailability of some core subsystem. For this reason the topics of understanding and managing failure of an API gateway are vitally important to learn.

### **API Gateway as a Single Point of Failure**

The first essential step in identifying single points of failure in a system is to ensure that you have an accurate understanding of the current system. This is typically accomplished by investigation and the continual update of associated documentation and diagrams. Assembling a diagram that traces a user-initiated request for each core journey or use case, all the way from ingress to data store or external system and back to egress, that shows all the core components involved can be extremely eye opening. This is especially the case in large organizations, where ownership can be unclear and core technologies can accidentally become abandoned.

In a typical web-based system, the first obvious single point of failure is typically DNS. Although this is often externally managed, there is no escaping the fact that if this fails, then your site will be unavailable. The next single points of failure will typically then be the global and regional layer 4 load balancers, and depending on the deployment location and configuration, the security edge components, such as the firewall or WAF.

## CHALLENGE ASSUMPTIONS WITH SECURITY SINGLE POINTS OF FAILURE

Depending on the product, deployment, and configuration, some security components may “fail open”, i.e. if the component fails then traffic will simply be passed through to upstream components or the backend. For some scenarios where availability is the most important goal this is desired, but for others (e.g. financial or government systems), this is most likely not. Be sure to challenge assumptions in your current security configuration.

After these core edge components, the next layer is typically the API gateway. The more functionality you are relying on within the gateway, the bigger the risk involved and bigger the impact of an outage. As an API gateway is often involved in a software release the configuration is also continually being updated. It is critical to be able to detect and resolve issues, and mitigate any risks.

## Detecting and Owning Problems

The first stage in detecting issues is ensuring that you are collecting and have access to appropriate signals from your monitoring system, i.e. data from metrics, logs, and traces. Any critical system should have a clearly defined team that owns it and is accountable for any issues. Teams should communicate **service level objectives (SLOs)**, which can be codified into service level agreements (SLAs) for both internal and external customers.

## ADDITIONAL READING: OBSERVABILITY, ALERTING, AND SRE

If you are new to the concept of observability, then we recommend learning more about Brendan Gregg’s [utilization, saturation, and errors \(USE\) method](#), Tom Wilkie’s [rate, errors, and duration \(RED\) method](#), and Google’s [four golden signals of monitoring](#). If you want to learn more about associated organisational goals and processes the Google Site Reliability Engineering (SRE) book, is highly recommended.

## Resolving Incidents and Issues

First and foremost, each API gateway operating within your system needs an owner that is accountable if anything goes wrong with the component. In a smaller organisation this may be the developers or SRE team that are also responsible for the underlying services. In a larger organisation this may be a dedicated infrastructure team. As an API gateway is on the critical path of requests, some portion of this owning team should be on-call as appropriate (this may be 24/7/365). The on-call team will then face the tricky task of fixing the issue as rapidly as possible, but also learning enough (or locating and quarantining systems and configuration) to learn what went wrong. After each incident you should strive to conduct blameless a postmortem, and document and share all of your learning. Not only can this information be used to trigger remediate action to prevent this issue reoccurring, but this knowledge can be very useful for engineers learning the system and for external teams dealing with similar technologies or challenges. If you are new to this space then the [Learning from Incidents](#) website is a fantastic jumping off point.

## Mitigating Risks

Any component that is on the critical path for handling user requests should be made as highly available as is practical in relation to cost and operational complexity. Software architects and technical leaders deal with tradeoffs; this type is one of the most challenging. In the world of API gateways, high availability typically starts with running multiple instances. With on-premise/co-lo instances this translates into operating multiple (redundant) hardware appliances, ideally spread across separate locations. In the cloud, this translates into designing and running the API gateway instance in multiple availability zones/data centers and regions. If a (global) load balancer is deployed in front of the API gateway instances, then this must be configured appropriately with health checks and failover processes that must be test regularly. This is especially important if the API gateways instances run in active/passive or leader/node modes of operation.

## LOAD BALANCING CHALLENGES

You must ensure that your load balancer to API gateway failover process meets all of your requirements in relation to continuity of service. Common problems experienced during failover events include:

- User client state management issues, such as backend state not being migrated correctly, which causes the failure of sticky sessions
- Poor performance, as clients are not redirected based on geographical considerations e.g. European users being redirected to the US west coast when an east coast data center is available
- Unintentional cascading failure, such as a faulty leader election component that results in deadlock, which causes all backend systems to become unavailable

Care should be taken with any high availability strategy to ensure that dependent components are also included. For example, many enterprise API gateways rely on a data store to store configuration and state in order to function correctly. This must also be run in a HA configuration, and it must be available to all active and passive components. It is becoming increasingly common to split the deployment and operation of the control plane and data plane, but you will need to understand what happens if the control plane has an outage over a considerable period of time.

Often the biggest period of change with an API gateway is during configuration update, for example, a release of a new API or service. **Progressive delivery** techniques should be used, such as rolling upgrades, localized rollouts, blue/green deployments, and canary releasing. These allow fast feedback and incremental testing in a production environment, and also limit the blast radius of any outage.

Finally, you can put in place reactive fallback measures to handle failure. Examples of this include building functionality to serve a degraded version of the website from static servers, caches, or via a CDN.

# Common API Gateway Implementation Pitfalls

We've already mentioned that no technology is a silver bullet, but (continuing on the theme of technology cliches) it can be the case that when you have a hammer, everything tends to look like a nail. There are several traffic management use cases where an API gateway is not an exact fit. There are also many examples of technology that became overly coupled with the underlying application and added increasing friction to the software delivery process being a key culprit. There are several common API gateway common pitfalls that you should always try and avoid.

## API Gateway Loopback: “Service Mesh Lite”

As with all common pitfalls, the implementation of this pattern often begins with good intentions. When an organization has only a few services this typically doesn't warrant the installation of a service mesh. However, a subset of service mesh functionality is often required, particularly service discovery. An easy implementation is to route all traffic through the edge or API gateway, which maintains the official directory of all service locations. At this stage the pattern looks somewhat like a “hub and spoke” networking diagram. The challenges present themselves in two forms: Firstly, when all of the service-to-service traffic is leaving the network before reentering via the gateway this can present performance, security, and cost concerns (cloud vendors often charge for egress). Secondly, this pattern doesn't scale beyond a handful of services, as the gateway becomes overloaded and a bottleneck, and it becomes a true single point of failure.

Looking at the current state of the Conference System with the two Mappings that you have configured, you can see the emergence of this issue. Any external traffic, such as user requests, are correctly being routed to their target services by the API gateway. However, how does the legacy application discover the location of the attendees service? Often the first approach is to route all requests back through the publicly addressable gateway e.g the legacy application makes calls to

[www.conferencesystems.com/attendees](http://www.conferencesystems.com/attendees) Instead the legacy application should use some form of internal service discovery mechanism and keep all of the internal requests within the internal network. You will learn more about how to use a service mesh to implement this in the next chapter.

## API Gateway as an ESB

The vast majority of API gateways support the extension of their out-of-the-box functionality via the creation of plugins or modules. NGINX supported Lua modules, which OpenResty and Kong capitalised on. Envoy Proxy originally supported extensions in C, and now WebAssembly filters. And we've already discussed how the original implementation of Netflix's Zuul API gateway "[“2015 Onwards: Second Generation API Gateways”](#)". Many of the use cases realised by these plugins are extremely useful, such as authn/z, filtering, and logging. However, it can be tempting to put business logic into these plugins, which is a way to highly-couple your gateway with your service or application. This leads to a potentially fragile system, where a change in a single plugin ripples throughout the organization, or additional friction during release where the target service and plugin have to be deployed in lockstep.

## Turtles (API Gateways) All the Way Down

If one API gateway is good, more must be better, right? It is common to find multiple API gateways deployed within the context of large organization, often in a hierarchical fashion, or in an attempt to segment networks or departments. The intentions are typically good: either for providing encapsulation for internal lines of business, or for a separation of concerns with each gateway (e.g. “this is the transport security gateway, this is the auth gateway, this is the logging gateway...”). The common pitfall rears its head when the cost of change is too high, e.g. you have to coordinate with a large number of gateway teams to release a simple service upgrade, there are understandability issues (“who owns the tracing functionality?”), or performance is impacted as every network hop naturally incurs a cost.

# Selecting an API Gateway

Now that you learned about the functionality provided by an API gateway, the history of the technology, and how an API gateway fits into to the overall system architecture, next is the \$1M question: how do you select an API gateway to included in your stack?

## Identifying Requirements

One of the first steps with any new software delivery or infrastructure project is identifying the related requirements. This may appear obvious, but it is all too easy to get distracted by shiny technology, magical marketing, or good sales documentation!

You can look back to the earlier section “[Why Use an API Gateway?](#)” of this chapter to explore in more detail the high-level requirements you should be considering during the selection process:

- Reducing coupling between front ends and back ends
- Simplifying client consumption of APIs by aggregating and/or translating back end services
- Protecting APIs from overuse and abuse, via threat detection and mitigation
- Understanding how APIs are being consumed and how the underlying systems are performing
- Managing APIs as products i.e. API Lifecycle Management requirements
- Monetizing APIs, including the needs for account management, billing, and payment

It is important to ask question both focused on current pain points and also your future roadmap. For example, do you need to simply expose REST-based backend APIs with a thin layer of load balancing and security, or do

you need comprehensive API aggregation and translation functionality combined with enterprise-grade security?

## Exploring Constraints: Team, Technologies, and Roadmap

In addition to identifying requirements it is also essential that you identify your organization's constraints. Broadly speaking, we have found there are three key areas that are well worth exploring when choosing your API gateway: team structure and dynamics, existing technology investments, and future roadmap.

The much discussed [Conway's law](#), stating that organizations design systems mirror their own communication structure, has almost become a cliche. But there is truth to this. As explored by Matthew Skelton and Manuel Pais in [Team Topologies](#) your organizational structure will constrain your solution space. For example, if your organization has a separate InfoSec team, there is a good chance that you will end up with a separate security gateway.

Another well discussed concept is the notion of the [sunk cost fallacy](#), but in the real world we see this ignored time and time again. Existing technology deployed often constrains future decisions. You will need to investigate the history associated with many of the big ticket technology decision making within an organization, such as API gateway choice, in order to avoid messy politics. For example, if a platform team has already invested millions of dollars into creating a bespoke API gateway, there is a good chance they will be reluctant to adopt another gateway, even if this is a better fit for the current problem space.

Finally, you will need to identify and explore your organization's roadmap. Often this will constrain the solution space. For example, some organizations are banning the use of certain cloud vendors or technology vendors, for reasons related to competition or ethical concerns. Other times, you will find that the roadmap indicates the leadership is "all-in" on a

specific technology, say, Java, and so this may restrict the deployment of a gateway that doesn't run on a JVM.

## Build Versus Buy

A common discussion when selecting an API gateway is the “build versus buy” dilemma. This is not unique to this component of a software system, but the functionality offered via an API gateway does lead to some engineers gravitating to this that they could build this “better” than existing vendors, or that their organization is somehow “special”, and would benefit from a custom implementation. In general, we believe that the API gateway component is sufficiently well-established that it typically best to adopt an open source implementation or commercial solution rather than build your own. Presenting the case for build versus buy with software delivery technology could take an entire book, and so in this section we only want to highlight some common challenges:

- Underestimating the Total Cost of Ownership (TCO): Many engineers discount the cost of engineering a solution, the continued maintenance costs, and the ongoing operational costs.
- Not thinking about opportunity cost: Unless you are a cloud or platform vendor, it is highly unlikely that a custom API gateway will provide you with a competitive advantage. You can delivery more value to your customers by building some functionality closer to your overall value proposition
- Not being aware of current technical solutions of products. Both the open source and commercial platform component space moves fast, and it can be challenging to keep up to date. This, however, is a core part of the role of being a technical leader.

## Radars, Quadrants, and Trend Reports

Although you should always perform your own experiments and proof of concept work, we recommend keeping up to date with technology trends

via reputable IT press. This type of content can be especially useful when you are struggling with a problem or have identified a solution and are in need of a specific piece of technology that many vendors offer.

We recommend the following sources of information for learning more about the state of the art of technology within the API gateway space:

- ThoughtWorks Technology Radar
- Gartner Magic Quadrant for Full Life Cycle API Management
- Cloud Native Computing Foundation (CNCF) Tech Radar
- InfoQ Trends Reports

Several organizations and individuals also publish periodic API gateway comparison spreadsheets, and these can be useful for simple “paper evaluations” in order to shortlist products to experiment with. It should go without saying that you will need to check for bias across these comparisons (vendors frequently sponsor such work), and also ensure the publication date is relatively recent. Enterprise API gateways solutions do not change much from month to month, but this is not true in the cloud native and Kubernetes space.

## API Gateway: A Type 1 Decision

Jeff Bezos, the CEO of Amazon, is famous for many things, and one of them is his discussion of [Type 1 decisions and Type 2 decisions](#). Type 1 decisions are not easily reversible, and you have to be very careful making them. Type 2 decisions are easy to change: “like walking through a door — if you don’t like the decision, you can always go back.” Usually this concept is presented in relation to confusing the two, and using Type 1 processes on Type 2 decisions: “The end result of this is slowness, unthoughtful risk aversion, failure to experiment sufficiently, and consequently diminished invention. We’ll have to figure out how to fight that tendency.” However, in the majority of cases — especially within a

large enterprise context — choosing an API gateway is very much a Type 1 decision. Ensure your organization acts accordingly!

## **Checklist: Selecting an API Gateway**

*T*

*a*

*b*

*l*

*e*

*4*

-

*2*

.

*S*

*e*

*l*

*e*

*c*

*t*

*i*

*n*

*g*

*a*

*n*

*A*

*P*

*I*

*G*

*a*

*t*

*e*

*w*

*a*

*y*

*C*

*h*

*e*

*c*

*k*  
*l*  
*i*  
*s*  
*t*

Decision	How should we approach selecting an API gateway for our organization?
Discussion Points	<p>Have we identified and prioritized all of our requirements associated with selecting an API gateway?</p> <p>Have we identified current technology solutions that have been deployed in this space within the organization?</p> <p>Do we know all of our team and organizational constraints?</p> <p>Have we explored our future roadmap in relation to this decision?</p> <p>Have we honestly calculated the “build versus buy” costs?</p> <p>Have we explored the current technology landscape and are we aware of all of the available solutions?</p> <p>Have we consulted and informed all involved stakeholders in our analysis and decision making?</p>
Recommendations	<p>Focus particularly on your requirement to reduce API/system coupling, simplify consumption, protect APIs from overuse and abuse, understand how APIs are being consumed, manage APIs as products, and monetize APIs</p> <p>Key questions to ask include: is there are existing API gateway in use? Has a collection of technologies been assembled to provide similar functionality (e.g. hardware load balancer combined with a monolithic app that performs authentication and application-level routing)? How many components currently make up your edge stack (e.g. WAF, LB, edge cache, etc.)</p> <p>Focus on technology skill levels within your team, availability of people to work on a API gateway project, and available resources and budget etc</p> <p>It is important to identify all planning changes, new features, and current goals that could impact traffic management and the other functionality that an API gateway provides</p> <p>Calculate the total cost of ownership (TCO) of all of the current API gateway-like implementations and potential future solutions.</p> <p>Consult with well known analysts, trend reports, and product reviews in order to understand all of the current solutions available.</p>

Selecting and deploying an API gateway will impact many teams and individuals. Be sure to consult with the developers, QA, the architecture review board, the platform team, InfoSec etc.

---

## Summary

In this chapter you have learned what an API gateway is, and also explored the historical context that led to evolution of the features currently provided by this essential component in any web-based software stack. You have learned how an API gateway is a very useful tool for migrating and evolving systems, and got hands-on with how to use an API gateway to route to the Attendees service that was extracted from the Conference System use case. You have also explored the current taxonomy of API gateways and their deployment models, which has equipped you to think about how to manage potential single points of failure in an architecture where all user traffic is routed through an edge gateway. Building on the concepts of managing traffic at the (ingress) edge of systems, you have also learned about service-to-service communication and how to avoid common pitfalls such as deploying an API gateway as a less-functional enterprise service bus (ESB).

The combination of all of this knowledge has equipped you with the key thinking points, constraints, and requirements necessary in order to make an effective choice when selecting an API gateway for your current use cases. As with most decisions a software architect or technical leader has to make, there is no distinct correct answer, but there often can be quite a few bad solutions to avoid.

Now that you have explored the functionality that API gateways provide for managing north-south ingress traffic and related APIs, the next chapter will explore the role of service meshes for managing east-west service-to-service traffic.

## **Prospective Table of Contents**

(Subject to Change)

## **Part I: API Fundamentals**

## **Part II: Designing, Building, and Testing APIs**

## **Part III: Traffic Patterns and Management**

## **Part IV: API Operations and Security**

## **Part V: Evolutionary Architecture with APIs**

## About the Authors

**James Gough** has worked extensively with financial systems and is the architectural lead for Client APIs at Morgan Stanley. He has a very pragmatic approach to building software and has been responsible for building API Gateways and applications to support a large enterprise API transformation.

**Daniel Bryant** works as a product architect at Datawire. His technical expertise focuses on DevOps tooling, cloud/container platforms, and microservice implementations. Daniel is a Java Champion, and contributes to several open source projects. He also writes for InfoQ, O'Reilly, and TheNewStack, and regularly presents at international conferences such as OSCON, QCon, and JavaOne. In his copious amounts of free time he enjoys running, reading, and traveling.

**Matthew Auburn** has worked for Morgan Stanley on a variety of financial systems. Before working at Morgan Stanley he has built a variety of mobile and web applications. Matthew's Masters degree primarily focused on security and this has fed into working in the security space for building APIs.