

1. Introduction

Application programming interfaces (APIs) represent an abstraction of the underlying implementation of a digital enterprise — they describe what actions users are allowed to take. Given the spatial (e.g., on-premises, multiple clouds) and logical (e.g., microservices) nature of current enterprise applications, APIs are needed to integrate and establish communication pathways between internal and third-party services and applications. The growing prevalence of microservice-oriented architectures and software as a service (SaaS), which are nearly always delivered via APIs, has resulted in network-based APIs being utilized across organizations in every type of application, including server-based monolithic, microservices-based, browser-based client, and the Internet of Things (IoT).

An API is a collection of commands or *endpoints* that operate on data or *objects* via some protocol to define how two pieces of software communicate. At runtime, service instances send requests to a specific API endpoint. An API gateway hosts many APIs and is responsible for mapping each request to its target API endpoint, applying policy for that endpoint (e.g., authentication, rate limiting), and routing that request to a service instance, which implements that API endpoint, as shown in Fig. 1.

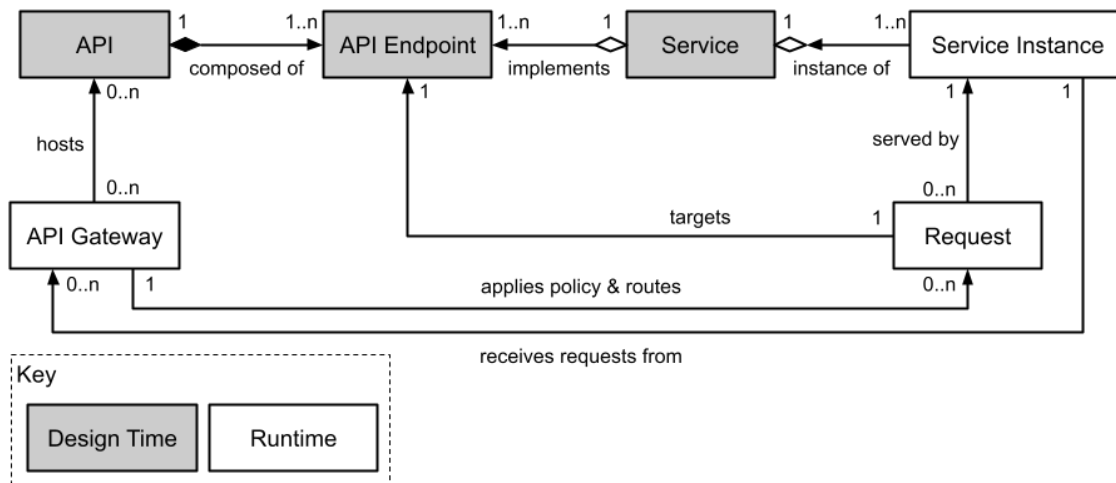


Fig. 1. API, API endpoint, service, and service instance

Network-based APIs are built to be consumed by remote applications over the network and present a unique set of challenges. Traditionally, network-based APIs are thought of as being customer-oriented, partner-oriented, or internal, which are referred to as “third-party,” “second-party,” and “first-party” APIs, respectively. First-party APIs can be exposed to callers inside of the organization on the same API gateway, but they are also often consumed directly by internal callers without traversing a dedicated API serving stack. Second- and third-party APIs are typically exposed to callers outside of the organization via an API gateway.

Most first-party API integrations occur via a service API (i.e., they map to a single service). However, APIs that are hosted by the API gateway typically have endpoints that map to many different services, especially for second- and third-party APIs. These are called *facade APIs*.

because they present a single facade to an outside caller over potentially many different service APIs. Multiple services are commonly grouped together into an *application* along organizational lines with an application mapping to a team. Fig. 2 shows a schematic diagram of a service API, a facade API, and an application (i.e., monolithic) API.

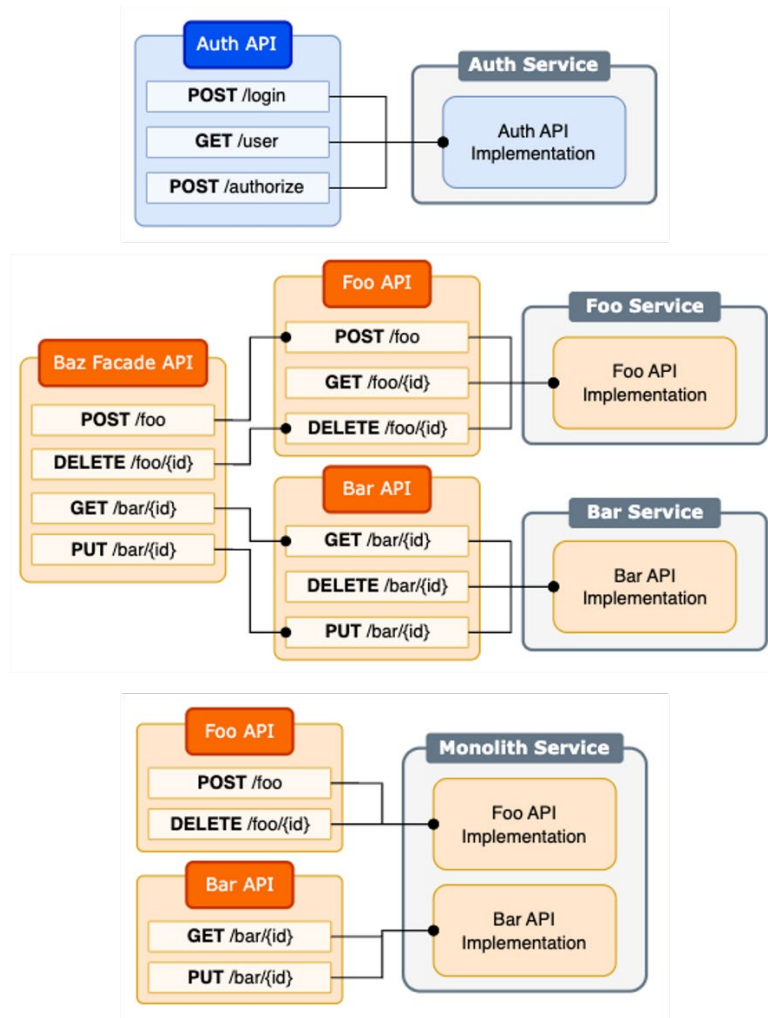


Fig. 2. Service API, facade API, and application (monolithic)

Whenever systems communicate, there is *some* API involved (e.g., Comma-Separated Values (CSV) over File Transfer Protocol (FTP)). While this Special Publication (SP) focuses on “modern” APIs that are exposed via mechanisms like Hypertext Transfer Protocol (HTTP)/Representational State Transfer (REST), gRPC Remote Procedure Calls (gRPC), or Simple Object Access Protocol (SOAP), the principles are universal and should be applied to *all* APIs and various communication styles (e.g., request/response, message-based/asynchronous).

1.1. Zero Trust and APIs: The Vanishing Perimeter

One of the most important implications of zero trust is that there is no meaningful distinction between an “internal” and “external” caller because the perimeter is the service instance itself.

That is, all callers are trusted if they are authorized to be trusted. This contrasts with traditional approaches to API security in which the only “APIs” are those exposed to “external” callers, and API-oriented controls are only enforced at the perimeter, typically via an API gateway.

SP 800-207A [6] discusses zero trust at runtime and the principle of shrinking the perimeter to the service instance using the five runtime controls of identity-based segmentation:

1. Encryption in transit — To ensure message authenticity and prevent eavesdropping, thus preserving confidentiality
2. Authenticate the calling service — Verify the identity of the software sending requests
3. Authorize the service — Using that authenticated identity, check that the action or communication being performed by the service is allowed
4. Authenticate the end user — Verify the identity of the entity triggering the software to send the request, often a non-person entity (NPE) (e.g., service account, system account)
5. Authorize the end user to access resources — Using the authenticated end-user identity, check that they are allowed to perform the requested action on the target resource

Achieving a zero-trust runtime requires applying these five controls to *all* API communications. Additional controls that are necessary for safe and secure API operations beyond identity-based segmentation should be enforced on all APIs in a system, including those exposed to the outside world (i.e., public APIs) and those intended only for other applications in a given infrastructure (i.e., internal APIs).

1.2. API Life Cycle

Like all software, APIs grow and change over time as requirements drift and usage patterns change. They also go through a continuous, iterative life cycle, including:

- Plan, Develop, Build, Test, Release — These “pre-runtime” life cycle phases lead to a service that can be deployed in production.
- Deploy, Operate, Monitor, Feedback — These “runtime” life cycle phases involve running and operating a service in production.

Department of Defense (DoD) Enterprise DevSecOps is an example of a software development life cycle paradigm. A detailed description of each phase of this paradigm is given in [1].

Application of the DevSecOps paradigm in the context of cloud-native applications can be found in [4][5].

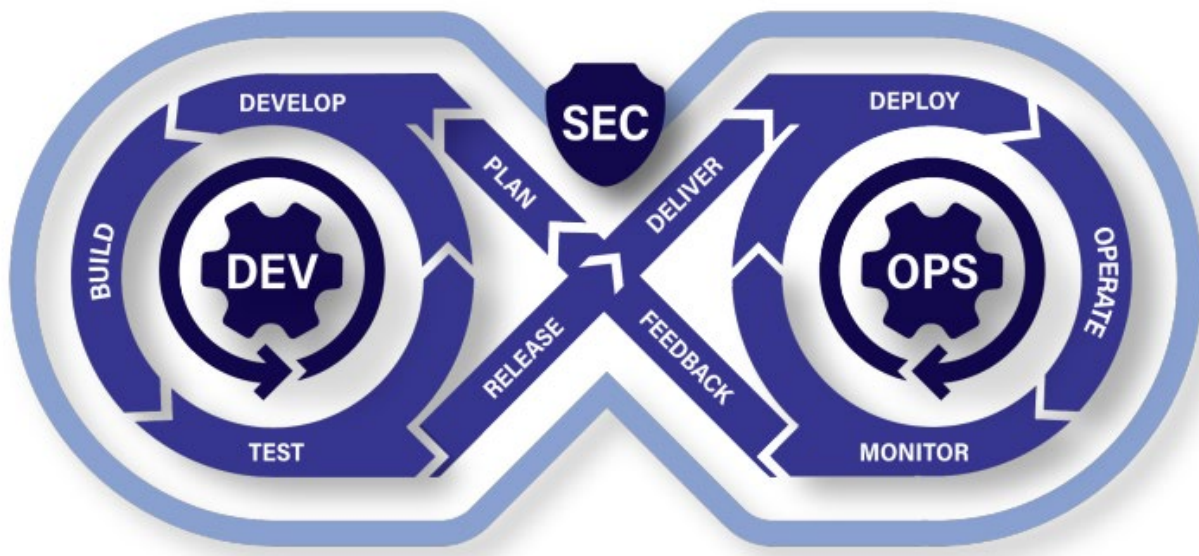


Fig. 3. DevSecOps life cycle phases

1.3. Document Goals

This document recommends controls for API protection except for API tools that embed one or more of the following control categories as part of its feature set:

1. Pre-runtime API protections — These controls need to be applied when designing and building APIs.
2. Runtime API protections — These controls need to be applied to every API request that an infrastructure serves, not just at the perimeter.

Each of these two categories is further divided into two subcategories based on organizational maturity (i.e., basic and advanced), which enables enterprises to adopt them using an incremental, risk-based approach.

A prerequisite for defining any API protection measure or policy irrespective of its category or sub-category is that the protections must be expressed in terms of nouns (e.g., resource) and verbs (e.g., “create customer record” using the verb POST/CR) that pertain to API components, API endpoint components, API requests, and API responses. These, in turn, contain references to resources (e.g., customer record [CR]), data, and operations on those resources.

1.4. Relationship to Other NIST Documents

Today, most software development and integration are based on APIs. Section 1.2 articulated the close relationship between software and APIs, demonstrated that API development and deployment follow the same iterative life cycle as the software, and provided NIST guidelines on DevSecOps.

Another distinguishing feature of the controls recommended for protecting APIs is the capacity to provide assurance for conforming to the principles of zero trust. This is because there is no distinction between internal and external API requests or calls due to the absence of an identifiable network perimeter and the distributed nature of applications on-premises and in multiple clouds. This security assurance can be achieved with authentication and authorization controls using identity-based segmentation [2]. Documents that provide recommendations on the configuration of authentication and authorization controls in the context of cloud-native applications (e.g., [2][3]) are also relevant in the context of configuring controls for API protection.

1.5. Document Structure

This document is organized as follows:

- Section 2 describes risk factors and vulnerabilities associated with APIs and the attack vectors that could exploit those vulnerabilities.
- Section 3 recommends controls to protect APIs and classifies them into basic and advanced categories that need to be applied prior to runtime or enforced during runtime.
- Section 4 provides a detailed analysis of implementation options or patterns for the controls described in Sec. 3 and outlines the advantages and disadvantages of each pattern.
- Section 5 provides the summary and conclusions.
- Appendix A provides the classification taxonomy for APIs.
- Appendix B illustrates the API controls related to each DevSecOps phase
- Appendix C provides a list of Limit Types

2. API Risks: Vulnerabilities and Exploits

This section considers some common risk factors that are associated with API deployments, including:

- Lack of visibility of APIs in the enterprise inventory [7]
- Missing, incorrect, or insufficient authorization [7]
- Broken authentication [7]
- Unrestricted resource consumption [7]
- Leaking of sensitive information
- Insufficient verification of input data

2.1. Lack of Visibility of APIs in the Enterprise Inventory

Most organizations have gaps in their API inventories, even if they otherwise have mature inventory management capabilities. The enterprise estate cannot be protected without an accurate API inventory, and unknown incidents may occur at the API level. Common reasons for the lack of visibility include:

- Organizational silos: APIs are built by many teams across an organization, deployed across cloud and on-premises environments, and inherited in mergers and acquisitions. Security concerns may not receive sufficient attention, and establishing accurate, up-to-date inventories may be difficult. Further, a lack of automation and integration with the API inventory management system exacerbates the challenge of maintaining an accurate inventory.
- Rogue or shadow APIs: APIs that are defined for internal use (e.g., debugging, testing, ad hoc solutions to business problems) may not be appropriately documented and often bypass standard security review practices.
- Zombie or deprecated APIs: APIs may have been replaced or superseded by newer systems but have not yet been entirely removed (e.g., because all callers have not yet migrated to the alternative, there no longer exists a team responsible for the system). They risk falling behind the latest security policies and protections.

2.2. Missing, Incorrect, or Insufficient Authorization

Authorization requires a high-reliability, low-latency system for making decisions about user access to resources at request time. Application developers must integrate their application with the same authorization system to keep it updated on users, resources, and permissions as the system changes over time (e.g., users create and delete resources, assign new permissions). Even then, developers may incorrectly enforce access decisions in their application code. In the industry-recognized catalogue of API risks, three of the top 10 (i.e., 1, 3, and 5) focus on authorization [7].

In line with identity-based segmentation, every service for an API endpoint should perform two levels of authorization: 1) service authorization and 2) end-user-to-resource authorization [6]. However, implementing both levels of authorization can still leave many APIs open to risk. Individual fields of a resource often need to be authorized independently of the resource itself. For example, if additional debug information is embedded in an internal field of the API object, that field should not be visible to external callers (i.e., callers not authorized to see privileged debug information).

There are at least three sources of authorization risks:

1. **Missing authorization:** There is no fine-grained, resource-level authorization present. For example, a legacy system may be operating under different access models (e.g., in a perimeter-based model, access *is* authorization), or there may be implementation bugs (i.e., an access check that should be enforced is not).
2. **Incorrect authorization:** The application performs an end-user-to-resource authorization check but fails because it checks the wrong end-user identity, the wrong permission, and/or the wrong target resource.
3. **Insufficient authorization:** The application performs a resource-level authorization that is successful, but the resource itself contains privileged information that is not intended for the level of access implied by access to the resource itself. This is often the root cause of leaking sensitive information (see Sec. 2.5).

2.3. Broken Authentication

Authentication is a prerequisite for authorization, particularly two aspects: the authentication system itself is robust, and the application uses the authenticated identities correctly. Risks that an authentication system needs to mitigate include [8]:

- Credential Stuffing is a type of brute force attack, where an attacker knows an account's name, and tries to brute force a ton of different passwords to unlock it. (They "stuff all the credentials they can find" into the victim system's authorization system in the hopes that something will stick). The attacker is able to carry out this since mitigation features such as rate limits, Captcha are absent.
- Brute-force attacks on a single account without mitigations, which is closely related to unrestricted resource consumption (see Sec. 2.4)
- Insecure practices, such as weak passwords, passing sensitive data in public channels (e.g., the URL), missing password validation for changes to sensitive account data, and using weak keys or poor algorithms to encrypt user data in transit and at rest
- Bad or incorrect token validation, including not validating at all, ignoring expiry, and using insecure signing schemes or weak signing keys

With a robust and secure authentication system in place, the application must use those credentials correctly. Risks to mitigate include:

- Missing authentication (e.g., tokens can be present but simply not checked), often due to a bug or misconfiguration in the application
- Weak or predictable tokens, default accounts, and default passwords (e.g., a hard-coded bootstrap account with the same username and password on all devices, test accounts with predictable names and weak/guessable passwords)

2.4. Unrestricted Resource Consumption

Services consume resources to serve APIs, many of which can affect external systems or the real world when serving an API call. The effects are an intended part of the business flow, but automation creates avenues for abuse by malicious users. Therefore, usage must be restricted to protect against malicious attackers abusing the system with a denial-of-service attack (DoS) or for its impact on external systems.

2.4.1. Unrestricted Compute Resource Consumption

Broadly, the risks associated with unrestricted compute resource consumption (e.g., memory, CPU, storage) are best mitigated via a combination of rate limiting, timeouts, circuit breaking (i.e., limits on the number of concurrent outstanding requests), bot/abuse detection, and application changes (e.g., reject file uploads over 20MB in size, return at most 10 items in response to a list request). These risks manifest as:

- DoS attacks via bandwidth saturation or resource starvation
- Unreliable performance due to resource utilization for one user or service that impacts others
- Cost amplification, in which an attacker can spend a small amount of resources (e.g., money, compute, bandwidth) to make requests that trigger a system to spend a much larger amount of resources to service the request

Even internal API consumption poses many of these risks. In most organizations, it is much easier for a developer to accidentally cause a DoS on an internal service than for an external attacker to maliciously cause such an attack. This is a potential security event that necessitates the need for a zero trust approach.

2.4.2. Unrestricted Physical Resource Consumption

Critical business operations can be impacted when an attacker targets software systems that control physical processes (e.g., SCADA systems). APIs may also result in text messages being sent to users, charges to credit cards, or the consumption of expensive third-party resources. For example, a common challenge seen by organizations that adopt AI is the accidental over-use of expensive AI APIs, which can result in large unplanned expenses for the business. These risks may manifest as:

- Impacts on business operations (e.g., damage to equipment and personnel, the creation of fake orders that require human effort to sort and remove)
- Impacts on customer relationships (e.g., scalpers automatically buying inventory to relist at a higher price elsewhere)
- Infrastructure co-opted for abuse or harassment (e.g., multi-factor authentication fatigue attacks, where an attacker triggers text spam to a user's phone via an SMS 2-factor authentication system [9])
- Unplanned expenses (e.g., consuming far more of a third-party service than planned by satisfying requests made by a malicious user)

These risks are best mitigated by a combination of rate limiting, quotas, spending policy controls in third-party software, bot/abuse detection, and application or business flow changes. Mitigations for both compute and physical resource consumption are similar. For compute resources, how users interact with a system should be limited. For physical resources, how the user interacts with a system *and* how a system interacts with external systems should be limited and considered early in the design phase. Mitigating these risks can sometimes require business flow changes.

2.5. Leaking Sensitive Information to Unauthorized Callers

Unintentionally leaking business data via APIs is closely related to missing, incorrect, or insufficient authorization (see Sec. 2.2). While correct and robust authorization should mitigate this risk, sensitive data can still be leaked from APIs via side channels. The two most common side channels exploited by attackers are response codes and error information, and common risks include:

- Enumeration of the resources (e.g., users, objects) in a system: This can have secondary impacts on the business, like revealing the customer set, information about product inventory, or the identity of employees in an organization. A common method of enumeration is enabled by services responding with "Not Found" status codes instead of "Permission Denied," allowing an attacker to distinguish between resources that exist (403) and those that do not (404).
- Revealing information about the internal implementation of the infrastructure to attackers: While security through obscurity is no security at all, it is still prudent to make it as hard as possible for attackers to discover an infrastructure's fine-grained specifics, which are often included in error messages (e.g., the exact versions of common software being run, internal names of systems for future pivot attacks).

2.6. Insufficient Verification of Input Data

Trusting unverified inputs is one of a major class of recurring security bugs in software. There are at least two levels of verification that APIs need to consider:

- Validating that the input is syntactically correct
- Ensuring that valid input is not malicious

2.6.1. Input Validation

A service must validate that each request (i.e., input) matches the API's definition, all expected fields are present and of the correct type, and no unexpected fields are present. For example, an API definition may say, "The 'name' field is required and must be a non-empty string less than 100 characters long," which must be verified at runtime on every request.

The lack of input validation results in a variety of risks, including:

- Impacting the availability of APIs
 - The "query of death" (QoD) [24] is a DoS attack via specially crafted requests that trigger pathological worst-case behavior in the server (i.e., the server itself may crash due to bad input handling).
- Invalid or malicious data being stored in the system, which can cause latent issues (e.g., failure to restart during recovery, crashes when accessing invalid records)
- Unanticipated error handling during request processing, which leaks internal information

2.6.2. Malicious Input Protection

While the input may satisfy "syntactic" validation, it also needs to be verified as non-malicious before it is used. Malicious input is any input that is syntactically valid but attempts to get the system to misbehave, potentially in a way that can be exploited to trigger an attack. Extending the "name" example above, a caller may send a request that contains a name field with a string less than 100 characters (i.e., valid), but that string may be a Structured Query Language (SQL) injection attack. Common risks include:

- Data leaks, which may lead to regulatory fines (General Data Protection Regulation (GDPR), California Consumer Privacy Act (CCPA)) or corruption (e.g., a SQL injection attack [7])
- Unanticipated or unrestricted resource utilization (e.g., an attacker automates account creation and uploads multi-gigabyte "profile pictures" to each account)
- Exposing a surface that attackers can use to pivot within the infrastructure or leverage to mount further attacks on others (e.g., by allowing servers to be used for server-side request forgery [SSRF])
- Cost amplification attacks, like the "billion laughs attack" (XML expansion) [10] or "zip bombs" (zip archive expansion) [11]

2.7. Credential Canonicalization: Preparatory Step for Controls

A common problem at the API gateway is handling the many different credentials that clients use to call APIs. For example, mobile apps use a certificate, clients use an API key and expect HTTPS, internal users may use JWTs, internal applications expect an mTLS connection with a SPIFFE identity, and others use HTTPS and a Kerberos ticket. All of them also need to convey the user's credential (e.g., OAuth Bearer token, a custom JWT, some trusted internal header). The combination is immense and challenging for application developers to perform correctly. As a result, organizations may only perform authentication and authorization at the edge via the API gateway. A solution to this problem is to standardize the credentials that an application sees at the API gateway — that is, to *canonicalize* them.

2.7.1. Gateways Straddle Boundaries

A gateway is something in an infrastructure that straddles a boundary and is typically the only way for traffic to cross that boundary. One of the most important policies that the API gateway enforces is authentication, ideally of both the user and the calling service.

Identity-based segmentation states that every server should authenticate and authorize both the calling service and the end user of every request and that those policies should be enforced at every hop in the infrastructure [6]. However, changing legacy systems to support new identities is often not possible. The challenge lies in implementing identity-based segmentation and support for both service and user identities without impacting other parts of the infrastructure.

API gateways can be used to draw a boundary around the parts of an infrastructure that perform identity-based segmentation. Within that boundary, all applications expect a standard set of credentials (e.g., user identity via a JWT in a specific header and service identity via a SPIFFE X.509 certificate). Common policies, practices, and tooling can then be used to ensure that all applications perform authentication and authorization correctly. Legacy schemes may continue to be used outside of the boundary. To reach inside, traffic must traverse a gateway that can canonicalize the incoming request's credentials into the expected form.

2.7.2. Requests With a Service Identity But No User Identity

Consider a batch job that runs nightly and touches data for many users. This is a risk because it requires special casing by the applications. For *some* service identities, end-user authorization is not required, but for all others, it is required. Any special casing increases the opportunity for incorrect or insufficient authorization.

The solution is to adopt service accounts that represent some system in a user identity domain. That service account can be for an internal system and, therefore, have permission to act on the data of many other users, or it can be for a user's applications with correspondingly fewer permissions. The API gateway can mediate with the user authentication system to exchange the service's runtime identity for a service account credential that represents the service in the user identity domain and attach that service account credential as the end-user credential to

requests that it forwards into the part of the infrastructure that supports identity-based segmentation.

Applications that perform identity-based segmentation will need to configure a policy for that service account user so that it can act on all of the data that the batch job previously used its service identity to access. At the same time, the application can remove any support for special data access without an end-user credential. Finally, the existing infrastructure can be leveraged to audit and manage both user and service access to data.

An implication of this is that all applications that attempt to implement identity-based segmentation without a user identity should adopt service accounts by changing their application code. This will simplify future migration into the identity segmentation domain and make the system more secure overall.

2.7.3. Requests With a User Identity But No Service Identity

Consider a cloud-provider API gateway that receives user traffic, terminates TLS, performs end-user authentication, and forwards requests to the infrastructure. The gateway enforces authentication, so some user credential is present. However, unless special care has been taken to communicate the service identity (e.g., via an API key or service account JWT), most notions of the calling workload will be lost at the external gateway provider.

Depending on the specifics of the setup, the only option may be to configure service identity-level policy via the external API gateway's controls and then implement fine-grained service-to-service policy for how requests can flow from that external gateway into the infrastructure. In other cases, the external gateway can be configured to pass some notion of the external workload (e.g., forwarding the client's certificate as a header) and then use that to create some canonical workload credential for internal communication (e.g., forwarding the client's certificate and creating a JWT that represents the external service identity from the certificate's common name).

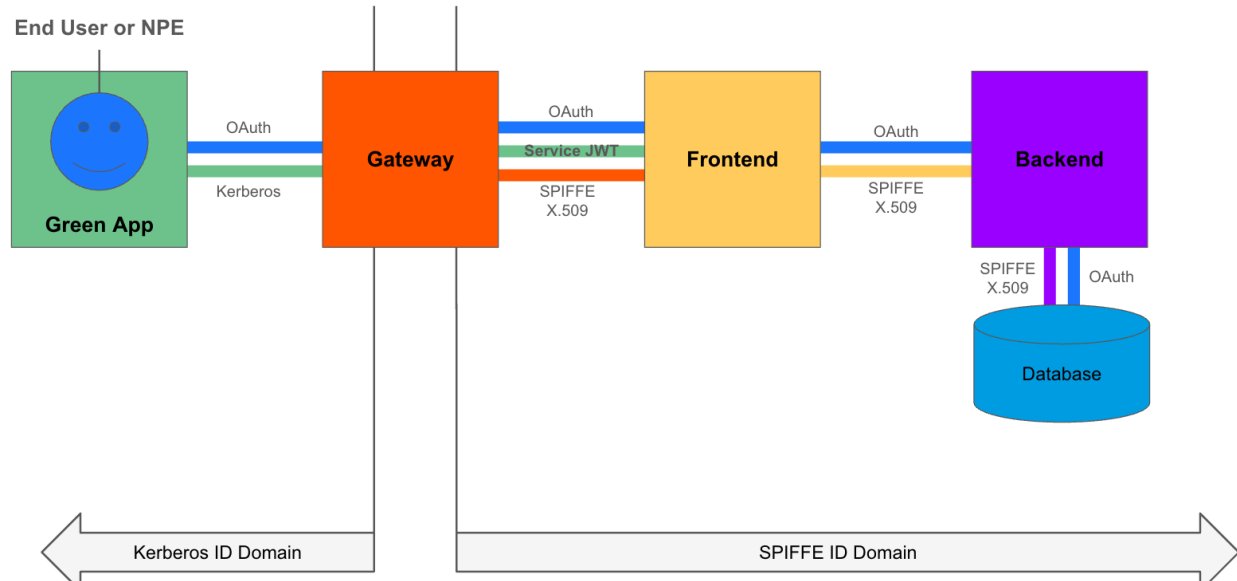


Fig. 4. Handling API calls with user identity but no service identity

However, the gateway's service identity is already in place between the gateway and the first service that performs identity-based segmentation. For that first hop, three identities need to be handled on the request: the gateway's service identity, the service identity of the external service, and the end user's identity. As before, external service authorization can be performed via the gateway and simply drop the external service identity. Services should support validating both the end user and a workload identity via metadata from the request in addition to validating workload identity via the transport (e.g., mTLS certificates).

For example, suppose that an organization A) uses a SPIFFE X.509 identity via mutual TLS for service identity as a service mesh does, B) uses a JWT bearer token for user identity, and C) chooses to represent an external service identity as a JWT token attached to the request. The mesh can then enforce that the gateway forward traffic to the service via (A), authenticate the service JWT and authorize the external service (C), and authenticate the end user (B) before forwarding a request to the application. This would fully support authenticating and authorizing all of the communicating parties, and the service in question would not need to be aware of the external service identity or credential. This is accomplished by Gateway bridging two identity domains (Kerberos and SPIFFE as shown in Fig.4). Gateways act as a policy enforcement point where we can "canonicalize" incoming credentials (e.g. a Kerberos ticket) into a standard form expected by internal systems (e.g. a SPIFFE Verifiable Identity Document -- "SVID"). The Gateways would simply need to manage a policy of "allowed external service callers" alongside their set of "allowed internal service callers."

2.7.4. Requests With Both User and Service Identities

In the best case, the legacy systems in question are already doing nearly the right thing in that they have both an end user and a service identity attached to requests. However, legacy system credentials likely do not fully conform to the credentials expected by the parts of the system

that implement identity-based segmentation. In that case, those other credentials will need to be translated into the canonical form expected by services that perform identity-based segmentation in the infrastructure. Essentially, the user's authenticated credential should be exchanged with an identity provider for the canonical form expected by the identity-based segmentation portion of the infrastructure (e.g., a JWT bearer token), and the external service's identity should be represented to the internal system as a token so that the policy can be enforced on all three identities in the first hop.

2.7.5. Reaching Out to Other Systems

A similar problem presents itself in reverse when a service that performs identity-based segmentation needs to reach out to legacy systems that expect legacy credentials. One option is to integrate modern applications with legacy credential systems so that those applications can fetch the legacy credential they need, which can significantly delay the sunsetting of those legacy systems. A better option is to perform a credential exchange on traffic leaving the identity-based segmentation subset of the infrastructure.

For example, an external SaaS API may expect a cloud provider service account as credentials. An egress gateway can be deployed to authenticate and authorize credentials that are used inside of the organization (i.e., identity-based segmentation) and exchange the internal identities for the external identities that are needed by the other system. In this way, services that perform modern identity-based segmentation can integrate with legacy systems with little impact and minimize any code dependencies on those legacy systems.

2.7.6. Mitigating the Confused Deputy

A 'confused deputy' is a type of privilege escalation where a privileged entity (the 'deputy') is tricked into using its authority on behalf of another, less privileged entity. One of the biggest risks in any scheme that involves credential exchange is a confused deputy [25], where one caller can trick the "deputy" responsible for handling credentials into using credentials that belong to another caller on its behalf, most often to escalate privileges. Any system that brokers multiple credentials needs more and better authentication and authorization before allowing credentials to be accessed.

An alternative approach is to break down the deputy into separate entities that hold only a single credential and map closely to a single application or service. This is the core idea behind the service mesh's sidecar presenting a service identity on behalf of the application: because the sidecar is one-to-one with a service instance, a service's identity cannot be confused for another at runtime. This same idea can be applied to API and egress gateways. Deploying them granularly — ideally per application — can minimize or eliminate any mixing of credentials, thereby mitigating any risk of a confused deputy. Section 4 discusses API gateway deployment patterns at length.

2.7.7. Identity Canonicalization

Canonicalizing credentials is really canonicalizing the identity domains for which one needs to write policy. Integrating identity providers to standardize credentials at the gateway inherently brings those identities into two identity domains: one for users and one for workloads. This allows for concise and consistent sets of policy that govern access to other services and user access to data. Having both policies in place implements identity-based segmentation and dramatically improves security posture.

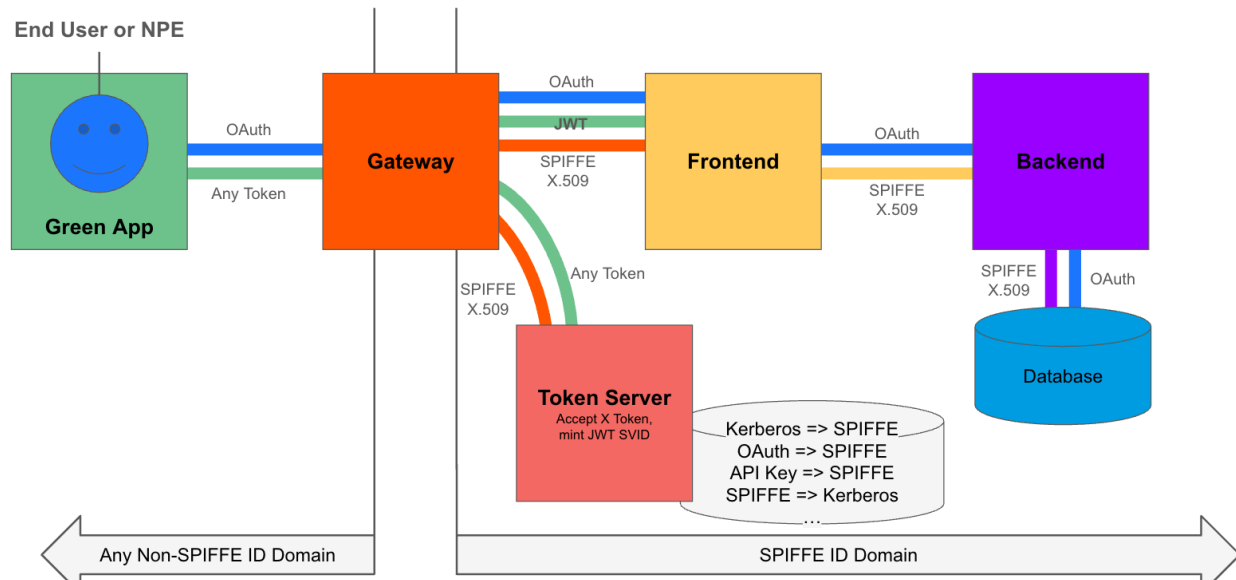


Fig. 5. Identity canonicalization for handling API calls

For most organizations, implementing credential canonicalization will require either adopting an identity provider wholesale and standardizing on that throughout (including working out legacy integration so that legacy credentials can be used to get credentials via the new provider) or performing identity exchanges, as described in this section. The API gateway is ideally situated to enforce either choice. Performing identity exchanges also requires a mapping of identities across domains as well as a “token server,” which uses that mapping to mint credentials. Fig.5 shows a Token Server that is a natural extension to the idea of canonicalizing credentials as they enter the system's boundary at the Gateway: rather than managing policy pair-wise for each potential type of credential that a Gateway may have to handle, we specialize and centralize that logic into a broker -- the Token Server -- that is responsible for mapping different types of identity tokens to (and from) the canonical credential the system expects (e.g. an SVID). An Egress Gateway (or the API Gateway itself) can also leverage such a token server to handle mapping from "internal" credentials (an SVID) to "external" credentials on behalf of an application (e.g. a Kerberos ticket or OAuth Bearer token).

3. Recommended Controls for APIs

In their earliest form, controls for APIs focused primarily on encryption in transit while delegating most other concerns to the application. Over time, a variety of challenges have emerged that necessitate the evolution of controls, including:

- The distributed nature of modern enterprise applications, which span multiple on-premises and cloud environments and communicate over the network using APIs
- The requirement to build robust systems that work around transient failures and handle large volumes of traffic
- An increasingly interconnected API with substantial system logic driven by business needs to integrate more deeply with partners and expose richer functionality to users
- Increasingly sophisticated attackers who have moved up the stack from low-level exploits and DoS attacks to application-level attacks that leverage the APIs that systems use to function

Controls for APIs should address all of the APIs in the organization, including those exposed to end users, those exposed to partners, and those that are only intended for internal consumption. This document's controls are structured into two primary sections based on the iterative API life cycle (see Sec. 1.2):

1. Pre-runtime protections, which should be applied during design, development, and testing. These include:
 - a. Creating a well-defined specification for the API's contract using some interface definition language (IDL) (e.g., OpenAPI, gRPC, Thrift)
 - b. Defining request and response schemas as part of that API specification
 - c. Defining valid ranges of values for fields of each request and response
 - d. Tagging the semantic type of each field of each request and response
 - e. Creating and maintaining an inventory of these API specifications across the organization, including ownership information
2. Runtime protections, which should be applied to each request and response to the API at runtime. These include:
 - a. Encryption in transit
 - b. End-user authentication and authorization
 - c. Service-to-service authentication and authorization
 - d. Request and response validation
 - e. Resource consumption mitigations, including rate limiting, timeouts, and circuit breaking

- f. Telemetry (e.g., logging and monitoring) to assess enforcement and detect attacks with minimum log data fields that include a timestamp, client ID, request/response status, and endpoint identifier

Within each section, the controls are grouped into “basic” and “advanced” categories:

- Basic protections should be pursued immediately with the goal of obtaining basic insight into the APIs that exist in an organization (*Identify* in NIST Cybersecurity Framework (CSF) [12]) and can be used to implement essential best practice controls (*Protect*). Generally, basic protections do not require deep introspection of the API’s request and response payloads but operate at the connection or request metadata level (i.e., on HTTP headers rather than the HTTP body).
- Advanced protections perform deeper analysis on requests and responses. Many of these policies require payload inspection, which is CPU- and latency-intensive. The goal is to enhance basic *Protection* and begin to cover the *Detect* and *Respond* functions in NIST CSF [12]. Addressing these concurrently with basic controls is recommended, but the basic protections may provide the most benefit for resource-constrained organizations.

All organizations should move immediately to act on basic controls, while advanced controls should be evaluated by the organization and applied to APIs based on risk profile.

3.1. Pre-Runtime Protections

All API controls must be well-defined and inventoried.

3.1.1. Basic Pre-Runtime Protections

REC-API-1: All APIs must have a specification in the form of a document that describes what endpoints the API exposes (“API spec” for short). To begin, the API spec can be a literal document, a set of internal wiki pages maintained by a team, or something similar. However, it should eventually migrate to a state-of-the-art IDL.

REC-API-2: API specifications should use a well-defined IDL (e.g., OpenAPI for HTTP/REST, gRPC for protobuf, Thrift, SOAP for XML).

- REC-API-2-1: API specs and implementations should conform to industry best practices (e.g., a Create-Read-Update-Delete [CRUD] API exposed as HTTP/REST should map the CRUD endpoints to the HTTP verbs POST, GET, PUT, and DELETE, respectively) for consistency [13]. For operations that do not cleanly map to CRUD, guidance should be provided to maintain consistency and avoid misuse of HTTP verbs.

REC-API-3: Request and response schema for each endpoint should be defined by the API specification, including validation guidelines for the values of each field of the request and response (e.g., “the name field is a string and must be shorter than 100 characters”). Additional information makes integration easier and less error-prone for clients and presents the

opportunity for automated enforcement, such as the maximum latency (e.g., “the server will drop requests that take longer than 5 seconds to process”) and rate limits (e.g., “by default, 5 calls per minute are allowed”).

REC-API-4: Establish a centralized API governance framework. The framework should incorporate the following functions at a minimum:

- REC-API-4.1: A centralized API platform team model should be created in which a central team provides a self-service platform for application teams to build APIs
- REC-API-4.2: Organizational API inventory of all internal (including shadow and Zombie APIs) and external APIs should be maintained. This is in line with the *Identify* directive of the CSF [12]. The documentation requirements for the inventory should include:
 - Each API’s specification, though the inventory does not need to be the *API documentation*
 - Ownership information about the API to simplify the translation of runtime problems to organizational response
 - Runtime information to enable operations and security teams to understand the impact of each API (e.g., service instances, instance IP addresses, runtime service ID, traffic volume, rate of requests and errors, the status of policy enforcement)
- REC-API-4.3: Use tools to improve visibility, such as **API discovery tools** and scheduled automated scans to detect all running APIs.
- REC-API-4.4: Establish strategies regarding API versioning, deprecation, and sunseting, including a focus on secure migration paths.

3.1.2. Advanced Pre-Runtime Protections

REC-API-5: Request and response validation in the schema should be included in the API’s specification (e.g., a string field must be non-empty and shorter than 255 characters, or an integer value must be non-negative and less than 2 million). This simplifies documentation and enables runtime tooling to validate request and response schema and syntax.

- Use primitive types in API schemas to reinforce this. For example, if a value is always semantically positive, model it in the schema as an unsigned integer rather than a regular integer (e.g., protobuf’s “uint” rather than “int”). Negative values are then disallowed by construction without any validation needed [14].
- This principle extends to zero or default values as well. Users (malicious or not) will frequently omit fields that the application expects. One approach to this is annotating fields as “required” or “optional” and rejecting requests with zero values for required fields. However, the application must handle missing optional fields. A second approach adopted by both Golang and protobuf/gRPC is to define “zero values” for each primitive type. Application code must either handle the zero value for each field or reject the request with a validation error.

REC-API-6: Annotate each field as public or internal for each request and response or with the level of trust or permission required for access. These annotations simplify documentation and enable runtime tooling to remove trusted data for untrusted callers as a cross-cutting policy rather than something that must be built into the business logic of each service. An in-application approach is much harder to implement correctly and to audit in practice.

- REC-API-6.1: Annotate endpoints and fields with required permissions to enable the use of tooling to automate fine-grained per-field authorization checks. Those authorization checks could then be performed by the API serving infrastructure on behalf of the application or via a common library in the application with standard logging and metrics to facilitate easy audits and ensure continuous enforcement. Once the annotations are present, a variety of runtime implementations are possible.

REC-API-7: Annotate each field with its semantic type to indicate fields that contain sensitive information, such as personally identifiable information (PII), protected health information (PHI), or payment card information (PCI). This enables runtime systems to track data flow through the system, trigger alerting, and apply cross-cutting policy to ensure that data does not leak across inappropriate boundaries.

REC-API-8: Include runtime information in the API inventory with ownership (REC-API-4). This becomes substantially more valuable when annotated with runtime information (e.g., service instances and their IP addresses, runtime identities of the service instances, metrics or health information for the service, runtime metrics for traffic between services). This information can help security identify the blast radius of an event, operations to identify problems and root causes, and application teams to understand their application's behavior. Correlating this information with the APIs being served makes it simple to link clients to servers as the problem is traced back to its root.

- API discovery tools must be deployed during runtime and reconciliation between declared specifications, and live traffic must be performed as part of maintaining an accurate API inventory. The objective of this type of discovery and reconciliation is to identify differences between what is deployed and what should be deployed in production, including shadow, orphan, and zombie API endpoints.

3.2. Runtime Protections

For runtime protections for APIs, apply zero trust principles as a baseline, and augment them with additional policy on requests and their payloads.

3.2.1. Basic Runtime Protections

REC-API-9: All runtime communication must be encrypted, even when the API is "public data" or otherwise unauthenticated. This is necessary to ensure that data has not been tampered with (integrity) and to prevent eavesdropping (confidentiality). Details on encryption in transit can be found in SP 800-53, control SC-8 [15] and SP 800-207A, control ID-SEG-REC-1 [6]. Details

on cryptographic algorithms and key lengths can be found in SP 800-57 [16] and Federal Information Processing Standards (FIPS) publication 140-3 [17].

REC-API-10: Perform general request and response validation policies (e.g., WAF, bot detection, DoS mitigation) to mitigate malicious payloads and unrestricted resource consumption. These can and should be executed early in the API serving stack to protect other components (e.g., authentication system) from DoS attacks. Since these protections are general and cross-cutting, there is little risk of unintentionally leaking sensitive information.

REC-API-11: Ensure the robust authentication and correct processing of authentication credentials.

- REC-API-11-1: Protect against credential stuffing and other brute-force attacks:
 - Implement rate limiting and account lockouts after repeated failed login attempts.
 - Enforce multi-factor authentication (MFA) to prevent account takeovers.
 - Use bot detection and CAPTCHAs to prevent automated attacks.
 - Implement adaptive authentication that adjusts security controls based on login behavior and risk level.
 - Leverage credential screening services (e.g., Have I Been Pwned API) to detect compromised credentials.
- REC-API-11-2: Authenticate the calling user and service, as described in SP 800-207A, controls ID-SEG-REC2 and ID-SEG-REC4 [6].
 - There are (at minimum) two identities in every API communication: the software calling the API and the end user of that software. For example, it is common to use an API key to identify calling software and an OAuth Bearer token to identify the end user. This is true even if the end-user identity is an NPE (i.e., internal software calling other internal software should use something like a service account to identify the user making the requests). The service identity may contain information (e.g., the device being used to access the system) in addition to a token from the software itself (e.g., an API key).
- REC-API-11-3: Identities (e.g., tokens) must be cryptographically verifiable and should not use weak signing algorithms (e.g., no JWTs with “alg: none,” weak algorithms, or short key lengths) or long expiration times (i.e., credentials are cycled regularly). SP 800-57 [16] discusses the strengths of cryptographic algorithms and the necessary key lengths for each. Token signing keys must also be rotated periodically to prevent token forgery attacks.
- REC-API-11-4: Authentication should use standard mechanisms whenever possible. For example, end-user authentication should use a mechanism such as OpenID Connect (OIDC), OAuth2, or SAML. Services should use a mechanism like SPIFFE SVIDs, JSON Web Tokens (JWTs), API keys, or similar.

- REC-API-11-5: Use opaque tokens for untrusted systems. Credential tokens commonly encode information about the internals of the system (e.g., minting a JWT to represent a user in the infrastructure that includes claims that represent the user's capabilities in the system) to simply and reliably enforce authorization per hop (e.g., validate the JWT, and check whether it contains the "claim" that represents the permission for an API endpoint). These claims encode all local operations that can be performed with data from the request and the local application. Returning a token with these details to an external user may risk leaking information about the internals of the system. This is where the following issues become critical to the safety of the API: how permissions are modeled, the set of internal permissions/claims that map to a given external API endpoint, and information about the path that the request traverses through the infrastructure.
- REC-API-11-6: Secure the storage and transmission of tokens using robust encryption algorithms (e.g., AES-256) during transit (e.g., TLS 1.2+). Authentication tokens should be sent using securing HTTP headers, not URLs.
- REC-API-11-7: Verify signatures, and check for expiry during token validation. For example, when processing JWTs, the "exp" claim RFC 7519 [18] must be checked. Similarly, when processing an X.509 SVID, check the validity period's "Not Before" and "Not After" [19].

REC-API-12: Authorize the calling user and service for each identity on the request, including whether the calling software system is allowed to access the API endpoint and whether the end user is authorized to take the action on the resource represented by the endpoint (see SP 800-207A [6], controls ID-SEG-REC2 and ID-SEG-REC4).

- REC-API-12.1: Use access control models (e.g., attribute-based access control) to achieve fine-grained (granular) service-to-service authorization.
- REC-API-12.2: Use standardized authorization schemes (e.g., OAuth 2.0 or JWT) for end-user-to-resource authorization.
- REC-API-12.3: Implement an authorization auditing tool to regularly check for missing or weak authorization mechanisms (see Sec. 2.2). Additionally, use unit and integration tests that identify exposed data to ensure that role permission assignments are consistent with data sensitivity levels.
- REC-API-6 discusses annotating each request or endpoint with the permission required by the end user to call that endpoint on a resource. Runtime tooling can then be implemented to ensure that those annotations are transformed into runtime permission checks against the authorization system. Combined with a robust DevOps process to ensure that annotations are present on APIs before they can be deployed, there can be a high degree of assurance that the correct authorization is being performed at the platform level. SP 800-204B discusses using the service mesh to achieve this [3].

REC-API-13: Validate each request and response per the API schema before it is processed by the business logic (e.g., ensure that the request has a “name” field that is a string and no other fields). This ensures that applications only receive well-formed input (called client-side validation) and minimizes a class of errors and data leaks due to inline validation in the business logic. Additionally, validate that each response from the server (called server-side validation) conforms to the expected response schema to help prevent a variety of data leaks, abuses, or mistakes. Responses after the repeated invocation of APIs are checked to ensure that implementation conforms to API specification.

REC-API-14: Authenticate, authorize, then validate in that order to minimize the risk of leaking data to attackers, since validation messages are at especially high risk of leaking information. For example, rejecting a request with a validation error for using a duplicate user-supplied name as another user may unintentionally leak information to callers regarding the existence of a resource. A likely mitigation may be an underlying per-user segregation of user-provided data, which often requires business logic changes in the application. Generic validations (REC-API-10) are exceptions to this because they are not business logic-aware and do not risk leaking information. They can be safely implemented by the platform ahead of authentication, which is often desirable to help protect the authentication and authorization systems from DoS and other attacks.

REC-API-15: Enforce limits on API and resource usage. API gateway teams should provide reasonable defaults for the organization, and application teams should be able to enforce more fine-grained limits in their application or leverage the platform. Those limits should include:

- REC-API-15-1: Rate-limit all API access for all callers to ensure fair utilization across users, help with capacity planning, and mitigate the risk of unrestricted resource consumption. See REC-API-16 for recommendations on specific rate-limiting implementations.
- REC-API-15-2: Apply timeouts to all requests, including the API gateway. This should be done at the TCP level, where connections are automatically timed out after a modest time (e.g., 5 minutes) rather than the operating system’s default of more than one hour per connection. Timeouts should also be configured at the application level. If a required operation should complete in five seconds as part of the API contract, set a 6-second timeout for it. This ensures that the resources in a service do not wait for a response that will never arrive.
- REC-API-15-3: Apply bandwidth and payload limits to enforce maximum request and response sizes. The “correct” limit is highly contextual and based on the organization and application (e.g., a bank will have very different expectations than a video streaming company). This helps avoid a variety of risks related to malicious input and DoS.
- REC-API-15-4: Validate and limit user-supplied query parameters (e.g., amount of processing done, size of their response based on user input), especially in the context of what the system can support and what is typical for users of the system. For example:

- The number of elements returned per page of a paginated list API. If a typical user has 100 items, cap the maximum number of elements per page to 1000.
- Time ranges in dynamic queries. If a system is intended for viewing recent events, and the user can provide a time range, limit that range to the last 30 days rather than allowing the user to query “from 1972 onward.”
- GraphQL and similar API facade systems that support query languages over many APIs. Enforce limits on the queries that users can execute (i.e., approved or predefined queries only) and caps on the number of outbound calls that are allowed in the execution of a single query.

REC-API-16: Rate limiting recommendations¹ are one of the most effective tools to mitigate unrestricted resource consumption and can increase the attacker’s cost of many attacks that try to leak sensitive information via data exfiltration from API calls (e.g., scraping all chat logs from an organization with a script impersonating a chat client). Most organizations apply some type of rate limit to external traffic, but it is equally important to rate-limit internal callers. Poorly conceived code can unintentionally cause a DoS on an internal system. It is equally critical to consider the limits placed on internal software that call out to external systems (see Sec. 2.6.2). The following recommendations on rate-limiting configuration address common pitfalls and misunderstandings:

- REC-API-16-1: Rate limits are not quotas. A quota is a usage limit on an API over an extended duration (e.g., per month) that is associated with a user’s payment or billing structure. Many organizations have “API usage tiers” that map prices to higher per-month limits. These quotas need to be strictly enforced and are typically used to generate billing reports that are sent to customers. In contrast, rate limits are intended to protect the system from overuse and help ensure fair usage across separate, concurrent callers. Rate limits do not need to be exact in the way that quotas must be.
- REC-API-16-2: Rate limits for total load provide little benefit and should be dimensioned by user (e.g., 83 requests per 5 minutes per user) using the source IP address or end-user credential as the key. Rate limits without a user dimension (e.g., service can receive 1,000 requests per 5 minutes total) are not particularly effective and allow some users to impact others (e.g., DoS risk). This is true even when total limits are dimensioned by service instance (e.g., a single instance cannot receive more than 100 requests per 5 minutes). Circuit-breaking functions must be used to provide protective limits on concurrency for a service instance. More information on circuit breaking and other resiliency and load-shedding techniques can be found in SP 800-204A, Sec. 2.3 [2].
- REC-API-16-3: Rate limits should be short in duration (e.g., per 60 seconds, per five minutes). A rate limit is defined as the number of calls allowed over a time period (e.g., 24,000 requests per 24 hours; 1,000 requests per hour; 16.5 requests per minute). Most systems allow for the configuration of both the number of calls and the amount of time over which they are allowed. Shorter time limits allow clients to experience a few

¹ Appendix C provides a non-exhaustive list of potential limit types.

intermittent failures every one or five minutes as their traffic grows organically rather than total failure with 24-hour limits. Additionally, the system will experience smoother traffic overall because a single client must pace their consumption over a longer duration, resulting in less load from each client at any given time.

- REC-API-16.4: Rate limits should be defined per user, service, or network parameter (e.g., IP).
- REC-API-16-5: Adopt robust threshold-setting strategies to avoid arbitrary or ineffective limits, such as token bucket or leaky bucket algorithms or adaptive limits based on user behavior or service tier.
- REC-API-16-6: Set limits for concurrent requests.
- REC-API-16-7: Deploy real-time monitoring tools (e.g., Open Telemetry) to detect sudden traffic spikes and send alerts. Set hard limits on third-party API consumption.

REC-API-17: Fine-grained request and user blocking allows the API serving stack to block individual users via their end-user credential and/or network address, which enables an effective response during an ongoing incident (see the *Respond* function in the CSF [12]). The actual enforcement can be handled by separate components (e.g., network-level blocking implemented by a firewall or the load balancer; credential-level blocking implemented by the API gateway, bot/abuse detection systems, or the authorization system). For relevant information on these techniques, refer to SP 800-53, AC-3 [15] and SP 800-204B, Sec. 4.6 [3].

REC-API-18: API access must be monitored to ensure that the API serving stack provides sufficient telemetry to assess the availability of APIs and ensure that policies are being enforced. The traditional triad of logging, metrics, and distributed traces is recommended. All three should be tagged with information about the API being accessed in addition to the runtime service so that service calls can be traced back to APIs. For the API gateway itself, a range of signals should be produced to enable the identification of:

- Basic communication information, like the information included in the Common Log Format [20] (e.g., who called, what method, from what origin)
- Health (e.g., rate of requests, rate of errors, latency) per API and API endpoint
- Enforcement results per policy class (e.g., requests allowed or denied due to missing or incorrect authentication or authorization checks, requests blocked due to rate limiting) to assess the aggregate enforcement of each policy
- The health of the services behind the API gateway

Fine-grained blocking by network addresses and user credentials (REC-API-17) that are augmented with blocking based on monitored data (REC-API-18) (e.g., payload pattern or anomaly score) is essential for real-time incident response and containment, especially for APIs that handle sensitive data. General information on audit and logging controls can be found in SP 800-53 [15], AU-2 Event Logging, AU-3 Content of Audit Records, and AU-12 Audit Record Generation. Information on service mesh telemetry, which can be used for audit and logging, can be found in SP 800-204A [2], SM-DR21 through SM-DR24.

3.2.2. Advanced Runtime Protections

REC-API-19: Field-level validation using API schema annotations can be used to validate the values of requests and responses at runtime. This is beyond the basic syntactic validation of REC-API-13 (e.g., “there is a name field, and it is a string”) and more like semantic validation (e.g., “the name field must not be longer than 100 characters,” or “the amount field must be positive and less than 2 million”). This can be implemented by the API gateway as part of a cross-cutting policy. An API spec is required (REC-API-2) and should be in a central inventory (REC-API-4). The API gateway team can then enforce the validation of all requests traversing an API gateway. This reduces the risks of insufficient input verification and leaking sensitive information compared to ad hoc, error-prone implementations in each application or standard implementations embedded in the application itself via SDK, which tend to be difficult to update. A timely update is an imperative for infrastructure that enforces security policy. The Controls REC-API-13 and REC-API-19 (Syntactic and Semantic Validation of API Schemas) should occur at every hop within the infrastructure, including internal service-to-service communications. This ensures that requests and responses are continuously validated, regardless of their origin, and helps prevent lateral movement and bypasses inside the environment.

REC-API-20: Authorization and filtering using API schema annotations enforce access to resources and fields per caller. The API gateway itself is the policy enforcement point, and it defers to an authorization system to make decisions. The information from the API schema is enough to extract credentials from the request, identify the target endpoint and its associated tags/permissions, and use them to form a call to the authorization service (e.g., “is the request’s end user allowed to perform the endpoint’s permission on the object targeted by the request?”). The API gateway can then enforce the result of the call at runtime. There are at least three levels of assurance that can be achieved, and each build on the previous one to further mitigate risks at increased runtime or development-time cost:

- REC-API-20-1: Resource-level authorization as a cross-cutting policy should be enforced on all requests using endpoint-level annotations that define the permissions required to call the endpoint (REC-API-6.1). This can be done at the platform level by leveraging the API gateway. When combined with a distributed gateway pattern (Sec. 4.3), this implements ID-SEG-REC-4 [6] at every hop.² This also helps prevent and potentially eliminate missing authorizations (Sec. 2.2), depending on the organizational guardrails in place. For example, an organization can build an API inventory by mandating an API spec with endpoint-level permission annotations as part of each app’s “ticket to the platform” (i.e., the data that an app team needs to submit to run their application on the organization’s infrastructure and platform). Combined with standard patterns for authentication (REC-API-11), this can ensure that the correct authentication and authorization are performed. However, additional organizational controls are required to ensure that the permissions are correct and sufficient to fully mitigate the risks around authorization (see Sec. 2.2).

² Other patterns have a wider perimeter and are susceptible to the API gateway being bypassed. Therefore, they do not satisfy ID-SEG-REC-4.

Achieving correct and sufficient authorization at the resource level is likely all that most organizations need to achieve. It mitigates the predominant risks identified by the OWASP API Security Top 10 [7] with respect to authorization. Moving beyond this level of assurance into REC-API-20.2 and REC-API-20.3 shifts the focus to mitigating the risk of leaking sensitive information.

- REC-API-20-2: Field-level visibility as a cross-cutting policy can leverage basic “public” and “private” annotations on each field. The authorization check effectively asks whether data should be visible to external callers.³ These coarse-grained public/private annotations are particularly effective on common types that are shared across many APIs in the organization. For example, a standard error reporting pattern used by all APIs can leverage field-level annotations to differentiate user-facing errors versus developer-facing errors to mitigate the risk of leaking sensitive information via errors. The gRPC Status proto [21] is an example of a consistent error reporting pattern. In the gRPC case, field-level annotations would reside in the message used for the status’s “details.”
- REC-API-20-3: Field-level authorization can be leveraged as a cross-cutting policy (REC-API-6.1). This extends the idea of REC-API-19.1 down to the level of each individual field of the response and allows for the filtering of API objects per use to implement sophisticated access control schemes.
 - While this kind of approach offers a very high level of data security, it causes a sharp increase in the number of policy checks that the authorization system must perform and requires active participation by application developers to keep permissions per field up to date as the application evolves. For example, a resource-level authorization check requires one authorization decision per request. A field-level authorization check requires one authorization decision for the request plus an additional decision for each field of the response. Even an object with a modest number of fields (e.g., 5) results in whole-number multiples more policy decisions made by the authorization system. For developers, the purpose and permission of an endpoint rarely change, but the fields of the request and response objects for that endpoint regularly evolve over time. This makes upkeep for permissions at the field level more expensive for application developers versus endpoint-level annotations (REC-API-19.1). As a result of the cost and load on the authorization system, this level of fine-grained checking is typically only used in the most high-risk situations and only by sophisticated organizations.

SP 800-204B [3] discusses the advantages of using a distributed API gateway architecture when implementing fine-grained authorization checks. When choosing to implement these authorization policy checks under the centralized and hybrid patterns, care must be taken to ensure that the gateways are not bypassed. For example, a service-level authorization policy

³ REC-API-19.1 focuses on requests, while this control focuses on the data that an application returns to callers in responses. They are complementary controls.

could disallow any traffic except from the API gateway as a means of defeating an attempt to bypass gateway checks via pivoting inside the infrastructure.

REC-API-21: Traffic monitoring and policy using semantic field labels can log and monitor the flow of sensitive data in a system. Further, the API gateway can be used as a policy enforcement point to control the flow of that data, potentially blocking traffic flows that transit significant amounts of data. Ultimately, with annotations and enforcement in place, the flow of sensitive data in the organization can be governed by mandatory access control (MAC) policies. A MAC policy is enforced by the authorization system, regardless of the user or resource in question. For example, a MAC policy may require PCI data to be isolated from systems that do not implement Payment Card Industry Data Security Standard (PCI DSS) controls to maintain security and prevent potential breaches. Such a MAC policy can be enforced with a combination of PCI-compliant services in the infrastructure and data tags on the semantic types of data that flow through the system.

REC-API-22: Non-signature payload scanning for generative AI APIs analyzes request and response data for sensitive information that may not be a literal attack signature. Tools typically analyze (e.g., via regression, AI, simple matching, and word filtering) the responses returned by servers to score the risk that they contain sensitive information and take action to block that traffic. Increasingly, AI agents are being deployed to assess the risk of data generated by other agents. At a high level, this technique is like a web application firewall (WAF), but WAFs are fundamentally signature-based, while these analyses are fundamentally content-based.

This is a general category of data egress analysis that is relevant across all APIs, but it has become increasingly important with the growth of generative AI. Generative agents are frequently trained on business-sensitive data, have insight into sensitive business operations and operational data, and are increasingly exposed to the organization and externally as APIs. Since the inception of generative AI agents, a variety of prompt injection attacks [22] have been created to exfiltrate data.

Tools for performing non-signature payload inspection should be used whenever an organization is handling data returned by their system, especially when that data is generated on demand (e.g., by AI agents). In most cases outside of dynamically generated output, implementing simple semantic and syntactic validations (REC-API-13, REC-API-18) will typically provide an organization with more risk mitigation for a lower runtime and operational cost.

- REC-API-22-1: Semantic data discovery tools are typically very good for identifying the type of information flowing through a system (e.g., string, email address). Building the inventory of APIs and adopting well-defined API schemas with meaningful annotations takes time. Such runtime tools are helpful for initial discovery and ensure that rollout is complete across all services and that services stay in compliance after the policy is rolled out. When it is reasonable to leverage due to compute and latency constraints, an organization benefits from inspecting traffic for sensitive data flow, even beyond field-level annotations.

API payload scanning/inspection (REC-API-22) must be extended to include behavioral analysis of API sessions, particularly for APIs that expose generative AI or dynamic content.

REC-API-23: Design error codes such that they do not provide a means for resource enumeration (e.g., return an error code 403 Forbidden instead of distinguishing between missing and unauthorized resources).

REC-API-24: Block resource enumeration attacks through rate limiting, including anomaly detection.

REC-API-25: Limit the exposure of sensitive data using data masking in responses and logs.

REC-API-26: Fine-grained blocking for specific requests can prevent a DoS or service crash. These bad inputs can often trigger a cascading failure [23], but the queries may not be malicious in nature (e.g., users using the system in ways that it was not intended or designed for, such as QoD) [24]. These tools help mitigate the risks of unrestricted resource consumption and malicious input validation. For APIs that expose generative AI or dynamic content, the entire session must be blocked. Depending on the complexity of the query and environment, it may be possible to leverage a WAF or non-signature payload scanning tools to block some types of QoDs. However, application code changes may be required — sometimes even rearchitecting the application itself — to mitigate the impact of these kinds of queries.

The detailed controls in this section fit into broad classes, and their association with the DevSecOps phases is discussed in Appendix B. This emphasizes the observation that APIs should be treated as any other software and go through an iterative, continuous life cycle.

4. Implementation Patterns and Trade-Offs for API Protections

Regardless of the mechanism or architecture of an API and its services, the following capabilities are required to realize the controls outlined in this document:

- Authentication and authorization
- Request and response validation
- Rate limiting
- Circuit breaking
- Error handling
- Logging and monitoring

In addition to these core capabilities for security, APIs that serve infrastructure typically deal with other common concerns, such as:

- Service discovery
- Routing
- Protocol conversion
- Caching

The following components are often deployed to provide the above functionality to an API:

1. A *gateway* to implement the API-oriented policy
2. The *service* itself to implement the API's business logic
3. A method to get traffic to gateway instances (e.g., DNS and a network load balancer) to facilitate service discovery, load balancing, and network reachability to horizontally scaled instances of the gateway itself

For example, if the gateway functionality is implemented via a Kubernetes ingress routing to a pod (i.e., the service instance), then callers outside of the network will require the cloud provider or data center network team to provision a network load balancer in front to route network traffic to the Kubernetes load balancer service.

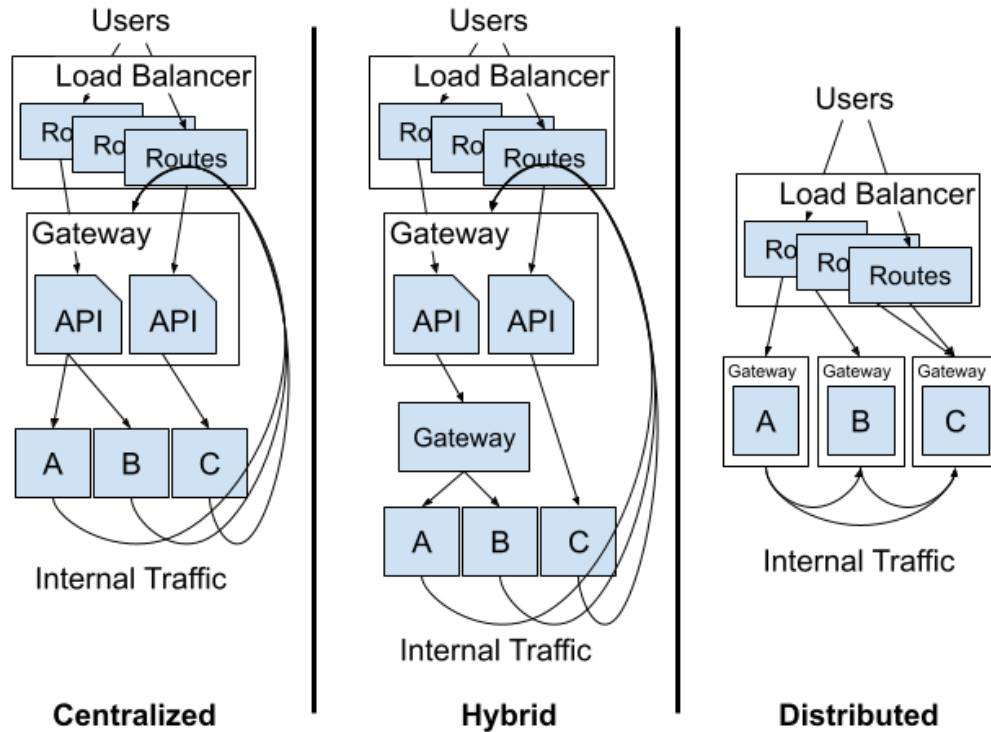


Fig. 6. API gateway patterns

Three patterns have been developed by industry to implement these capabilities, as shown in Fig. 6:

1. Centralized gateway — Protections for all APIs in the enterprise are implemented by a single shared component: an API gateway.
2. Hybrid gateway — Cross-cutting policies (e.g., authentication) are implemented in the centralized shared gateway, but application-specific policies (e.g., authorization) are implemented in the application itself or by components that are owned by the application team.
3. Distributed gateway — All policy checks are performed by gateways that are dedicated to each application, often deployed beside each service instance.

All three patterns can achieve all of the controls outlined in this document and be used by organizations to operate their APIs safely and confidently. Further, many of these patterns may be in use within a single organization. This section explores the engineering design trade-offs that each pattern provides in terms of risks and operational overhead.

Many API gateway products provide management capabilities, such as API key issuance, discovery documentation (i.e., API definition) hosting, documentation for client developers, and support for quotas and billing tiers. These are all valuable features in the enterprise setting, but all of them can be supported across any implementation pattern and are therefore not addressed in this section.

4.1. Centralized API Gateway

The centralized API gateway pattern implements protections for all APIs with a single component (i.e., the API gateway) that is often deployed close to the perimeter of the system. External traffic enters through the gateway, typically via a load balancer. Internal traffic “hairpins” through the gateway as well, which facilitates service-to-service communication inside the infrastructure. That internal, service-to-service traffic may also have to traverse the load balancer for some service instances. Fig. 7 shows a common configuration for a centralized API gateway pattern.

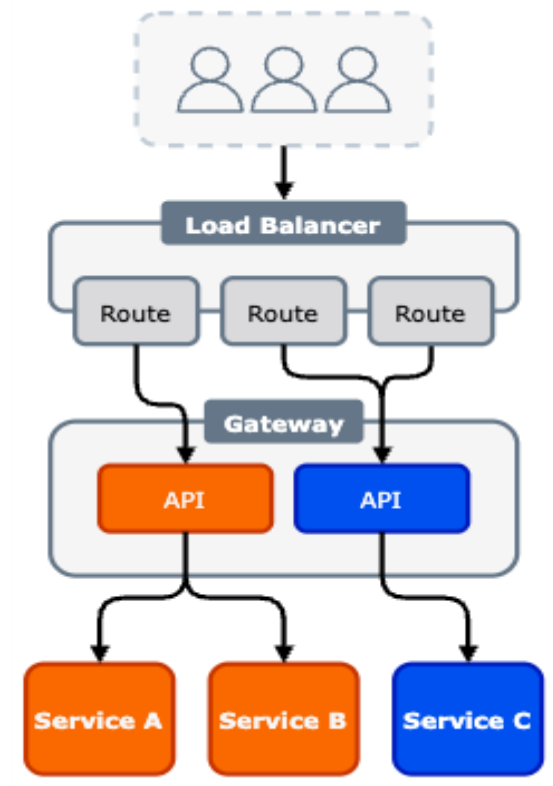


Fig. 7. Centralized API gateway pattern

An API gateway is typically a software application that can be scaled horizontally (i.e., more instances can be deployed side by side). This is one of the reasons why an API gateway often sits behind a load balancer, even for internal service-to-service traffic use cases. Advantages of this pattern include:

- A single policy enforcement point that is easy to monitor and audit, making it simple to verify that policy is enforced for all traffic that traverses the gateway.
- Implementation that matches the organizational structure. Typically, large organizations have a single API team that owns the centralized gateway component and is responsible for when an API is available, which API endpoints are failing, whether policies are being enforced, whether the configuration is up to date, and other issues.

- Streamlined setup for application developers who need to “onboard” their API but do not need to deploy or maintain any additional runtime components.

Disadvantages of this pattern include:

- Shared fate outages. Because there is a single component, an outage of that component causes an outage for all APIs, which can be problematic for mission-critical APIs that need to operate continuously.
- Noisy neighbors, where traffic consumes resources for some APIs and increases latency for all APIs. In the worst case, one application team may submit invalid configuration parameters for a service that may crash or cause a DoS on the API gateway, triggering a shared fate outage for other APIs.
- Long change lead times due to managing how the changes to an individual team’s API configuration impact the shared gateway. This is a frequent side effect of controls that are added to mitigate shared fate outages and noisy neighbors.
- Cost attribution. All requests are handled by the central gateways, and resources spent per request per API (e.g., on payload validation) are uneven. Therefore, it can be difficult to attribute API gateway runtime costs to internal application teams. This can be a problem for companies that implement an internal resource economy for planning by assigning cost centers for each application team.
- Caching the results of policy decisions at runtime becomes critical when implementing the policies outlined in these guidelines due to the sheer number of policy checks required. Caching both increases client-perceived availability and reduces the load on key systems, like authentication and authorization. However, two layers of load balancing (i.e., network load balancer to API gateway and API gateway to service instance) tend to result in poor cache hit rates across policies enforced by the API gateway and for user data in the application layer itself. While some techniques can mitigate this (e.g., distributed caches or streaming connections), they generally add additional development or operational overhead for the application team, API gateway team, or both.
- Because a shared gateway is located at the perimeter, it can be bypassed (e.g., via an attacker pivoting inside the perimeter), which in turn bypasses the policy checks that are enforced by that gateway. This can be mitigated with techniques like service-to-service access policies that ensure that applications only receive traffic via the centralized gateway or by attaching proofs (i.e., credentials) to the request that allow an application to authenticate that the request was handled by the gateway.

4.2. Hybrid Deployments

Hybrid gateway deployments split policy enforcement responsibilities between a centralized gateway and the applications themselves. Cross-cutting policies (e.g., authentication, service discovery, routing, rate limiting, caching) are handled by the centralized gateway. Application-

specific policies (e.g., authorization, request and response validation, protocol conversion, error handling, logging, monitoring) are handled by the application team. This can manifest in the application itself (e.g., gRPC) or as a separate deployment that handles traffic before the application (e.g., GraphQL, Spring Cloud Gateway). As with the centralized pattern, all internal and external traffic between applications must first go through the centralized gateway and, in some instances, through the load balancer. Fig. 8 shows the schematic diagram of a Hybrid gateway pattern.

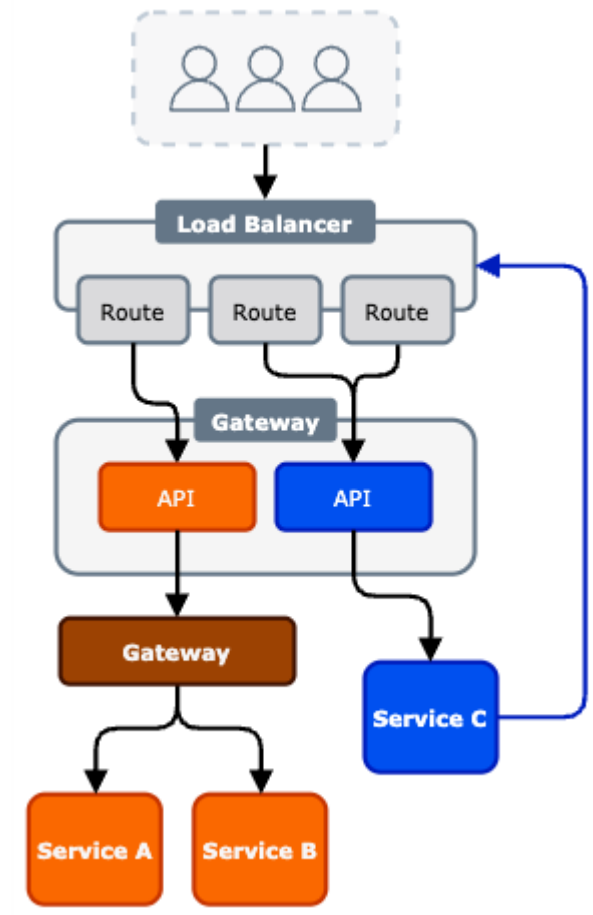


Fig. 8. Hybrid gateway pattern

Overall, this pattern behaves similarly to the centralized API gateway pattern, except that some of the most failure-prone parts of the centralized pattern are delegated to the application teams. This streamlines API gateway operations and enables app teams to move at their own pace. However, it also shifts the responsibility for some runtime operational and security concerns from the API gateway team to those application teams. The exact split of these responsibilities (e.g., sidecar in a service mesh architecture) can vary greatly across different organizations based on their risk profiles and past experiences. Typically, the gateway takes responsibility for:

- Authentication

- Rate limiting
- Circuit breaking
- Service discovery
- Routing
- Caching
- Network-level load balancing

The application or dedicated gateway is responsible for:

- Authorization
- Request/response validation
- Protocol conversion
- Error handling
- Application-instance load balancing

Both are responsible for logging and monitoring to enable visibility into the state of the system and to ensure that policies are being enforced at runtime.

There are similar advantages as the centralized gateway pattern that also include:

- Mitigation of most shared-fate outages and noisy neighbors by moving the most error-prone processing (e.g., request validation) out of the shared gateway and delegating to the application or dedicated gateway.
- Increased iteration speed due to the ability to update configurations with less process overhead and quicken the time involved. This is possible due to the reduced risk of a shared fate outage.

Disadvantages include:

- The enforcement of policies is split across the API gateway and many service instances, which makes it more challenging to ensure that the policy is being enforced consistently and correctly.
- There is increased operational burden on application teams compared to the centralized API gateway pattern, as they are now responsible for ensuring that some policies are enforced in their application.
- Not all classes of shared fate outages and noisy neighbors can be eliminated because the shared central gateway is doing at least some application layer processing.
- Cost attribution is significantly improved compared to the centralized pattern because the most expensive runtime policies are implemented by the application teams. However, the centralized gateway can still be very expensive to operate at high scales and is as difficult to attribute costs as in the centralized pattern.

- Caching hit rates also suffer similarly to the centralized pattern for the same reasons.
- Bypassability/pivot – the module that enforces policies may be bypassed

4.3. Distributed Gateway Pattern

In a distributed approach, the gateway is directly associated with the application, which is owned by a single team. This ensures that changes are isolated to services owned by that team and that the potential for shared fate outages does not arise. Changes to each gateway are “safe” from the organization’s perspective: a bad change will not cause additional problems for other teams, and the team that caused the outage to occur can fix the problem.

External traffic must still enter through a load balancer (see Fig. 9), which does not enforce any policy and only performs routing. Internal traffic may use the same load balancer but may be routed directly peer-to-peer, removing the central gateway from internal traffic as desired, since enforcement of policy happens at the service instance.

This leaves two key challenges that the implementation must address:

1. Ensuring that the remaining shared configuration (i.e., the load balancer) is safe for each team’s changes
2. Ensuring that both cross-cutting and application-specific policies are enforced consistently across the organization

Keeping the load balancer’s configuration safe is a universal problem across all three implementations. However, it is most acute in the distributed pattern because the load balancer must cope with configuration for many applications, while only the API gateway’s configuration needs to be present in the other patterns. Regardless of implementation pattern, this is most often handled at the business process level. Organizations decide on a fixed naming scheme that is enforced by the continuous integration and continuous delivery (CI/CD) process or is otherwise hidden by the organization’s platform (e.g., subdomains-per-service, such as `foo.api.example.com` and `bar.api.example.com`; paths-per-service, such as `api.example.com/v1/foo` and `api.example.com/v2/bar`).

The challenge of a cross-cutting policy is unique to this pattern. In recent years, it has been solved robustly in open source via the service mesh, which can provide a single point for policy management and use its proxies to enforce those policies (i.e., API protections) at each service instance. The service’s properties [2] and use for security [3][6] have been covered in other NIST guidelines.

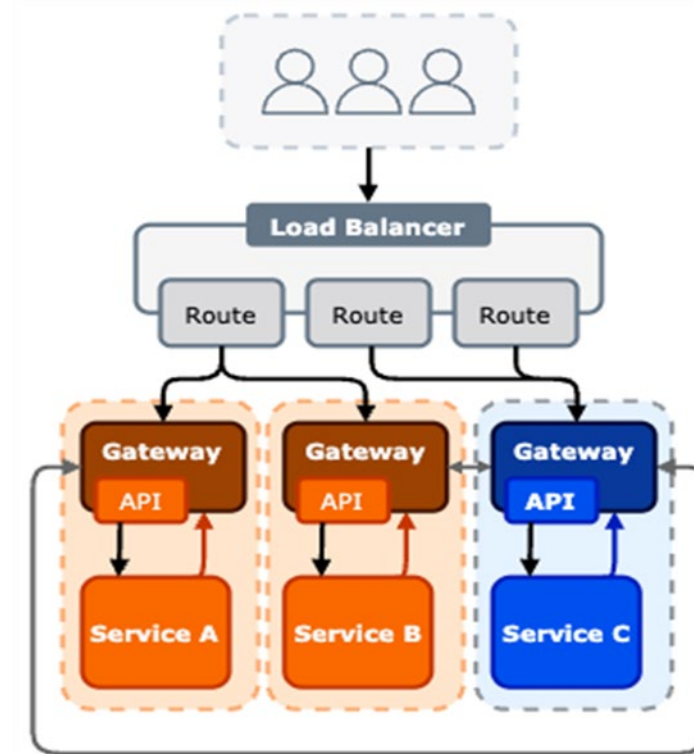


Fig. 9. Distributed API gateway pattern

The advantages of a distributed gateway pattern include:

- All processing is done per application team (i.e., no noisy neighbors), and the risk of a shared fate outage is only present on the load balancer (see Fig. 10), which is a risk shared across all implementation patterns.
- It has the highest rate of change for app teams because they have no external dependencies and little chance of causing outages for other teams.
- A cross-cutting policy can be managed by the central API gateway team via the gateway's control plane (e.g., with the service mesh). This pattern can be adopted harmoniously in a mixed environment, where some APIs are implemented via any of the three patterns in a single organization.
- Cost attribution is straightforward and no more or less challenging than attributing any compute resource spent by teams in the organization.
- Cache locality is typically better than in the other patterns because there is only a single layer of load balancing, and the gateway is co-located with the application. This means that gateway policy checks for a given user are cached alongside the application instance that caches business logic data for that user. However, if a user's request is load-balanced across multiple service instances, then duplicate policy checks have to be performed that would not be required in the other patterns.

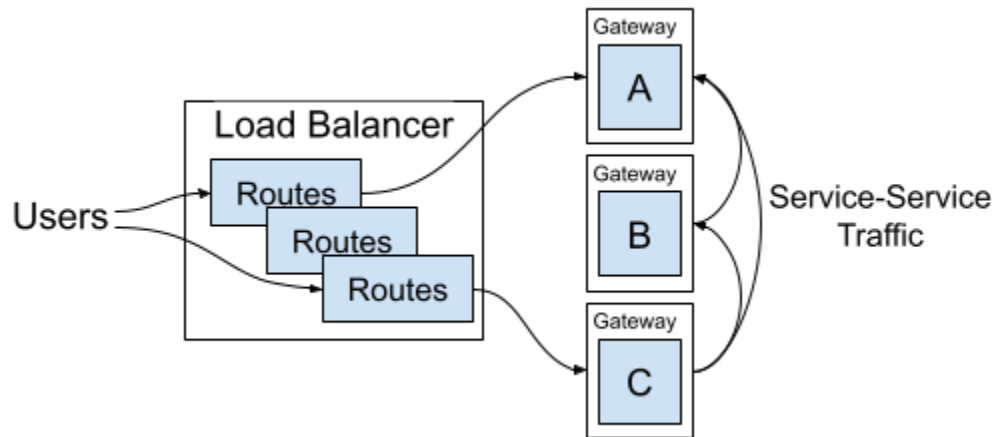


Fig. 10. Service-to-service traffic flows in distributed API gateway pattern

Disadvantages include:

- Because the policy is checked and easily cached per application instance, there can be many more policy checks in the system overall. Any time a user's request is load-balanced to a new service instance, it is highly likely that a new policy check has to be performed. This is an inherent problem in any zero-trust system, which pushes enforcement to the application instance and likely necessitates the adoption of a distributed cache that is managed alongside or as part of the API-serving infrastructure.
- The pattern puts the most burden on application teams. Those teams have to interact with the team managing the load balancer for each API they expose and need to operate at least some of the API-serving infrastructure (e.g., making sure that they have a gateway deployed and routing). Technology like a service mesh can help simplify this, but a burden remains.
- Auditing and verifying policy enforcement can be challenging as enforcement is distributed across all application instances. A robust, distributed gateway implementation (e.g., a service mesh) can help mitigate this via centralized configuration control combined with distributed enforcement and consistent telemetry. If an organization can audit and verify a hybrid gateway pattern, a distributed gateway pattern can be supported with little additional effort.

4.4. Related Technologies

Other technologies fit in and overlap with simplified API gateway patterns and architectures. Notable companion technologies include:

- Web application firewalls (WAFs)
- Bot detection
- DoS and distributed denial of service (DDoS) mitigation

- API endpoint protection
- Web application and API protection (WAAP)

4.4.1. Web Application Firewalls

WAFs mitigate risks related to a request's metadata and payload without needing the application to be involved. In other words, they can be treated as a cross-cutting policy and managed by a central team. WAFs work at the application level and operate on parsed HTTP requests (i.e., they can implement policy per header and on request bodies). However, WAFs generally do not work at the API level. A WAF can scan a request for a payload that looks like a SQL injection attack, but it cannot assert, for example, that a request has a "name" field that is a string less than 100 characters long. As such, a WAF is an excellent first step for organizations to implement the policies outlined in this document, but it is not a complete solution.

The Open Worldwide Application Security Project (OWASP) publishes research on vulnerabilities based on data from its partners [7] as well as a generic set of WAF rules — the Core Rule Set (CRS) [25] — that aim to mitigate many common attacks. The CRS should be treated as a starting point for any organization's WAF policy. Deploying a WAF with at least the CRS enabled helps mitigate risks, including malicious inputs (see Sec. 2.6.2), unrestricted resource consumption (see Sec. 2.4), and the leaking of sensitive information (see Sec. 2.5).

There are two primary downsides with WAFs:

1. WAFs are relatively expensive to run in terms of both latency and compute. They need to parse every request, perform a variety of scans to identify attack signatures (the number of scans depends on the policy configured), and either block or forward the request. While this overlaps heavily with the functionality of an API gateway, a WAF is typically deployed and operated by a separate team in isolation from the API gateway, often as part of the load balancer. This is convenient because the load balancer is the first place where requests are decrypted in the infrastructure. A secondary consequence is that WAF policies are typically only enforced at the perimeter.
2. WAFs are fundamentally reactive. They operate based on matching requests to known attack signatures. As a result, they are largely ineffective at mitigating novel attacks, and attackers can leverage a variety of obfuscation techniques to hide known attacks behind novel signatures. Care must be taken to ensure that the WAF is running with the latest attack signature configurations, and custom rules must often be written for the organization.

In line with a zero trust posture, WAF policies should be enforced as close to the application as possible. This helps mitigate a variety of mechanisms that attackers might use to pivot within or otherwise compromise an infrastructure. As a practical matter, it can be cost-prohibitive to run a full suite of WAF mitigations on every internal and external request. This cost can be mitigated in two ways, which can be combined:

1. Incorporate the WAF as part of the overall API-serving infrastructure, deploy the WAF itself in a hybrid model (i.e., keep a centralized WAF at the load balancer with a full suite of policies to protect against untrusted traffic), and enforce a minimum set of app-specific WAF policies near each of the applications (e.g., in the distributed gateway). This minimizes policies run on east-west (i.e., generally assumed to be trusted) traffic while still sanitizing less trusted external traffic and tends to result in a good compromise of risk versus cost.
2. Deploy the WAF as part of the API gateway implementation itself, which can avoid parsing the request multiple times (i.e., reduce the latency and compute costs of WAF policies), regardless of the API-serving implementation pattern chosen. If the API gateway is hybrid or distributed, then this technique can also be incorporated for further performance improvement.

4.4.2. Bot Detection

Bot detection typically involves evaluating risk signals, including origin (e.g., source IP, user credentials) and API usage patterns, over time to determine whether a seemingly legitimate user is likely to be a bot (i.e., an automated script acting maliciously). In response to flagging a high-risk user, bot detection systems will either block traffic or serve some kind of bot-defeating measure (e.g., CAPTCHA) before allowing the system to continue to be used. These tools primarily mitigate the risks of unrestricted resource access (see Sec. 2.4) (e.g., maliciously automating account creation in an email system) and leaking sensitive information (see Sec. 2.5), especially data exfiltration by repeated calls.

Bot detection is frequently deployed in user-facing applications. It can be more challenging with a purely machine-to-machine API because legitimate and malicious traffic patterns are even harder to differentiate. Many APIs are *intended* for use by scripts or non-user-facing applications, so human versus computer checks are irrelevant.

4.4.3. Distributed Denial of Service (DDoS) Mitigation

A DDoS attack is a DoS that originates from many different locations or users. This makes it more challenging to mitigate than a traditional DoS attack, which can often be prevented by blocking a small set of users. While DoS attacks may be targeted and application-level, DDoS attacks are often network-oriented in nature and seek to saturate the server's bandwidth or ability to establish new connections. Because of the primarily network-oriented nature of DDoS attacks, most DDoS mitigation tools are deployed at the network edge as part of the load balancer or even before the load balancer as part of the CDN and DNS system (often called "Global Traffic Management"). Predictably, DDoS mitigation tools help mitigate unrestricted resource consumption (see Sec. 2.4).

4.4.4. API Endpoint Protection

“API protection” or “API endpoint protection” are nebulous terms for describing a set of capabilities around API inventory, authentication, rate limiting, and data analysis. The exact set of capabilities tends to vary with the implementation. For example, sophisticated implementations can scan requests and responses to tag suspect data on the wire (e.g., to help tag sensitive data and pinpoint possible leaks or exfiltration).

API protection products are typically packaged with the API gateway. API gateway vendors primarily deliver their products in the centralized API gateway pattern, so these controls are often only enforced at the perimeter. Like a WAF, the policies they enforce are typically cross-cutting and do not require an in-depth understanding at the API payload level. As such, the two products are often marketed in a similar niche. The exact set of risks mitigated by these tools depends on the feature set, but they usually attempt to mitigate a lack of API visibility (see Sec. 2.1), broken authentication (see Sec. 2.3), some aspects of unrestricted compute consumption (see Sec. 2.4), and the leaking of sensitive information (see Sec. 2.5).

There is value in any tool that helps organizations inventory and manage their APIs and traffic. However, policy enforcement should be as close to the individual service instance as possible in order to achieve robust API security assurance. In the use case of data classification, these tools can be especially helpful when building an initial inventory. As API definitions are rolled out across the organization, data tagging should be implemented as part of the API schema, and the data flow policy should be enforced via explicit policy (e.g., with an authorization system). The runtime discovery of data flow is particularly important for protecting against exfiltration.

4.4.5. Web Application and API Protection (WAAP)

Gartner coined the term “web application and API protection” (WAAP) [27] to describe the trend of packaging these technologies (i.e., WAF, bot detection, DDoS mitigation, API protection) into a single product. Regardless of how the capabilities are implemented, organizations must understand the risks that they are trying to mitigate in the context of their existing security posture.

4.5. Summary of Implementation Patterns

Combining the three patterns in API gateway architecture with the companion technologies discussed Sec. 4.4 provides a comprehensive set of enterprise security solutions for API protection. The key point in each pattern is identifying where to enforce each policy. These decisions result in trade-offs in runtime, architecture, and operations for the application teams utilizing the API-serving infrastructure. Many organizations use a combination of all three patterns deployed in production precisely because of those trade-offs. While all three patterns can be used to successfully implement the controls outlined in this document, the distributed gateway pattern and its companion technologies best align with the principles of zero trust and are strongly recommended for organizations that want to adopt a security-forward approach.

5. Conclusions and Summary

APIs are critical for integrating applications into the digital infrastructure of an enterprise. Given the highly distributed nature of both physical and logical applications, NIST recommends that APIs be operated under zero trust principles, irrespective of whether they are exposed to the outside world or meant to be consumed by other applications within the enterprise's infrastructure. Like all software, APIs go through an iterative life cycle with phases (i.e., Develop, Build, Deploy, Operate) that can be broadly classified into pre-runtime and runtime stages.

The sheer proliferation of API deployments, the heterogeneous infrastructures under which they operate, and the access to valuable corporate data that they enable make them targets for exploitation. A detailed analysis of their vulnerabilities and the potential attack vectors that can exploit them is a prerequisite for identifying the appropriate set of protection measures or controls to ensure API security. This document analyzes a spectrum of risk factors that give rise to vulnerabilities, such as the lack of a formal schema, improper inventorying, the lack of robust authentication and authorization support, the improper monitoring of resource consumption, and inadequate control over the leakage of sensitive information.

The recommended controls in this document are classified into pre-runtime and runtime protections. They are further subdivided into basic and advanced protections to enable enterprises to use a risk-based and incremental approach to securing their digital assets. Pre-runtime protections focus on API specification parameters (i.e., syntactic and semantic aspects), while runtime protections focus on API request and response operations (e.g., encrypted communication channels, proper authentication and authorization).

These guidelines present a landscape of real-world and state-of-practice implementation options to configure and enforce the recommended controls by describing the advantages and disadvantages of each type of protection deployment or pattern. This will enable practitioners to make an informed decision to realize a robust and cost-effective API security infrastructure for their enterprises.