



NIST Special Publication 800
NIST SP 800-228 ipd

Guidelines for API Protection for Cloud-Native Systems

Initial Public Draft

Ramaswamy Chandramouli
Zack Butcher

This publication is available free of charge from:
<https://doi.org/10.6028/NIST.SP.800-228.ipd>

NIST Special Publication 800
NIST SP 800-228 ipd

Guidelines for API Protection for Cloud-Native Systems

Initial Public Draft

Ramaswamy Chandramouli
*Computer Security Division
Information Technology Laboratory*

Zack Butcher
Tetrade, Inc.

This publication is available free of charge from:
<https://doi.org/10.6028/NIST.SP.800-228.ipd>

March 2025



U.S. Department of Commerce
Howard Lutnick, Secretary

National Institute of Standards and Technology
Craig Burkhardt, Acting Under Secretary of Commerce for Standards and Technology and Acting NIST Director

Certain commercial equipment, instruments, software, or materials, commercial or non-commercial, are identified in this paper in order to specify the experimental procedure adequately. Such identification does not imply recommendation or endorsement of any product or service by NIST, nor does it imply that the materials or equipment identified are necessarily the best available for the purpose.

There may be references in this publication to other publications currently under development by NIST in accordance with its assigned statutory responsibilities. The information in this publication, including concepts and methodologies, may be used by federal agencies even before the completion of such companion publications. Thus, until each publication is completed, current requirements, guidelines, and procedures, where they exist, remain operative. For planning and transition purposes, federal agencies may wish to closely follow the development of these new publications by NIST.

Organizations are encouraged to review all draft publications during public comment periods and provide feedback to NIST. Many NIST cybersecurity publications, other than the ones noted above, are available at <https://csrc.nist.gov/publications>.

Authority

This publication has been developed by NIST in accordance with its statutory responsibilities under the Federal Information Security Modernization Act (FISMA) of 2014, 44 U.S.C. § 3551 et seq., Public Law (P.L.) 113-283. NIST is responsible for developing information security standards and guidelines, including minimum requirements for federal information systems, but such standards and guidelines shall not apply to national security systems without the express approval of appropriate federal officials exercising policy authority over such systems. This guideline is consistent with the requirements of the Office of Management and Budget (OMB) Circular A-130.

Nothing in this publication should be taken to contradict the standards and guidelines made mandatory and binding on federal agencies by the Secretary of Commerce under statutory authority. Nor should these guidelines be interpreted as altering or superseding the existing authorities of the Secretary of Commerce, Director of the OMB, or any other federal official. This publication may be used by nongovernmental organizations on a voluntary basis and is not subject to copyright in the United States. Attribution would, however, be appreciated by NIST.

NIST Technical Series Policies

[Copyright, Use, and Licensing Statements](#)

[NIST Technical Series Publication Identifier Syntax](#)

Publication History

Approved by the NIST Editorial Review Board on YYYY-MM-DD [Will be added to final publication.]

How to Cite this NIST Technical Series Publication:

Chandramouli R, Butcher Z (2025) Guidelines for API Protection for Cloud-Native Systems. (National Institute of Standards and Technology, Gaithersburg, MD), NIST Special Publication (SP) NIST SP 800-228 ipd.
<https://doi.org/10.6028/NIST.SP.800-228.ipd>

Author ORCID iDs

Ramaswamy Chandramouli: 0000-0002-7387-5858

NIST SP 800-228 ipd (Initial Public Draft)
March 2025

Guidelines for API Protection
for Cloud-Native Systems

Public Comment Period

March 25, 2025 – May 12, 2025

Submit Comments

sp800-228-comments@nist.gov

National Institute of Standards and Technology
Attn: Computer Security Division, Information Technology Laboratory
100 Bureau Drive (Mail Stop 8930) Gaithersburg, MD 20899-8930

Additional Information

Additional information about this publication is available at <https://csrc.nist.gov/pubs/sp/800/228/ipd>, including related content, potential updates, and document history.

All comments are subject to release under the Freedom of Information Act (FOIA).

1 **Abstract**

2 Modern enterprise IT systems rely on a family of application programming interfaces (APIs) for
3 integration to support organizational business processes. Hence, a secure deployment of APIs is
4 critical for overall enterprise security. This, in turn, requires the identification of risk factors or
5 vulnerabilities in various phases of the API life cycle and the development of controls or
6 protection measures. This document addresses the following aspects of achieving that goal: (a)
7 the identification and analysis of risk factors or vulnerabilities during various activities of API
8 development and runtime, (b) recommended basic and advanced controls and protection
9 measures during pre-runtime and runtime stages of APIs, and (c) an analysis of the advantages
10 and disadvantages of various implementation options for those controls to enable security
11 practitioners to adopt an incremental, risk-based approach to securing their APIs.

12 **Keywords**

13 API; API endpoint; API gateway; API key; API schema; web application firewall.

14 **Reports on Computer Systems Technology**

15 The Information Technology Laboratory (ITL) at the National Institute of Standards and
16 Technology (NIST) promotes the U.S. economy and public welfare by providing technical
17 leadership for the Nation’s measurement and standards infrastructure. ITL develops tests, test
18 methods, reference data, proof of concept implementations, and technical analyses to advance
19 the development and productive use of information technology. ITL’s responsibilities include
20 the development of management, administrative, technical, and physical standards and
21 guidelines for the cost-effective security and privacy of other than national security-related
22 information in federal information systems. The Special Publication 800-series reports on ITL’s
23 research, guidelines, and outreach efforts in information system security, and its collaborative
24 activities with industry, government, and academic organizations.

25 **Call for Patent Claims**

26 This public review includes a call for information on essential patent claims (claims whose use
27 would be required for compliance with the guidance or requirements in this Information
28 Technology Laboratory (ITL) draft publication). Such guidance and/or requirements may be
29 directly stated in this ITL Publication or by reference to another publication. This call also
30 includes disclosure, where known, of the existence of pending U.S. or foreign patent
31 applications relating to this ITL draft publication and of any relevant unexpired U.S. or foreign
32 patents.

33 ITL may require from the patent holder, or a party authorized to make assurances on its behalf,
34 in written or electronic form, either:

- 35 a) assurance in the form of a general disclaimer to the effect that such party does not hold
36 and does not currently intend holding any essential patent claim(s); or
- 37 b) assurance that a license to such essential patent claim(s) will be made available to
38 applicants desiring to utilize the license for the purpose of complying with the guidance
39 or requirements in this ITL draft publication either:
 - 40 i. under reasonable terms and conditions that are demonstrably free of any unfair
41 discrimination; or
 - 42 ii. without compensation and under reasonable terms and conditions that are
43 demonstrably free of any unfair discrimination.

44 Such assurance shall indicate that the patent holder (or third party authorized to make
45 assurances on its behalf) will include in any documents transferring ownership of patents
46 subject to the assurance, provisions sufficient to ensure that the commitments in the assurance
47 are binding on the transferee, and that the transferee will similarly include appropriate
48 provisions in the event of future transfers with the goal of binding each successor-in-interest.

49 The assurance shall also indicate that it is intended to be binding on successors-in-interest
50 regardless of whether such provisions are included in the relevant transfer documents.

51 Such statements should be addressed to: sp800-228-comments@nist.gov

52

53	Table of Contents	
54	Executive Summary	1
55	1. Introduction	2
56	1.1. Building Blocks and Structures.....	2
57	1.2. Zero Trust and APIs: The Vanishing Perimeter	4
58	1.3. API Life Cycle	5
59	1.4. Document Goals.....	5
60	1.5. Relationship to Other NIST Documents	6
61	1.6. Document Structure.....	6
62	2. API Risks — Vulnerabilities and Exploits	8
63	2.1. Lack of Visibility of APIs in the Enterprise Inventory	8
64	2.2. Missing, Incorrect, or Insufficient Authorization	8
65	2.3. Broken Authentication	9
66	2.4. Unrestricted Resource Consumption.....	10
67	2.4.1. Unrestricted Compute Resource Consumption	10
68	2.4.2. Unrestricted Physical Resource Consumption	10
69	2.5. Leaking Sensitive Information to Unauthorized Callers.....	11
70	2.6. Insufficient Verification of Input Data.....	12
71	2.6.1. Input Validation	12
72	2.6.2. Malicious Input Protection	12
73	2.7. Credential Canonicalization— Preparatory Step for Controls	13
74	2.7.1. Gateways Straddle Boundaries	13
75	2.7.2. Requests With a Service Identity but No User Identity.....	13
76	2.7.3. Requests With a User Identity But No Service Identity.....	14
77	2.7.4. Requests With Both User and Service Identities.....	15
78	2.7.5. Reaching Out to Other Systems	16
79	2.7.6. Mitigating the Confused Deputy	16
80	2.7.7. Identity Canonicalization.....	16
81	3. Recommended Controls for APIs	18
82	3.1. Pre-Runtime Protections.....	19
83	3.1.1. Basic Pre-Runtime Protections.....	19
84	3.1.2. Advanced Pre-Runtime Protections	20
85	3.2. Runtime Protections	21
86	3.2.1. Basic Runtime Protections.....	21

87	3.2.2. Advanced Runtime Protections	26
88	4. Implementation Patterns and Trade-Offs for API Protections	30
89	4.1. Centralized API Gateway	32
90	4.2. Hybrid Deployments	34
91	4.3. Decentralized Gateways	36
92	4.4. Related Technologies	38
93	4.4.1. Web Application Firewalls	39
94	4.4.2. Bot Detection	40
95	4.4.3. Distributed Denial of Service (DDoS) Mitigation	40
96	4.4.4. API Endpoint Protection	41
97	4.4.5. Web Application and API Protection (WAAP)	41
98	4.5. Summary of Implementation Patterns	41
99	5. Conclusions and Summary	43
100	References	44
101	Appendix A. API Classification Taxonomy	46
102	A.1. API Classification Based on Degree of Exposure	46
103	A.2. API Classification Based on Communication Patterns	46
104	A.3. API Classification Based on Architectural Style or Pattern (API Types)	47
105	Appendix B. DevSecOps Phase and Associated Class of API Controls	48
106		
107		

108 **Acknowledgments**

109 The authors would like to thank Orion Letizi, technical writer at Tetrade for providing
110 continuous, ongoing edits during the development of this document. We would also like to
111 thank Erica Hughberg, an engineer at Tetrade and James Gough, a Distinguished Engineer at
112 Morgan Stanley for their feedback on the initial outline for controls. Their extensive hands-on
113 experience in running API security programs in large enterprises helped us to address the
114 current API security issues and incorporate state of practice API security controls in our
115 recommendations. Last but not the least, the authors would also like to express their thanks to
116 Isabel Van Wyk of NIST for her detailed and extensive editorial review.

117 **Executive Summary**

118 Application programming interfaces (APIs) provide the means to integrate and communicate
119 with the modern enterprise IT application systems that support business processes. However, a
120 lack of due diligence can introduce vulnerabilities and risk factors that exploit the connectivity
121 and accessibility features of APIs. If these vulnerabilities are not identified, analyzed and
122 addressed through control measures, attack vectors could threaten the security posture of the
123 application systems spanned by these APIs. A systematic and effective means of identifying and
124 addressing these vulnerabilities is only possible by treating the development and deployment of
125 APIs as an iterative life cycle using paradigms like DevSecOps.

126 This document provides guidance and recommendations on controls and protection measures
127 for secure API deployments in the enterprise. In addition, an analysis of the advantages and
128 disadvantages of various implementation options (called patterns) for those controls enable
129 security practitioners to choose the most effective option for their IT ecosystem.

130 Developing these controls and analyzing their implementation options should be guided by
131 several overarching principles:

132 The guidance for controls should cover all APIs, regardless of whether they are exposed to
133 customers/partners or used internally within the enterprise.

134 With the vanishing of perimeters in modern enterprise IT applications, all controls should
135 incorporate the concept of zero trust.

136 The controls should span the entire API life cycle and be classified into (a) pre-runtime
137 protections and (b) runtime protections that are then subdivided into basic and advanced
138 protections to enable incremental risk-based adoption.

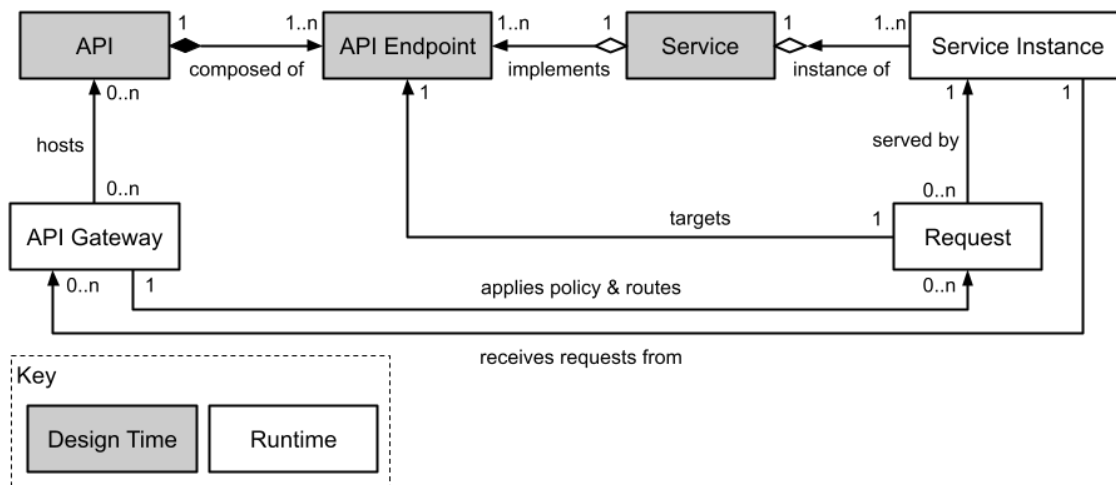
139 **1. Introduction**

140 Application programming interfaces (APIs) represent an abstraction of the underlying
141 implementation of a digital enterprise. Given the spatial (e.g., on-premises, multiple clouds)
142 and logical (e.g., microservices) nature of current enterprise applications, APIs are needed to
143 integrate and establish communication pathways between internal and third-party services and
144 applications. Informally, APIs are the lingua franca of modern IT systems: they describe what
145 actions users are allowed to take. They are also used in every type of application, including
146 server-based monolithic, microservices-based, browser-based client, and IoT.

147 **1.1. Building Blocks and Structures**

148 An Application Programming Interface (API) defines how any two pieces of software
149 communicate – they are ubiquitous in software. An API is a collection of commands or
150 *endpoints* that operate on data or *objects* via some protocol. *Network-based APIs* are APIs built
151 to be consumed by remote applications over the network. Because they're exposed and
152 consumed over the network, they present a unique set of challenges. The growth of (micro-)
153 service-oriented architectures, coupled with Software-as-a-Service (SaaS) becoming
154 commonplace – which are nearly always delivered via APIs – has resulted in an explosion in
155 network-based APIs across organizations. This document focuses on controls for network-based
156 APIs.

157 Before we can discuss API controls, we need a common understanding and language for the
158 building blocks, and how they relate to each other. The taxonomy is: an API is composed of a
159 set of API Endpoints; API Endpoints are implemented by Services; at runtime, Requests to a
160 specific API Endpoint are served by Service Instances. An API Gateway hosts many APIs and is
161 responsible for mapping each Request to its target API Endpoint, applying policy for that
162 Endpoint (e.g. authentication and rate limiting), then routing that Request to a Service Instance
163 which implements that API Endpoint.



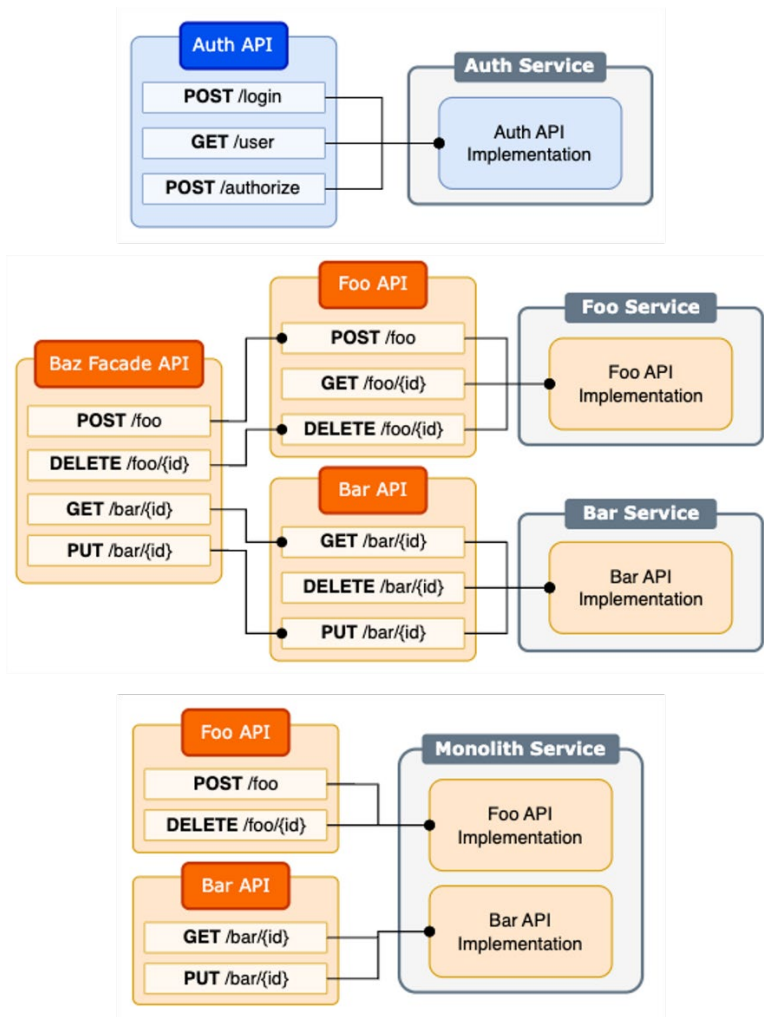
164

165

Fig. 1. API, API Endpoint, Service and Service Instance

166 Traditionally, we think of network-based APIs as being customer-oriented, partner-oriented, or
167 internal – often called “third-party”, “second-party”, and “first-party” APIs, respectively.
168 Second- and third-party APIs are typically exposed to callers outside of the organization via an
169 API gateway. First-party APIs can be exposed to callers inside of the organization on the same
170 API gateway, but they are also often consumed directly by internal callers without traversing a
171 dedicated API serving stack.

172 An API is a set of API Endpoints, and a Service implements a set of API Endpoints – so every
173 Service implements *some* API. We call these *Service APIs*. Most first-party API integrations
174 happen via the Service API, i.e. they map to a single service. On the other hand, APIs hosted by
175 the API gateway typically have Endpoints that map to many different Services. This is especially
176 common for second- and third-party APIs. We call these *Facade APIs*, because they present a
177 single facade to an outside caller over (potentially many) different Service APIs. Finally, it’s
178 common that multiple Services are grouped together into an *Application* – typically along
179 organizational lines (often an Application maps to a team). Schematic diagrams of a Service API,
180 Façade API and an Application (Monolithic) API are given below:



181
182

Fig. 2. (Top to Bottom) Service API, Façade API, Service and Application (Monolithic)

183

184 Less formally: we can think of the APIs we expose outside the organization as a facade over a
185 set of Services. Those Services implement internal APIs (*Service APIs*). Services in the
186 organization communicate with each other via those internal APIs – sometimes directly, and
187 sometimes via an API gateway. The API Gateway is responsible for some policies, like
188 authentication and rate limiting, as well as being responsible for mapping the facade APIs for
189 external clients to internal APIs. Then, to get a handle on things organizationally, we often
190 group related Services into a bucket called an Application.

191 While we tend to think of APIs in the context of exposing functionality to clients or partners,
192 APIs don't exist solely at the edge of our infrastructure. Any time systems communicate, there's
193 *some* API involved. Even if that API is something like CSV over FTP. The examples in SP focus
194 primarily on "modern" APIs exposed via mechanisms like HTTP/REST, gRPC, or SOAP, but we
195 believe the principals in this SP are universal and should be applied to *all* APIs.

196 **1.2. Zero Trust and APIs: The Vanishing Perimeter**

197 APIs are built out of services that communicate with each other via APIs, similar to how the
198 internet is a "network of networks." One of the most important implications of zero trust is that
199 there is no meaningful distinction between an "internal" and "external" caller because the
200 perimeter is the service instance itself. Rather, all callers are trusted if they are authorized to be
201 trusted. This contrasts with traditional approaches to API security in which the only "APIs" are
202 those exposed to "external" callers, and API-oriented controls are only enforced at the
203 perimeter, typically via an API gateway.

204 NIST Special Publication (SP) 800-207A [6] discusses zero trust at runtime and the principle of
205 shrinking the perimeter to the service instance using the five runtime controls of identity-based
206 segmentation:

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

217 Achieving a zero-trust runtime requires applying these five controls to *all* API communications.
218 This guidance further describes additional controls that are necessary for safe and secure API
219 operations beyond identity-based segmentation. These controls should be enforced on all APIs

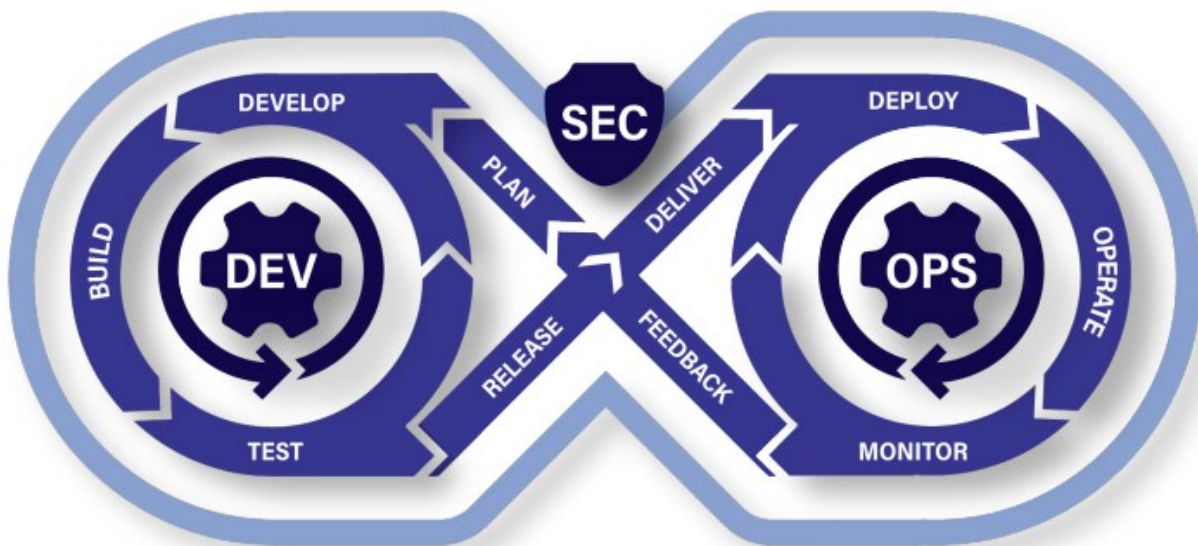
220 in a system, including those exposed to the outside world (i.e., public APIs) and those intended
221 only for other applications in a given infrastructure (i.e., internal APIs).

222 1.3. API Life Cycle

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

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

229 DoD Enterprise DevSecOps Fundamentals [1] provides a detailed description of each phase of
230 the software development life cycle. Application of the DevSecOps paradigm in the context of
231 cloud-native applications can be found in [4][5].



232

233

Fig. 3. DevSecOps life cycle phases

234 1.4. Document Goals

235 The goal of this document is to recommend guidance or controls for API protection. These
236 controls are classified into two categories:

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

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

244 A prerequisite for defining any API protection measure or policy irrespective of its category or
245 sub-category is that the protections must be expressed in terms of nouns and verbs that pertain
246 to API components, API endpoint components, API requests, and API responses that in turn
247 contain references to resources/data and operations on those resources. These nouns and
248 verbs form the fundamental surface that is exposed to the consumers of APIs and API
249 endpoints.

250 **1.5. Relationship to Other NIST Documents**

251 Today, most enterprise software development and integration are based on APIs. Section 1.3
252 articulated the close relationship between software and APIs, demonstrated that API
253 development and deployment follow the same iterative life cycle as the software, and provided
254 NIST guidance on DevSecOps.

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

264 **1.6. Document Structure**

265 This document is organized as follows:

- 266 • Section 2 looks at the risk factors and vulnerabilities associated with APIs and the attack
267 vectors that could exploit those vulnerabilities.
- 268 • Section 3 recommends controls to protect APIs and classifies them into basic and
269 advanced categories that need to be applied prior to runtime or enforced during
270 runtime.
- 271 • Section 4 provides a detailed analysis of implementation options or patterns for the
272 controls described in Sec. 3 and outlines the advantages and disadvantages of each
273 pattern.
- 274 • Section 5 provides the summary and conclusions.
- 275 • Appendix A provides the classification taxonomy for APIs.

- 276
- Appendix B illustrates the API controls related to each DevSecOps phase

277 **2. API Risks — Vulnerabilities and Exploits**

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

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

285 **2.1. Lack of Visibility of APIs in the Enterprise Inventory**

286 Most organizations have gaps in their API inventories, even if they otherwise have mature
287 inventory management capabilities. Without an accurate API inventory, one cannot begin to
288 protect the enterprise estate, and there may be incidents occurring at the API level that the
289 security organization is entirely unaware of. Common reasons for lack of visibility are:

- 290 • Organizational silos: APIs are built by many teams across the organization, deployed
291 across cloud and on-premises environments, and inherited in mergers and acquisitions.
292 This results in uneven attention to security concerns and difficulty establishing accurate,
293 up-to-date inventories. This is worsened in organizations that achieve a high degree of
294 developer agility: the faster development happens; the faster inventories grow stale.
- 295 • Rogue or shadow APIs: APIs defined for internal use (e.g., debugging, testing, ad hoc
296 solutions to business problems) may not be appropriately documented and often bypass
297 standard security review practices.
- 298 • Zombie or deprecated APIs: APIs may have been replaced or superseded by newer
299 systems but have not yet been entirely removed (e.g., because all callers have not yet
300 migrated to the alternative, there no longer exists a team responsible for the system).
301 They risk falling behind the latest security policies and protections.

302 **2.2. Missing, Incorrect, or Insufficient Authorization**

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

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

317 Authorization risks can be categorized in three ways:

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

329 **2.3. Broken Authentication**

330 Authentication is a prerequisite for authorization, particularly two aspects: the authentication
331 system itself is robust, and the application uses the authenticated identities correctly.

332 Risks that an authentication system needs to mitigate include [8]:

- 333 • Credential stuffing, where an attacker brute forces usernames and passwords without
334 mitigation (e.g., rate limits and Captcha).
- 335 • Brute force attacks on a single account without mitigations. This and the previous bullet
336 are closely related to unrestricted resource consumption (see Sec. 2.4).
- 337 • Insecure practices like weak passwords, passing sensitive data in public channels (e.g.,
338 the URL), missing password validation for changes to sensitive account data, and using
339 weak keys or poor algorithms to encrypt user data in transit and at rest.
- 340 • Bad or incorrect token validation, including not validating at all, ignoring expiry, and
341 using insecure signing schemes or weak signing keys.

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

- 344 • Missing authentication. Tokens can be present but simply not checked. This is often due
345 to a bug or misconfiguration in the application.

- 346 • Weak or predictable tokens, default accounts, and default passwords (e.g., a hard-coded
347 bootstrap account with the same username and password on all devices, test accounts
348 with predictable names and weak/guessable passwords).

349 **2.4. Unrestricted Resource Consumption**

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

355 **2.4.1. Unrestricted Compute Resource Consumption**

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

- 361 • DoS attacks via bandwidth saturation or resource starvation
- 362 • Unreliable performance due to resource utilization for one user or service impacting
363 others
- 364 • Cost amplification, where an attacker can spend a small amount of resources (e.g.,
365 money, compute, bandwidth) to make requests that trigger a system to spend a much
366 larger amount of resources servicing the request

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

371 **2.4.2. Unrestricted Physical Resource Consumption**

372 The risks associated with the unrestricted consumption of physical resources are often ignored
373 by software engineers, who tend to be better versed in the threat landscape of the virtual
374 world. Critical business operations can be impacted when an attacker targets software systems
375 that control physical processes (e.g., SCADA systems).

376 APIs may also result in text messages being sent to users, charges to credit cards, or the
377 consumption of expensive third-party resources. For example, a common challenge seen by
378 organizations that adopt AI is the accidental over-use of expensive AI APIs, resulting in large
379 unplanned expenses for the business.

380 These risks are best mitigated by a combination of rate limiting, quotas, spending policy
381 controls in third-party software, bot/abuse detection, and application or business flow changes.
382 These risks manifest as:

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

392 Mitigations for both compute and physical resource consumption are similar. For compute
393 resources, how users interact with a system should be limited. For physical resources, how the
394 user interacts with a system *and* how a system interacts with external systems should be
395 limited and considered early in the design phase. Mitigating these risks can sometimes require
396 business flow changes.

397 **2.5. Leaking Sensitive Information to Unauthorized Callers**

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

402 Risks include:

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

414 **2.6. Insufficient Verification of Input Data**

415 Trusting unverified input is one of the largest classes of recurring security bugs in software.
416 There are at least two levels of verification that APIs need to be concerned with:

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

419 **2.6.1. Input Validation**

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

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

- 425 • Impacting the availability of APIs
 - 426 ○ The "Query of Death" [24] is a DoS attack via specially crafted requests that
 - 427 trigger pathological worst-case behavior in the server.
 - 428 ○ In the worst case, the server itself may crash due to bad input handling, which
 - 429 can be exploited by an attacker to cause DoS on systems.
- 430 • Invalid or malicious data being stored in the system, which can cause latent issues (e.g.,
431 failure to restart during recovery, crashes when accessing invalid records)
- 432 • Unanticipated error handling during request processing, which leaks internal
433 information

434 **2.6.2. Malicious Input Protection**

435 While the input may satisfy "syntactic" validation, it also needs to be verified as non-malicious
436 before it is used. Extending the "name" example above, a caller may send a request that
437 contains a name field with a string less than 100 characters (i.e., valid), but that string may be a
438 SQL injection attack. Common risks include:

- 439 • Data leaks or corruption (e.g., a SQL injection attack)
- 440 • Unanticipated or unrestricted resource utilization (e.g., an attacker automates account
441 creation and uploads multi-gigabyte "profile pictures" to each account)
- 442 • Exposing a surface that attackers can use to pivot within the infrastructure or leverage
443 to mount further attacks on others (e.g., by allowing servers to be used for server-side
444 request forgery [SSRF])
- 445 • Cost amplification attacks, like the "billion laughs attack" (XML expansion) [10] or "zip
446 bombs" (zip archive expansion) [11]

447 **2.7. Credential Canonicalization— Preparatory Step for Controls**

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

457 **2.7.1. Gateways Straddle Boundaries**

458 A gateway is something in an infrastructure that straddles a boundary and is typically the only
459 way for traffic to cross that boundary. As a result, the gateway is uniquely positioned to enforce
460 policy. One of the most important policies that the API gateway enforces is authentication,
461 ideally of both the user and the calling service.

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

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

475 **2.7.2. Requests With a Service Identity but No User Identity**

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

480 The solution here is to adopt service accounts that present some system in a user identity
481 domain. That service account can be for an internal system and, therefore, have permission to
482 act on the data of many other users, or it can be for a user’s applications with correspondingly
483 fewer permissions. The API gateway can mediate with the user authentication system to
484 exchange the service’s runtime identity for a service account credential that represents the

485 service in the user identity domain and attach that service account credential as the end user
486 credential to requests that it forwards into the part of the infrastructure that supports identity-
487 based segmentation.

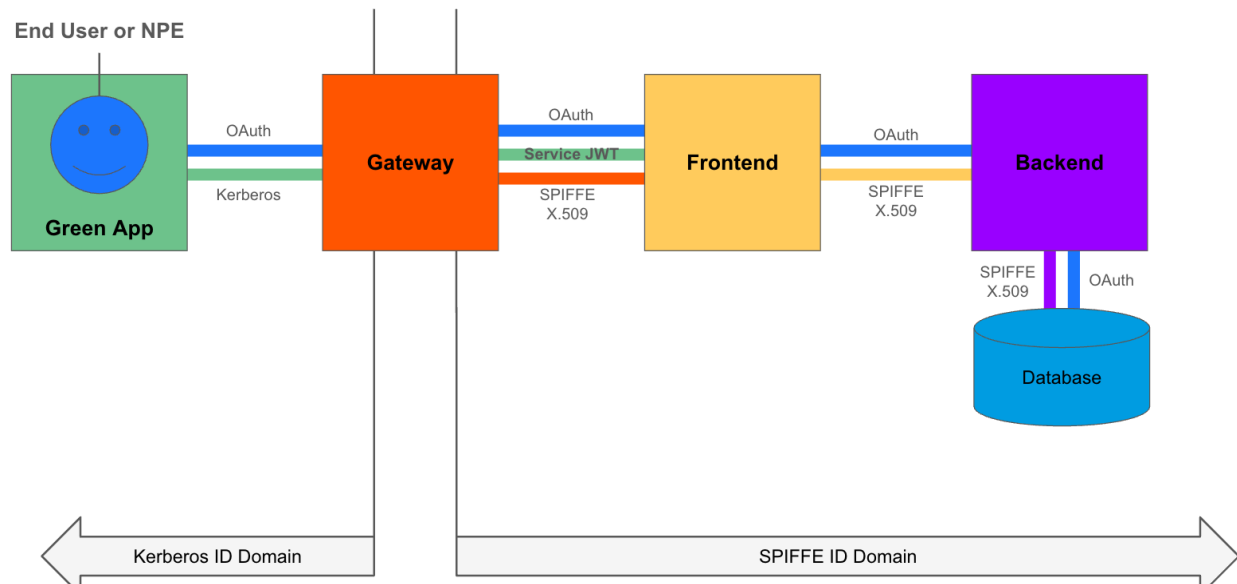
488 Applications that perform identity-based segmentation will need to configure policy for that
489 service account user so that it can act on all of the data that the batch job previously used its
490 service identity to access. At the same time, the application can remove any support for special
491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

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

497 **2.7.3. Requests With a User Identity But No Service Identity**

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

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



511

512

Fig. 4. Handling API Calls with User Identity & No Service Identity

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

520 For example, suppose that an organization A) uses a SPIFFE X.509 identity via mutual TLS for
521 service identity as a service mesh does, B) uses a JWT bearer token for user identity, and C)
522 chooses to represent external service identity as a JWT token attached to the request. The
523 mesh can then enforce that the gateway forward traffic to the service via (A), authenticate the
524 service JWT and authorize the external service (C), and authenticate the end user (B) before
525 forwarding a request to the application. This would fully support authenticating and authorizing
526 all of the communicating parties, and the service in question would not need to be aware of the
527 external service identity or credential. They would simply need to manage a policy of “allowed
528 external service callers” alongside their set of “allowed internal service callers.”

529 **2.7.4. Requests With Both User and Service Identities**

530 In the best case, the legacy systems in question are already doing nearly the right thing in that
531 they have both an end user and a service identity attached to requests. However, since they are
532 a legacy system, those credentials likely do not fully conform to the credentials expected by the
533 parts of the system implementing identity-based segmentation. In that case, those other
534 credentials will need to be translated into the canonical form expected by services performing
535 identity-based segmentation in the infrastructure. Essentially, the user’s authenticated
536 credential should be exchanged with an identity provider for the canonical form expected by

537 the identity-based segmentation portion of the infrastructure (e.g., a JWT bearer token), and
538 the external service's identity should be represented to the internal system as a token so that
539 the policy can be enforced on all three identities in the first hop.

540 **2.7.5. Reaching Out to Other Systems**

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

547 Rather than integrating the API gateway with a variety of identity providers to canonicalize
548 inbound credentials, a gateway with a variety of identity providers to fetch outbound
549 credentials can be integrated instead. For example, an external SaaS API may expect a cloud
550 provider service account as credentials. An egress gateway can be deployed to authenticate
551 and authorize credentials used inside of the organization (i.e., identity-based segmentation)
552 and then exchange the internal identities for the external identities needed by the other
553 system. In this way, services that perform modern identity-based segmentation can integrate
554 with legacy systems with little impact and minimize any code dependencies on those legacy
555 systems.

556 **2.7.6. Mitigating the Confused Deputy**

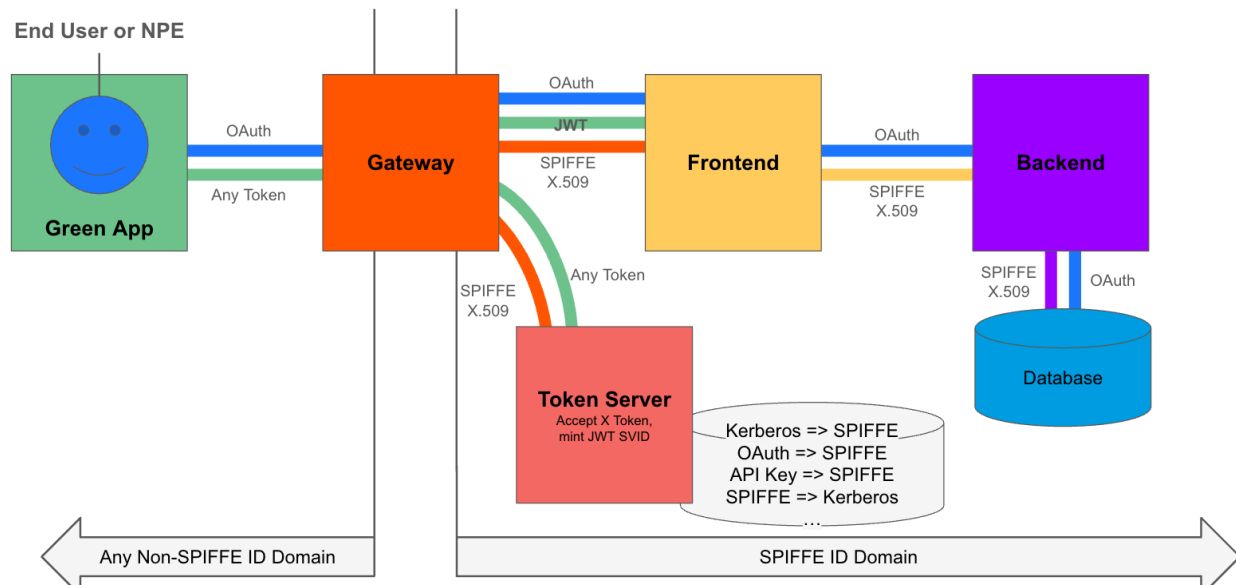
557 One of the biggest risks in any scheme that involves credential exchange is a confused deputy
558 [26], where one caller can trick the "deputy" responsible for handling credentials into using
559 credentials that belong to another caller on its behalf, most often to escalate privileges. Any
560 system that brokers multiple credentials needs more and better authentication and
561 authorization before allowing credentials to be accessed.

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

570 **2.7.7. Identity Canonicalization**

571 Canonicalizing credentials is really canonicalizing the identity domains for which one needs to
572 write policy. Integrating identity providers to standardize credentials at the gateway inherently
573 brings those identities into two identity domains: one for users and one for workloads. This

574 allows for concise and consistent sets of policy that govern access to other services and user
575 access to data. Having both policies in place implements identity-based segmentation and
576 dramatically improves security posture.



577

578

Fig. 5. Identity Canonicalization for Handling API calls

579 For most organizations, implementing credential canonicalization will require either adopting
580 an identity provider wholesale and standardizing on that throughout (including working out
581 legacy integration so that legacy credentials can be used to get credentials via the new
582 provider) or performing identity exchanges, as described in this section. The API gateway is
583 ideally situated to enforce either choice. Performing identity exchanges also requires a mapping
584 of identities across domains as well as a “token server,” which uses that mapping to mint
585 credentials.

586

587 3. Recommended Controls for APIs

588 APIs are the language of the systems we build. They define the “grammar” as well as the
589 “nouns” and “verbs” any user of your system works within. As attackers move up the stack, it’s
590 more important than ever to elevate our security posture to think and manage policy in terms
591 of those “nouns” and “verbs” our APIs expose.

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

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

604 Controls for APIs should cover all of the APIs in the organization, including those exposed to end
605 users, those exposed to partners, and those that are only intended for “internal” consumption.
606 This document’s controls are structured into two primary sections based on the iterative API
607 life cycle discussed in Sec. 1:

- 608 1. Pre-runtime protections, which should be applied during design, development, and
609 testing. These include:
 - 610 a. Creating a well-defined specification for the API’s contract using some interface
611 definition language (IDL) (e.g., OpenAPI, gRPC, Thrift)
 - 612 b. Defining request and response schemas as part of that API specification
 - 613 c. Defining valid ranges of values for fields of each request and response
 - 614 d. Tagging the semantic type of each field of each request and response
 - 615 e. Creating and maintaining an inventory of these API specifications across the
616 organization, including ownership information
- 617 2. Runtime protections, which should be applied to each request and response to the API
618 at runtime. These include:
 - 619 a. Encryption in transit
 - 620 b. End-user authentication and authorization
 - 621 c. Service-to-service authentication and authorization

- 622 d. Request and response validation
- 623 e. Resource-consumption mitigations, including rate limiting, timeouts, and circuit
- 624 breaking
- 625 f. Telemetry (e.g., logging and monitoring) to assess enforcement and detect
- 626 attacks

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

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

640 All organizations should move immediately to act on basic controls, while advanced controls

641 should be evaluated by the organization and applied to APIs based on risk profile.

642 3.1. Pre-Runtime Protections

643 All API controls must be well-defined and inventoried.

644 3.1.1. Basic Pre-Runtime Protections

645 REC-API-1: All APIs must have a specification in the form of a document that describes what

646 endpoints the API exposes (“API spec” for short). To begin, the API spec can be a literal

647 document, a set of internal wiki pages maintained by a team, or something similar. However, it

648 should eventually migrate to a state-of-the-art IDL.

649 REC-API-2: API specifications should use a well-defined IDL (e.g., OpenAPI for HTTP/REST, gRPC

650 for protobuf, Thrift, SOAP for XML).

- 651 • REC-API-2-1: API specs and implementations should conform to industry best practices
- 652 (e.g., a Create-Read-Update-Delete [CRUD] API exposed as HTTP/REST should map the
- 653 CRUD endpoints to the HTTP verbs POST, GET, PUT, and DELETE, respectively) for
- 654 consistency [13].

655 REC-API-3: Request and response schema for each endpoint should be defined by the API

656 specification, including validation guidelines for the values of each field of the request and

657 response (e.g., “the name field is a string and must be shorter than 100 characters”). Additional

658 information makes integration easier and less error-prone for clients and presents the
659 opportunity for automated enforcement, such as the maximum latency (e.g., “the server will
660 drop requests that take longer than 5 seconds to process”) and rate limits (e.g., “by default, 5
661 calls per minute are allowed”).

662 REC-API-4: Organizational API inventory of all internal and external APIs should be maintained.
663 This is in line with the *Identify* directive of the CSF [12]. That inventory should include:

- 664 • Each API’s specification, though the inventory does not need to be the *API*
665 *documentation*
- 666 • Ownership information about the API to simplify the translation of runtime problems to
667 organizational response
- 668 • Runtime information to enable operations and security teams to understand the impact
669 of each API (e.g., service instances, instance IP addresses, runtime service ID, traffic
670 volume, rate of requests and errors, the status of policy enforcement)

671 **3.1.2. Advanced Pre-Runtime Protections**

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

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

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

- 692 • REC-API-6.1: Annotate endpoints and fields with permissions required to enable the use
693 of tooling to automate fine-grained per-field authorization checks. Those authorization
694 checks could then be performed by the API serving infrastructure on behalf of the
695 application or via a common library in the application with standard logging and metrics

696 to facilitate easy audit and ensure continuous enforcement. Once the annotations are
697 present, a variety of runtime implementations are possible.

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

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

711 **3.2. Runtime Protections**

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

714 **3.2.1. Basic Runtime Protections**

715 REC-API-9: All runtime communication must be encrypted, even when the API is "public data"
716 or otherwise unauthenticated. This is necessary to ensure that data has not been tampered
717 with (integrity) and to prevent eavesdropping (confidentiality). Details on encryption in transit
718 can be found in SP 800-53, control SC-8 [15] and SP 800-207A, control ID-SEG-REC-1 [6]. Details
719 on cryptographic algorithms and key lengths can be found in SP 800-57 [16] and FIPS 140-3
720 [17].

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

726 REC-API-11: Authenticate the calling user and service, as described in SP 800-207A, controls ID-
727 SEG-REC2 and ID-SEG-REC4 [6].

728 There are (at minimum) two identities in every API communication: the software calling the API
729 and the end user of that software. For example, it is common to use an API key to identify
730 calling software and an OAuth Bearer token to identify the end user. This is true even if the
731 end-user identity is an NPE (i.e., internal software calling other internal software should use
732 something like a service account to identify the user making the requests). The service identity

733 may contain information (e.g., the device being used to access the system) in addition to a
734 token from the software itself (e.g., an API key).

- 735 • REC-API-11-1: Identities must be cryptographically verifiable and should not use weak
736 signing algorithms (e.g., no JWTs with “alg: none,” weak algorithms, or short key-
737 lengths). SP 800-57 [16] discusses the strengths of cryptographic algorithms and the
738 necessary key lengths for each.
- 739 • REC-API-11-2: Authentication should use standard mechanisms whenever possible. For
740 example, end user authentication should use a mechanism such as OpenID Connect
741 (OIDC), OAuth2, or SAML. Services should use a mechanism like SPIFFE SVIDs, JSON Web
742 Tokens (JWTs), API keys, or similar.
- 743 • REC-API-11-3: Tokens must support expiry so that credentials are cycled regularly.
744 Checking for expiry must be an inherent part of token validation. For example, when
745 processing JWTs, the “exp” claim RFC 7519 [18] must be checked. Similarly, when
746 processing an X.509 SVID, check the validity period’s “Not Before” and “Not After” [19].
- 747 • REC-API-11-4: Return opaque tokens to untrusted systems. It is common for credential
748 tokens to encode information about the internals of the system (e.g., minting a JWT to
749 represent a user in the infrastructure that includes claims that represent the user’s
750 capabilities in the system). This is a common scheme to simply and reliably enforce
751 authorization per hop: validate the JWT, and check whether it contains the “claim” that
752 represents the permission for an API endpoint. These claims encode all local operations
753 that can be performed with data from the request and the local application.

754 Returning a token with these details to an external user may risk leaking information
755 about the internals of the system. This is where the following issues become critical to
756 the safety of the API: how permissions are modeled, the set of internal
757 permissions/claims that map to a given external API endpoint, and information about
758 the path that the request traverses through the infrastructure.

759 REC-API-12: Authorize the calling user and service for each identity on the request, including
760 whether the calling software system is allowed to access the API endpoint and whether the end
761 user is authorized to take the action on the resource represented by the endpoint. See SP 800-
762 207A [6], controls ID-SEG-REC2 and ID-SEG-REC4.

763 Getting these authorization checks correct is one of the most common mistakes in API security
764 [7]. REC-API-6 discusses annotating each request or endpoint with the permission required by
765 the end user to call that endpoint on a resource. With annotations like those in place, runtime
766 tooling can be implemented to ensure that those annotations are transformed into runtime
767 permission checks against the authorization system. Combined with a robust DevOps process to
768 ensure that annotations are present on APIs before they can be deployed, there can be a high
769 degree of assurance that the correct authorization is being performed at the platform level. The
770 idea of using the service mesh to achieve this is discussed in SP 800-204B [3].

771 REC-API-13: Validate each request and response per the API schema before it is processed by
772 the business logic (e.g., ensure that the request has a “name” field that is a string and no other

773 fields). This ensures that applications only receive well-formed input and minimizes a class of
774 errors and data leaks due to validation inline in the business logic. Additionally, validate that
775 each response from the server conforms to the expected response schema to help prevent a
776 variety of data leaks, abuses, or mistakes.

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

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

- 791 • REC-API-15-1: Rate limit all API access for all callers to ensure fair utilization across
792 users, help with capacity planning, and mitigate the risk of unrestricted resource
793 consumption. See REC-API-16 for recommendations on specific rate-limiting
794 implementations.
- 795 • REC-API-15-2: Apply timeouts to all requests, including the API gateway. This should be
796 done at the TCP level, where connections are automatically timed out after a modest
797 time (e.g., 5 minutes) rather than the kernel's default of more than one hour per
798 connection. Timeouts should also be configured at the application level. If a required
799 operation should complete in five seconds as part of the API contract, set a 6-second
800 timeout for it. This ensures that the resources in a service do not wait for a response
801 that will never arrive.
- 802 • REC-API-15-3: Apply bandwidth and payload limits to enforce maximum request and
803 response sizes. The "correct" limit is highly contextual and based on the organization
804 and application (e.g., a bank will have very different expectations than a video streaming
805 company). This helps avoid a variety of risks related to malicious input and DoS.
- 806 • REC-API-15-4: Validate and limit user-supplied query parameters (e.g., amount of
807 processing done, size of their response based on user input), especially in the context of
808 what the system can support and what is typical for users of the system. For example:
 - 809 ○ The number of elements returned per page of a paginated list API. If a typical
810 user has 100 items, cap the maximum number of elements per page to 1000.

- 811 ○ Time ranges in dynamic queries. If a system is intended for viewing recent
812 events, and the user can provide a time range, limit that range to the last 30 days
813 rather than allowing the user to query “from 1972 onward.”
- 814 ○ GraphQL and similar API facade systems that support query languages over many
815 APIs should have limits on the queries that users can execute (i.e., approved or
816 predefined queries only) and caps on the number of outbound calls allowed in
817 the execution of a single query.

818 REC-API-16: Rate limiting recommendations are one of the most effective tools to mitigate
819 unrestricted resource consumption and can increase the challenge and discoverability of many
820 attacks with a goal of leaking sensitive information via data exfiltration from API calls (e.g.,
821 scraping all chat logs from an organization with a script impersonating a chat client). Most
822 organizations apply some type of rate limit to “external” traffic, but it is equally important to
823 rate-limit internal callers. It is very easy to unintentionally cause a DoS on an internal system
824 with poorly conceived code. It is equally critical to consider the limits placed on internal
825 software that call out to external systems (see Sec. 2.6.2).

826 The following recommendations on rate-limiting configuration address common pitfalls and
827 misunderstandings:

- 828 ● REC-API-16-1: Rate limits are not quotas. A quota is a usage limit on an API over an
829 extended duration (e.g., per month) that is associated with a user’s payment or billing
830 structure. Many organizations have “API usage tiers” that map prices to higher per-
831 month limits. These quotas need to be strictly enforced and are typically used to
832 generate billing reports that are sent to customers. In contrast, rate limits are intended
833 to protect the system from overuse and help ensure fair usage across separate,
834 concurrent callers. Rate limits do not need to be exact in the way that quotas must be.
- 835 ● REC-API-16-2: Rate limits for total load provide little benefit and should be dimensioned
836 by user (e.g., 83 requests per 5 minutes per user) using the source IP address or end-
837 user credential as the key. Rate limits without a user dimension (e.g., service can receive
838 1,000 requests per 5 minutes total) are not particularly effective and allow some users
839 to impact others (e.g., DoS risk). This is true even when total limits are dimensioned by
840 service instance (e.g., a single instance cannot receive more than 100 requests per 5
841 minutes). Circuit breaking functions must be used to provide protective limits on
842 concurrency for a service instance. More information on circuit breaking and other
843 resiliency and load-shedding techniques can be found in 800-204A, Sec. 2.3 [2].
- 844 ● REC-API-16-3: Rate limits should be short in duration (e.g., per 60-seconds, per 5-
845 minutes). A rate limit is defined as the number of calls allowed over a time period (e.g.,
846 24,000 requests per 24 hours; 1,000 requests per hour; 16.5 requests per minute). Most
847 systems allow for the configuration of both the number of calls and the amount of time
848 over which they are allowed.

849 However, there are two problems with per 24-hour rate limits. First, they cause outages
850 for callers that resolve themselves when the rate limit server resets for the next 24-hour

851 period, even if the rate limit was originally set correctly based on the client's expected
852 usage. A successful API ecosystem will see the increased usage of APIs over time, which
853 results in the increased usage of their dependencies and those APIs. This is the typical
854 organic growth of API usage. Adjusting rate limits before they caused outages is almost
855 never a priority for application teams, so over time, clients may see the API begin to
856 randomly fail with 400 errors. Second, per 24-hour rate limits can result in spiky traffic
857 for the service, where a client consumes the entire 24-hour limit over a very short time
858 and causes a heavy load on the services.

859 Shorter time limits allow clients to experience a few intermittent failures every minute
860 or five as their traffic grows organically rather than total failure with per 24-hours.
861 Additionally, the system will experience smoother traffic overall because a single client
862 must pace their consumption over a longer duration, resulting in less load from each
863 client at any given time.

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

872 REC-API-18: API access must be monitored to ensure that the API serving stack provides
873 sufficient telemetry to assess the availability of APIs and to ensure that policies are being
874 enforced. The traditional triad of logging, metrics, and distributed traces is recommended. All
875 three should be tagged with information about the API being accessed in addition to the
876 runtime service so that service calls can be traced back to APIs.

877 For the API gateway itself, a range of signals should be produced to enable the identification of:

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

885 General information on audit and logging requirements can be found in SP 800-53 [15], AU-2
886 Event Logging, AU-3 Content of Audit Records, and AU-12 Audit Record Generation.
887 Information on service mesh telemetry, which can be used for audit and logging, can be found
888 in SP 800-204A [2], SM-DR21 through SM-DR24.

889 3.2.2. Advanced Runtime Protections

890 REC-API-18: Field-level validation using API schema annotations can be used to validate the
891 values of requests and responses at runtime. This is beyond the basic syntactic validation of
892 REC-API-13 (e.g., “there is a name field, and it is a string”) and more like semantic validation
893 (e.g., “the name field must not be longer than 100 characters,” or “the amount field must be
894 positive and less than 2 million”). This can be implemented by the API gateway as part of a
895 cross-cutting policy. An API spec is required (REC-API-2) and should be in a central inventory
896 (REC-API-4). The API gateway team can then enforce the validation of all requests traversing an
897 API gateway. This reduces the risks of insufficient input verification and leaking sensitive
898 information compared to ad hoc, error-prone implementations in each application or standard
899 implementations embedded in the application itself via SDK, which tend to be difficult to
900 update. A timely update is an imperative for infrastructure that enforces security policy.

901 REC-API-19: Authorization and filtering using API schema annotations enforce access to
902 resources and fields per caller. In this case, the API gateway itself is the policy enforcement
903 point, and it defers to an authorization system to make decisions. The information from the API
904 schema is enough to extract credentials from the request, identify the target endpoint and its
905 associated tags/permissions, and use those to form a call to the authorization service (e.g., “is
906 the request’s end user allowed to perform the endpoint’s permission on the object targeted by
907 the request?”). The API gateway can then enforce the result of the call at runtime. There are at
908 least three levels of assurance that can be achieved, and each build on the previous one to
909 further mitigate risks at increased runtime or development-time cost:

- 910 • REC-API-19.1: Resource-level authorization as a cross-cutting policy should be enforced
911 on all requests using endpoint-level annotations that define the permissions required to
912 call the endpoint (REC-API-6.1). This can be done at the platform level leveraging the API
913 gateway. When combined with a decentralized gateway pattern (Sec. 4.3), this
914 implements ID-SEG-REC-4 [6] at every hop.¹ This also helps prevent and potentially
915 eliminate missing authorization (Sec. 2.2), depending on the organizational guardrails in
916 place. For example, an organization can build an API inventory by mandating an API spec
917 with endpoint-level permission annotations as part of each app’s “ticket to the
918 platform” (i.e., the data that an app team needs to submit to run their application on
919 the organization’s infrastructure and platform). Combined with standard patterns for
920 authentication (REC-API-11), this can ensure that the correct authentication and some
921 authorization are performed. However, additional organizational controls are required
922 to ensure that the permissions are correct and sufficient in order to fully mitigate risks
923 around authorization (see Sec. 2.2).

924 Achieving correct and sufficient authorization at the resource level is likely all that most
925 organizations need to achieve. It mitigates the predominant risks identified by the
926 OWASP API Security Top 10 [7] with respect to authorization. Moving beyond this level

¹ Other patterns have a wider perimeter and are susceptible to the API gateway being bypassed. Therefore, they do not satisfy ID-SEG-REC-4.

927 of assurance into REC-API-19.2 and REC-API-19.3 shifts the focus to mitigating the risk of
928 leaking sensitive information.

929 • REC-API-19.2: Field-level visibility as a cross-cutting policy can leverage basic “Public”
930 and “Private” annotations on each field. The authorization check effectively asks
931 whether data should be visible to “external” callers.² These coarse-grained
932 Public/Private annotations are particularly effective on common types shared across
933 many APIs in the organization. For example, a standard error reporting pattern used by
934 all APIs can leverage field-level annotations to differentiate “user” facing errors versus
935 “developer” facing errors, mitigating the risk of leaking sensitive information via errors.
936 The gRPC Status proto [21] is an example of a consistent error reporting pattern. In the
937 gRPC case, field-level annotations would reside in the message used for the status’s
938 “details.”

939 • REC-API-19.3: Field-level authorization as a cross-cutting policy can be leveraged to
940 perform fine-grained field-level authorization (REC-API-6.1). This extends the idea of
941 REC-API-19.1 down to the level of each individual field of the response and allows for
942 the filtering of API objects per-use to implement sophisticated access control schemes.

943 While this kind of approach offers a very high level of data security, it causes a sharp
944 increase in the number of policy checks that the authorization system must perform and
945 requires active participation by application developers to keep permissions per field up
946 to date as the application evolves. For example, a resource-level authorization check
947 requires one authorization decision per request. A field-level authorization check
948 requires one authorization decision for the request plus an additional decision for each
949 field of the response. Even an object with a modest number of fields (e.g., 5) results in
950 whole-number multiples more policy decisions made by the authorization system. For
951 developers, the purpose and therefore permission of an endpoint rarely changes, but
952 the fields of the request and response objects for that endpoint regularly evolve over
953 time. This makes upkeep for permissions at the field level more expensive for
954 application developers versus endpoint-level annotations (REC-API-19.1).

955 As a result of the cost and load on the authorization system, this level of fine-grained
956 checking is typically only used in the most high-risk situations and only by sophisticated
957 organizations.

958 SP 800-204B [3] discusses the advantages of using a decentralized API gateway architecture
959 when implementing fine-grained authorization checks. When choosing to implement these
960 authorization policy checks under the centralized and hybrid patterns, care must be taken to
961 ensure that the gateways are not bypassed. For example, a service-level authorization policy
962 could disallow any traffic except from the API gateway as a means of defeating an attempt to
963 bypass gateway checks via pivoting inside the infrastructure.

964 REC-API-20: Traffic monitoring and policy using semantic field labels can log and monitor the
965 flow of sensitive data in a system. Further, the API Gateway can be used as a policy

² REC-API-19.1 i focuses on requests, while this control focuses on the data that an application returns to callers in responses. They are complementary controls.

966 enforcement point to control the flow of that data, potentially blocking traffic flows that transit
967 significant amounts of data. Ultimately, with annotations and enforcement in place, the flow of
968 sensitive data in the organization can be governed by mandatory access control (MAC) policies.
969 A MAC policy is enforced by the authorization system, regardless of the user or resource in
970 question. For example, while not explicitly stated as a hard rule in PCI DSS, a MAC policy
971 followed in implementation of systems handling PCI data is that they should be isolated from
972 systems that do not implement PCI DSS controls to maintain security and prevent potential
973 breaches. Such a MAC policy can be enforced with a combination of understanding PCI-
974 compliant services in the infrastructure and data tags on the semantic types of data that flow
975 through the system.

976 REC-API-21: Non-signature payload scanning (for generative AI APIs) analyzes request and
977 response data for sensitive information that may not be a literal attack signature. Tools typically
978 analyze (e.g., via regression, AI, simple matching and word filtering) the responses returned by
979 servers to score the risk that they contain sensitive information and take action to block that
980 traffic. Increasingly, AI agents are being deployed to assess the risk of data generated by other
981 agents. At a high level, this technique is like a web application firewall (WAF), but WAFs are
982 fundamentally signature-based, while these analyses are fundamentally content-based.

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

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

- 994 • REC-API-21.1: Semantic data discovery tools are typically very good for identifying the
995 type of information flowing through a system (e.g., string, email address). Building the
996 inventory of APIs and the developers adopting well-defined API schemas with
997 meaningful annotations takes time. Runtime tools such as these are very helpful for
998 initial discovery, ensuring that rollout is complete across all services, and ensuring that
999 services stay in compliance after the policy is rolled out. When it is reasonable to
1000 leverage due to compute and latency constraints, an organization benefits from
1001 inspecting traffic for sensitive data flow, even beyond field-level annotations.

1002 REC-API-22: Fine-grained blocking for specific requests can prevent a DoS or service crash.
1003 These bad inputs can often trigger a cascading failure [23], but the queries may not be
1004 malicious in nature (e.g., users using the system in ways that it was not intended or designed
1005 for). In cybersecurity, this is sometimes called the “query of death” (QoD) [24]. These tools help
1006 mitigate the risks of unrestricted resource consumption and malicious input validation.

1007 As a system grows in size and complexity, it is necessary to be able to pin-point block these
1008 kinds of queries to keep the system stable and available. Depending on the complexity of the
1009 query and environment, it may be possible to leverage a WAF or non-signature payload
1010 scanning tools to block some types of QoDs. However, application code changes may be
1011 required — sometimes even rearchitecting the application itself — to mitigate the impact of
1012 these kinds of queries.

1013 The detailed controls in this section fit into broad classes, and their association with the
1014 DevSecOps phases is discussed in Appendix B. This emphasizes the observation that APIs should
1015 be treated as any other software and go through an iterative, continuous life cycle.

1016

1017 **4. Implementation Patterns and Trade-Offs for API Protections**

1018 Regardless of the mechanism or architecture of an API and its services, there is a core set of
1019 capabilities required to realize the controls outlined in this document:

- 1020 • Authentication and authorization
- 1021 • Request and response validation
- 1022 • Rate limiting
- 1023 • Circuit breaking
- 1024 • Error handling
- 1025 • Logging and Monitoring

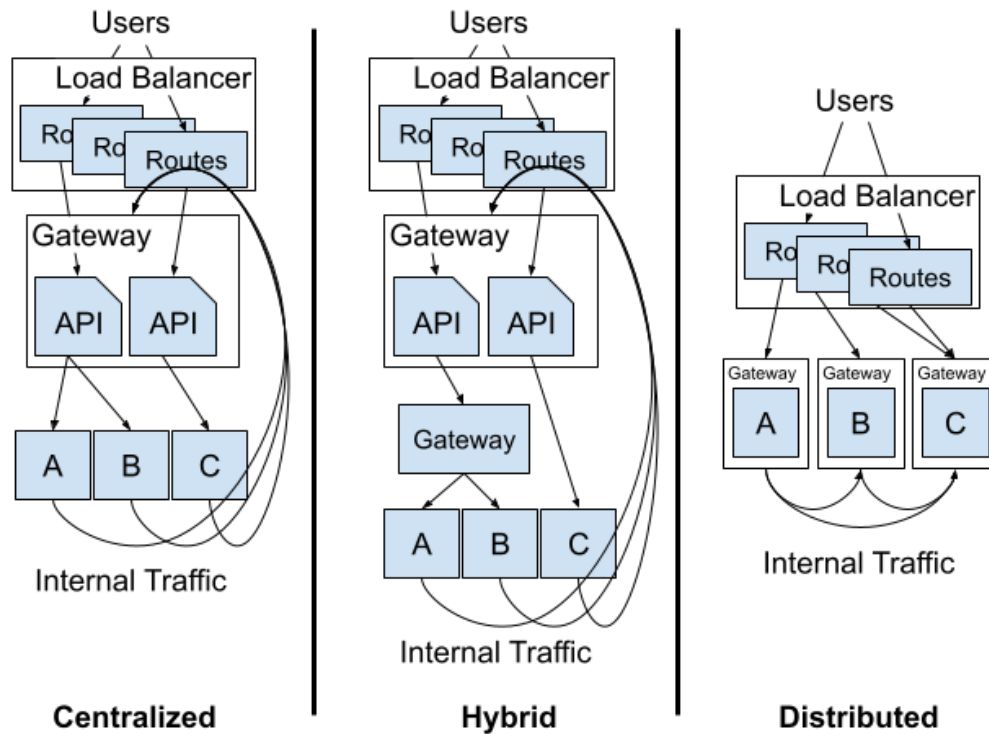
1026 In addition to these core capabilities for security, APIs that serve infrastructure typically deal
1027 with other common concerns:

- 1028 • Service discovery
- 1029 • Routing
- 1030 • Protocol conversion
- 1031 • Caching

1032 Three components are needed to provide this functionality to serve an API:

- 1033 1. A *gateway* to implement the API-oriented policy
- 1034 2. The *service* itself to implement the API's business logic
- 1035 3. A method to get traffic to gateway instances (e.g., DNS and a network load balancer) to
1036 facilitate service discovery, load balancing, and network reachability to horizontally
1037 scaled instances of the gateway itself

1038 For example, if the Gateway functionality is implemented via a Kubernetes ingress routing to a
1039 pod (i.e., the service instance), then callers outside of the network to reach the gateway will
1040 require the cloud provider or data center network team to provision a network load balancer in
1041 front to route network traffic to the Kubernetes load balancer service.



1042

1043

Fig. 6. API gateway patterns

1044 Three patterns have been developed by industry to implement these capabilities:

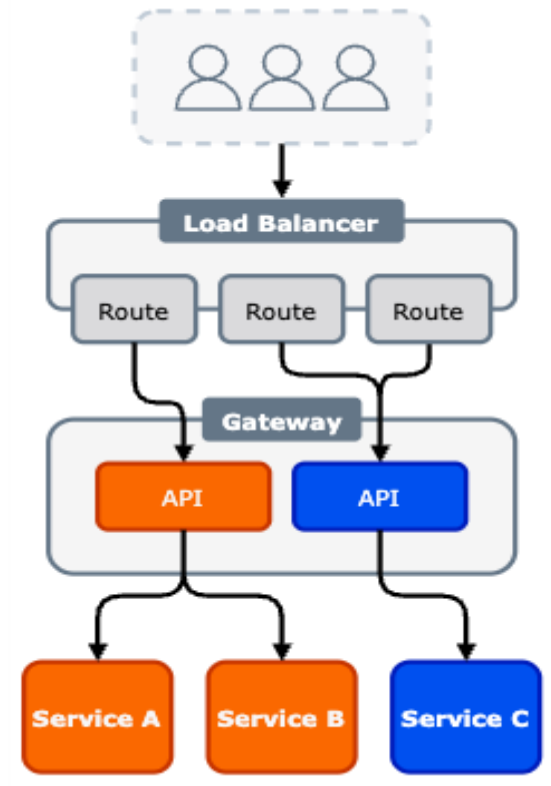
- 1045 1. Centralized gateway — Protections for all APIs in the enterprise are implemented by a
1046 single shared component: an API gateway.
- 1047 2. Hybrid deployment — Cross-cutting policies (e.g., authentication) are implemented in
1048 the centralized shared gateway, but application-specific policies (e.g., authorization) are
1049 implemented in the application itself or by components owned by the application team.
- 1050 3. Decentralized gateways — All policy checks are performed by gateways dedicated to
1051 each application, often deployed beside each service instance.

1052 All three patterns can achieve all of the controls outlined in this document and be used by
1053 organizations to operate their APIs safely and confidently. Further, many of these patterns may
1054 be in use within a single organization. This section explores the engineering design trade-offs
1055 that each pattern provides in terms of risks and operational overhead.

1056 Many API gateway products provide management capabilities, such as API key issuance,
1057 discovery documentation (i.e., API definition) hosting, documentation for client developers, and
1058 support for quotas and billing tiers. These are all valuable features in the enterprise setting, but
1059 all of them can be supported across any implementation pattern and are therefore not
1060 addressed in this section.

1061 **4.1. Centralized API Gateway**

1062 The centralized API gateway pattern implements protections for all APIs with a single
1063 component: an “API gateway” that is often deployed close to the perimeter of the system.
1064 External traffic enters through the gateway, typically via a load balancer. Internal traffic
1065 “hairpins” through the gateway as well, which facilitates service-to-service communication
1066 inside the infrastructure. That internal, service-to-service traffic may also have to traverse the
1067 load balancer for some service instances. Fig. 5 shows a common configuration for a centralized
1068 API gateway pattern.



1069
1070 **Fig. 7. Centralized API gateway pattern**

1071 An API gateway is typically a software application that can be scaled horizontally (i.e., more
1072 instances can be deployed side by side). This is one of the reasons why an API gateway often
1073 sits behind a load balancer, even for internal service-to-service traffic use cases.

1074 Advantages of this pattern include:

- 1075 • A single policy enforcement point that is easy to monitor and audit, making it simple to
1076 verify that policy is enforced for all traffic that traverses the gateway.
- 1077 • Implementation matches the organizational structure. Typically, large organizations
1078 have a single API team that owns the centralized gateway component. That team is
1079 responsible for and able to execute on when an API is available, which API endpoints are

1080 failing, whether policies are being enforced, whether the configuration up to date, and
1081 other issues.

- 1082 • Streamlined setup for application developers who need to “onboard” their API but do
1083 not need to deploy or maintain any additional runtime components.

1084 Disadvantages of this pattern include:

- 1085 • Shared fate outages. Because there is a single component, an outage of that component
1086 causes an outage for all APIs, which can be problematic for mission-critical APIs that
1087 need to operate continuously.

- 1088 • Noisy neighbors, where traffic consumes resources for some APIs and increases latency
1089 for all APIs. In the worst case, one application team may submit invalid configuration
1090 parameters for a service that may crash or cause DoS on the API gateway, triggering a
1091 shared fate outage for other APIs.

- 1092 • Long change lead times due to managing how the changes to an individual team’s API
1093 configuration impact the shared gateway. This is a frequent side-effect of controls
1094 added to mitigate shared fate outages and noisy neighbors.

- 1095 • Cost attribution. All requests are handled by the central gateways, and resources spent
1096 per request per API (e.g., on payload validation) are uneven. Therefore, it can be difficult
1097 to attribute API gateway runtime costs to internal application teams. This can be a
1098 problem for companies that implement an internal resource economy for planning by
1099 assigning cost centers for each application team.

- 1100 • Caching the results of policy decisions at runtime becomes critical when implementing
1101 the policies outlined in this SP due to the sheer number of policy checks required.
1102 Caching both increases client-perceived availability and reduces the load on key
1103 systems, like authentication and authorization. However, two layers of load balancing
1104 (i.e., network load balancer to API gateway and API gateway to service instance) tend to
1105 result in poor cache hit rates across policies enforced by the API gateway and for user
1106 data in the application layer itself. While some techniques can be used to mitigate this
1107 (e.g., distributed caches or streaming connections), they generally add additional
1108 development or operational overhead for the application team, API gateway team, or
1109 both.

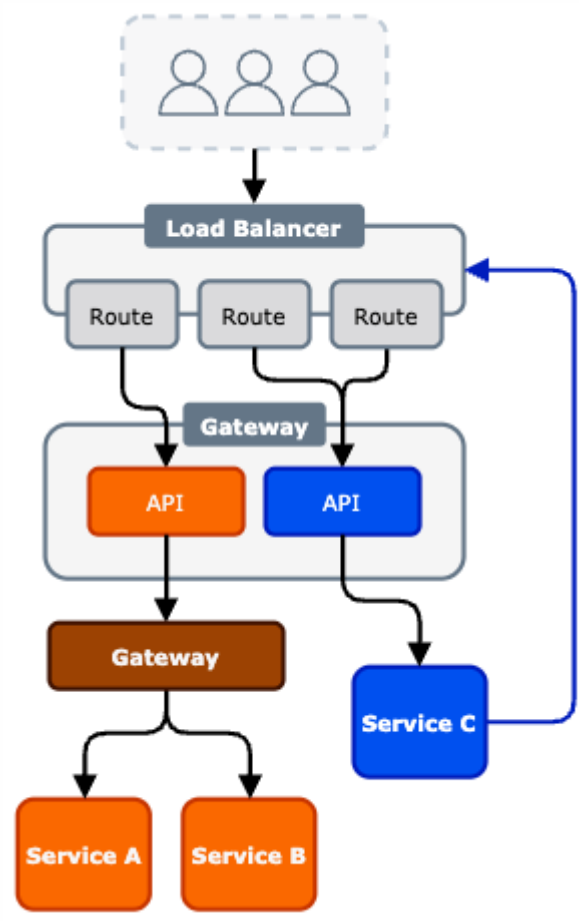
- 1110 • Because a shared gateway is located at the perimeter, it can be bypassed (e.g., via an
1111 attacker pivoting inside the perimeter), which in turn bypasses the policy checks
1112 enforced by that gateway. This can be mitigated with techniques like service-to-service
1113 access policies that ensure that applications only receive traffic via the centralized

1114 gateway or by attaching proofs (i.e., credentials) to the request that allow an application
1115 to authenticate that the request was handled by the gateway.

1116 4.2. Hybrid Deployments

1117 Hybrid gateway deployments split policy enforcement responsibilities between a centralized
1118 gateway and the applications themselves. Cross-cutting policies (e.g., authentication, service
1119 discovery, routing, rate limiting, caching) are handled by the centralized gateway. Application-
1120 specific policies (e.g., authorization, request and response validation, protocol conversion, error
1121 handling, logging, monitoring) are handled by the application team. This can manifest in the
1122 application itself (e.g., gRPC) or as a separate deployment that handles traffic before the
1123 application (e.g., GraphQL or Spring Cloud Gateway). As with the centralized pattern, all
1124 internal and external traffic between applications must first go through the centralized gateway
1125 and, in some instances, through the load balancer. Fig. 6 shows the schematic diagram of a
1126 distributed gateway pattern.

1127



1128

1129

Fig. 8. Distributed gateway pattern (hybrid deployment)

1130 Overall, this pattern behaves similarly to the centralized API gateway pattern, except that some
1131 of the most failure-prone parts of the centralized pattern are delegated to the application
1132 teams. This streamlines API gateway operations and enables app teams to move at their own
1133 pace. However, it also shifts the responsibility for some runtime operational and security
1134 concerns from the API gateway team to those application teams.

1135 The exact split of responsibilities between gateway and application (e.g., sidecar in a service
1136 mesh architecture) can vary greatly across different organizations based on their risk profiles
1137 and past experiences. Typically, the gateway takes responsibility for:

- 1138 • Authentication
- 1139 • Rate limiting
- 1140 • Circuit breaking
- 1141 • Service discovery
- 1142 • Routing
- 1143 • Caching
- 1144 • Network-level load balancing

1145 The application or dedicated gateway is responsible for:

- 1146 • Authorization
- 1147 • Request/response validation
- 1148 • Protocol conversion
- 1149 • Error handling
- 1150 • Application-instance load balancing

1151 Both are responsible for logging and monitoring to enable visibility into the state of the system
1152 and to ensure that policies are being enforced at runtime.

1153 There are similar advantages as the centralized gateway pattern that also include:

- 1154 • Mitigation of most shared-fate outages and noisy neighbors by moving the most error-
1155 prone processing like request validation out of the shared gateway and delegating to
1156 the application or dedicated gateway.
- 1157 • Increased iteration speed due to the ability to update configurations with less process
1158 overhead and hence quickening the time involved. This is possible due to reduced risk of
1159 shared fate outage.

1160 Disadvantages include:

- 1161 • The enforcement of policies is split across the API gateway and many service instances,
1162 which makes it more challenging to ensure that the policy is being enforced consistently
1163 and correctly.

- 1164 • There is increased operational burden on application teams compared to the centralized
1165 API gateway pattern, as they are now responsible for ensuring that some policies are
1166 enforced in their application.
- 1167 • Not all classes of shared fate outages and noisy neighbors can be eliminated because
1168 the shared central gateway is doing at least some application layer processing.
- 1169 • Cost attribution is significantly improved compared to the centralized pattern because
1170 the most expensive runtime policies are implemented by the application teams.
1171 However, the centralized gateway can still be very expensive to operate at high scales
1172 and is as difficult to attribute costs as in the centralized pattern.
- 1173 • Caching hit rates also suffer similarly to the centralized pattern for the same reasons.
- 1174 • Bypassability/pivot

1175 **4.3. Decentralized Gateways**

1176 In a decentralized approach, the gateway is directly associated with the application, which is
1177 owned by a single team. This ensures that changes are isolated to services owned by that team
1178 and that the potential for shared fate outages does not arise. Changes to each gateway are
1179 “safe” from the organization’s perspective: a bad change will not cause additional problems for
1180 other teams, and the team that caused the outage to occur can fix the problem.

1181 External traffic must still enter through a load balancer, which does not enforce any policy and
1182 only performs routing. Internal traffic may use the same load balancer but may be routed
1183 directly peer-to-peer, removing the central gateway from internal traffic as desired, since
1184 enforcement of policy happens at the service instance.

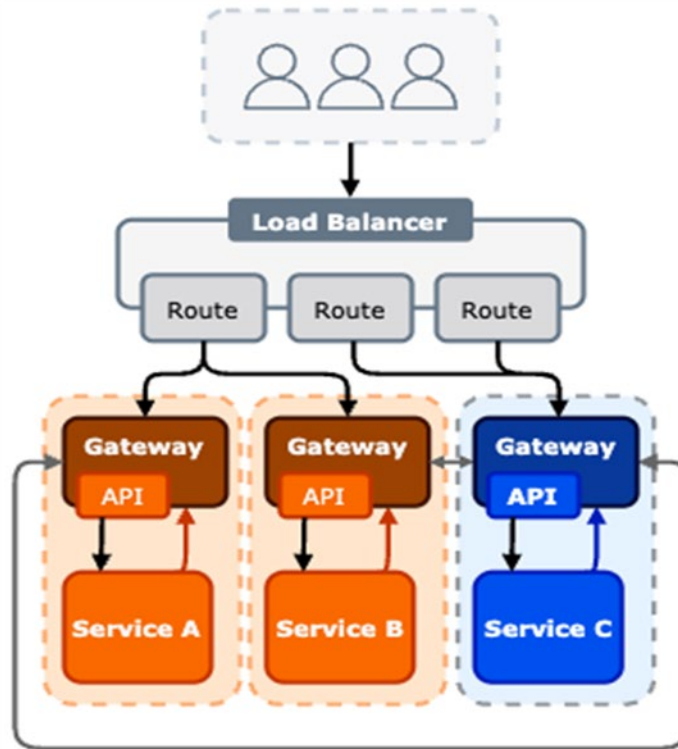
1185 This leaves two key challenges that the implementation must address:

- 1186 1. Ensuring that the remaining shared configuration (i.e., the load balancer) is safe for each
1187 team’s changes
- 1188 2. Ensuring that both cross-cutting and application-specific policies are enforced
1189 consistently across the organization

1190 Keeping the load balancer’s configuration safe is a universal problem across all three
1191 implementations. However, it is most acute in the decentralized pattern because the load
1192 balancer must cope with configuration for many applications, while only the API gateway’s
1193 configuration needs to be present in the other patterns. Regardless of implementation pattern,
1194 this is most often handled at the business process level. Organizations decide on a fixed naming
1195 scheme that is enforced by and during by the CI/CD process or is otherwise hidden by the
1196 organization’s platform (e.g., subdomains-per-service, such as foo.api.example.com,
1197 bar.api.example.com; paths-per-service, such as api.example.com/v1/foo,
1198 api.example.com/v2/bar).

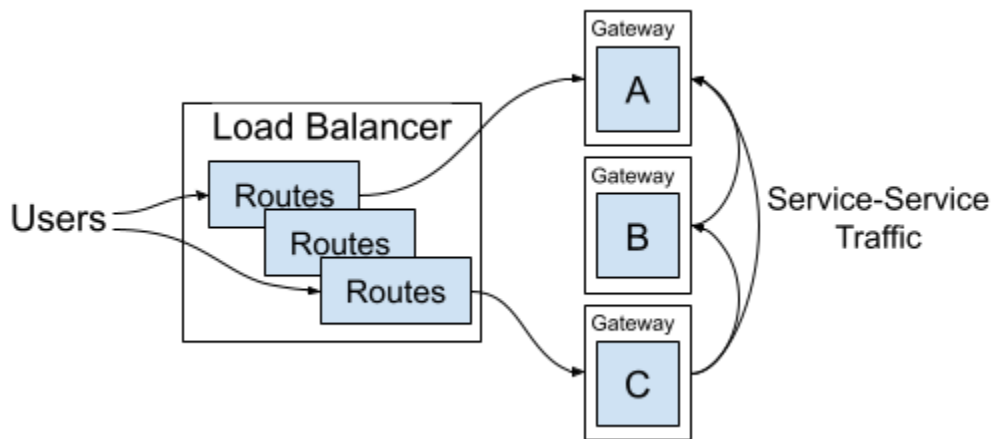
1199 The challenge of cross-cutting policy is unique to this pattern. In recent years, it has been solved
1200 robustly in open source via the service mesh, which can provide a single point for policy

1201 management and use its proxies to enforce those policies (i.e., API protections) at each service
1202 instance. The service's properties [2] and use for security [3][6] have been covered in other
1203 NIST guidance documents.



1204
1205

Fig. 9. Decentralized API gateway pattern



1206
1207

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

1208 The advantages of a decentralized gateway pattern include:

- 1209 • All processing is done per application team (i.e., no noisy neighbors), and the risk of a
1210 shared fate outage is only present on the load balancer, which is a risk shared across all
1211 implementation patterns.
- 1212 • It has the highest rate of change for app teams because they have no external
1213 dependencies and little chance of causing outages for other teams.
- 1214 • A cross-cutting policy can be managed by the central API gateway team via the
1215 gateway's control plane (e.g., with the service mesh). This pattern can be adopted
1216 harmoniously in a mixed environment, where some APIs are implemented via any of the
1217 three patterns in a single organization.
- 1218 • Cost attribution is straightforward and no more or less challenging than attributing any
1219 compute resource spent by teams in the organization.
- 1220 • Cache locality is typically better than in the other patterns because there is only a single
1221 layer of load balancing, and the gateway is co-located with the application. This means
1222 that gateway policy checks for a given user are cached alongside the application
1223 instance caching business logic data for that user. However, if a user's request is load-
1224 balanced across multiple service instances, then "duplicate" policy checks have to be
1225 performed that would not be required in the other patterns.
- 1226 Disadvantages include:
 - 1227 • Because the policy is checked and easily cached per application instance, there can be
1228 many more policy checks in the system overall. Any time a user's request is load-
1229 balanced to a new service instance, it is highly likely that a new policy check has to be
1230 performed. This is an inherent problem in any zero-trust system, which pushes
1231 enforcement to the application instance and likely necessitates the adoption of a
1232 distributed cache managed alongside or as part of the API-serving infrastructure.
 - 1233 • The pattern puts the most burden on application teams. Those teams have to interact
1234 with the team managing the load balancer for each API they expose and need to
1235 operate at least some of the API-serving infrastructure (e.g., making sure that they have
1236 a gateway deployed and routing). Technology like a service mesh can help simplify this,
1237 but a burden remains.
 - 1238 • Auditing and verifying policy enforcement can be challenging as enforcement is
1239 distributed across all application instances. A robust, distributed gateway
1240 implementation (e.g., a service mesh) can help mitigate this via centralized
1241 configuration control combined with distributed enforcement and consistent telemetry.
1242 If an organization can audit and verify a hybrid gateway pattern, a distributed gateway
1243 pattern can be supported with little additional effort.

1244 4.4. Related Technologies

1245 Other technologies fit in and overlap with simplified API gateway patterns and architectures.
1246 Notable companion technologies include:

- 1247 • Web application firewalls (WAFs)
- 1248 • Bot detection
- 1249 • DoS/DDoS mitigation
- 1250 • API endpoint protection
- 1251 • Web application and API protection (WAAP)

1252 **4.4.1. Web Application Firewalls**

1253 Web application firewalls (WAFs) mitigate risks related to a request’s metadata and payload
1254 without needing the application to be involved. In other words, they can be treated as a cross-
1255 cutting policy and managed by a central team.

1256 WAFs work at the application level and operate on parsed HTTP requests (i.e., they can
1257 implement policy per header and on request bodies). However, WAFs generally do not work at
1258 the API level. A WAF can scan a request for a payload that looks like a SQL injection attack, but
1259 it cannot assert, for example, that a request has a “name” field that is a string less than 100
1260 characters long. As such, a WAF is an excellent first step for organizations to implement the
1261 policies outlined in this document, but it is not a complete solution.

1262 The Open Worldwide Application Security Project (OWASP) publishes research on
1263 vulnerabilities based on data from its partners. The OWASP API Security Top 10 [7] was
1264 consulted extensively in the preparation of this document. OWASP *also* publishes a generic set
1265 of WAF rules — the Core Rule Set (CRS) [25] — that aim to mitigate many common attacks. The
1266 CRS should be treated as a starting point for any organization’s WAF policy. Deploying a WAF
1267 with at least the CRS enabled helps mitigate risks, including malicious input (see Sec. 2.6.2),
1268 unrestricted resource consumption (see Sec. 2.4), and leaking sensitive information (see Sec.
1269 2.5).

1270 There are two primary downsides with WAFs:

- 1271 1. WAFs are relatively expensive to run in terms of both latency and compute. They need
1272 to parse every request, perform a variety of scans to identify attack signatures (the
1273 number of scans depends on the policy configured), and either block or forward the
1274 request. While this overlaps heavily with the functionality of an API gateway, a WAF is
1275 typically deployed and operated by a separate team in isolation from the API gateway,
1276 often as part of the load balancer. This is convenient because the load balancer is the
1277 first place where requests are decrypted in the infrastructure. A secondary consequence
1278 is that WAF policies are typically only enforced at the perimeter.
- 1279 2. WAFs are fundamentally reactive. They operate based on matching requests to known
1280 attack signatures. As a result, they are largely ineffective at mitigating novel attacks, and
1281 attackers can leverage a variety of obfuscation techniques to hide known attacks behind
1282 novel signatures. Care must be taken to ensure that the WAF is running with the latest
1283 attack signature configurations, and custom rules must often be written for the
1284 organization.

1285 In line with a zero-trust posture, WAF policies should be enforced as close to the application as
1286 possible. This helps mitigate a variety of mechanisms that attackers might use to pivot within or
1287 otherwise compromise an infrastructure. As a practical matter, it can be cost-prohibitive to run
1288 a full suite of WAF mitigations on every internal and external request. This cost can be
1289 mitigated in two ways, which can be combined:

- 1290 1. Incorporate the WAF as part of the overall API-serving infrastructure and deploy the
1291 WAF itself in a “hybrid” model (i.e., keep a centralized WAF at the load balancer with a
1292 full suite of policies to protect against untrusted traffic). Then enforce a minimum set of
1293 app-specific WAF policies near each of the applications (e.g., in the distributed
1294 gateway). This minimizes policies run on east-west (i.e., more trusted) traffic while still
1295 sanitizing less trusted external traffic and tends to result in a good compromise of risk
1296 versus cost.
- 1297 2. Deploy the WAF as part of the API gateway implementation itself, which can avoid
1298 parsing the request multiple times (i.e., reduce the latency and compute costs of WAF
1299 policies), regardless of the API-serving implementation pattern chosen. If the API
1300 gateway is hybrid or distributed, then this technique can also be incorporated for
1301 further performance improvement.

1302 **4.4.2. Bot Detection**

1303 Bot detection typically involves evaluating risk signals, including origin (e.g., source IP, user
1304 credentials) and API usage patterns, over time to determine whether a seemingly legitimate
1305 user is likely to be a bot (i.e., an automated script acting maliciously). In response to flagging a
1306 high-risk user, bot detection systems will either block traffic or serve some kind of bot-
1307 defeating measure (e.g., CAPTCHA) before allowing the system to continue to be used. These
1308 tools primarily mitigate the risks of unrestricted resource access (see Sec. 2.4) (e.g., maliciously
1309 automating account creation in an email system) and leaking sensitive information (see Sec.
1310 2.5), especially data exfiltration by repeated calls.

1311 Bot detection is frequently deployed in user-facing applications. It can be more challenging with
1312 a purely machine-to-machine API because legitimate and malicious traffic patterns are even
1313 harder to differentiate. Many APIs are *intended* for use by scripts or non-user-facing
1314 applications, so human versus computer checks are irrelevant.

1315 **4.4.3. Distributed Denial of Service (DDoS) Mitigation**

1316 A distributed denial of service (DDoS) attack is a DoS that originates from many different
1317 locations or users. This makes it more challenging to mitigate than a traditional DoS attack,
1318 which can often be prevented by blocking a small set of users. While DoS attacks may be
1319 targeted and application-level, DDoS attacks are often network-oriented in nature and seek to
1320 saturate the server’s bandwidth or ability to establish new connections. When a determined
1321 attacker is able to build and execute an application-level DDoS attack, it is one of the most
1322 challenging attacks to mitigate. Because of the primarily network-oriented nature of DDoS

1323 attacks, most DDoS mitigation tools are deployed at the network edge as part of the load
1324 balancer or even before the load balancer as part of the CDN and DNS system (often called
1325 “Global Traffic Management”). Predictably, DDoS mitigation tools help mitigate unrestricted
1326 resource consumption (see Sec. 2.4).

1327 **4.4.4. API Endpoint Protection**

1328 “API protection” or “API endpoint protection” are nebulous terms for describing a set of
1329 capabilities around API inventory, authentication, rate limiting, and data analysis. The exact set
1330 of capabilities tends to vary with the implementation. For example, sophisticated
1331 implementations can scan requests and responses to tag suspect data on the wire (e.g., to help
1332 tag sensitive data and pinpoint possible leaks or exfiltration).

1333 API protection products are typically packaged with the API gateway. API gateway vendors
1334 primarily deliver their products in the centralized API gateway pattern, so these controls are
1335 often only enforced at the perimeter. Like a WAF, the policies they enforce are typically cross-
1336 cutting and do not require an in-depth understanding at the API payload level. As such, the two
1337 products are often marketed in a similar niche.

1338 The exact set of risks mitigated by these tools depends on the feature set, but they typically
1339 attempt to mitigate lack of API visibility (see Sec. 2.1), broken authentication (see Sec. 2.3),
1340 some aspects of unrestricted compute consumption (see Sec. 2.4), and leaking sensitive
1341 information (see Sec. 2.5).

1342 There is value in any tool that helps organizations inventory and manage their APIs and traffic.
1343 However, the enforcement of any policies should be as close to the individual service instance
1344 as possible in order to achieve robust API security assurance. In the use case of data
1345 classification, these tools can be especially useful for building an initial inventory. However, as
1346 API definitions are rolled out across the organization, data tagging should be implemented as
1347 part of the API schema, and the data flow policy should be enforced via explicit policy (e.g., with
1348 an authorization system). The runtime discovery of data flow is primarily important in
1349 protecting against exfiltration.

1350 **4.4.5. Web Application and API Protection (WAAP)**

1351 Gartner coined the term “web application and API protection” (WAAP) [27] to describe the
1352 trend of packaging the technologies listed here (i.e., WAF, bot detection, DDoS mitigation, and
1353 API protection) into a single product. Whether these capabilities are implemented with a single
1354 product or a range of technologies, the key is understanding what risk each capability is helping
1355 to mitigate and evaluating how it fits into the organization’s existing security posture.

1356 **4.5. Summary of Implementation Patterns**

1357 Combining the three patterns in API gateway architecture with the companion technologies
1358 discussed Sec. 4.4 provides a comprehensive set of enterprise security solutions for API

1359 protection. The key point in each pattern is identifying where to enforce each policy. These
1360 decisions result in trade-offs in runtime, architecture, and operations for the application teams
1361 utilizing the API-serving infrastructure. Many organizations use a mix of all three patterns
1362 deployed in production precisely because of those trade-offs. All three patterns can be used to
1363 successfully implement all of the controls outlined in this document. That said, the distributed
1364 gateway pattern and its companion technologies best align with the principals of zero trust and
1365 are strongly recommended for organizations that want to adopt a security-forward approach.

1366

1367 **5. Conclusions and Summary**

1368 No business-critical enterprise application can be integrated into the digital infrastructure of an
1369 enterprise without the use of APIs. With the highly distributed nature of applications (both
1370 physically and logically), APIs must be operated under zero trust principles irrespective of
1371 whether they are exposed to the outside world or meant to be consumed by other applications
1372 within the enterprise infrastructure. Like all software, APIs go through an iterative life cycle
1373 whose phases (i.e., Develop, Build, Deploy, Operate) can be broadly classified into pre-runtime
1374 and runtime stages.

1375 The sheer proliferation of API deployments, the heterogeneous infrastructures under which
1376 they operate, and the access to valuable corporate data that they enable make them targets for
1377 exploitation. A detailed analysis of their vulnerabilities and the potential attack vectors that can
1378 exploit them is a prerequisite for identifying the appropriate set of protection measures or
1379 controls to ensure API security. This document analyzes a spectrum of risk factors that give rise
1380 to vulnerabilities, such as the lack of a formal schema, improper inventorying, the lack of robust
1381 authentication and authorization support, improper monitoring of resource consumption, and
1382 the least leakage of sensitive information.

1383 The recommended controls in this document are classified into pre-runtime and runtime
1384 protections. They are further subdivided into basic and advanced protections to enable
1385 enterprises to use a risk-based and incremental approach to securing their digital assets. Pre-
1386 runtime protections focus on API specification parameters (i.e., syntactic and semantic aspects),
1387 while runtime API protections focus on protections during API request and response operations
1388 (e.g., encrypted communication channels, proper authentication and authorization).

1389 This document presents a landscape of real-world and state-of-practice implementation
1390 options to configure and enforce the recommended controls by describing the advantages and
1391 disadvantages of each type of protection deployment or pattern. This will enable practitioners
1392 to make an informed decision to realize a robust and cost-effective API security infrastructure
1393 for their enterprises.

1394

1395 **References**

- 1396 [1] U.S. Department of Defense Chief Information Officer (2024) DoD Enterprise DevSecOps
1397 Fundamentals. Version 2.5, October 2024. Available at
1398 <https://dodcio.defense.gov/Portals/0/Documents/Library/DoD%20Enterprise%20DevSecOps%20Fundamentals%20v2.5.pdf>
1399
- 1400 [2] Chandramouli R, Butcher Z (2020) Building Secure Microservices-based Applications
1401 Using Service-Mesh Architecture. (National Institute of Standards and Technology,
1402 Gaithersburg, MD), NIST Special Publication (SP) NIST SP 800-204A.
1403 <https://doi.org/10.6028/NIST.SP.800-204A>
- 1404 [3] Chandramouli R, Butcher Z, Aradhna C (2021) Attribute-based Access Control for
1405 Microservices-based Applications using a Service Mesh. (National Institute of Standards
1406 and Technology, Gaithersburg, MD), NIST Special Publication (SP) NIST SP 800-204B.
1407 <https://doi.org/10.6028/NIST.SP.800-204B>
- 1408 [4] Chandramouli R (2022) Implementation of DevSecOps for a Microservices-based
1409 Application with Service Mesh. (National Institute of Standards and Technology,
1410 Gaithersburg, MD), NIST Special Publication (SP) NIST SP 800-204C.
1411 <https://doi.org/10.6028/NIST.SP.800-204C>
- 1412 [5] Chandramouli R, Kautz F, Torres-Arias S (2024) Strategies for the Integration of Software
1413 Supply Chain Security in DevSecOps CI/CD Pipelines. (National Institute of Standards and
1414 Technology, Gaithersburg, MD), NIST Special Publication (SP) NIST SP 800-204D.
1415 <https://doi.org/10.6028/NIST.SP.800-204D>
- 1416 [6] Chandramouli R, Butcher Z (2023) A Zero Trust Architecture Model for Access Control in
1417 Cloud-Native Applications in Multi-Cloud Environments. (National Institute of Standards
1418 and Technology, Gaithersburg, MD), NIST Special Publication (SP) NIST SP 800-207A.
1419 <https://doi.org/10.6028/NIST.SP.800-207A>
- 1420 [7] OWASP (2023) OWASP Top 10 API Security Risks. Available at <https://owasp.org/API-Security/editions/2023/en/0x11-t10/>
1421
- 1422 [8] OWASP (2023) API2:2023 Broken Authentication. Available at <https://owasp.org/API-Security/editions/2023/en/0xa2-broken-authentication/>
1423
- 1424 [9] Wikipedia (2025) Fatigue Attack. Available at https://en.wikipedia.org/wiki/Multi-factor_authentication_fatigue_attack
1425
- 1426 [10] Wikipedia (2024) Billion laughs attack. Available at
1427 <https://en.wikipedia.org/wiki/BillionLaughsAttack>
- 1428 [11] Wikipedia (2025) Zip bomb. Available at https://en.wikipedia.org/wiki/Zip_bomb
- 1429 [12] National Institute of Standards and Technology (2024) The NIST Cybersecurity
1430 Framework (CSF) 2.0. (National Institute of Standards and Technology, Gaithersburg,
1431 MD), NIST Cybersecurity White Paper (CSWP) NIST CSWP 29.
1432 <https://doi.org/10.6028/NIST.CSWP.29>
- 1433 [13] Wikipedia (2025) Principle of least astonishment. Available at
1434 https://en.wikipedia.org/wiki/Principle_of_least_astonishment
- 1435 [14] F# for fun and profit (2013) Designing with types: Making illegal types unrepresentable.
1436 Available at <https://fsharpforfunandprofit.com/posts/designing-with-types-making-illegal-states-unrepresentable/>
1437

- 1438 [15] Joint Task Force (2020) Security and Privacy Controls for Information Systems and
1439 Organizations. (National Institute of Standards and Technology, Gaithersburg, MD), NIST
1440 Special Publication (SP) NIST SP 800-53r5. Includes updates as of December 10, 2020.
1441 <https://doi.org/10.6028/NIST.SP.800-53r5>
- 1442 [16] Barker E (2020) Recommendation for Key Management: Part 1 – General. (National
1443 Institute of Standards and Technology, Gaithersburg, MD), NIST Special Publication (SP)
1444 NIST SP 800-57pt1r5. <https://doi.org/10.6028/NIST.SP.800-57pt1r5>
- 1445 [17] National Institute of Standards and Technology (2019) Security Requirements for
1446 Cryptographic Modules. (Department of Commerce, Washington, D.C.), Federal
1447 Information Processing Standards Publications (FIPS) NIST FIPS 140-3.
1448 <https://doi.org/10.6028/NIST.FIPS.140-3>
- 1449 [18] rfc7519 (2015) JSON Web Token (JWT). Available at
1450 <https://datatracker.ietf.org/doc/html/rfc7519>
- 1451 [19] rfc5280 (2008) Internet X.509 Public Key Infrastructure Certificate and Certificate
1452 Revocation List (CRL) Profile. Available at <https://datatracker.ietf.org/doc/html/rfc5280>
- 1453 [20] Wikipedia (2023) Common Log Format. Available at
1454 https://en.wikipedia.org/wiki/Common_Log_Format
- 1455 [21] Github (2016) grpc. Available at
1456 <https://github.com/grpc/grpc/blob/master/src/proto/grpc/status/status.proto>
- 1457 [22] Wikipedia (2025) Prompt injection. Available at
1458 https://en.wikipedia.org/wiki/Prompt_injection
- 1459 [23] Wikipedia (2024) Cascading failure. Available at
1460 https://en.wikipedia.org/wiki/Cascading_failure
- 1461 [24] Infoq.com (2020) How to Avoid Cascading Failures in Distributed Systems. Available at
1462 <https://www.infoq.com/articles/anatomy-cascading-failure/>
- 1463 [25] Coreruleset.org (2025) OWASP CRS PROJECT. Available at <https://coreruleset.org>
- 1464 [26] Wikipedia (2025) *Confused deputy problem*. Available at
1465 https://en.wikipedia.org/wiki/Confused_deputy_problem
- 1466 [27] Gartner.com (2025) *Cloud Web Application and API Protection*. Available at
1467 <https://www.gartner.com/reviews/market/cloud-web-application-and-api-protection>
- 1468
1469
1470
1471

1472 **Appendix A. API Classification Taxonomy**

1473 **A.1. API Classification Based on Degree of Exposure**

1474 Since APIs are interfaces that are exposed to relevant stakeholders, they should be classified
1475 based on their degree of exposure. Three kinds of APIs are prevalent:

- 1476 1. Open/Public APIs are exposed to a broader and wider audience (i.e., customers) and
1477 used with external partnerships or services. These are also called “facade APIs,” as they
1478 may provide limited access to certain functionalities.
- 1479 2. Private APIs are used to link various systems within an enterprise and are closely
1480 guarded, such as a contract between microservices that are internal to an organization.
1481 Variations of private APIs are:
 - 1482 a. Internal APIs (service APIs): Used by enterprises to streamline their internal
1483 workflows and create flexible systems that can adapt to changing business needs
 - 1484 b. Composite APIs: Allow multiple data and service calls to be combined to realize
1485 efficiency in system design [10]
- 1486 3. Partner APIs are used in the context of collaborative ventures between enterprises, as
1487 both rely on shared services or data to deliver value to their end users. In terms of
1488 exposure, they represent a middle ground between public and private APIs since access
1489 is restricted based on collaborative agreements [10].

1490 **A.2. API Classification Based on Communication Patterns**

1491 There are two fundamental API communication patterns that govern how information flows
1492 between the components involved in system interactions:

- 1493 1. Request-response APIs: A communication pattern in which a client sends a request to a
1494 server and awaits a corresponding response. It operates synchronously with stateless,
1495 independent requests. This pattern is widely employed in diverse API architectures
1496 (such as RESTful APIs, GraphQL, and various web services [11]) and is appropriate for
1497 immediate data retrieval or any instant action (e.g., downloading a user’s profile in a
1498 social media app). For example, requests are made with verbs that are appropriate for
1499 the API architecture (e.g., HTTP method GET in RESTful architecture, a structured query
1500 specifying the exact data needed in response in GraphQL architecture).
- 1501 2. Event-driven APIs: A better choice for receiving real-time updates (e.g., user’s activities
1502 in the same app).

1503 **A.3. API Classification Based on Architectural Style or Pattern (API Types)**

1504 **Table 1. API classification based on Architectural Patterns**

API Name	Network Protocol	Data Format
REST	HTTP/1.1	Text-based JSON
gRPC	HTTP/2	Binary — Protocol Buffers (Protobuf)
GraphQL	HTTP – POST only	JSON
WebSocket	WebSocket	JSON

1505

1506 **Appendix B. DevSecOps Phase and Associated Class of API Controls**

1507 The detailed controls in Sec.3 fit into several broad classes, and Table 2 shows their associations
1508 with the DevSecOps phases (see Sec. 1).

1509 **Table 2. DevSecOps phase and associated class of API controls**

DevSecOps Phase	Class of API Controls
Coding	Well-defined API schema definition that calls to routines for annotating schema definitions
Build	Generate routines that validate on-field values in the request and response payloads of API calls and responses, respectively
Test	Ensure that validation routines perform as intended in various runs of API requests and responses
Deployment	Ensure that the deployment package contains all of the runtime policy enforcement routines, API schema definitions, and APIs and is signed off by the right authorities
Observe and Monitor	Ensure that certain security incidents (e.g., data leakage) do not occur due to (a) inherent flaws in API design, (b) the lack of input data validation, or (c) engineered attacks realized through a sequence of requests that each pass all validation tests

1510