

E6156 – Topics in SW Engineering: Cloud Computing

Lecture 5: APIGW, Lambda, SSO, Events



Contents

Agenda

- Level set – a microservice walkthrough
 - What do you actually get when receiving a REST call?
 - Implementing REST paths
 - Dependency injection
 - CORS
- REST/HATEOAS:
 - Concepts
 - Deep dive on “stateless”
- Middleware
- Event Driven Processing
- Lambda functions
- Webhooks

Code to Show Today

Code to Show

- /demo path and rest_utils (UserProfileNew)
- oauth-python project, including setup on GCP.
- UserProfileNew
 - Pagination
 - Fields
 - Links
 - Middleware
- Lambda function.
- WebHook

Microservice Walkthrough

REST Method

- When your handler receives a REST operation, there are the following elements. Not all methods have all elements:
 - Path
 - Path parameters
 - Query string parameters
 - Headers
 - Body
 - Flask provides access through the
 - Request object.
 - Response object.
 - Walkthrough
- The REST handler's response has the following elements:
 - Status code
 - Body
 - Headers
 - The route handler maps between the REST representation of the elements and the language/runtime representation.

REST Paths

<https://restful-api-design.readthedocs.io/en/latest/>

- There are many patterns. I typically follow something like:

URL Structure

Each collection and resource in the API has its own URL. URLs should never be constructed by an API client. Instead, the client should only follow links that are generated by the API itself.

The recommended convention for URLs is to use alternate collection / resource path segments, relative to the API entry point. This is best described by example. The table below uses the “:name” URL variable style from Rail’s “Routes” implementation.

URL	Description
/api	The API entry point
/api/:coll	A top-level collection named “coll”
/api/:coll/:id	The resource “id” inside collection “coll”
/api/:coll/:id/:subcoll	Sub-collection “subcoll” under resource “id”
/api/:coll/:id/:subcoll/:subid	The resource “subid” inside “subcoll”

Even though sub-collections may be arbitrarily nested, in my experience, you want to keep the depth limited to 2, if possible. Longer URLs are more difficult to work with when using simple command-line tools like [curl](#).

- We will not be too rigorous about status codes, but some common examples are:
 - 200: OK (success)
 - 201: CREATED (for POST)
 - 404: NOT FOUND
 - 500: INTERNAL SERVICE ERROR
 - 418: I’M A TEAPOT

Method	Scope	Semantics
GET	collection	Retrieve all resources in a collection
GET	resource	Retrieve a single resource
HEAD	collection	Retrieve all resources in a collection (header only)
HEAD	resource	Retrieve a single resource (header only)
POST	collection	Create a new resource in a collection
PUT	resource	Update a resource
PATCH	resource	Update a resource
DELETE	resource	Delete a resource
OPTIONS	any	Return available HTTP methods and other options

Success Response Codes

Operation	HTTP Request	HTTP Response Codes Supported
READ	GET	200 - OK with message body 204 - OK no message body 206 - OK with partial message body
CREATE	POST	201 - Resource created (Operation Complete) 202 - Resource accepted (Operation Pending)
UPDATE	PUT	202 - Accepted (Operation Pending) 204 - Success (Operation Complete)
DELETE	DELETE	202 - Accepted (Operation Pending) 204 - Success (Operation Complete)

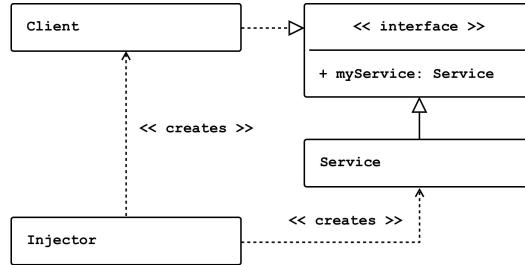
202 means
Your request went asynch.
The HTTP header Link
is where to poll for rsp.

Examples of Link Headers in HTTP response:

Link: <<http://api/jobs/j1>>;rel=monitor;title="update profile"
Link: <<http://api/reports/r1>>;rel=summary;title="access report"

Dependency Injection

- Concepts (https://en.wikipedia.org/wiki/Dependency_injection):
 - In software engineering, dependency injection is a technique in which an object receives other objects that it depends on, called dependencies.
 - The intent behind dependency injection is to achieve separation of concerns of construction and use of objects. This can increase readability and code reuse.



```
class UserResource(BaseApplicationResource):  
    def __init__(self, config_info):  
        super().__init__(config_info)
```

- There is no single approach. There are several frameworks, tools, etc.
- Like many other things, it is easy to get carried away and become dogmatic.
- I follow a single approach:
 - A Context class converts environment and other configuration and provides to application.
 - The top-level application injects a config_info object into services.

- “For security reasons, browsers restrict cross-origin HTTP requests initiated from scripts. For example, XMLHttpRequest and the [Fetch API](#) follow the [same-origin policy](#). This means that a web application using those APIs can only request resources from the same origin the application was loaded from unless the response from other origins includes the right CORS headers.”

Untitled Request

OPTIONS http://54.242.71.165:5000/users/11

Send Save Cookies Code

Params Authorization Headers (7) Body Pre-request Script Tests Settings

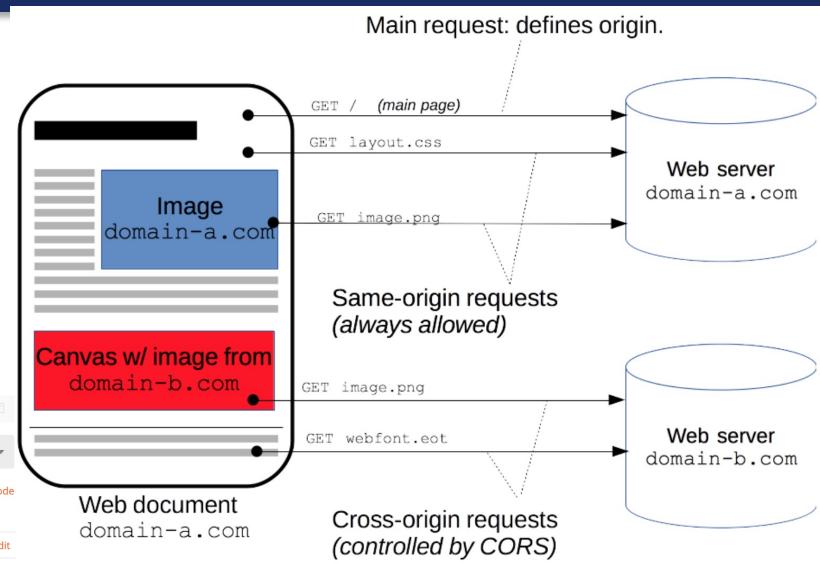
Query Params

KEY	VALUE	DESCRIPTION	Bulk Edit
Key	Value	Description	

Body Cookies Headers (6) Test Results

Status: 200 OK Time: 78 ms Size: 225 B Save Response

KEY	VALUE
Content-Type	text/html; charset=utf-8
Allow	PUT, OPTIONS, GET, HEAD, DELETE
Access-Control-Allow-Origin	*
Content-Length	0
Server	Werkzeug/2.0.1 Python/3.7.10
Date	Wed, 06 Oct 2021 16:39:32 GMT



```
app = Flask(__name__)
CORS(app)
```

REST Principles

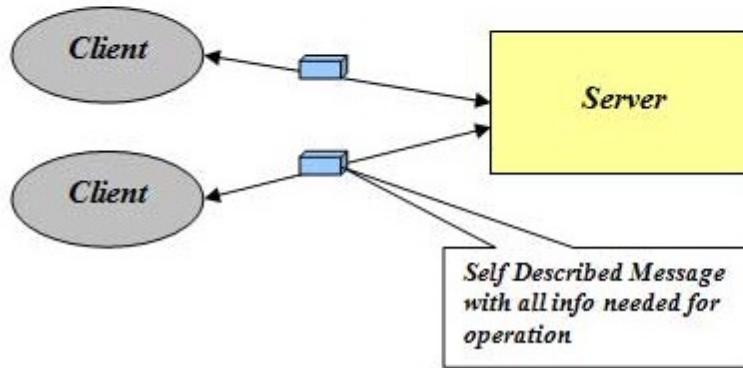
HATEOAS

REST Principles

- REST Principles:
 - Uniform interface
 - Client–server
 - Stateless
 - Cacheable
 - Layered system
 - Code on demand (optional)
 - Many of these are straightforward. We will explore some of the concepts in various lectures.
 - Today, we cover “statelessness.”
- HATEOAS: Hypertext as the Engine of Application State --*
The principle is that a client interacts with a network application entirely through [hypermedia](#) provided dynamically by application servers. A REST client needs no prior knowledge about how to interact with any particular application or server beyond a generic understanding of hypermedia.

Stateless

The notion of statelessness is defined from the perspective of the server. The constraint says that the server should not remember the state of the application. As a consequence, the client should send all information necessary for execution along with each request, because the server cannot reuse information from previous requests as it didn't memorize them. All info needed is in message.



Statelessness

2. Application State vs Resource State

It is important to understand the difference between the application state and the resource state. Both are completely different things.

Application state is server-side data that servers store to identify incoming client requests, their previous interaction details, and current context information.

Resource state is the current state of a resource on a server at any point in time – and it has nothing to do with the interaction between client and server. It is what we get as a response from the server as the API response. We refer to it as resource representation.

REST statelessness means being free from the application state.

- This concept can be confusing with experience with other approach to session state.
- Some application frameworks did/do keep session state on the server.
- Application servers and frameworks have various support for “session state,” e.g. <https://flask-session.readthedocs.io/en/latest/>

3. Advantages of Stateless APIs

There are some very noticeable advantages of having REST APIs **stateless**.

1. Statelessness helps in **scaling the APIs to millions of concurrent users** by deploying it to multiple servers. Any server can handle any request because there is no session related dependency.
2. Being stateless makes REST APIs **less complex** – by removing all server-side state synchronization logic.
3. A stateless API is also **easy to cache** as well. Specific softwares can decide whether or not to cache the result of an HTTP request just by looking at that one request. There's no nagging uncertainty that state from a previous request might affect the cacheability of this one. It **improves the performance** of applications.
4. The server never loses track of “where” each client is in the application because the client sends all necessary information with each request.

Pagination and Field Selection

Pagination

Most endpoints that returns a list of entities will need to have some sort of pagination.

Without pagination, a simple search could return millions or even billions of hits causing extraneous network traffic.

Paging requires an implied ordering. By default this may be the item's unique identifier, but can be other ordered fields such as a created date.

Offset Pagination

This is the simplest form of paging. Limit/Offset became popular with apps using SQL databases which already have LIMIT and OFFSET as part of the SQL SELECT Syntax. Very little business logic is required to implement Limit/Offset paging.

Limit/Offset Paging would look like `GET /items?limit=20&offset=100`. This query would return the 20 rows starting with the 100th row.

Example

(Assume the query is ordered by created date descending)

1. Client makes request for most recent items: `GET /items?limit=20`
2. On scroll/next page, client makes second request `GET /items?limit=20&offset=20`
3. On scroll/next page, client makes third request `GET /items?limit=20&offset=40`

<https://www.moesif.com/blog/technical/api-design/REST-API-Design-Filtering-Sorting-and-Pagination/>

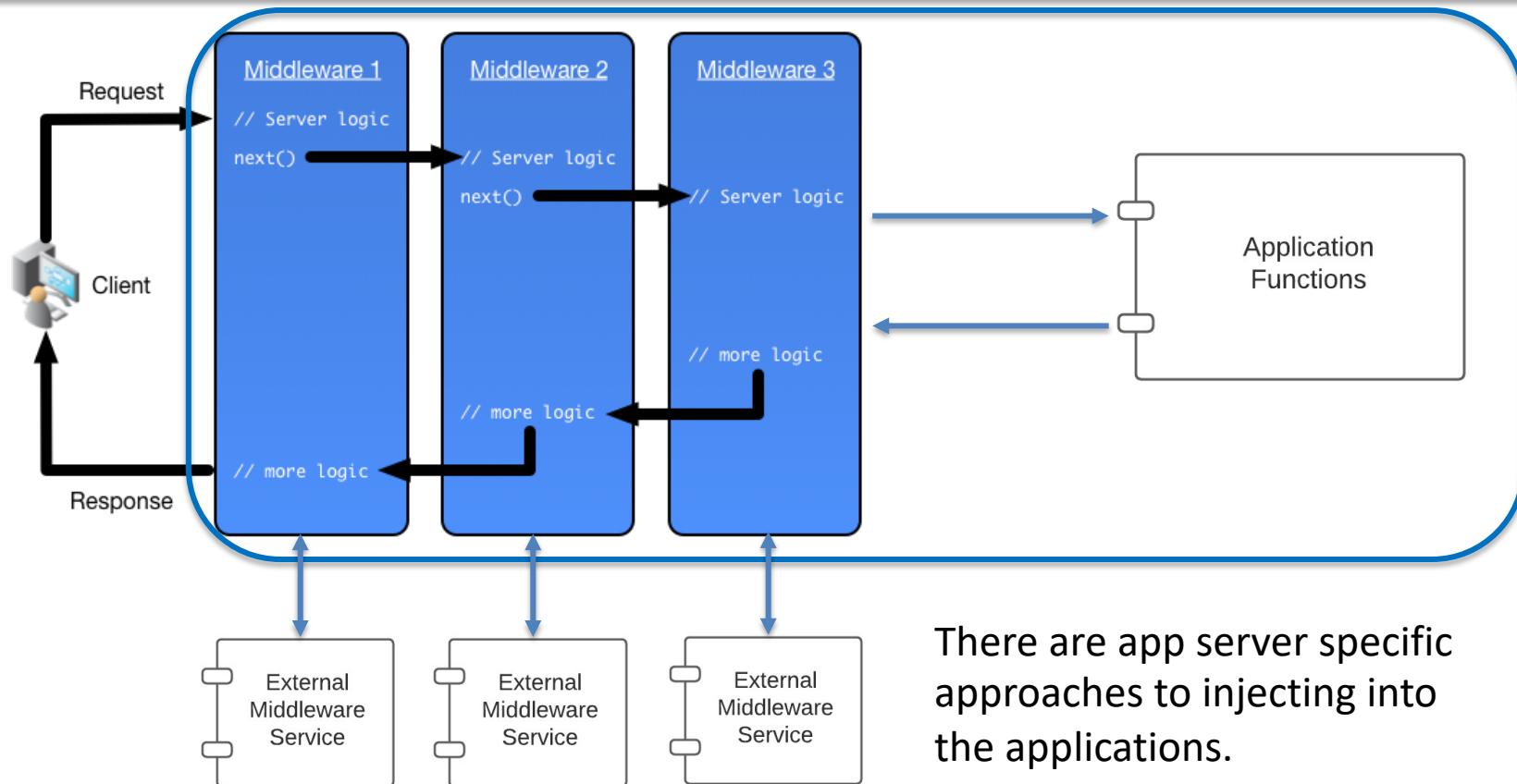
Middleware

What is middleware?

Middleware is software that lies between an operating system and the applications running on it. Essentially functioning as hidden translation layer, middleware enables communication and data management for distributed applications. It's sometimes called plumbing, as it connects two applications together so data and databases can be easily passed between the "pipe." Using middleware allows users to perform such requests as submitting forms on a web browser, or allowing the web server to return dynamic web pages based on a user's profile.

Common middleware examples include database middleware, application server middleware, message-oriented middleware, web middleware, and transaction-processing monitors. Each program typically provides messaging services so that different applications can communicate using messaging frameworks like simple object access protocol (SOAP), web services, representational state transfer (REST), and JavaScript object notation (JSON). While all middleware performs communication functions, the type a company chooses to use will depend on what service is being used and what type of information needs to be communicated. This can include security authentication, transaction management, message queues, applications servers, web servers, and directories. Middleware can also be used for distributed processing with actions occurring in real time rather than sending data back and forth.

Middleware – Conceptual Model



Middleware and Flask

```
# These methods get called before/after. You can check security information for all requests.
```

```
@application.before_request
```

```
def before_decorator():
```

```
    a_ok = sec.check_authentication(request)
```

```
    print("a_ok = ", a_ok)
```

```
    if a_ok[0] != 200:
```

```
        handle_error(a_ok[0], a_ok[1], a_ok[2])
```

- White list of URLs
- If not on white list:
 - Verify token
 - Verify authorization

```
@application.after_request
```

```
def after_decorator(rsp):
```

```
    print("... In after decorator ...")
```

```
    notify.notify(request, rsp)
```

```
    return rsp
```

- List of filters of the form:
 - URL
 - HTTP Method
 - SNS topic to use
- Send an event to SNS if matches filter

```
@application.route("/api/users", methods=["GET", "POST"])
```

```
def users_get(): ...
```

Show some demo code.

So, How Will We Use This?

- Add before_request and after_request to microservices.
(Note: There are other approaches).
- Implement two middleware interceptors:
 - Security
 - Configured with a JSON file.
 - File lists which {path, method} require being logged on.
 - While IDs can perform which methods on which paths.
 - Automates driving login redirect (401 – Unauthorized)
 - Checks authorization (403– Forbidden)
 - Notification:
 - Configured with JSON file.
 - File lists which operations emit a before and/or after event to an SNS topic.
- We will use notification to drive some other things like webhooks.

Event Drive Processing

Publish/Subscribe

Pub/Sub – Event Driven Processing

"In [software architecture](#), **publish–subscribe** is a [messaging pattern](#) where senders of [messages](#), called publishers, do not program the messages to be sent directly to specific receivers, called subscribers, but instead categorize published messages into classes without knowledge of which subscribers, if any, there may be. Similarly, subscribers express interest in one or more classes and only receive messages that are of interest, without knowledge of which publishers, if any, there are.

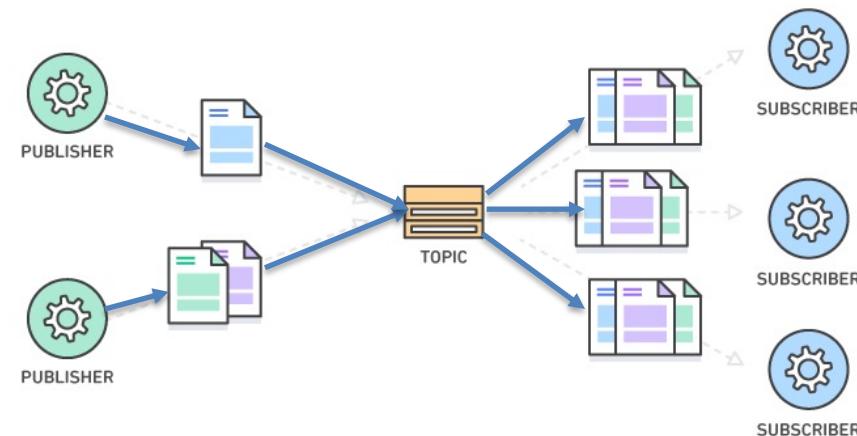
Publish–subscribe is a sibling of the [message queue](#) paradigm, and is typically one part of a larger [message-oriented middleware](#) system. Most messaging systems support both the pub/sub and message queue models in their [API](#), e.g. [Java Message Service](#) (JMS).

This pattern provides greater network [scalability](#) and a more dynamic [network topology](#), with a resulting decreased flexibility to modify the publisher and the structure of the published data."

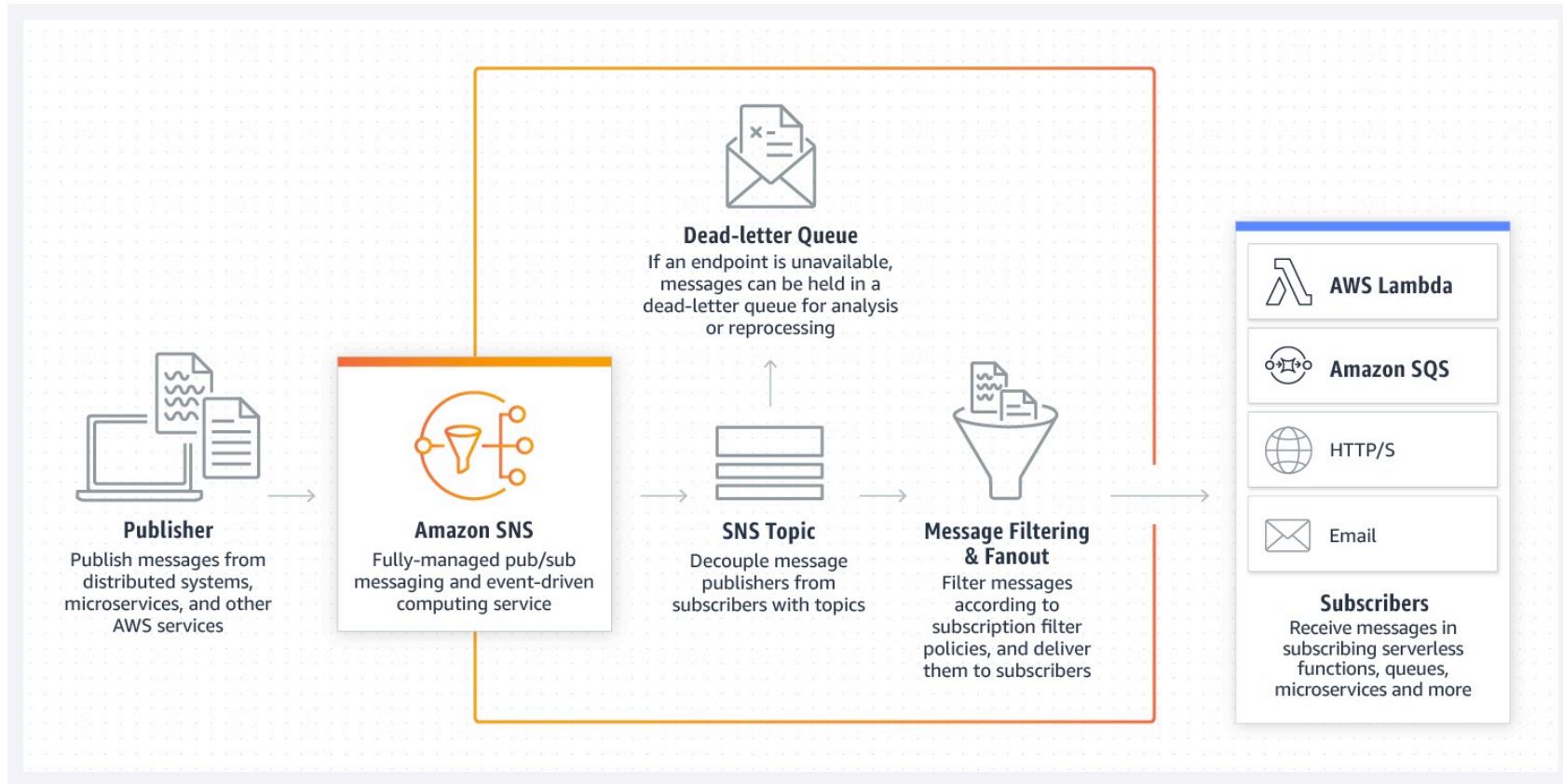
https://en.wikipedia.org/wiki/Publish%E2%80%93subscribe_pattern

It also supports loose coupling for application evolution and agility.

I can add and evolve services without modifying the services.

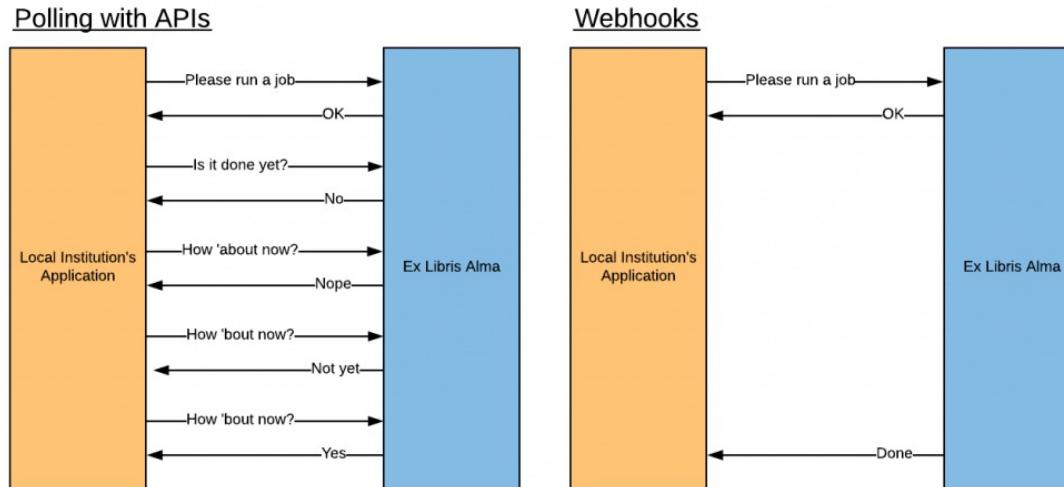


Pub/Sub – Simple Notification Service



Webhooks

- “A webhook in web development is a method of augmenting or altering the behavior of a web page or web application with custom callbacks. These callbacks may be maintained, modified, and managed by third-party users and developers who may not necessarily be affiliated with the originating website or application. The term “webhook” was coined by Jeff Lindsay in 2007 from the computer programming term hook.”



Function-as-a-Service

Serverless and Function-as-a-Service

- **Serverless computing** is a [cloud computing execution model](#) in which the cloud provider runs the [server](#), and dynamically manages the allocation of machine resources. Pricing is based on the actual amount of resources consumed by an application, rather than on pre-purchased units of capacity.^[1] It can be a form of [utility computing](#).

Serverless computing can simplify the process of [deploying code](#) into production. Scaling, capacity planning and maintenance operations may be hidden from the developer or operator. Serverless code can be used in conjunction with code deployed in traditional styles, such as [microservices](#). Alternatively, applications can be written to be purely serverless and use no provisioned servers at all.^[2] This should not be confused with computing or networking models that do not require an actual server to function, such as [peer-to-peer](#) (P2P)."
- **Function as a service (FaaS)** is a category of [cloud computing services](#) that provides a [platform](#) allowing customers to develop, run, and manage application functionalities without the complexity of building and maintaining the infrastructure typically associated with developing and launching an app.^[1] Building an application following this model is one way of achieving a "[serverless](#)" architecture, and is typically used when building [microservices](#) applications."
- That was baffling:
 - IaaS, CaaS – You control the SW stack and the cloud provides (virtual) HW.
 - PaaS – The cloud hides the lower layer SW and provides an application container (e.g. Flask) with “Your code goes here.” You are aware of container.
 - FaaS – The cloud provides the container, and you implement functions (corresponding to routes in REST).

Serverless and Function-as-a-Service

Private Cloud	IaaS	PaaS	FaaS	SaaS
	Infrastructure as a Service	Platform as a Service	Function as a Service	Software as a Service
Function	Function	Function	Function	Function
Application	Application	Application	Application	Application
Runtime	Runtime	Runtime	Runtime	Runtime
Operating System	Operating System	Operating System	Operating System	Operating System
Virtualization	Virtualization	Virtualization	Virtualization	Virtualization
Server	Server	Server	Server	Server
Storage	Storage	Storage	Storage	Storage
Networking	Networking	Networking	Networking	Networking

<https://medium.com/@tanmayct/serverless-architecture-function-as-a-service-19e127b8c990>

Managed by the customer



Managed by the provider



What is Serverless Good/Not Good For ... ?

Serverless is **good** for
*short-running
stateless
event-driven*



- Microservices
- Mobile Backends
- Bots, ML Inferencing
- IoT
- Modest Stream Processing
- Service integration

Serverless is **not good** for
*long-running
stateful
number crunching*

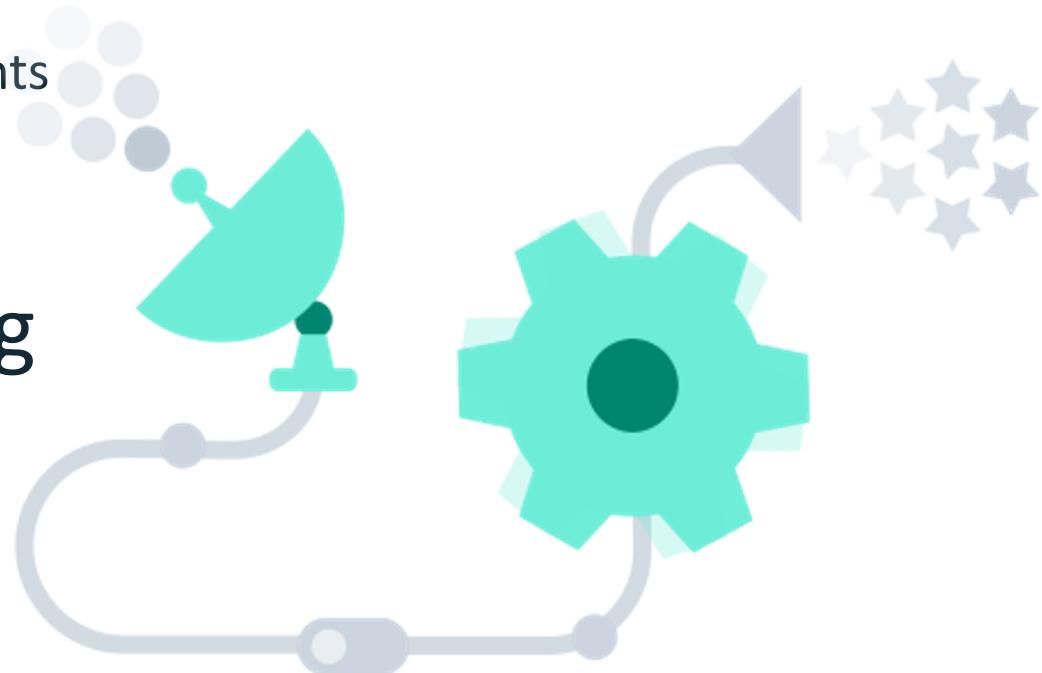


- Databases
- Deep Learning Training
- Heavy-Duty Stream Analytics
- Numerical Simulation
- Video Streaming

What triggers code execution?

Runs code **in response** to events

Event-programming
model



(Some) Current Platforms for Serverless



IBM Cloud
Functions



Azure
Functions



Red-Hat



Google
Functions



Kubernetes



API Gateway

API Management

API management is the process of creating and publishing web [application programming interfaces](#) (APIs), enforcing their usage policies, controlling access, nurturing the subscriber community, collecting and analyzing usage statistics, and reporting on performance. API Management components provide mechanisms and tools to support developer and subscriber community.^[1]

Components [edit]

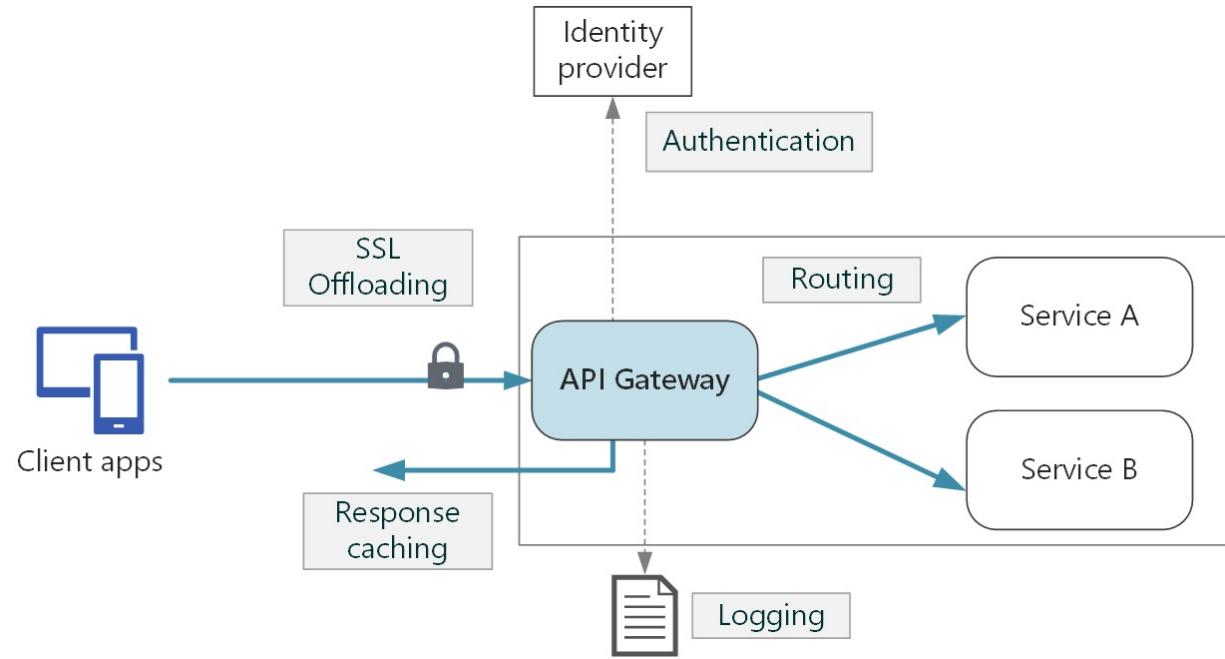
While solutions vary, components that provide the following functionality are typically found in API management products:

- **Gateway:** a server that acts as an API front-end, receives API requests, enforces throttling and security policies, passes requests to the back-end service and then passes the response back to the requester.^[2] A gateway often includes a transformation engine to [orchestrate](#) and modify the requests and responses on the fly. A gateway can also provide functionality such as collecting analytics data and providing caching. The gateway can provide functionality to support authentication, authorization, security, audit and regulatory compliance.^[3]
It can be implemented by a [Reverse proxy](#).
- **Publishing tools:** a collection of tools that API providers use to define APIs, for instance using the [OpenAPI](#) or [RAML](#) specifications, generate API documentation, govern API usage through access and usage policies for APIs, test and debug the execution of API, including security testing and automated generation of tests and test suites, deploy APIs into production, staging, and quality assurance environments, and coordinate the overall API lifecycle.
- **Developer portal/API store:** community site, typically branded by an API provider, that can encapsulate for API users in a single convenient source information and functionality including documentation, tutorials, sample code, software development kits, an interactive API console and sandbox to trial APIs, the ability to subscribe to the APIs and manage subscription keys such as [OAuth2](#) Client ID and Client Secret, and obtain support from the API provider and user and community.
- **Reporting and analytics:** functionality to monitor API usage and load (overall hits, completed transactions, number of data objects returned, amount of compute time and other internal resources consumed, volume of data transferred). This can include real-time monitoring of the API with alerts being raised directly or via a higher-level [network management system](#), for instance, if the load on an API has become too great, as well as functionality to analyze historical data, such as transaction logs, to detect usage trends. Functionality can also be provided to create synthetic transactions that can be used to test the performance and behavior of API endpoints. The information gathered by the reporting and analytics functionality can be used by the API provider to optimize the API offering within an organization's overall [continuous improvement process](#) and for defining software [Service-Level Agreements](#) for APIs.
- **Monetization:** functionality to support charging for access to commercial APIs. This functionality can include support for setting up pricing rules, based on usage, load and functionality, issuing invoices and collecting payments including multiple types of credit card payments.

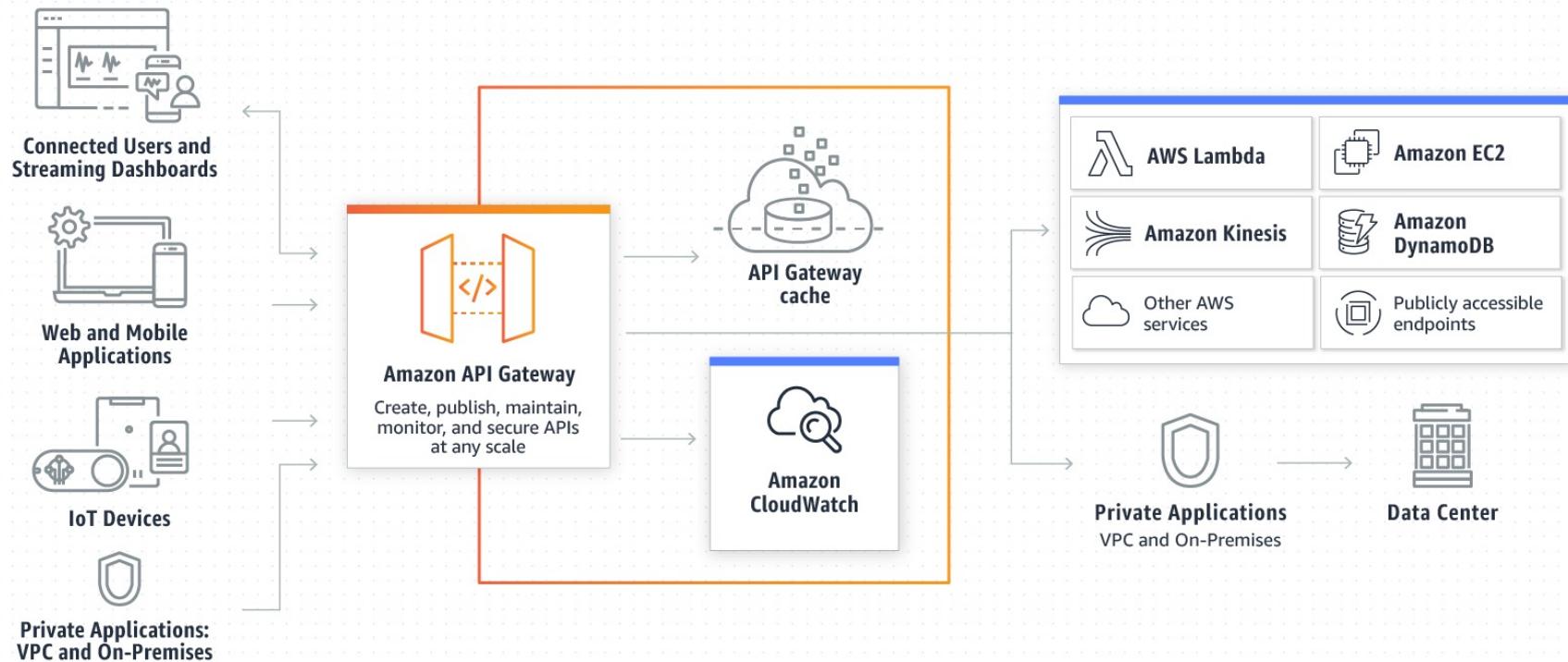
API Gateway (<https://docs.microsoft.com/en-us/azure/architecture/microservices/design/gateway>)

Some functions:

- SSL termination
- Authentication
- IP allow/block list
- Client rate limiting (throttling)
- Logging and monitoring
- Response caching
- Web application firewall
- GZIP compression
- Servicing static content



AWS API Gateway



Function-as-a-Service

Putting Some Pieces Together

- API Gateway:
 - We will encapsulate the 3 microservices behind the API Gateway.
 - Single, logical endpoint and locations for security, etc.
- Events/SNS:
 - Setup SNS topics.
 - State changes in resources, e.g. user, triggers event publication.
 - Lambda functions handle the events and integrate with things like webhooks.
- Middleware:
 - Add middleware to microservices to automate generating events.
 - Use middleware to enforce security and login, which we will later replace/complement with API GW.
- Implement