

Serverless Computing

Author: Richard Clauß

Docent: Prof. Dr. Martin Leischner

Date: 06.12.2022



Hochschule
Bonn-Rhein-Sieg

Agenda

1. Introduction
2. Triggers
3. Interfaces to functions
4. Runtimes
5. Cold-Starts
6. Keeping state
7. Serverless Trilemma
8. Workflows
9. Programming for serverless computing

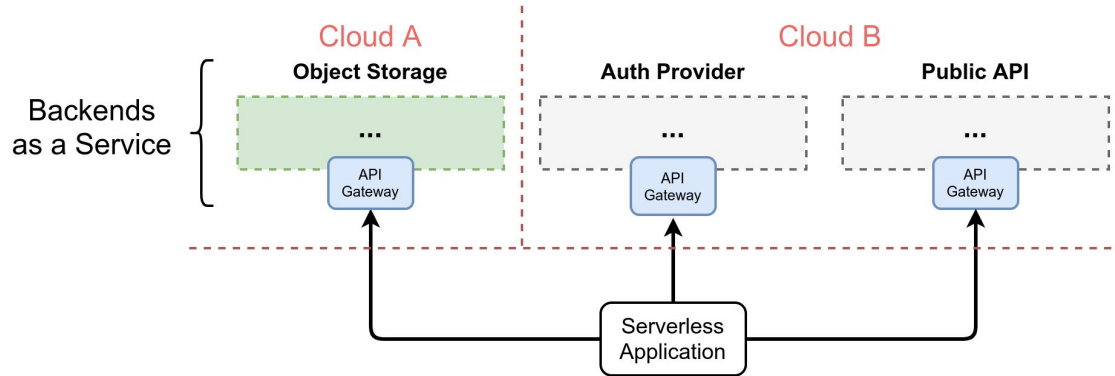
Agenda

10. Suitable workloads and use cases
11. Advantages, disadvantages, learnings

1. Introduction

1.1 What does “serverless” mean?

a) Serverless Architecture



Serverless Architecture

- multiple backends
- might span different cloud providers
- transparent and automatic scaling
- specific backend servers unknown

1.1 What does “serverless” mean?

b) Serverless as cloud-computing execution model

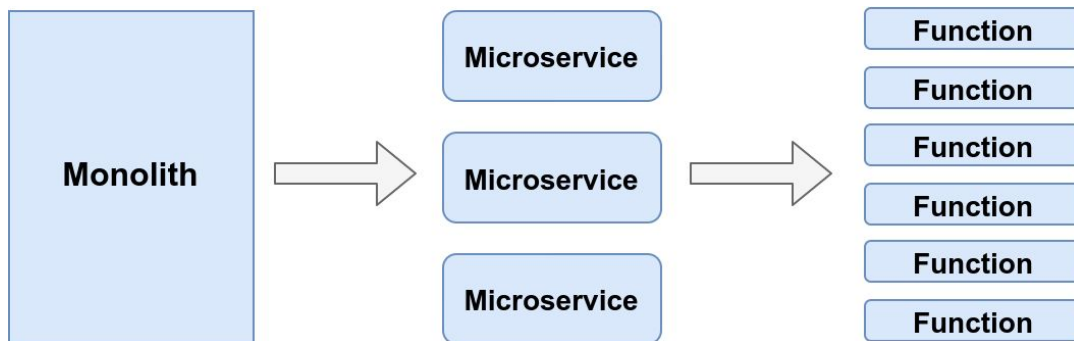
called “**serverless computing**” and
commodified as “**Function as a Service**”

1.1 What does “serverless” mean?

b) Serverless as cloud-computing execution model

called “**serverless computing**” and commodified as “**Function as a Service**”

Are serverless functions just smaller pieces of Microservices?

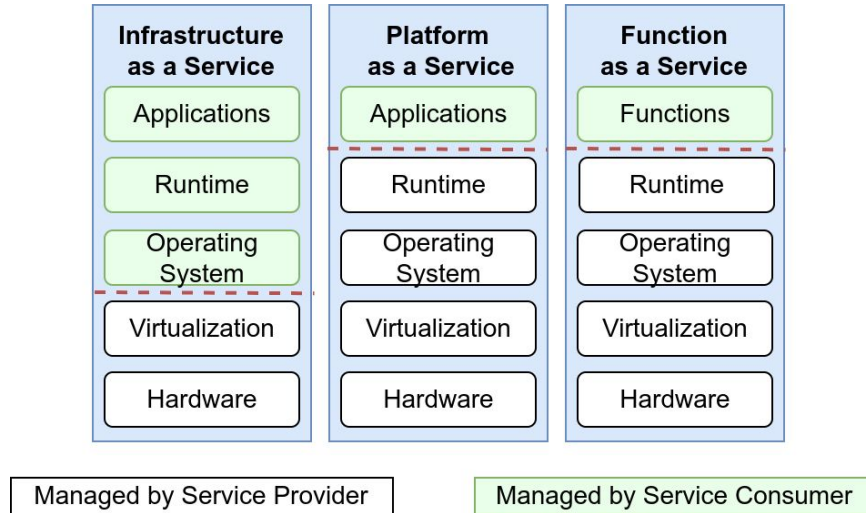


1.1 What does “serverless” mean?

b) Serverless as cloud-computing execution model

called “**serverless computing**” and commodified as “**Function as a Service**”

Which are then run on a runtime like usual PaaS Containers?



1.1 What does “serverless” mean?

Are serverless functions just smaller pieces of Microservices?

Which are then run on a runtime like usual PaaS Containers?

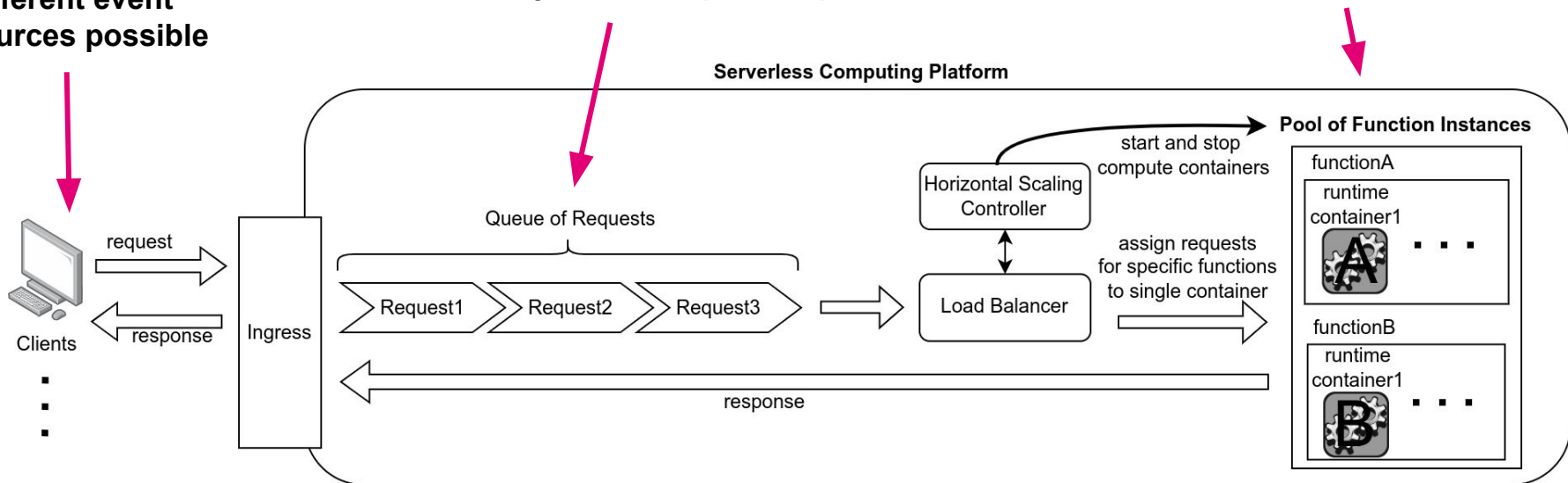
Yes, serverless functions are usually small fast starting CRI containers,
but the technology in general does not demand that

1.2 Attempt to find the main properties of serverless computing

Scales transparently, horizontally, automatically and often down to zero

Different event sources possible

May be kernel packet queue



Broadest thinkable serverless computing definition:

“Horizontally automatically scaling programs are executed in response to events”

1.3 FaaS vs. PaaS

How are PaaS and FaaS different if both are containers?



adrian cockcroft

@adrianco



If your PaaS can efficiently start instances in 20ms that run for half a second, then call it serverless.

[twitter.com/doctor_julz/st...](https://twitter.com/doctor_julz/status/1111111111111111111)

1.3 FaaS vs. PaaS

Platform as a Service

- Runs containers (or on other runtimes)
- Long running (usually)
- Stateless or stateful
- Scales by configuration
- Event-driven or permanently running
- Can have side effects

Function as a Service

- Runs containers (usually)
- Short lived = ephemeral = transient
- Stateless
- Scales automatically
- Event-driven → executed when triggered
- Can have side effects

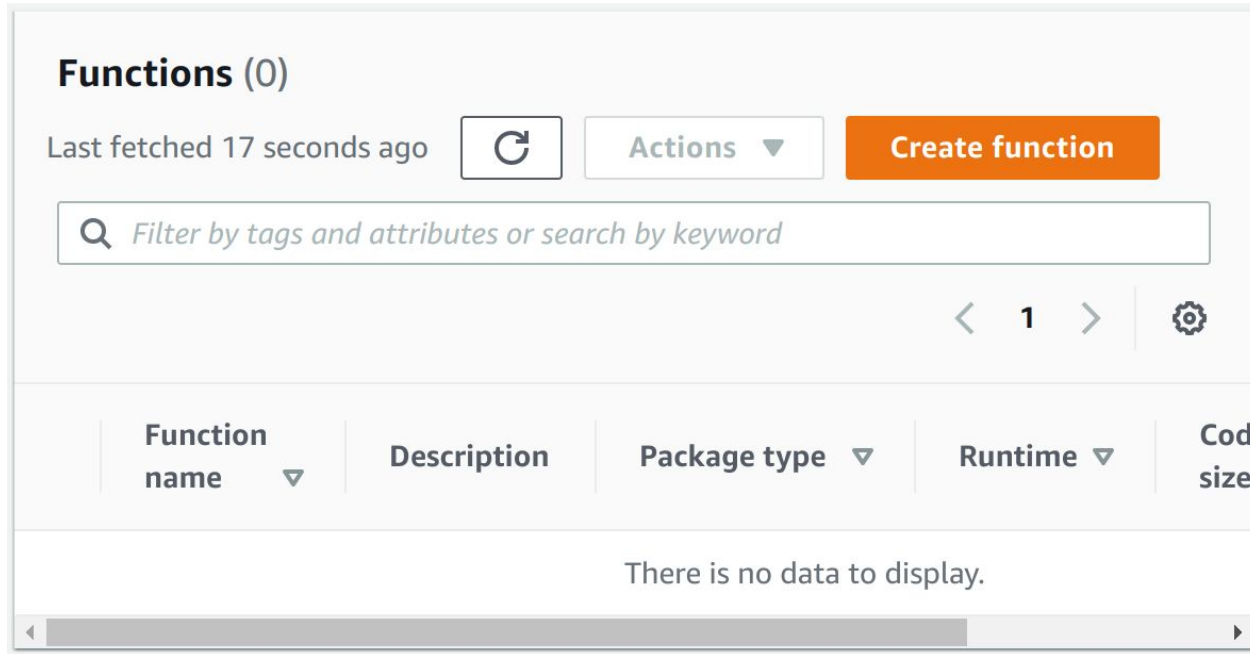
implicit
need for
small size



Main benefits of serverless computing:

- More tasks automatically shifted to the cloud compared to PaaS
- Transparent automatic horizontal scaling makes the technology appear “serverless”
- In public clouds finely granular on-demand billing based on milliseconds execution time
- In public clouds no cost because of scale-to-zero

1.4 AWS Lambda Example



The screenshot shows the AWS Lambda console interface. At the top, it displays "Functions (0)" and "Last fetched 17 seconds ago". There is a refresh button, an "Actions" dropdown menu, and a prominent orange "Create function" button. Below this is a search bar with the placeholder text "Filter by tags and attributes or search by keyword". To the right of the search bar are navigation arrows, a page number "1", and a settings gear icon. A table header is visible with columns for "Function name", "Description", "Package type", "Runtime", and "Code size". The table area is currently empty, displaying the message "There is no data to display." at the bottom.

1.4 AWS Lambda Example

Create function [Info](#)

Choose one of the following options to create your function.

Author from scratch

Start with a simple Hello World example.

Use a blueprint

Build a Lambda application from sample code and configuration presets for common use cases.

Container image

Select a container image to deploy for your function.

Browse serverless app repository

Deploy a sample Lambda application from the AWS Serverless Application Repository.

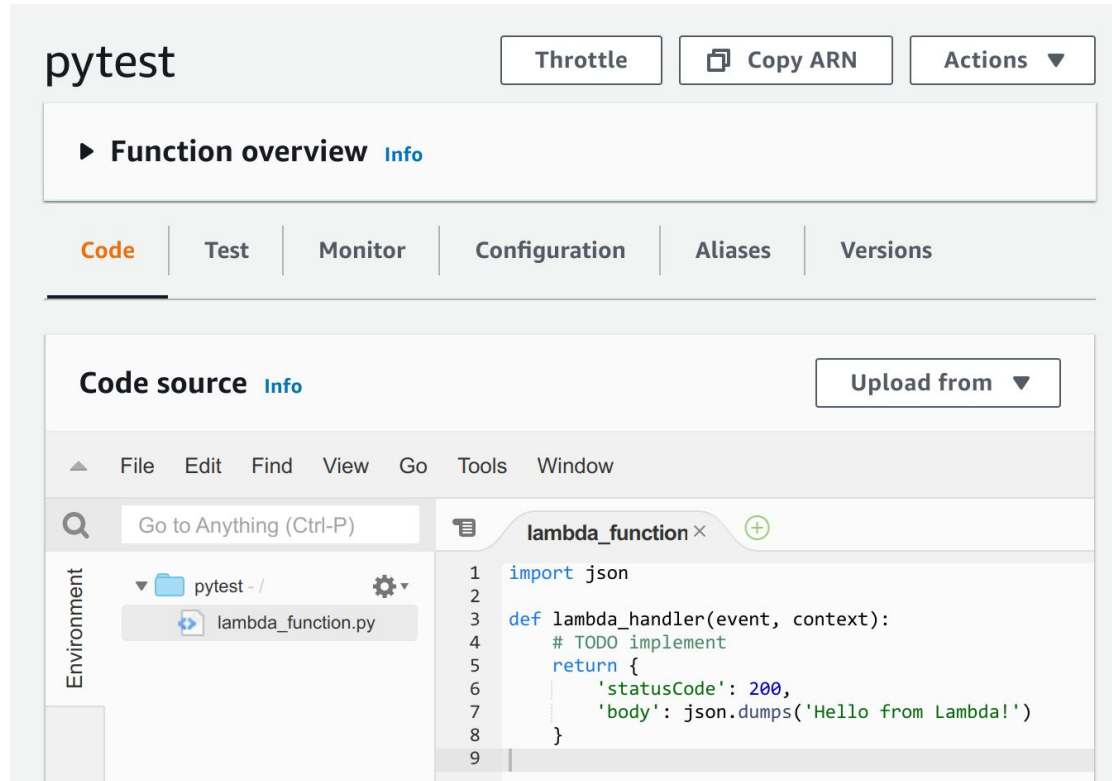
Basic information

Function name
Enter a name that describes the purpose of your function.

Use only letters, numbers, hyphens, or underscores with no spaces.

Runtime [Info](#)
Choose the language to use to write your function. Note that the console code editor supports only Node.js, Python, and Ruby.

1.4 AWS Lambda Example



The screenshot displays the AWS Lambda console interface for a function named "pytest". At the top, there are buttons for "Throttle", "Copy ARN", and "Actions". Below this is a "Function overview" section with an "Info" link. A navigation bar contains tabs for "Code", "Test", "Monitor", "Configuration", "Aliases", and "Versions". The "Code source" section is active, showing a menu with "Upload from" and an "Info" link. Below the navigation is a menu bar with "File", "Edit", "Find", "View", "Go", "Tools", and "Window". A search bar contains "Go to Anything (Ctrl-P)". The "Environment" sidebar shows a folder "pytest -/" and a file "lambda_function.py". The main editor window, titled "lambda_function", displays the following Python code:

```
1 import json
2
3 def lambda_handler(event, context):
4     # TODO implement
5     return {
6         'statusCode': 200,
7         'body': json.dumps('Hello from Lambda!')
8     }
9
```






1.4 AWS Lambda Example

Add trigger

Trigger configuration

Select a trigger

Q

-  **API Gateway**
api application-services aws serverless
-  **AWS IoT**
aws devices iot
-  **Alexa Skills Kit**
alexa iot
-  **Alexa Smart Home**
alexa iot
-  **Apache Kafka**
aws stream



API Gateway: test-API

arn:aws:execute-api:us-east-1:033762176954:██████████/*/*/pytest

▼ Details

API endpoint: <https://██████████.execute-api.us-east-1.amazonaws.com/default/pytest>

API type: **HTTP**

Authorization: **NONE**

Cross-origin resource sharing (CORS): **No**

Enable detailed metrics: **No**

Method: **ANY**

Resource path: **/pytest**



Stage: **default**


1.4 AWS Lambda Example

pytest Throttle Copy ARN Actions ▾

✓ The trigger test-API was successfully added to function pytest. The function is now receiving events from the trigger. ✕

▼ **Function overview** [Info](#)


 **pytest**
 Layers (0)

 **API Gateway** + Add destination

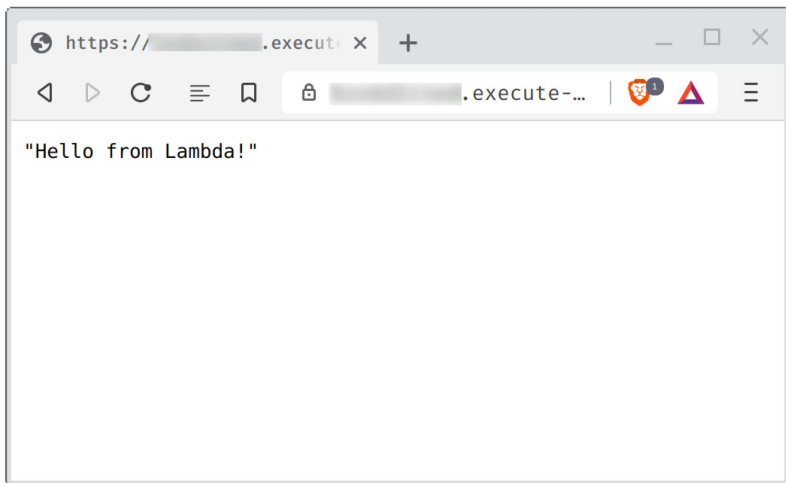
+ Add trigger

Description
-

Last modified
3 minutes ago

Function ARN
 `arn:aws:lambda:us-east-1:762176954:function:pytest`

1.4 AWS Lambda Example



REPORT for **second invocation**:

Duration: 0.89 ms

Billed Duration: 1 ms

("Init Duration" disappeared)

Function Logs **first invocation**

START RequestId: df51350a-dd12-46dd-95d0-d23ba7c524cd Version: \$LATEST

END RequestId: df51350a-dd12-46dd-95d0-d23ba7c524cd

REPORT

RequestId: df51350a-dd12-46dd-95d0-d23ba7c524cd

Duration: 1.32 ms

Billed Duration: **2 ms** ← billed

Memory Size: **128 MB** ← billed

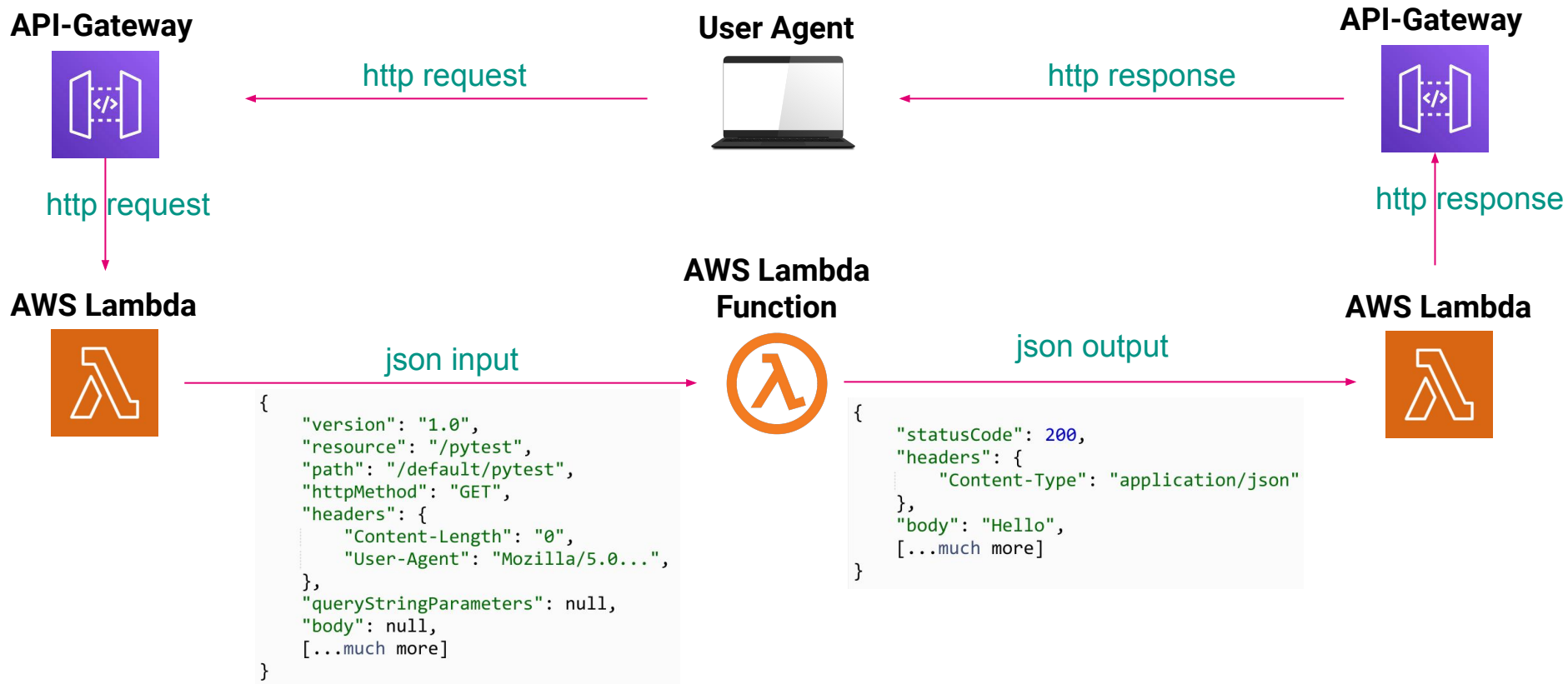
Max Memory Used: **50 MB** ← not billed

Init Duration: **132.34 ms** ← not billed

Request ID

df51350a-dd12-46dd-95d0-d23ba7c524cd

1.4 AWS Lambda Example



1.5 Monoliths as serverless function

Example request to AWS Lambda:

```
{
  "version": "1.0",
  "resource": "/pytest",
  "path": "/default/pytest",
  "httpMethod": "GET",
  "headers": {
    "Content-Length": "0",
    "User-Agent": "Mozilla/5.0...",
  },
  "queryStringParameters": null,
  "body": null,
  [...much more]
}
```

One function can handle multiple subpaths
→ monolith possible in just one function

1.6 AWS Lambda Function

Event = json Input

```
{  
  "key1": "value1",  
  "key2": "value2",  
  "key3": "value3"  
}
```

Context













- App Name
- Function Name
- Memory Limit
- Amazon Features (Log Group, etc.)
- Custom Environment variables

```
import json  
  
def lambda_handler(event, context):  
    # TODO implement  
    return {  
        'statusCode': 200,  
        'body': json.dumps('Hello from Lambda!')  
    }
```













← json output

1.7 Software and Providers

Commercial

- Amazon AWS Lambda  
- IBM Cloud Functions  
- Oracle Cloud Functions 
- Google Cloud Functions  
- Microsoft Azure Functions  
- Cloudflare Workers 
- Vercel Cloud Functions 
- Tencent Cloud Functions  Tencent Cloud

Open Source

- OpenWhisk (Apache-2.0 License)  
- Fn (Apache-2.0 License)  
- Knative (Apache-2.0 License)  
- OpenFaaS (MIT License)  
- Kubeless (Apache-2.0 License)  
- Fission (Apache-2.0 License)  

2. Triggers

2.1 Triggers in Amazon AWS Lambda

Amazon AWS Lambda is highly integrated into the AWS Service Portfolio, Event Sources include:

Invoke functions on Database updates



Trigger on File manipulations in Object Storage



Consume real time data streams



Amazon MSK

Message Queues and Work Queues



Amazon MQ ->



External Events via Amazon EventBridge



Amazon EventBridge

Scheduled Events (Cronjob) via CloudWatch Events

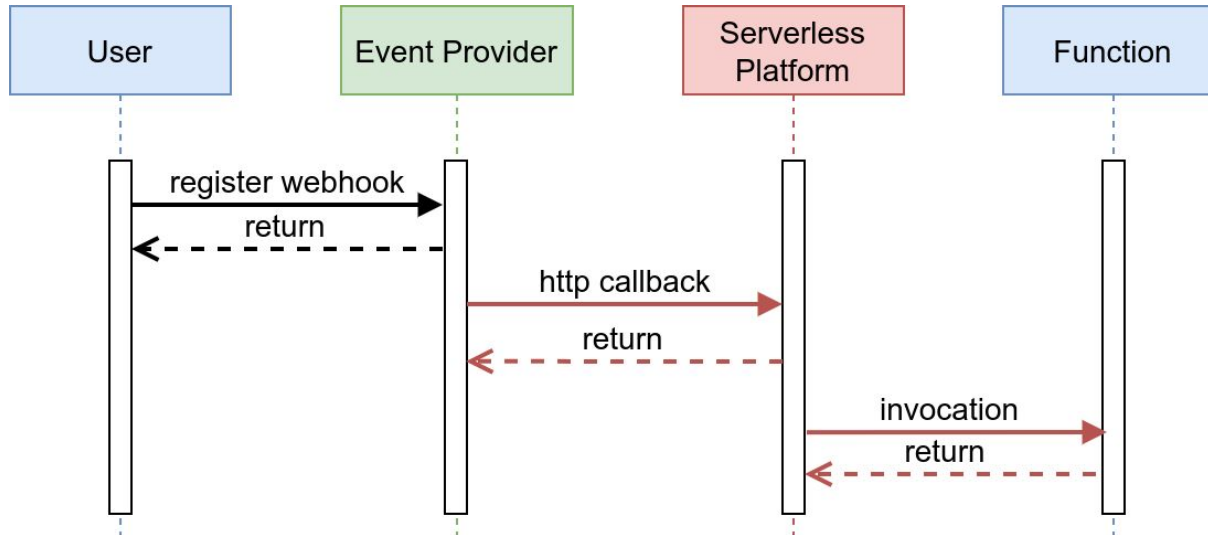


Amazon EventBridge CloudWatch

2.2 Webhooks as event source

a) Webhooks

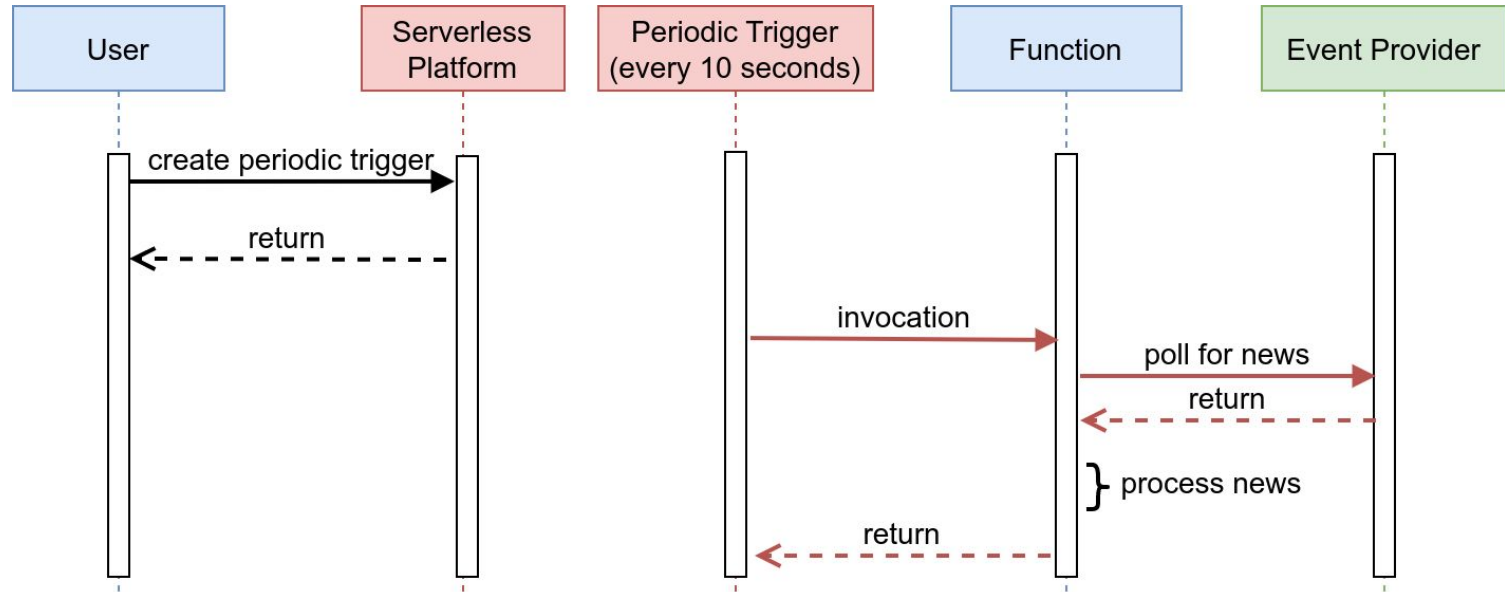
“Telegram, send me a http request, if my bot received a new message”



2.3 Polling as event source

for example polling an RSS feed every 10 minutes

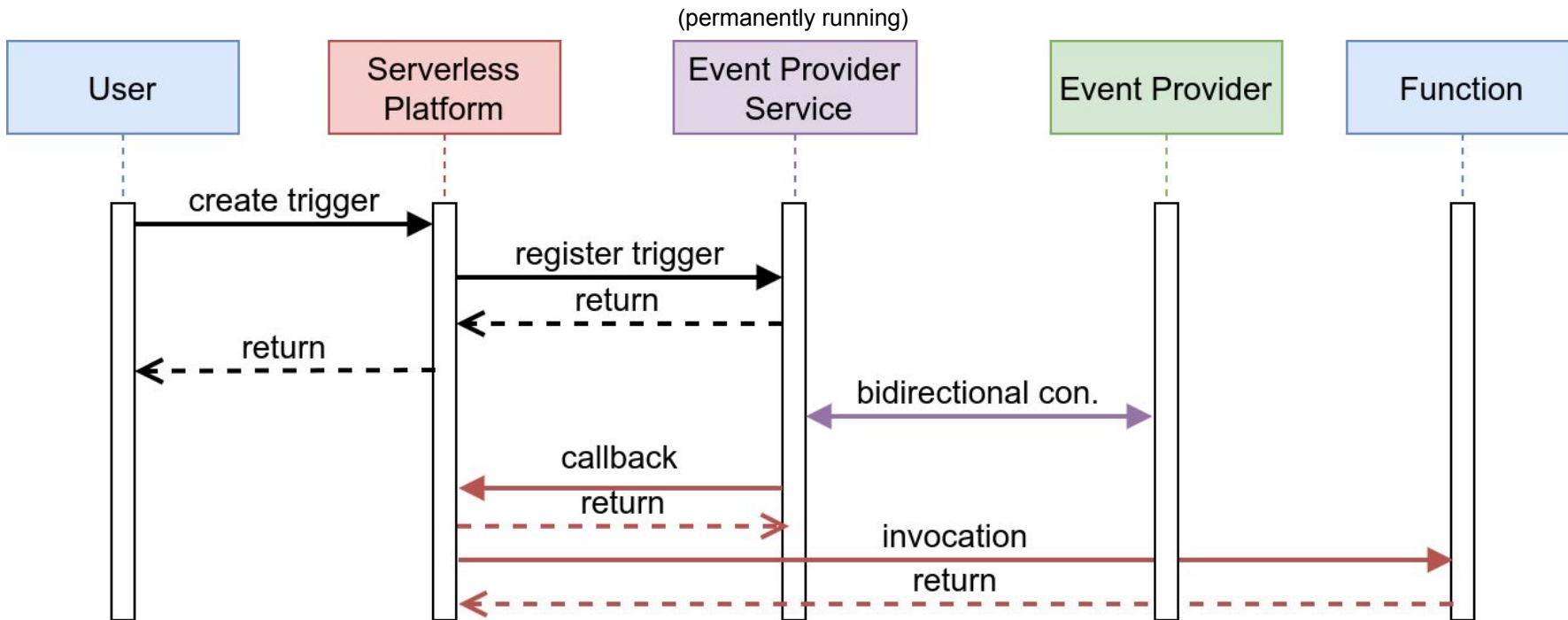
b) Polling



2.4 Permanent connections as event source

c) Connections Pattern

for example subscribing to an mqtt topic



2.5 Any protocol can be connected to serverless computing

Conclusion

- There are patterns to connect any stateless and stateful protocol to serverless computing
- But this may result in performance penalties or the necessity of a permanently running component

2.6 Synchronous and asynchronous invocations

Synchronous Invocation

- Blocks until result is available

Sync Example (OpenWhisk)

```
wsk action invoke \  
  /whisk.system/samples/greeting --result --blocking  
{  
  "payload": "Hello, World!"  
}
```

Asynchronous Invocation

- Returns immediately and gives an Activation ID
- Client can retrieve result later

Async Example (OpenWhisk)

```
wsk action invoke \  
  /whisk.system/samples/greeting  
ok: invoked /_/pythonfunction with id  
733404104295414ab404104295c14ae2
```

3. Interfaces to functions

(Programmers point of view)

3.1 Interfaces to functions

Question: How do programmers receive function invocations?

→ Answer depends on serverless computing implementation

3.2 Language-native map-like data structures

Used in: OpenWhisk, AWS Lambda

a) Platform translates incoming events to json objects

AWS Lambda



json input

```
{
  "version": "1.0",
  "resource": "/pytest",
  "path": "/default/pytest",
  "httpMethod": "GET",
  "headers": {
    "Content-Length": "0",
    "User-Agent": "Mozilla/5.0...",
  },
  "queryStringParameters": null,
  "body": null,
  [...much more]
}
```



json output

```
{
  "statusCode": 200,
  "headers": {
    "Content-Type": "application/json"
  },
  "body": "Hello",
  [...much more]
}
```

AWS Lambda



3.2 Language-native map-like data structures

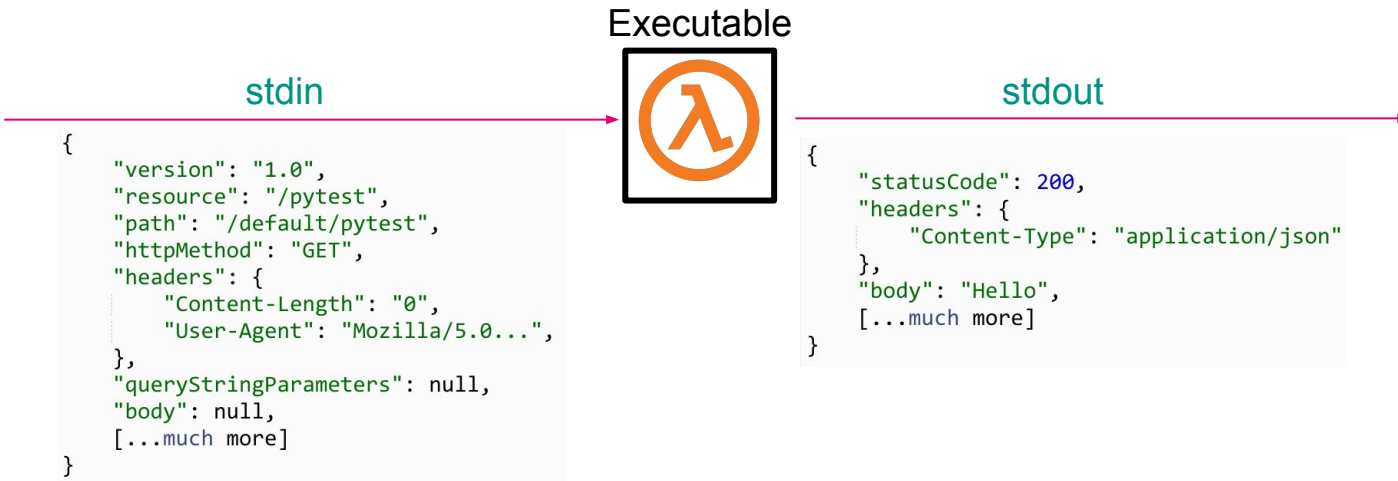
b) json objects are then deserialized to language-native map-like data structures

Go example from OpenWhisk:

```
1 func Main(inobj map[string]interface{}) map[string]interface{} {
2     outobj := make(map[string]interface{})
3     out_header_map := make(map[string]string)
4     out_header_map["headername"] = "headerval"
5     outobj["headers"] = out_header_map
6     outobj["statusCode"] = 200
7     outobj["body"] = "body"
8     return outobj
9 }
```

3.3 json over file descriptors (for example stdin and stdout)

Used in: OpenWhisk, OpenFaaS



→ **Arbitrary executables and bash scripts can be used as serverless functions**

3.4 Language-specific runtime SDK

Used in: Fn, OpenFaas

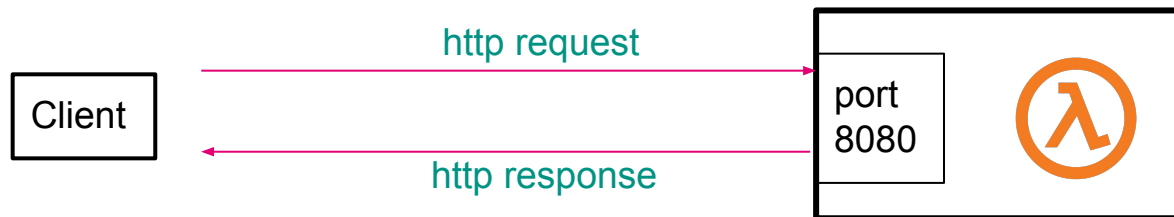
Example from fn project:

```
1 def handler(context, data: io.BytesIO=None):
2     body = json.loads(data.getvalue())
3
4     return response.Response(
5         context,
6         response_data=json.dumps({"message": "Hello {0}".format(name)}),
7         headers={"Content-Type": "application/json"}
8     )
```

3.5 Function receives http-requests on tcp port

Used in: Knative

Function is containerized web-server



- Developers can use their favourite web-server libraries
(In my opinion the most modern approach)
- Function invocation is no additional layer of complexity

3.6 Software abstraction to well-known webserver library

Used in: Fission, OpenFaas, Knative functions

```
~> faas-cli new --lang python3-flask myfunction
Folder: myfunction created.
```



```
Function created in folder: myfunction
Stack file written: myfunction.yml
```

```
~> tree myfunction
myfunction
├── handler.py
├── handler_test.py
├── __init__.py
├── requirements.txt
└── tox.ini
```

- Developers can use their favourite web-server libraries (In my opinion the most modern approach)
- abstraction needs to implement runtime specification, which is sometimes not as trivial as with knative

```
~> cat myfunction/handler.py
def handle(req):
    """handle a request to the function
    Args:
        req (str): request body
    """
    return req
```

3.7 Summary for the function invocation interfaces

Summary

- There are many different ways how programmers receive function invocations
- Letting programmers use well-known web-server libraries in their favourite language is the most modern approach used in fission, knative and openfaas

4. Runtimes

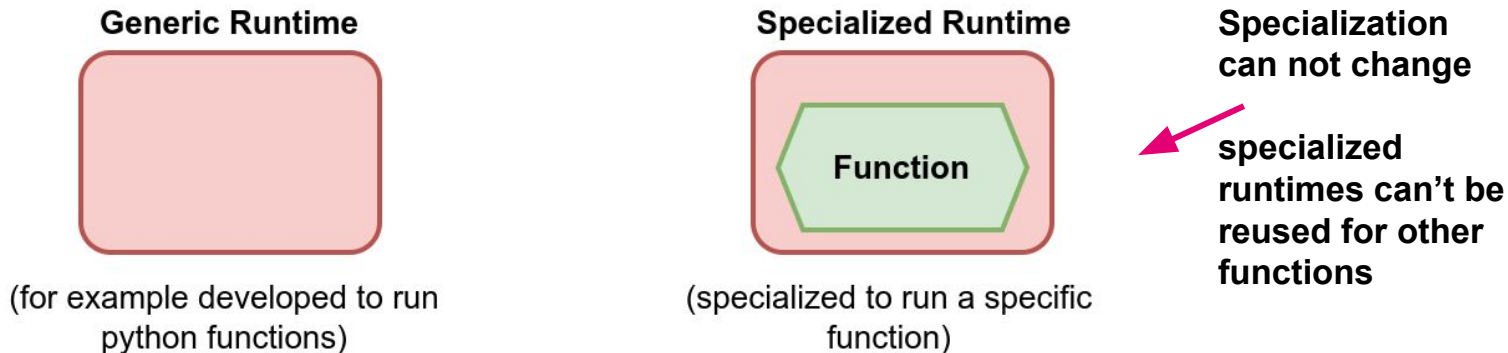
4.1 Deployment unit

There are two observed possibilities:

- A) Function is deployed as a container image
- B) Function is deployed as source-code file or archive containing source-code files

In case B, generic runtime containers are specialized with function code at runtime





















- Initialization effort contributes to cold-start time
- Building a container image in the cluster after deployment would also be possible



4.2 Runtimes for programming languages

Runtimes are usually created at the level of programming languages

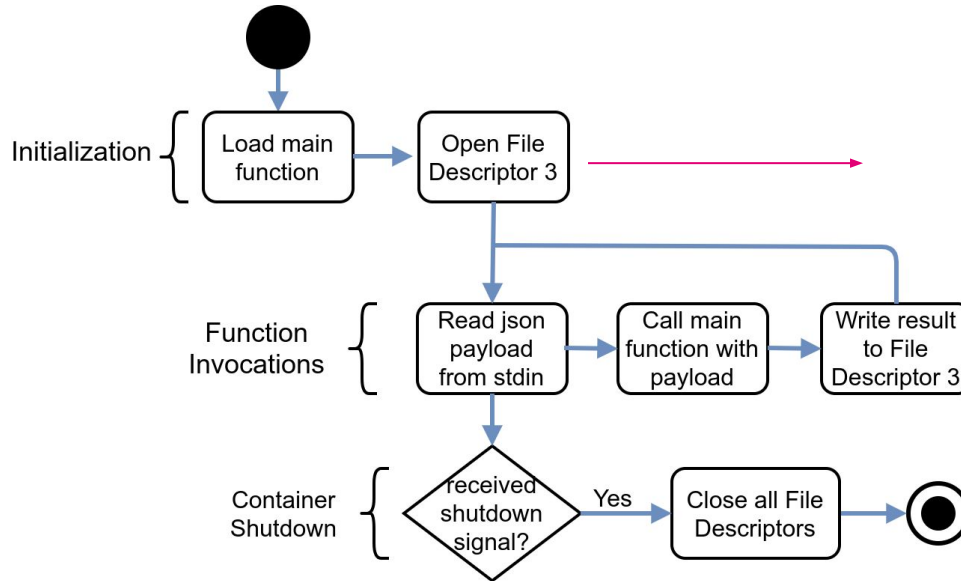
Official Fission “Environments”

 Node.js 	 Go 
 Python 	 PHP 
 Java 	 Java (JVM-Jersey) 
 Perl 	 Ruby 
 TensorFlow 	 Binary 

→ Runtimes usually provide a few common libraries

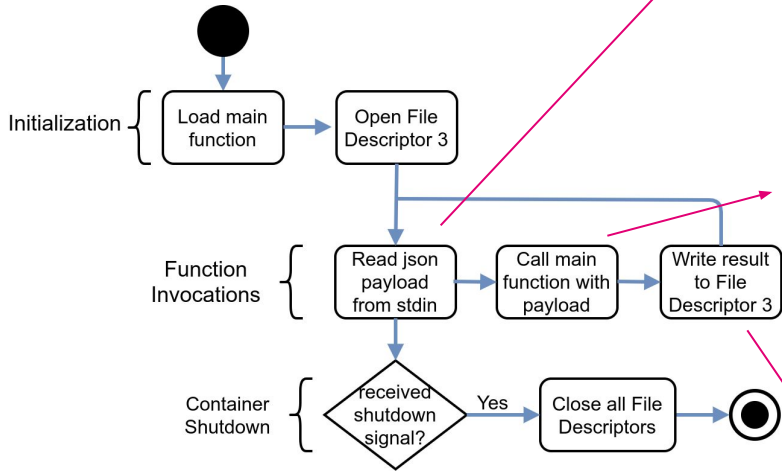
4.3 OpenWhisk Python Runtime

(Code slightly shortened version of the original OpenWhisk Python Runtime Implementation)



```
1 # import the action == function
2 from main_ import main as main
3
4 # Open File Descriptor 3 for output
5 out = fdopen(3, "wb")
6
```

4.3 OpenWhisk Python Runtime



```
8 while True:
9     # read line on each invocation
10    line = stdin.readline()
11    if not line: break
12
13    # Parse json input
14    args = json.loads(line)
15    payload = {}
16    for key in args:
17        if key == "value":
18            payload = args["value"]
19        else:
20            env["__OW_%s" % key.upper()] = args[key]
21    res = {}
22
23    # execute the function
24    try:
25        res = main(payload)
26    except Exception as ex:
27        print(traceback.format_exc(), file=stderr)
28        res = {"error": str(ex)}
29
30    # write result to fd 3
31    out.write(json.dumps(res, ensure_ascii=False).encode('utf-8'))
32    out.write(b'\n')
33    stdout.flush()
34    stderr.flush()
35    out.flush()
```

4.4 Consequences of OpenWhisks Implementation

Consequences of OpenWhisks Implementation

- Json is already parsed to native data structures
- Deserialization and Serialization means overhead
- One runtime serves exactly one action
- Runtimes can't be reused to serve other actions
- Reexecuting same function is fast (warm start)
- First execution is delayed (cold start)

4.5 Limitations

Introduced Problem:

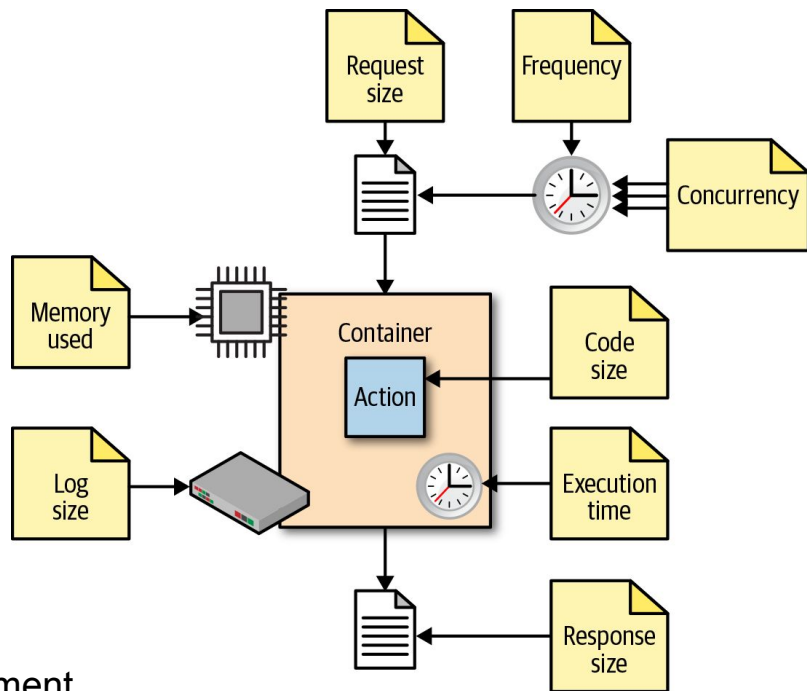
Programmers need to take limitations into account which are specific to the used runtime and software

Limitations in AWS Lambda:

- **unlimited** concurrency (function scale first to region dependent **500-3000** concurrent instances, then **500** more each minute)
- **10240 MiB** max memory usage (def: **128 MiB**)
- **900 sec** max execution time (def: **3 sec**)
- **6 MiB** synchronous invocation payload
- **256 KiB** asynchronous invocation payload

→ payload limits make some use cases hard to implement

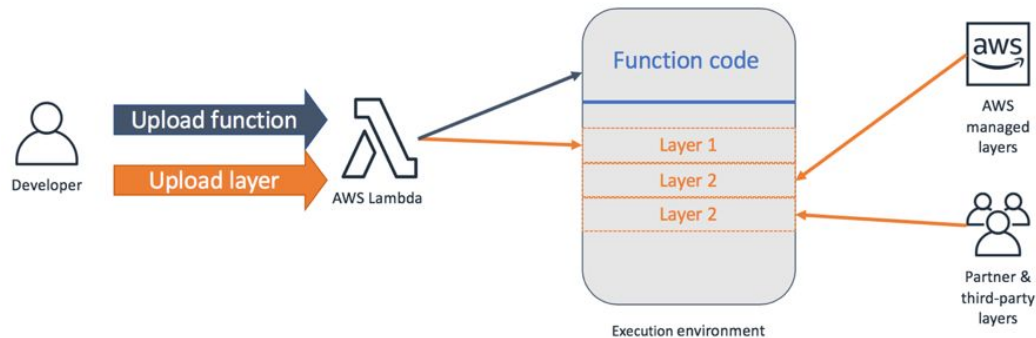
Limitations in OpenWhisk:



4.6 Extending Runtimes

Runtimes can be extended

- Add new container layers
- Provide libraries in deployments




Extend Python Runtimes in OpenWhisk

- 1) `virtualenv virtualenv`
- 2) `source virtualenv/bin/activate`
- 3) `pip install <dependency>`
- 4) `zip -r helloPython.zip virtualenv __main__.py`
- 5) `wsk action create helloPython --kind python:3 helloPython.zip`

```
faas_project
├── hello.py
└── virtualenv
```

4.7 Relationship between serverless computing and container orchestration

- Scaling the number of runtime containers requires interfacing with container orchestration
- Most Self-Hosted serverless software supports  **kubernetes** natively

FaaS Software	Kubernetes Support	Helm Chart available
OpenWhisk	✓	✓
Fn	✓	✓
Kubeless	✓	✓
OpenFaas	✓	✓

→ Other supported platforms are highly implementation dependent

→ Fission and Knative are integrated into Kubernetes with **Custom Resource Definitions**

4.8 Interface to runtime and whether it supports concurrent invocations

	OpenWhisk	OpenFaas	Knative	Fn	Fission
Foundation	OCI containers	OCI containers	OCI containers	OCI containers	OCI containers
Interface to container	tcp-port	tcp-port	tcp-port	unix socket	tcp-port
Protocol to runtime	http (platform API)	http (reverse proxy like)	http (own webserver)	http (own webserver)	http (extendable API)
Concurrency (inside runtime)	no, but activatable	yes → classic yes → of-watchdog (http mode)	yes	no	yes

The serverless platform can of course handle more than one concurrent request.

So why does it matter if a function can be invoked concurrently multiple times for a single runtime?

- *Soft limits make it possible to completely avoid cold-starts during scaling out*
- *Requests can be immediately answered without waiting for a new pod to become available, if runtimes can always tolerate one more request*

5. Cold starts

5.1 Cold- and warm-starts

Cold-start

- Before a function invocation can be processed, a new function instance needs to be started
→ cold-start delay = initialization time + execution time

Warm-start

- The invocation can be forwarded to an already existing function instance with free capacity
→ warm-start delay = execution time

Occurrence of cold-starts

- a) When no function instances are running
 - after scale-to-zero, failure or deployment
- b) During scaling out if all present function instances can't serve more requests.

5.2 Avoiding cold-starts

Occurrence of cold-starts

- a) When no function instances are running
 - after scale-to-zero, failure or deployment
- b) During scaling out if all present function instances can't serve more requests.

Avoid case a)

- disable scale-to zero
and
- set minimum/initial number of function instances ≥ 1

5.2 Avoiding cold-starts

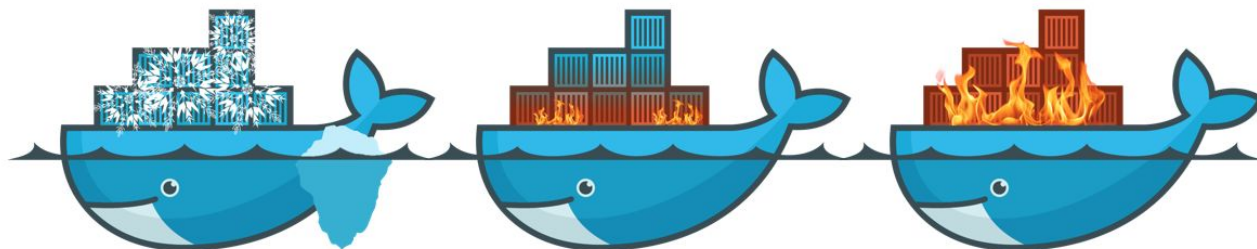
Occurrence of cold-starts

- a) When no function instances are running
 - after scale-to-zero, failure or deployment
- b) During scaling out if all present function instances can't serve more requests.

Avoid case b)

- always have more function instances than needed
 - so there is always an idling function instance available to process an event immediately
- or**
- have function instances, which support concurrent function invocations and do not set a hard-limit for concurrency
 - by using a soft-limit always one more invocation can be processed immediately
 - scaling is detached from function invocation and happens in the background ("eventual scaling")

5.3 Cold- vs. Prewarm- vs. Warm-start



What needs to be done?

cold

prewarm

warm

Starting the container

Initializing the action

Running the action






-> prewarming only possible when generic containers are specialized

5.4 Prewarm containers in OpenWhisk

OpenWhisks invoker launches prewarm containers to accelerate cold starts

Pods

Name	Labels	Node	Status
 wskowdev-invoker-00-6-prewarm-nodejs10	<p>invoker: invoker0</p> <p>name: wskowdev-invoker-00-6-prewarm-nodejs10</p> <p>release: owdev user-action-pod: true</p>	minikube	Running
 wskowdev-invoker-00-1-prewarm-nodejs10	<p>invoker: invoker0</p> <p>name: wskowdev-invoker-00-1-prewarm-nodejs10</p> <p>release: owdev user-action-pod: true</p>	minikube	Running
 owdev-invoker-0	<p>app: owdev-openwhisk</p> <p>chart: openwhisk-1.0.0</p> <p>controller-revision-hash: owdev-invoker-5cd96b87db</p> <p>heritage: Helm name: owdev-invoker</p> <p>Alles anzeigen</p>	minikube	Running

6. Keeping state

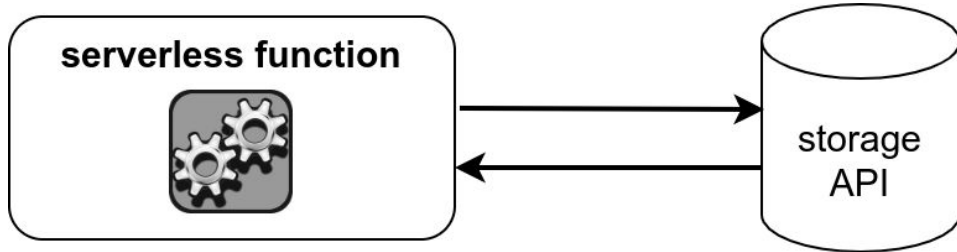
6.1 Why serverless functions need to be stateless

Platform can scale down automatically and transparently

- Content of memory can be removed any time
 - State needs to be externalized
 - Memory is still usable for temporary caching

6.2 Externalize state

External storage API



→ **Suitable solution**

Advantages

- Modular development of storage API
- Separated administration
- Separated scalability

Disadvantages

- Might need additional complexity like locking mechanisms or ACID transactions

6.3 Suitable database and storage APIs for FaaS

Properties of database APIs suitable for use with serverless computing

- Pay-per-use
 - Stateless
 - Fast (AWS bills idling)
 - No expensive handshakes
 - Sacrifice database normalization for speed
 - Automatic Scaling
- *Same properties like those of FaaS functions themselves*

6.4 AWS Aurora Provisioned vs. Serverless



Amazon AWS Aurora

- Relational Database
- Drop-In Replacement for MySQL and PostgreSQL

AWS Aurora - Provisioned

Billing:

- storage (**per GiB/month**)
- I/O rate (**per 1 Mio req**)
- instance size (**per hour**)
(e.g. db.t3.medium = 0,082 USD/h)

AWS Aurora - Serverless

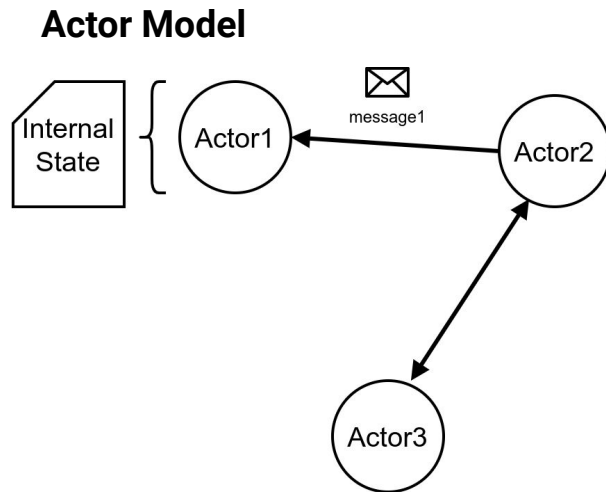
Billing for

- storage (**per GiB/month**)
- I/O rate (**per 1 Mio req**)
- Aurora Capacity Unit (**per hour**)
(1 ACU ~ 2GiB Memory usage)

Differences in Serverless

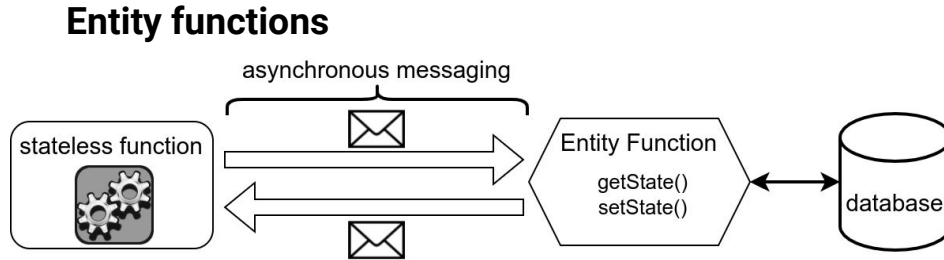
*Different pricing
+ Autoscales up and down
+ Stateless http "Data API"*

6.5 Actors model



- Actor1 has an internal state
- Access to that internal state only via messages

6.6 Entity functions



- Entity functions have a persistent state
- Access to internal state only via event-driven messages
- Backend can be normal storage API

(Practical example: Azure Functions Durable Entities)

Ways to communicate

Signaling: fire-and-forget

→ write operations without confirmation

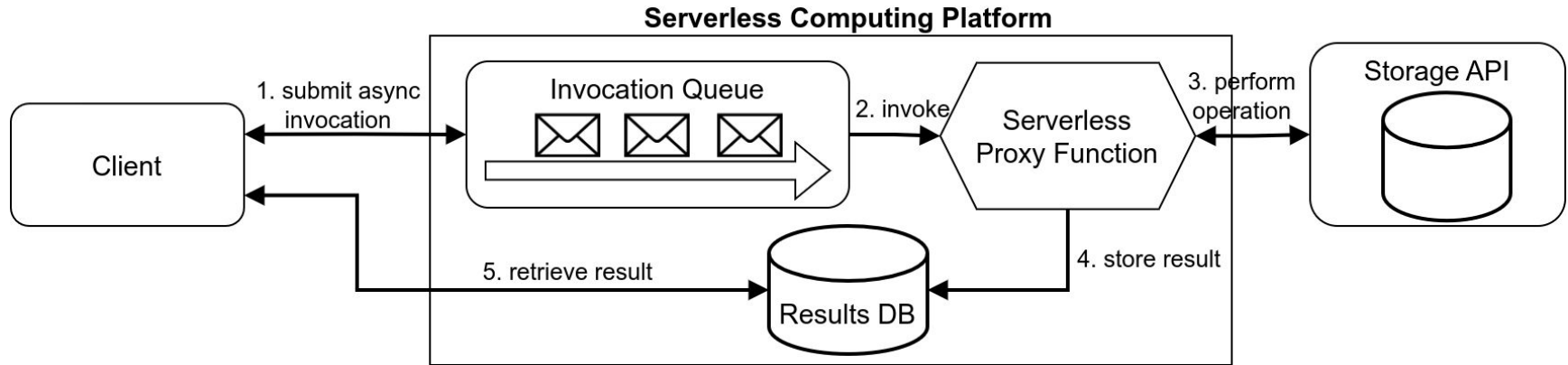
Call: blocking call, waiting for response

→ read and write operations

(identical to normal synchronous APIs)

6.7 Serverless computing platforms as entity functions

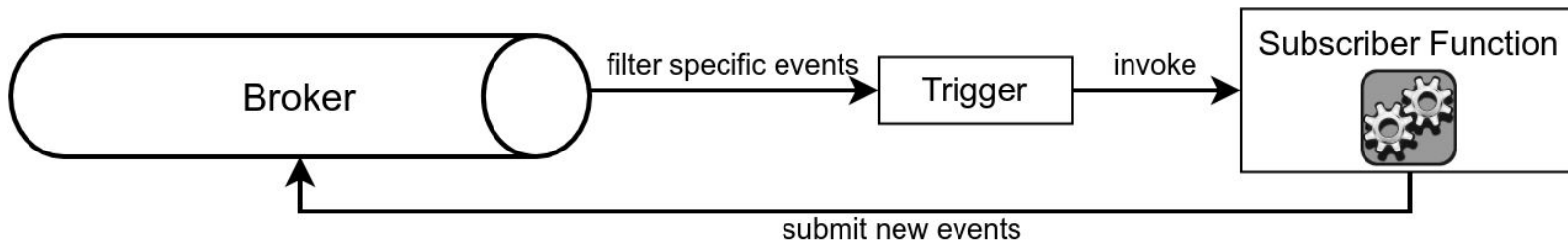
Serverless platforms together with asynchronous invocations can be programmed to be an implementation of entity functions



6.8 Message-broker for asynchronous invocations

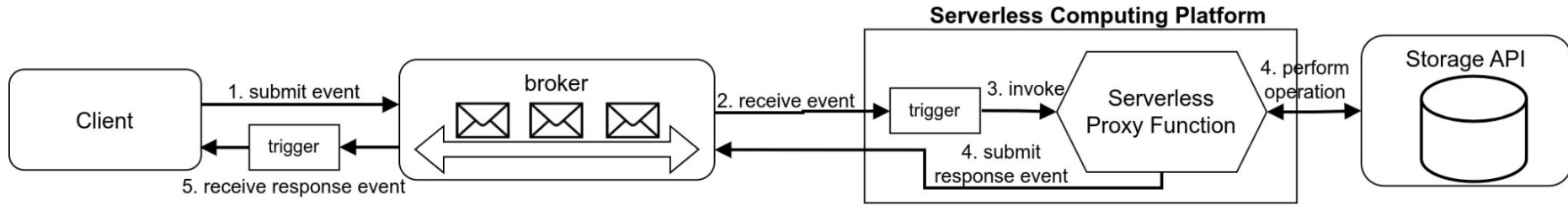
Not all platforms support asynchronous invocations

- But all platforms can be configured to receive events from a message broker



6.8 Message-broker for asynchronous invocations

Connecting a message broker to add the feature of asynchronous invocations to a serverless platform

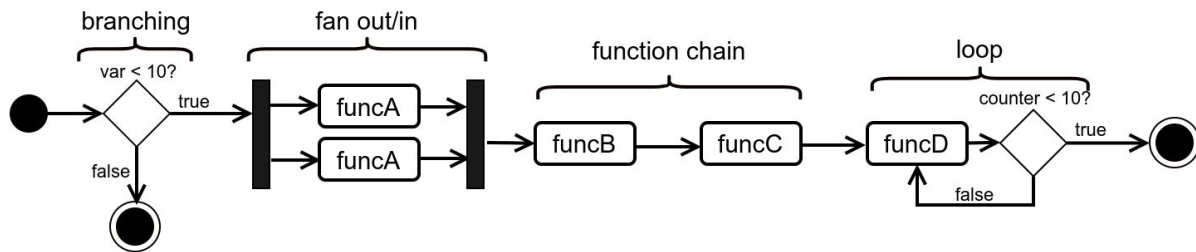


→ **Also a possible implementation of entity functions**

6.9 Store data using in a platform-native workflow system

Storing data in workflows

(some platforms provide a workflow system to connect multiple functions together)



Locations to store data:

- Current workflow state (as in a state machine)
- Parameters forwarded from function to function
- Platform managed workflow context

Advantages

- Function instances use less memory
- Implicit garbage collection
- No double-billing if workflows are ST-safe (explained later)

Disadvantages

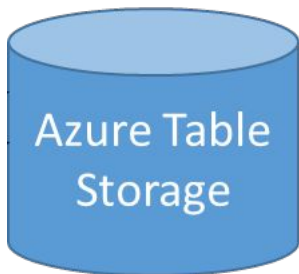
- Code needs to be platform aware
- Impacts modularity (data moved out of API specification to workflows)

6.10 Every method to externalize state uses an external storage

All observed practical implementations of the described methods under the hood externalize state in an external API

Example: “Azure Functions” and “Azure Functions Durable Entities”

Current state is saved in:



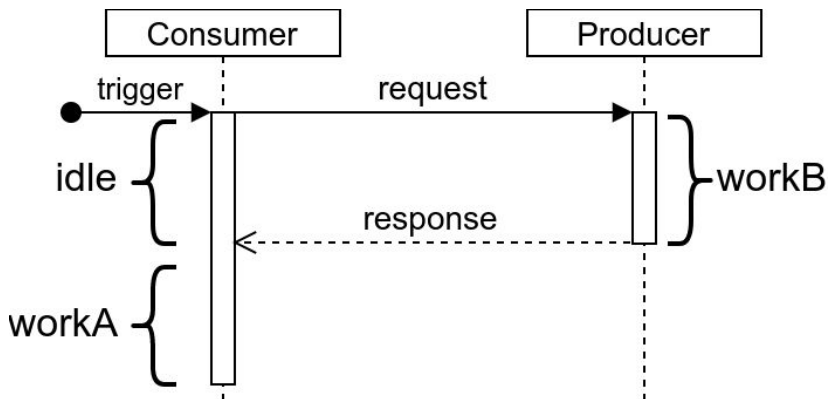
Example: Invocations and their parameters inside “OpenWhisk Sequences” are stored in CouchDB.



7. Serverless Trilemma

7.1 Double billing

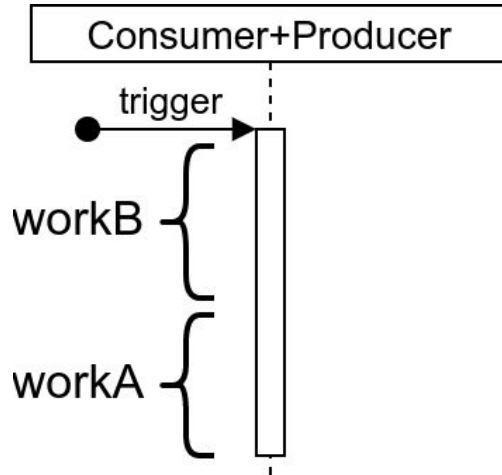
If functions invoke other functions synchronously, the time is billed twice



$$\text{billed_duration} = \text{idle} + \text{workA} + \text{workB}$$

- applies only if billing is based on function execution time
- consumer also blocks a concurrency slot in the function instance

7.2 Solution A: Inline producer code



billed_duration = workA + workB

Problems

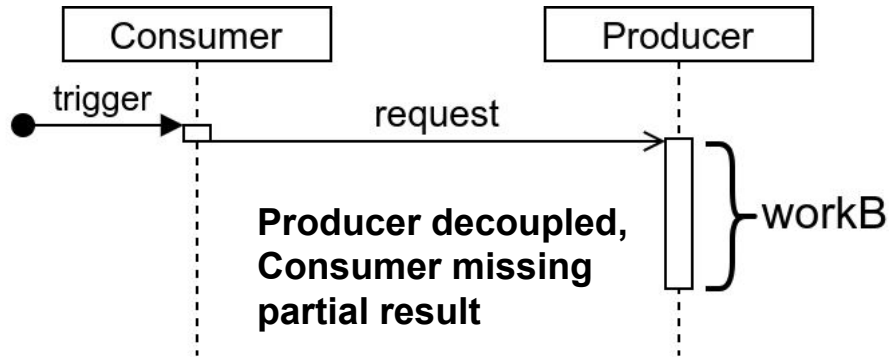
- Knowledge of internal programming of Producer required
- API of Producer not used anymore

→ **Double billing avoided**

→ **Black-Box principle of producer violated**

7.3 Solution B: Asynchronous split off

Invoke other functions asynchronously



→ **Double billing avoided**

→ **Consumer can't produce same result as before**

→ **Consumer can't be a substitutable part of a composition anymore, because it doesn't behave identical to an atomic function**

→ **If intermediate result is not required, solution can be still suitable**
(for example when sending an email notification)

Note:

- Waiting for a result would be synchronous operation with double billing

7.4 Serverless Trilemma

Desired properties of serverless computing architectures:

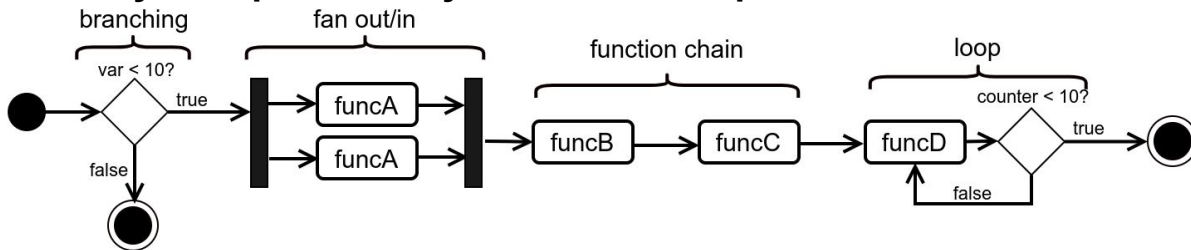
1. **Treat functions as black-box**
2. **No Double-Billing**
3. **Substitutable functions**

Serverless Trilemma: In an event-driven system one of the properties must be violated

8. Workflows

8.1 ST-Safe workflows avoid serverless trilemma

Workflow system provided by the serverless platform



ST-Safe

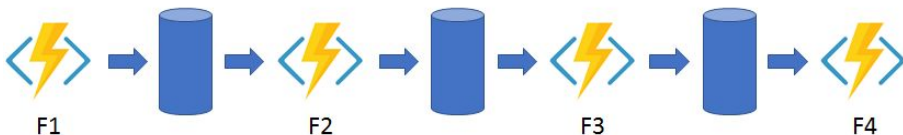
- A serverless function orchestration system that satisfies all three conditions of the serverless trilemma is called ST-Safe

Serverless trilemma can be applied “backwards”

- If a program doesn't use ST-Safe workflows, it must violate one of the desired properties of the serverless trilemma

8.2 Workflows in Azure Functions

Function Chaining



→ Azure Durable functions uses syntax known from async programming

```
[FunctionName("Chaining")]
public static async Task<object> Run(
    [OrchestrationTrigger] IDurableOrchestrationContext context)
{
    try
    {
        var x = await context.CallActivityAsync<object>("F1", null);
        var y = await context.CallActivityAsync<object>("F2", x);
        var z = await context.CallActivityAsync<object>("F3", y);
        return await context.CallActivityAsync<object>("F4", z);
    }
    catch (Exception)
    {
        // Error handling or compensation goes here.
    }
}
```

“await” pauses billed execution time

activities

orchestrator function

8.3 Workflows in Azure Functions - Replaying and Checkpoints

```
[FunctionName("Chaining")]
public static async Task<object> Run(
    [OrchestrationTrigger] IDurableOrchestrationContext context)
{
    try
    {
        var x = await context.CallActivityAsync<object>("F1", null);
        var y = await context.CallActivityAsync<object>("F2", x);
        var z = await context.CallActivityAsync<object>("F3", y);
        return await context.CallActivityAsync<object>("F4", z);
    }
    catch (Exception)
    {
        // Error handling or compensation goes here.
    }
}
```

orchestrator function

activities

How it works:

- orchestrator function is invoked multiple times → replay
- if an activity already happened result is returned immediately → checkpoint from Azure Table Storage

9. Programming for serverless computing

9.1 Consequences of formal foundations paper

Paper “Formal Foundations of Serverless Computing”

- Contains theoretical model of serverless computing based on operational semantics
- does not fit perfectly to recent developments
- but introduces relevant properties programmers need to account for.

Usually given guarantee:

At-Most Once execution

(Retries of failed function invocations are usually a separated feature, which is turned off by default)

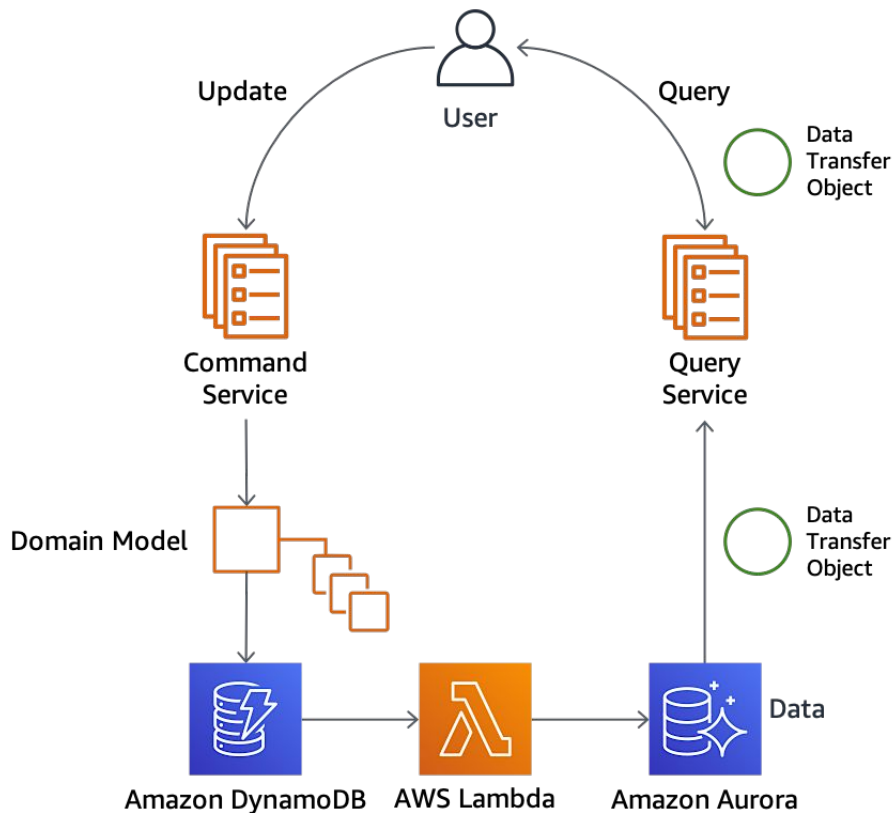
Function as a Service

- Execution of requests
 - not necessarily in order
 - possibly in parallel
 - may happen twice

Solution in short:

ACID transactions in storage backend
+ locking/mutex mechanisms
+ idempotent functions

9.2 Command-Query-Responsibility-Segregation



Idea

- Different models for reading and writing data

Command Model

- Scalable and schema-less database for creating and updating data fast

Query Model

- Relational database, normalized with schema for complex queries

Lambda function is triggered on each update and translates between the two data models!

10. Suitable workloads and use cases

10.1 Relationship between Edge Computing and FaaS

Edge Computing

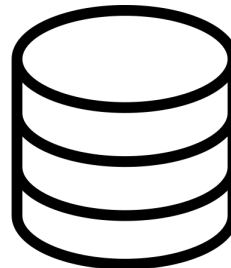
- Process data where it is generated and/or needed

The data a function operates on may be a lot bigger

Functions have small size



(Small source-code files or small docker image layers)



Conclusions

- Distributing serverless functions world-wide is comparatively easy
 - Moving data world-wide may be hard
- Small deployment size makes serverless computing especially suitable for edge computing

10.2 Bursty traffic patterns

Deployment with permanently provisioned resources

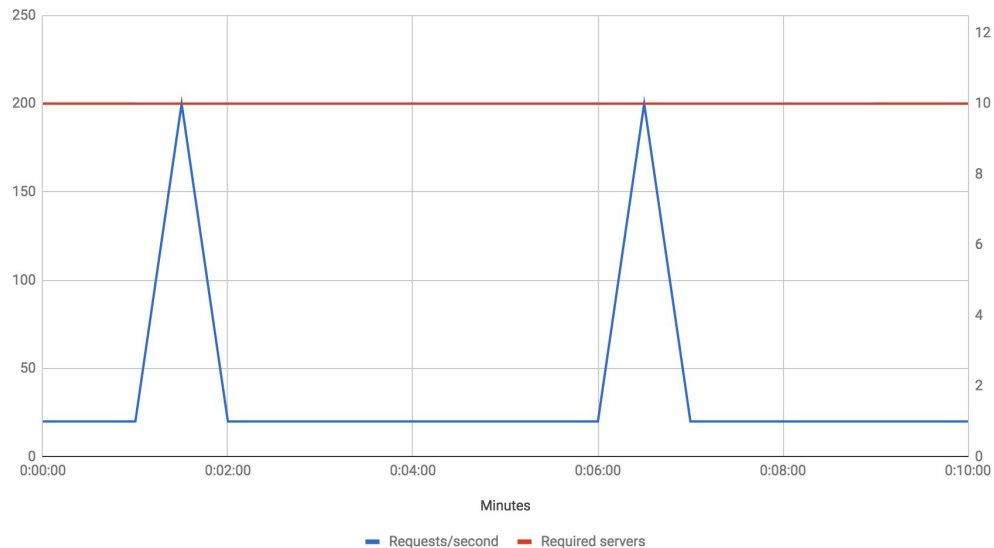
- permanent need for 10 servers to account for all possible traffic spikes

Deployment with on-demand resource usage

- scaling down avoids idling servers

→ Billing based on milliseconds execution time makes serverless computing very suitable for bursty traffic patterns

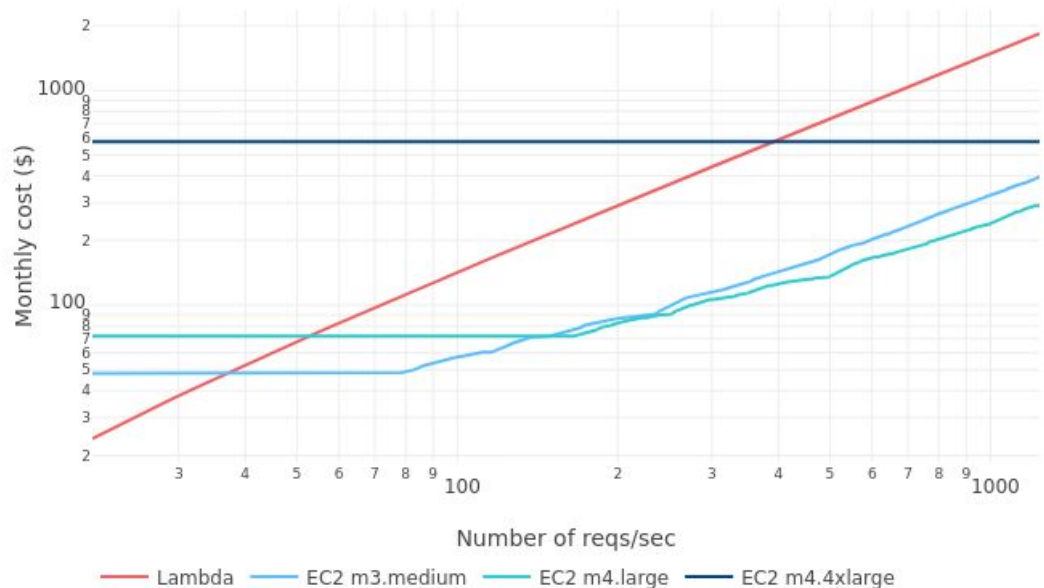
Bursty traffic pattern (example with provisioned resources)



10.3 Exemplary cost calculation - provisioned vs serverless

AWS Lambda vs AWS EC2 Instances

Monthly cost by number of requests per second



In this specific example

→ A high number of reqs/sec makes FaaS uneconomic in any case

In general

Serverless functions good for

- Low total number of reqs/sec
 - High idle times (or bursty traffic)
 - Unpredictable scenarios
- no sizing necessary beforehand

10.4 Predicting cost is hard in AWS Lambda

(Prices for EU Frankfurt)

Memory (MB)	Price per 1ms
128	\$0.0000000021
512	\$0.0000000083
1024	\$0.0000000167
1536	\$0.0000000250
2048	\$0.0000000333

execution_time_price

+ price for number of requests

+ traffic price inbound and outbound a region

+ other services used

For example: S3 Object Storage

- GB stored per month
- per 1000 requests
- Replication

Price	
Requests	\$0.20 per 1M requests

Pricing	
Data Transfer IN To Amazon EC2 From Internet	
All data transfer in	\$0.00 per GB
Data Transfer OUT From Amazon EC2 To Internet	
Up to 1 GB / Month	\$0.00 per GB
Next 9.999 TB / Month	\$0.09 per GB
Next 40 TB / Month	\$0.085 per GB
Next 100 TB / Month	\$0.07 per GB
Greater than 150 TB / Month	\$0.05 per GB

AWS Lambda can easily cost more

11. Advantages, disadvantages, learnings

11.1 Complexity prevented from developers

Complexity prevented from developers

- Managing horizontal scaling **(K)**
- Achieving high redundancy and availability **K**
- Evenly distributed load balancing **K**
- Performing rolling updates **(K)**
- No need for system upgrades **K**
- No infrastructure management **K**

K = Advantage already given by kubernetes

(K) = Advantage already given, but limited

Why using serverless computing at all?

→ small config file implicitly gives many features

→ new features

11.2 Using serverless computing enables features implicitly

```
1  apiVersion: apps/v1
2  kind: Deployment
3  metadata:
4    name: my-nginx
5  spec:
6    selector:
7      matchLabels:
8        run: my-nginx
9    replicas: 2
10   template:
11     metadata:
12       labels:
13         run: my-nginx
14     spec:
15       containers:
16         - name: my-nginx
17           image: nginx
18           ports:
19             - containerPort: 80
20   ---
21  apiVersion: v1
22  kind: Service
23  metadata:
24    name: my-nginx
25    labels:
26      run: my-nginx
27  spec:
28    ports:
29      - port: 80
30      protocol: TCP
31    selector:
32      run: my-nginx
```

← Kubernetes

Knative

```
1  apiVersion: serving.knative.dev/v1
2  kind: Service
3  metadata:
4    name: my-nginx
5  spec:
6    template:
7      spec:
8        containers:
9          - image: nginx
10         ports:
11           - containerPort: 80
```

Less yaml, but implicitly more features like

- scale-to-zero
- on-demand resource usage
- revisions
- etc.

→ **new projects can leverage these features without any effort**

11.3 New features of serverless computing platforms

Many new features introduced.

Example: Sophisticated Autoscaling in Knative

```
10 # scaling targets 10 requests per 30s for each replica
11 autoscaling.knative.dev/target: "10"
12 autoscaling.knative.dev/metric: "rps"
13 autoscaling.knative.dev/window: "30s"
14
15 # panic mode starts when 200% requests arrive,
16 # relative to what current replicas can handle
17 autoscaling.knative.dev/panic-threshold-percentage: "200.0"
18
19 # panic mode decreases window to 10% size
20 autoscaling.knative.dev/panic-window-percentage: "10.0"
```

11.4 Complexity added for developers

Complexity added for developers

- New software and cli tools to learn
- Need to consider unique properties
 - Stateless
 - Cold-starts
 - Time and resource limitations
 - Serverless-trilemma
- Additional problems
 - Double-billing
 - Serverless trilemma
 - How and why to use workflows
 - New security vulnerabilities <- *there are papers*
 - Need to program event generators
- Less comfort
 - Harder debugging of applications

11.5 Operations when using FaaS

Does using serverless computing in a public cloud mean we don't need an operations team? -> **No**

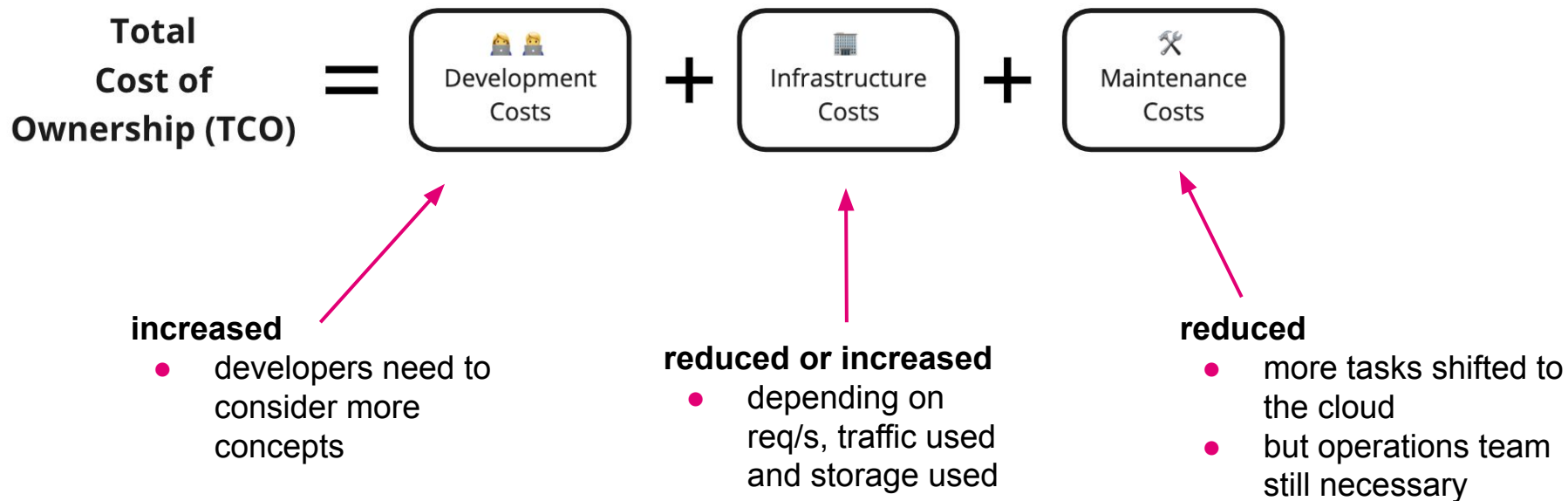
FaaS software still requires operations:

- Monitoring
- Logging
- Backup and Recovery strategy
- IT-Security (Intrusion detection, etc.)
- ...

But usually cloud provided:

- Firewall
- Reverse Proxy
- Load Balancers
- Cluster and Nodes
- ...

11.6 Total cost of ownership



11.7 New features of serverless computing platforms

My conclusion:

- Making apps stateless is a good idea to automatically achieve horizontal scalability, but this is independent from serverless computing
- The technology serverless computing doesn't add enough advantages compared to kubernetes to justify another complexity layer
- New and small projects can benefit from effortlessly usable features like automatic horizontal scaling

End of presentation

**Thank you for your
participation**

Feel free to ask questions

Sources (1 / 4)

01. Sources on the state of research (you of course do not need to read everything, this is just for interested people)

H. B. Hassan, S. A. Barakat, and Q. I. Sarhan, "Survey on serverless computing," J Cloud Comp, vol. 10, no. 1, p. 39, Dec. 2021, doi: 10.1186/s13677-021-00253-7.

(starting on page 16 in the section "Serverless computing challenges and issues (RQ7)" there is a good and concise overview of the open questions of the current state of research on serverless computing. The paper is from 2021 and therefore very current. Starting on page 19, there are also research questions)

A. Bocci, S. Forti, G.-L. Ferrari, and A. Brogi, "Secure FaaS orchestration in the fog: how far are we?," Computing, vol. 103, no. 5, pp. 1025–1056, May 2021, doi: 10.1007/s00607-021-00924-y.

(Very good overview of selected current papers around the topic of serverless computing. Short paragraphs describe the content of each paper. I recommend reading chapters 4 and 5)

A. Jangda, D. Pinckney, Y. Brun, and A. Guha, "Formal Foundations of Serverless Computing," Proc. ACM Program. Lang., vol. 3, no. OOPSLA, pp. 1–26, Oct. 2019, doi: 10.1145/3360575.

(The paper defines Operational Semantics for Serverless Computing, highlighting a few specifics that derive from this theory for FaaS platforms. First of interest are pages 4 and 5. Page 4 shows a sketch and a long source code. The long source code compensates for the problems of FaaS by writing more code. Which problems these are can be read in the following text.)

G. C. Fox, V. Ishakian, V. Muthusamy, and A. Slominski, "Status of Serverless Computing and Function-as-a-Service(FaaS) in Industry and Research," 2017, doi: 10.13140/RG.2.2.15007.87206.

(The results of a conference on serverless computing. The text provides a few more practical ideas on the direction the technology might continue to take.)

A. Brogi, S. Forti, C. Guerrero, and I. Lera, "How to place your apps in the fog: State of the art and open challenges," Softw: Pract Exper, vol. 50, no. 5, pp. 719–740, May 2020, doi: 10.1002/spe.2766.

(In edge or fog computing, there are a number of scattered servers to which one or more applications are to be deployed. The resulting "Fog Application Placement Problem (FAPP)" is a heavily researched topic. The paper shows on page 5 a review of the various papers that attempt to solve the FAPP problem.)

Sources (2 / 4)

E. Jonas et al., "Cloud Programming Simplified: A Berkeley View on Serverless Computing." arXiv, Feb. 09, 2019. Accessed: Dec. 05, 2022. [Online]. Available: <http://arxiv.org/abs/1902.03383>

(General evaluation of the technology and its current problems)

P. Aditya et al., "Will Serverless Computing Revolutionize NFV?," in Proceedings of the IEEE, vol. 107, no. 4, pp. 667-678, April 2019, doi: 10.1109/JPROC.2019.2898101.

(The paper is about running virtualized network functions like those of the 5G core on serverless computing.)

H. Shafiei, A. Khonsari, and P. Mousavi, "Serverless Computing: A Survey of Opportunities, Challenges, and Applications," ACM Comput. Surv., vol. 54, no. 11s, pp. 1–32, Jan. 2022, doi: 10.1145/3510611.

(Another overview)

**02. Amazon (Hrsg.). (.). aws. Abgerufen am 9. Juni 2021
von <https://docs.aws.amazon.com/lambda/latest/dg/invoke-eventsourcemapping.html>**

**03. GitHub. (17. Mai 2010). Abgerufen am 9. Juni 2021
von <https://github.com/apache/openwhisk/blob/master/docs/feeds.md>.**

**04. GitHub. (17. Mai 2010). Abgerufen am 9. Juni 2021
von <https://github.com/apache/openwhisk/blob/master/docs/actions.md>**

Sources (3 / 4)

05. GitHub. (19. Februar 2020). Abgerufen am 9. Juni 2021 von

<https://github.com/apache/openwhisk-runtime-python/blob/master/core/python3ActionLoop/lib/launcher.py>

06. Sciabarrà, M. (2021). O'Reilly. Abgerufen am 9. Juni 2021 von

<https://learning.oreilly.com/library/view/learning-apache-openwhisk/9781492046158/ch01.html>

07. Amazon (Hrsg.). (kein Datum). aws. Abgerufen am 9. Juni 2021 von

<https://aws.amazon.com/de/blogs/compute/working-with-aws-lambda-and-lambda-layers-in-aws-sam/>

08. Amazon (Hrsg.). (kein Datum). aws. Abgerufen am 9. Juni 2021 von

<https://github.com/apache/openwhisk/blob/master/docs/actions-python.md>

09. Thömmes, M. (20. April 2017). Apache OpenWhisk. Abgerufen am 9. Juni 2021 von

<https://medium.com/openwhisk/squeezing-the-milliseconds-how-to-make-serverless-platforms-blazing-fast-aea0e9951bd0>

10. Storti, B. (09. Juli 2015). brianstorti.com. Abgerufen am 23. Juni 2021 von

<https://www.brianstorti.com/the-actor-model/>

11. Microsoft (Hrsg.). (kein Datum). aws. Abgerufen am 9. Juni 2021 von

<https://learn.microsoft.com/en-us/azure/azure-functions/durable/durable-functions-entities?tabs=csharp>

12. Microsoft (Hrsg.). (kein Datum). aws. Abgerufen am 9. Juni 2021 von

<https://github.com/Azure/durabletask/wiki/Core-Concepts>

Sources (4 / 4)

13. Agarwal, A., Choudhary, C., & Bhagat, S. (2018). *The Serverless Trilemma - Function Composition for Serverless Computing*. Washington: University of Washington .
14. GitHub (02. Juli 2020). Abgerufen am 23. Juni 2021 von <https://github.com/Azure/durabletask/wiki/Writing-Task-Orchestrations>
15. Abhinav Jangda, Donald Pinckney, Yuriy Brun, and Arjun Guha. 2019. Formal foundations of serverless computing. Proc. ACM Program. Lang. 3, OOPSLA, Article 149 (October 2019), 26 pages. <https://doi.org/10.1145/3360575>
16. Amazon (Hrsg.). (kein Datum). aws. Abgerufen am 23. Juni 2021 von <https://docs.aws.amazon.com/whitepapers/latest/modern-application-development-on-aws/command-query-responsibility-segregation.html>
17. Rodríguez, Á. A. (9. Dezember 2020). BBVA. Abgerufen am 9. Juni 2021 von <https://www.bbva.com/en/economics-of-serverless/>
18. Amazon (Hrsg.). (kein Datum). aws. Abgerufen am 23. Juni 2021 von https://aws.amazon.com/lambda/pricing/?nc1=h_ls
19. Amazon (Hrsg.). (kein Datum). aws. Abgerufen am 23. Juni 2021 von <https://aws.amazon.com/de/s3/pricing/>