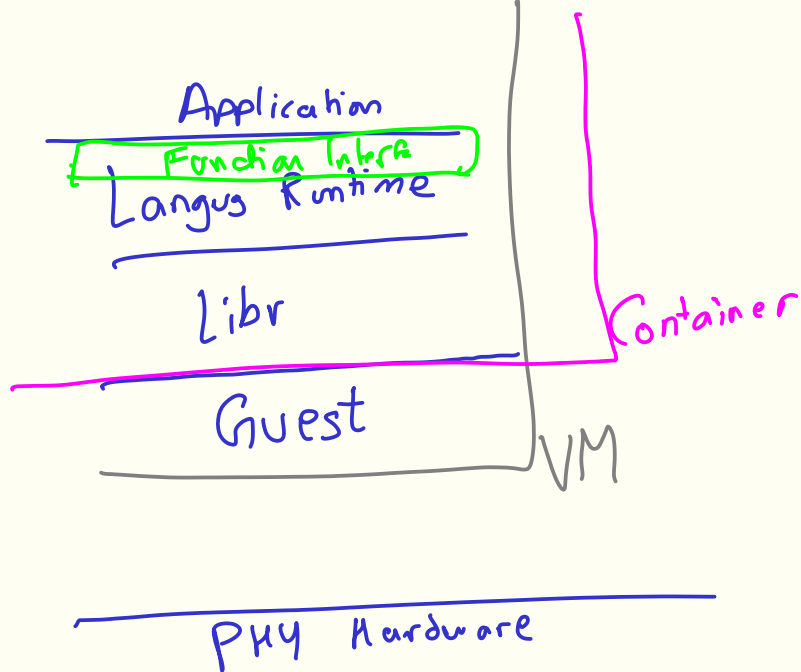# Functions as a Service

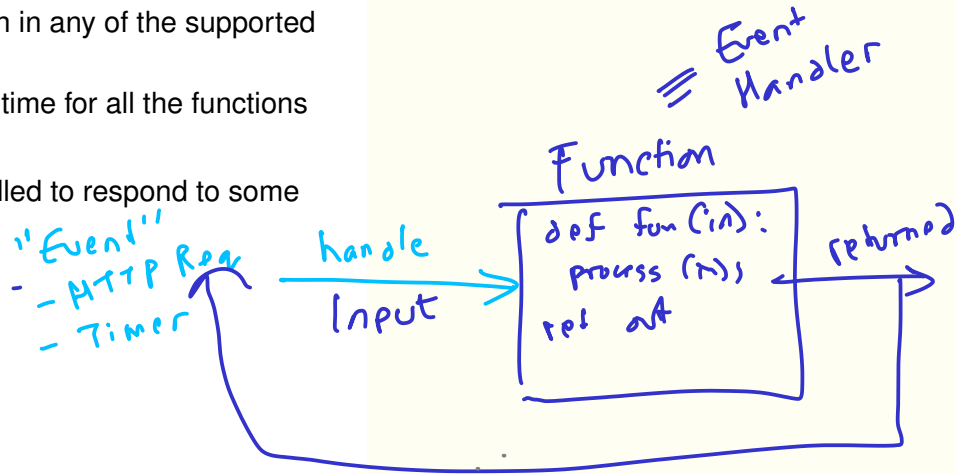# Infrastructure as a Service Pitfalls

Infrastructure is now software, but even that is too hard...

- Configuring and managing complete software stack
- OS and software upgrades, security patches, . . .
- Monitoring and logging
- Auto-scaling, redundancy, geo-replication, . . .

Application

Function Interfa

Langus Runtime

Libr

Container

Guest

VM

PHY Hardware

# Cloud Functions

- Cloud platform runs *functions* on behalf of user
- `ret-val function-name(arguments) {...}`
- Users provide the function implementation in any of the supported languages
- Cloud platform provides the language runtime for all the functions
- Python, Javascript, Java, Go, ...
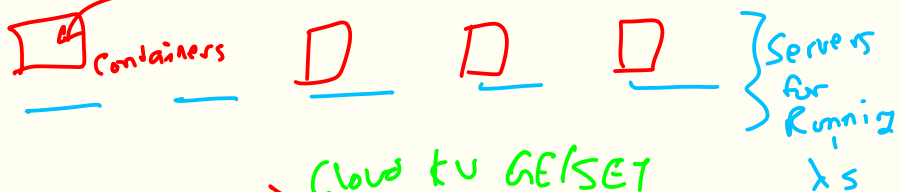- Event-driven programming: function is called to respond to some event



≡ Event Handler

Function

"Event"
- HTTP Req
- Timer

handle
Input

```
def fun (in):
    process (in);
    ret out
```

returned

# Functions as a Service

- AWS Lambda ($\sim$ 2015)
- FaaS $\equiv$ "Serverless"
- No need to manage VMs explicitly
- Functions run inside sandboxed runtimes

**FaaS is a new programming model :**

- Functions are "pure functions" and stateless
- Each function invocation is in a new sandboxed environment
- All state is stored in cloud storage services (like S3 buckets)

$\lambda$: functions in FP

evt trigger $\longrightarrow \lambda$

containers

Servers for Running $\lambda$s

Cloud to GET/SET

$X$ count = 0 $\longrightarrow$
```
def handle (fn):
    if (count)
    X += count;
    return
```

# Pricing

- Function invocations are charged
- Pay in proportion to usage
- Application is not charged if not used!
- Resource limits on CPU/memory utilization
- Current pricing is very low: $ 10^{-7}$ per request

*NOT allocation*

*Rare / Sporadic / Bursty events*

# Function triggers

1. Explicit HTTP requests
2. Changes in cloud storage state (new bucket is added etc)
3. Queuing service (new items added in queue)
4. Publish/Subscribe changes

.

1 Calling a function:

```
curl -X POST
"https://region-project-id.cloudfunctions.net/function-name"
-H "Content-Type:application/json"
--data '{"name":''foo''}'
```

2 Register function: gcloud functions deploy function-name
-runtime nodejs8 -trigger-http

3 NodeJS function to be called:

```
exports.helloHttp = (req, res) => {
  res.send(Hello ${escapeHtml(req.query.name ||
    req.body.name || 'World')}!);
};
```

*(handwritten annotations)*

≡ Simple Web Server

- Select a VM
- Configure/update OS
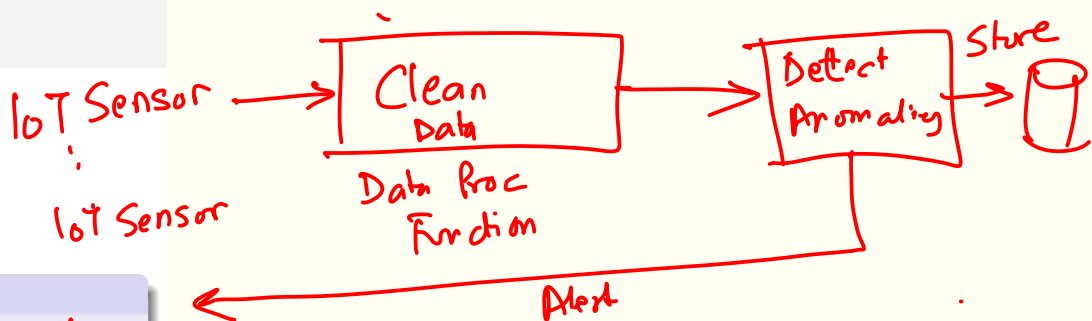- Setup Apache
- Deploy core

# Use-cases

- Web services (static web hosting, blogs, etc.)
- Data processing
- API serving

**New use cases:**

- Highly parallel video processing
- Linear algebra
- ML training

*low current costs*

$\infty$ scalabili/ Paralleliz

$10^6$ $\lambda$'s in parallel

IoT Sensor → Clean Data

Data Proc Function

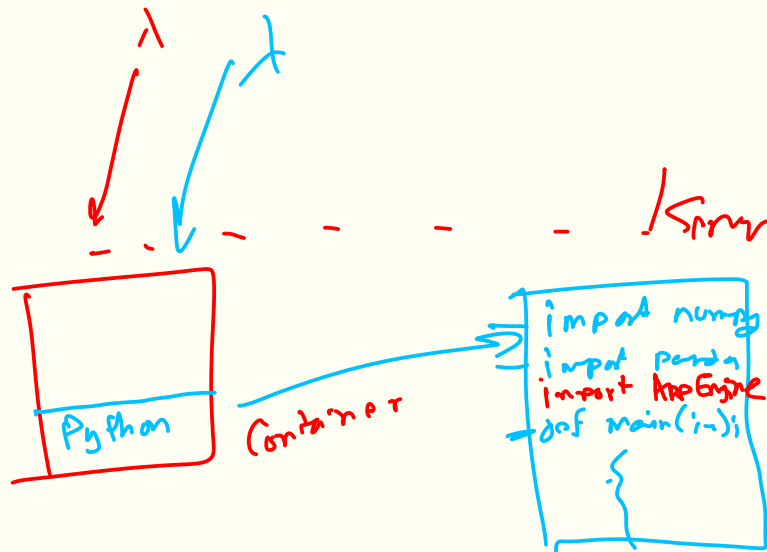IoT Sensor

Detect Anomalies → Store

Alert

# Implementation

1. On trigger, cloud finds a free server with enough resources
2. Instantiates the language runtime on server
3. Example: If user provides Python code, then launch the Python interpreter, import all dependencies, and run function
4. Once function exits, destroy the execution environment

**Cold-start problem: Function invocation latency can be high**

- Frameworks like OpenWhisk, OpenFaaS can be used to setup FaaS on local environments.

keep alive
container i-mem
future invocations "warm start"

λ

5/10 ms

Srrv

Python    Container

import numpy
import pandas
import the Engine
def main(in);
{

# Performance Issues

- Limited lifetimes ($\sim$ 15–30 minutes)
- Storage-bound: excess communication and storage overhead
- Startup-latency: $\sim$ 100 ms . Bad for response times
- No state: Can't use any caching or batching

# History

- CGI (Common Gateway Interface)
- Google AppEngine

# Pitfalls and Challenges

- Vendor lock-in.
- High costs for applications with steady workloads
- Restricted programming and manageability flexibility
- Heterogenous hardware: using GPUs?