# The State of Serverless Computing
## or, Fixing Dysfunction-as-a-Service

Chenggang Wu

RISE Lab, UC Berkeley

QCon New York 06/24/2019

UC Berkeley: CS 61A
Fall 2018

# Making Programmers Productive

Key Question: Where will code be run? In the ☁️!

# Background: Serverless Computing

# What is serverless computing?

*Serverless computing is a programming abstraction that enables users to upload programs, run them at any scale, and pay only for resources used.*

# Functions-as-a-Service (FaaS)

- AWS Lambda, Google Cloud Functions, OpenWhisk (IBM), Azure Functions, OpenLambda, OpenFaaS, kNative...
- Optimized for simplicity – register functions, enable triggers, and scale transparently



**Function code** Info

Code entry type
Edit code inline ▼

Runtime
Python 3.6 ▼

```
File   Edit   Find   View   Go   Tools   Window

▲     ▼ 📁 election        ⚙ ▾   🗐   lambda_function ×   ⊕

Environment       ◆ lambda_function.py   1  |
                                         2  import boto3
                                         3  import pickle
                                         4  import random
                                         5  import time
                                         6
                                         7  client = boto3.resource('dynamodb')
                                         8  table = client.Table('vsreekanti')
                                         9  candidates = client.Table('candidates')
                                        10  THRESHOLD = 50
                                        11
                                        12  def main(thisid):
                                        13      print("My invocation's id is: " + thisid)
                                        14      thisid = int(thisid)
                                        15      vote_round_count = 0
                                        16      am_leader = False
                                        17
                                        18      while True:
                                        19          time.sleep(.25)
```

## Add triggers

Choose a trigger from the list below to add it to your function.

API Gateway

AWS IoT

Alexa Skills Kit

Alexa Smart Home

Application Load Balancer

CloudFront

🔑

Add triggers from the list on the left

# Academic Interest in Serverless

**Occupy the Cloud: Distrib...**

Eric Jonas, Qifan Pu, Shivaram Ve...
University of ...
{jonas, qifan, shivaram, ist...

**ABSTRACT**

Distributed computing remains inaccessible to a large number of users, in spite of many open source platforms and extensive commercial offerings. While distributed computation frameworks have moved beyond a simple map-reduce model, many users are still left to struggle with complex cluster management and configuration tools, even for running simple embarrassingly parallel jobs. We argue that stateless functions represent a viable platform for these users, eliminating cluster management overhead, fulfilling the promise of elasticity. Furthermore, using our prototype implementation, PyWren, we show that this model is general enough to implement a number of distributed computing models, such as BSP efficiently. Extrapolating from recent ...
vent of disaggregated stor...
a natural fit for data proc...

**CCS CONCEPTS**

• Computer systems org...
puting methodologies —...

**KEYWORDS**

Serverless, Distributed Computing, AWS Lambda, PyWren

**ACM Reference Format:**
Eric Jonas, Qifan Pu, Shivaram Venkataraman, Ion Stoica, Benjamin Recht University of California, Berkeley {jonas, qifan, shivaram, istoica, brecht} @eecs.berkeley.edu. 2017. Occupy the Cloud: Distributed Computing for the 99%. In *Proceedings of SoCC '17, Santa Clara, CA, USA, September 24–27, 2017*, 7 pages.
https://doi.org/10.1145/3127479.3128601

**1 INTRODUCTION**

Despite a decade of availability, the twin promises of scale and elasticity [2] remain out of reach for a large number of cloud computing users. Academic and commercially-successful platforms (Apache Hadoop, Apache Spark) with tremendous corporate backing (Amazon, Microsoft, Google) still present high barriers to entry for the average data scientist or scientific computing user. In fact, taking advantage of elasticity remains challenging for even sophisticated users, as the majority of these frameworks were designed to first...

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
*SoCC '17, September 24–27, 2017, Santa Clara, CA, USA*
© 2017 Copyright held by the owner/author(s). Publication rights licensed to Association for Computing Machinery.
ACM ISBN 978-1-4503-5028-0/17/09...$15.00

---

**Encoding, Fa...**
**Low-Latency Video Processing U...**

Sadjad Fouladi §, Riad S. Wa...
Karthikeyan Vasuki Balasubramaniam ...
Anirudh Sivaraman ▥, Georg...
*Stanford University §, University of California San ...*

**Abstract**

We describe ExCamera, a system that can edit, transform, and encode a video, including 4K and VR material, with low latency. The system makes two major contributions.

First, we designed a framework to run general-purpose parallel computations on a commercial "cloud function" service. The system starts up thousands of threads in seconds and manages inter-thread communication.

**1 Introduction**

The pace of data analysis and processing has advanced rapidly, enabling new applications over large data sets. Providers use data-parallel frameworks such as MapReduce [8], Hadoop [12], and Spark [32] to analyze a variety of data streams: click logs, user ratings, medical records, sensor histories, error logs, and financial transactions.

Yet video, the largest source of data transiting the Internet [6], has proved one of the most vexing to analyze and manipulate. Users increasingly seek to apply complex computational pipelines to video content. Examples include video editing, scene understanding, object recognition and classification, and compositing. Today, these jobs often take hours, even for a short movie.

There are several reasons that interactive video-processing applications have yet to arrive. First, video jobs take a lot of CPU. In formats like 4K or virtual reality, an hour of video will typically take more than 30 CPU-hours to process. A user who desires results in a few seconds would need to invoke thousands of threads of execution—even assuming the job can be parallelized into thousands of pieces.

Second, existing video encoders do not permit fine-grained parallelism. Video is generally stored in compressed format, but unlike the per-record compression used by data-parallel frameworks [2], video compres-

---

## Cloud Prog
## A Berkeley View

Eric Jonas          Johann Schl...
Anurag Khandelwal   Qifan ...
Karl Krauth         Neeraja Ya...
                    Ion St...

server

predict these issues are solvable and that...
computing.

## Contents

1  Introduction to Serverless Com...
2  Emergence of Serverless Compu...
   2.1  Contextualizing Serverless Con...
   2.2  Attractiveness of Serverless Co...
3  Limitations of Today's Serverle...
   3.1  Inadequate storage for fine-gra...
   3.2  Lack of fine-grained coordinati...
   3.3  Poor performance for standard...
   3.4  Predictable Performance . . .
4  What Serverless Computing Sh...
   4.1  Abstraction challenges . . . .
   4.2  System challenges . . . . . . .
   4.3  Networking challenges . . . .
   4.4  Security challenges . . . . . .
   4.5  Computer architecture challeng...
5  Fallacies and Pitfalls

This article is published under a Creative Commons (http://creativecommons.org/licenses/by/3.0/), which permits distribution and reproduction in any medium as well as allowing derivative works, provided that you attribute the original work to the author(s) and CIDR 2019.

---

**Serverless Compu...**

Joseph M. Hellerstein, Jose Fa...

{hellerstein,jmfa...

**ABSTRACT**

Serverless computing offers the potential to...
an autoscaling, pay-as-you go manner. In t...
critical gaps in first-generation serverless co...
its autoscaling potential at odds with domi...
computing: notably data-centric and distri...
also open source and custom hardware. Pu...
make current serverless offerings a bad fit...
and particularly bad for data systems inno...
pinpointing some of the main shortfalls of...
chitectures, we raise a set of challenges we...

available to the general public, managed as...

Despite that potential, we have yet to ha...
in radical ways. The cloud today is largely...
platform for standard enterprise data servic...
creative developers need programming fra...
them to leverage the cloud's power.

New computing platforms have typically...
programming languages and environments...
is difficult to identify the new programming...
cloud. And whether cause or effect, the res...
in practice: the majority of cloud services ar...
easier-to-administer clones of legacy enter...
object storage, databases, queueing systems...
Multitenancy and administrative simplicity a...
able goals, and some of the new services hav...
in their own right. But this is, at best, only...
offered by millions of cores and exabytes of...

Recently, public cloud vendors have beg...
gramming interfaces under the banner of *ser...
interest is growing. Google search trends sho...
term "serverless" recently matched the histo...
of the phrase "Map Reduce" or "MapReduce...
also been a significant uptick in attention to t...
from the research community [13, 6, 27, 14].

with highlights including that AWS Lambda adopts a bin-packing-like strategy to maximize VM memory utilization, that severe contention between functions can arise in AWS and Azure, and that Google had bugs that allow customers to use resources for free.

**1 Introduction**

Cloud computing has allowed backend infrastructure maintenance to become increasingly decoupled from application development. Serverless computing (or function-as-a-service, FaaS) is an emerging application deployment architecture that completely hides server management from tenants (hence the name). Tenants receive minimal access to an application's runtime configuration. This allows tenants to focus on developing their functions — small applications dedicated to specific tasks. A function usually executes in a dedicated *function instance* (a container or other kind of sandbox) with restricted resources such as CPU time and memory. Unlike virtual machines (VMs) in more traditional infrastructure-as-a-service (IaaS) platforms, a function instance will be launched only when the function is invoked and is put to sleep immediately after handling a request. Tenants are charged on a per-invocation basis, without paying for unused and idle resources.

---

**Peeking Behind the Curtains of Serverless Platforms**

Liang Wang [1], Mengyuan Li [2], Yinqian Zhang [2], Thomas Ristenpart[3], Michael Swift[1]

[1]UW-Madison, [2]Ohio State University, [3]Cornell Tech

**Abstract**

Serverless computing is an emerging paradigm in which an application's resource provisioning and scaling are managed by third-party services. Examples include AWS Lambda, Azure Functions, and Google Cloud Functions. Behind these services' easy-to-use APIs are opaque, complex infrastructure and management ecosystems. Taking on the viewpoint of a serverless...

Serverless computing originated as a design pattern for handling low duty-cycle workloads, such as processing in response to infrequent changes to files stored on the cloud. Now it is used as a simple programming model for a variety of applications [14, 22, 42]. Hiding resource management from tenants enables this programming model, but the resulting opacity hinders adoption for many potential users, who have expressed concerns about: security in terms of the quality of isolation, ...23, 35, 37, 40]; the need to ...ent to improve application ...28, 40]; and the ability ...performance [10–12, 29– ...es made to shed light on ...ent and security [33, 34], ...es, as we will show, fail to ...

We therefore perform the most in-depth study of resource management and performance isolation to date in three popular serverless computing providers: AWS Lambda, Azure Functions, and Google Cloud Functions (GCF). We first use measurement-driven approaches to partially reverse-engineer the architectures of Lambda and Azure Functions, uncovering many undocumented details. Then, we systematically examine a series of issues related to resource management: how quickly function instances can be launched, function instance placement strategies, function instance reuse, and more. Several security issues are identified and discussed.[1] We further explore how CPU, I/O and network bandwidth are allocated among functions and the ensuing performance implications. Last but not least, we explore whether all resources are properly accounted for, and report on two resource accounting bugs that allow tenants to use extra resources for free. Some highlights of our results include:

- AWS Lambda achieved the best scalability and the lowest coldstart latency (the time to provision a new function instance), followed by GCF. But

[1]We responsibly disclosed our findings to related parties before this paper was made public.

---

# Industrial Interest in Serverless
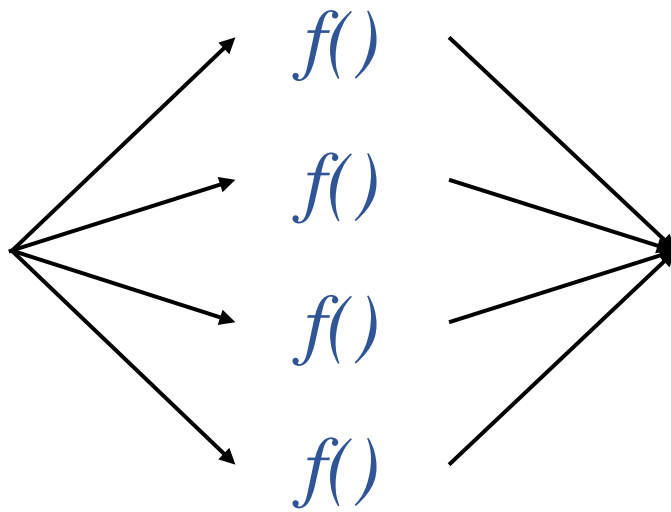
# What is FaaS good at today?



Embarrassingly parallel tasks
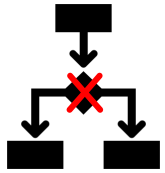
Workflow orchestration

Wait! What about...

- There are other serverless services, too!
  - e.g., Google Cloud Dataflow, AWS Athena, Snowflake...
- Often referred to as Backend-as-a-Service (BaaS)

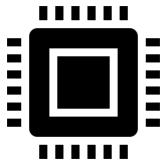We're primarily interested in generality.

# Limitations on FaaS Today

Limited execution lifetimes

No inbound network connections

IO is a bottleneck

No specialized hardware

# But that's ~~okay~~ NOT: Everything is ~~functional~~ NOT!

- Functional programs don't have side effects or mutable state!



- And it *is called* AWS Lambda

# Dysfunction-as-a-Service

- FaaS is not designed for functional programming because real applications share state

$$f(g(x))$$

- FaaS is poorly suited for all of these

# Quantifying The Pain of FaaS

How FaaS Disappoints Famous Computer Scientists

# Even Functional Programming is Slow!

$$f(g(x))$$



Median and 99th percentile latencies for composing two arithmetic functions on AWS Lambda.

# Even Functional Programming is Slow!

$$f(g(x))$$



Median and 99th percentile latencies for composing two arithmetic functions on AWS Lambda.

# Even Functional Programming is Slow!

$$f(g(x))$$



Median and 99th percentile latencies for composing two arithmetic functions on AWS Lambda.

# Shared Mutable State

# Shared Mutable Storage

# ~~Serverless~~ Serverless Storage

Autoscaling

Tradeoff!

Low Latency

# (In)Consistency Guarantees



Shared
Counter

# No Inbound Network Connections

Enables Process
Migration

Easy Fault
Tolerance

# Indirect Communication



Write      Read

We can
fix that!

# A Platform for Stateful Serverless Computing

# Background: Anna

- High performance across orders of magnitude in scale
  - ✓ 10x faster than Redis/Cassandra in a geo-distributed deployment
- Autoscaling & cost-efficient
  - ✓ 500x faster than Amazon DynamoDB for the same cost

Chenggang Wu, Jose Faleiro, Yihan Lin, and Joseph M. Hellerstein. "Anna: A KVS for Any Scale." *IEEE Transactions on Knowledge and Data Engineering* (2019).
Chenggang Wu, Vikram Sreekanti, and Joseph M. Hellerstein. "Autoscaling Tiered Cloud Storage in Anna." *Proceedings of the VLDB Endowment* 12, no. 6 (2019): 624-638.

# Fluent: FaaS-over-Anna

- Maintain disaggregation of compute & state
- Make serverless a viable option for stateful applications
- Use Anna for both storage and communication



Anna

# Fluent: FaaS-over-Anna



Network Boundary

Anna

# Logical disaggregation with physical colocation

# Fluent: FaaS-over-Anna

Network
Boundary

Anna

# Key Idea: Caching

- Enable low-latency data access by caching data close to code execution

- Communication (and composition) is achieved via a fast-path on top of KVS puts and gets

# Challenge: Cache Consistency

- tl;dr: we can provide a variety of coordination-avoiding consistency modes – which is better than S3 or DynamoDB!
- This is done by encapsulating program state in lattices

# Lattice

- Data structure that accepts incoming update in a way that is associative, commutative, and idempotent (ACI).

- Achieves eventual replica convergence

# Causal Consistency

- Strongest consistency level that doesn't require coordination

- Causally-related updates will be revealed in an order that respects causality

- In addition, guarantee
  - Repeatable read
  - Atomic visibility

# Function Composition, Revisited

$$f(g(x))$$



Median and 99th percentile latencies for composing two arithmetic functions on AWS Lambda and Fluent.

# Case Study: Prediction Serving

# Prediction Serving

- Generate predictions from pretrained machine learning models



- At first blush, a great fit for serverless infrastructure

# Background: SqueezeNet

- State-of-the-art image classification model (developed at Berkeley!)



INPUT     CONVOLUTION + RELU     POOLING     CONVOLUTION + RELU    POOLING     FLATTEN     FULLY CONNECTED    SOFTMAX

CAR — TRUCK — VAN — BICYCLE

FEATURE LEARNING         CLASSIFICATION

# Prediction Serving



Median and 99th percentile latencies for SqueezeNet on Fluent and AWS SageMaker.

# The Future of Cloud Programming

# Looking Back: Disappointed Computer Scientists

- Functional programming is slow
- Communication through slow storage
- Poor consistency guarantees

# Making FaaS Functional

- Embrace state

- Easy things become better

- Hard things become easy

$$f(g(x))$$

A step on our road towards a programmable cloud.

# Our Vision

- Serverless will change the way that we write software and the way that programming infrastructure works

- Cloud-native programming models
  - Enable users to take advantage of millions of cores and petabytes of RAM

# Moving Forward from FaaS



Building
Developer
Tools

UC Berkeley: CS 61A
Fall 2018

The State of Serverless Computing
QCon New York 06/24/2019

# Moving Forward from FaaS

**Building**
Developer
Tools

**Developing**
Autoscaling
Policy

# Moving Forward from FaaS

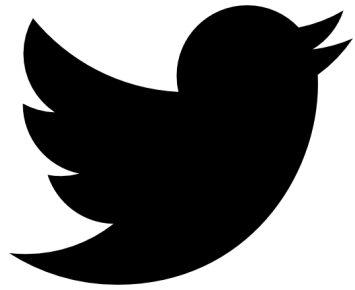

**Building**
Developer
Tools

**Developing**
Autoscaling
Policy

**Designing**
SLOs & SLAs

# Thanks!

@cgwu0530

cgwu@berkeley.edu

fluent-project/fluent