## Making Serverless Computing More Serverless

Zaid Al-Ali, Sepideh Goodarzy, Ethan Hunter, Sangtae Ha, Richard Han, Eric Keller University of Colorado Boulder

> Eric Rozner IBM Research



University of Colorado Boulder



\* Views not representative of IBM policy, products, or strategies

# Serverless Background

• Serverless offerings are widespread today





**IBM Cloud Functions** 





Amazon Lambda

Apache OpenWhisk



OpenFaaS



fission



### Today's serverless abstraction

#### Making Serverless Computing More Serverless

#### FaaS, Lamda, OpenWhisk, ...

#### Break limitations of single server

Zaid Al-Ali<sup>†</sup>, Sepideh Goodarzy<sup>†</sup>, Ethan Hunter<sup>†</sup>, Sangtae Ha<sup>†</sup>, Richard Han<sup>†</sup>, Eric Keller<sup>†</sup>, Eric Rozner<sup>\*</sup> <sup>†</sup>University of Colorado Boulder; \*IBM Research

Abstract—In serverless computing, developers define a function to handle an event, and the serverless framework horizontally scales the application as needed. The downside of this function-based abstraction is it limits the type of application supported and places a bound on the function to be within the physical resource limitations of the server the function executes on. In this paper we propose a new abstraction through memory or storage. The challenge, of course, is realizing a process-based serverless framework which can map our serverless process abstraction to an underlying, physically distributed infrastructure. To that end, we propose a new architecture called ServerlessOS to enable our vision and argue three key components are necessary to make our



This talk: propose new serverless abstraction and overview its design

## A new serverless abstraction

- Expand serverless beyond the bounds of FaaS
- Goals of our new abstraction:
  - Flexible enough to support general set of applications
  - Familiar to developers, operating systems, and admins
    - Easy to transition existing codebases to serverless
    - Same simplicity and scale-out as FaaS

## A new serverless abstraction

Flexible enough to support general set of applications
Familiar to developers, operating systems, and admins
Easy to transition existing codebases to serverless
Same simplicity and scale-out as FaaS





Flexible enough to support general set of applications

Multiple threads, I/O via sockets, persist state, ...

Familiar to developers, operating systems, and admins

Abstraction already used today in non-serverless

Easy to transition existing codebases to serverless

Pool of CPU, I/O, memory, storage: server  $\rightarrow$  datacenter

Same simplicity and scale-out as FaaS

Challenge: map serverless process abstraction to underlying physically distributed architecture

- Goal: provide seamless, scale-out process abstraction
- This talk: high-level outline of our ServerlessOS vision



- Goal: provide seamless, scale-out process abstraction
- This talk: high-level outline of our ServerlessOS vision



 Break coupling between process & underlying physical server resources

### Fluid multi-resource disaggregation



Disaggregation: decouple process's resources from single server

 Break coupling between process & underlying physical server resources

### Fluid multi-resource disaggregation



Decouple memory, compute, I/O to increase flexibility

 Break coupling between process & underlying physical server resources

#### Fluid multi-resource disaggregation



Acessing remote memory incurs much higher overhead than local memory

 Break coupling between process & underlying physical server resources



Fluidity: allow process to move to data when more efficient

 Break coupling between process & underlying physical server resources

#### Fluid multi-resource disaggregation



Fluidity: enable process to exploit locality to improve performance

### Fluidity over multiple resources



Already provided by prior works (RamCloud, DSM, InfiniSwap, ...)





Move processing to data or other server with more compute resources (Initial results show 2-3x speedup over a DSM scheme)

Decouple device that captured I/O from device that will process I/O. Additionally, move I/O to more bandwidth. (CPU fluidity can move processing with socket)

- Goal: provide seamless, scale-out process abstraction
- This talk: high-level outline of our ServerlessOS vision



### Fine-grained live orchestration layer

 Monitor, allocate, and optimize run-time performance by automatically assigning, migrating, or scaling workloads

![](_page_15_Figure_2.jpeg)

### Fine-grained live orchestration layer

 Monitor, allocate, and optimize run-time performance by automatically assigning, migrating, or scaling workloads

![](_page_16_Figure_2.jpeg)

- Goal: provide seamless, scale-out process abstraction
- This talk: high-level outline of our ServerlessOS vision

![](_page_17_Figure_3.jpeg)

## **Coordinated Isolation**

![](_page_18_Figure_1.jpeg)

ServerlessOS: extend isolation across multiple servers in coordinated fashion

# **Coordinated Isolation**

![](_page_19_Figure_1.jpeg)

- Goal: provide seamless, scale-out process abstraction
- This talk: high-level outline of our ServerlessOS vision

![](_page_20_Figure_3.jpeg)

## Conclusions

- New abstraction for serverless: a seamless, scale-out process
- High-level overview of ServerlessOS architecture
  - Fluid multi-resource disaggregation
  - Fine-grained live orchestration layer
  - Coordinated isolation
- Complementary to current serverless techniques
- Next steps: refine design, build prototype, conquer the world!
- Thanks! mailto: erozner@us.ibm.com