Hello!

# CS 744: PYWREN

Shivaram Venkataraman

Fall 2020

# ADMINISTRIVIA

Project checkins due Nov 20th → *Friday*

In-class project presentations

Dec 8th and Dec 10th → *5 min talks about your project*

Project grade breakdown → *Canvas soon!*

Intro: 5%

Mid-semester checkin: 5%

Presentation: 10%

Final Report: 10%

*Tonight deadline for submitting regrade requests! for Midterm 1*

# NEW HARDWARE MODELS

Implications → Society
of Big Data
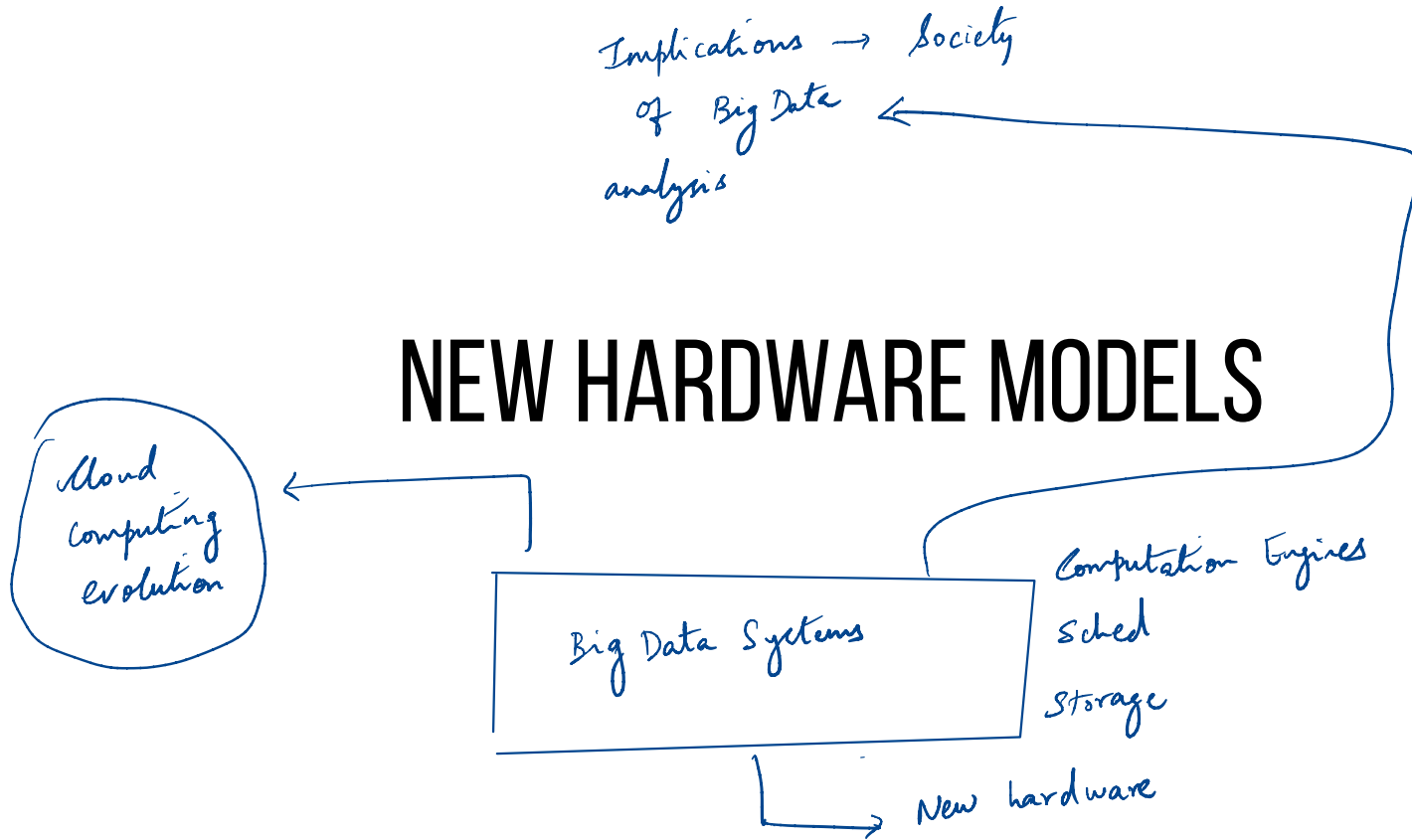analysis
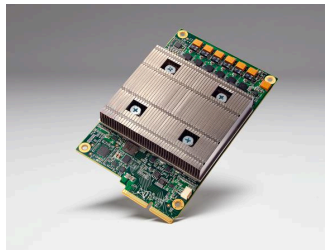
Cloud
computing
evolution

Big Data Systems

Computation Engines
Sched
Storage

New hardware

Serverless Computing

Compute Accelerators

Infiniband Networks

Non-Volatile Memory

# SERVERLESS COMPUTING

No servers ??

# MOTIVATION: USABILITY

Azure, Google etc.

Data Scientist ≡ Analysis

**What instance type?**

**What base image?**

**How many to spin up?**

**What price? Spot?**

Makes it difficult to use the cloud

EC2Instances.info Easy Amazon EC2 Instance Comparison

EC2  RDS

Region: US East (N. Virginia) ▾   Cost: Hourly ▾   Reserved: 1 yr - No Upfront ▾   Columns ▾   Compare Selected   Clear Filters

Filter: Min Memory (GB): [ ]   Compute Units: [ ]   Storage (GB): [ ]

| Name | API Name | Memory | Compute Units (ECU) | vCPUs | Storage | Arch | Network Performance | EBS Optimized: Max Bandwidth | VPC Only | Linux On Demand cost | Linux Reserved cost | Windows On Demand cost | Windows Reserved cost |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Cluster Compute Eight Extra Large | cc2.8xlarge | 60.5 GB | 88 units | 32 vCPUs | 3360.0 GB (4 * 840.0 GB) | 64-bit | 10 Gigabit | N/A | No | $2.000 hourly | $1.090 hourly | $2.570 hourly | $1.336 hourly |
| Cluster GPU Quadruple Extra Large | cg1.4xlarge | 22.5 GB | 33.5 units | 16 vCPUs | 1680.0 GB (2 * 840.0 GB) | 64-bit | 10 Gigabit | N/A | No | $2.100 hourly | unavailable | $2.600 hourly | unavailable |
| T2 Nano | t2.nano | 0.5 GB | Burstable | 1 vCPU | 0 GB (EBS only) | 64-bit | Low | N/A | Yes | $0.006 hourly | $0.005 hourly | $0.009 hourly | $0.007 hourly |
| T2 Micro | t2.micro | 1.0 GB | Burstable | 1 vCPUs | 0 GB (EBS only) | 32/64-bit | Low to Moderate | N/A | Yes | $0.013 hourly | $0.009 hourly | $0.018 hourly | $0.014 hourly |
| T2 Small | t2.small | 2.0 GB | Burstable | 1 vCPUs | 0 GB (EBS only) | 32/64-bit | Low to Moderate | N/A | Yes | $0.026 hourly | $0.018 hourly | $0.036 hourly | $0.032 hourly |
| T2 Medium | t2.medium | 4.0 GB | Burstable | 2 vCPUs | 0 GB (EBS only) | 64-bit | Low to Moderate | N/A | Yes | $0.052 hourly | $0.036 hourly | $0.072 hourly | $0.062 hourly |
| T2 Large | t2.large | 8.0 GB | Burstable | 2 vCPUs | 0 GB (EBS only) | 64-bit | Low to Moderate | N/A | Yes | $0.104 hourly | $0.072 hourly | $0.134 hourly | $0.106 hourly |
| M4 Large | m4.large | 8.0 GB | 6.5 units | 2 vCPUs | 0 GB (EBS only) | 64-bit | Moderate | 450.0 Mbps | Yes | $0.120 hourly | $0.083 hourly | $0.246 hourly | $0.184 hourly |
| M4 Extra Large | m4.xlarge | 16.0 GB | 13 units | 4 vCPUs | 0 GB (EBS only) | 64-bit | High | 750.0 Mbps | Yes | $0.239 hourly | $0.164 hourly | $0.491 hourly | $0.366 hourly |
| M4 Double Extra Large | m4.2xlarge | 32.0 GB | 26 units | 8 vCPUs | 0 GB (EBS only) | 64-bit | High | 1000.0 Mbps | Yes | $0.479 hourly | $0.329 hourly | $0.983 hourly | $0.735 hourly |
| M4 Quadruple Extra Large | m4.4xlarge | 64.0 GB | 53.5 units | 16 vCPUs | 0 GB (EBS only) | 64-bit | High | 2000.0 Mbps | Yes | $0.958 hourly | $0.658 hourly | $1.966 hourly | $1.469 hourly |
| M4 Deca Extra Large | m4.10xlarge | 160.0 GB | 124.5 units | 40 vCPUs | 0 GB (EBS only) | 64-bit | 10 Gigabit | 4000.0 Mbps | Yes | $2.394 hourly | $1.645 hourly | $4.914 hourly | $3.672 hourly |
| M4 16xlarge | m4.16xlarge | 256.0 GB | 188 units | 64 vCPUs | 0 GB (EBS only) | 64-bit | 20 Gigabit | 10000.0 Mbps | Yes | $3.830 hourly | $2.632 hourly | $7.862 hourly | $5.875 hourly |
| C4 High-CPU Large | c4.large | 3.75 GB | 8 units | 2 vCPUs | 0 GB (EBS only) | 64-bit | Moderate | 500.0 Mbps | Yes | $0.105 hourly | $0.078 hourly | $0.193 hourly | $0.170 hourly |
| C4 High-CPU Extra Large | c4.xlarge | 7.5 GB | 16 units | 4 vCPUs | 0 GB (EBS only) | 64-bit | High | 750.0 Mbps | Yes | $0.209 hourly | $0.155 hourly | $0.386 hourly | $0.339 hourly |
| C4 High-CPU Double Extra Large | c4.2xlarge | 15.0 GB | 31 units | 8 vCPUs | 0 GB (EBS only) | 64-bit | High | 1000.0 Mbps | Yes | $0.419 hourly | $0.311 hourly | $0.773 hourly | $0.679 hourly |
| C4 High-CPU Quadruple Extra Large | c4.4xlarge | 30.0 GB | 62 units | 16 vCPUs | 0 GB (EBS only) | 64-bit | High | 2000.0 Mbps | Yes | $0.838 hourly | $0.621 hourly | $1.546 hourly | $1.357 hourly |
| C4 High-CPU Eight Extra Large | c4.8xlarge | 60.0 GB | 132 units | 36 vCPUs | 0 GB (EBS only) | 64-bit | 10 Gigabit | 4000.0 Mbps | Yes | $1.675 hourly | $1.242 hourly | $3.091 hourly | $2.769 hourly |
| P2 Extra Large | p2.xlarge | 61.0 GB | 12 units | 4 vCPUs | 0 GB (EBS only) | 64-bit | High | 750.0 Mbps | No | $0.900 hourly | $0.684 hourly | $1.084 hourly | $0.868 hourly |
| P2 Eight Extra Large | p2.8xlarge | 488.0 GB | 94 units | 32 vCPUs | 0 GB (EBS only) | 64-bit | 10 Gigabit | 5000.0 Mbps | No | $7.200 hourly | $5.476 hourly | $8.672 hourly | $6.948 hourly |
| P2 16xlarge | p2.16xlarge | 732.0 GB | 188 units | 64 vCPUs | 0 GB (EBS only) | 64-bit | 20 Gigabit | 10000.0 Mbps | No | $14.400 hourly | $10.951 hourly | $17.344 hourly | $13.895 hourly |
| G2 Double Extra Large | g2.2xlarge | 15.0 GB | 26 units | 8 vCPUs | 60.0 GB SSD | 64-bit | High | 1000.0 Mbps | No | $0.650 hourly | $0.474 hourly | $0.767 hourly | $0.611 hourly |
| G2 Eight Extra Large | g2.8xlarge | 60.0 GB | 104 units | 32 vCPUs | 240.0 GB (2 * 120.0 GB SSD) | 64-bit | 10 Gigabit | N/A | No | $2.600 hourly | $1.896 hourly | $2.878 hourly | $1.979 hourly |
| X1 16xlarge | x1.16xlarge | 976.0 GB | 174.5 units | 64 vCPUs | 1920.0 GB SSD | 64-bit | 10 Gigabit | 5000.0 Mbps | No | $6.669 hourly | $4.579 hourly | $9.613 hourly | $7.523 hourly |
| X1 32xlarge | x1.32xlarge | 1952.0 GB | 349 units | 128 vCPUs | 3840.0 GB (2 * 1920.0 GB SSD) | 64-bit | 20 Gigabit | 10000.0 Mbps | No | $13.338 hourly | $9.158 hourly | $19.226 hourly | $15.046 hourly |
| R3 High-Memory Large | r3.large | 15.25 GB | 6.5 units | 2 vCPUs | 32.0 GB SSD | 64-bit | Moderate | N/A | No | $0.166 hourly | $0.105 hourly | $0.291 hourly | $0.238 hourly |
| R3 High-Memory Extra Large | r3.xlarge | 30.5 GB | 13 units | 4 vCPUs | 80.0 GB SSD | 64-bit | Moderate | 500.0 Mbps | No | $0.333 hourly | $0.209 hourly | $0.583 hourly | $0.428 hourly |
| R3 High-Memory Double Extra Large | r3.2xlarge | 61.0 GB | 26 units | 8 vCPUs | 160.0 GB SSD | 64-bit | High | 1000.0 Mbps | No | $0.665 hourly | $0.418 hourly | $1.045 hourly | $0.824 hourly |
| R3 High-Memory Quadruple Extra Large | r3.4xlarge | 122.0 GB | 52 units | 16 vCPUs | 320.0 GB SSD | 64-bit | High | 2000.0 Mbps | No | $1.330 hourly | $0.836 hourly | $1.944 hourly | $1.490 hourly |
| R3 High-Memory Eight Extra Large | r3.8xlarge | 244.0 GB | 104 units | 32 vCPUs | 640.0 GB (2 * 320.0 GB SSD) | 64-bit | 10 Gigabit | N/A | No | $2.660 hourly | $1.672 hourly | $3.500 hourly | $1.989 hourly |
| I2 Extra Large | i2.xlarge | 30.5 GB | 14 units | 4 vCPUs | 800.0 GB SSD | 64-bit | Moderate | 500.0 Mbps | No | $0.853 hourly | $0.424 hourly | $0.973 hourly | $0.565 hourly |
| I2 Double Extra Large | i2.2xlarge | 61.0 GB | 27 units | 8 vCPUs | 1600.0 GB (2 * 800.0 GB SSD) | 64-bit | High | 1000.0 Mbps | No | $1.705 hourly | $0.848 hourly | $1.946 hourly | $1.131 hourly |
| I2 Quadruple Extra Large | i2.4xlarge | 122.0 GB | 53 units | 16 vCPUs | 3200.0 GB (4 * 800.0 GB SSD) | 64-bit | High | 2000.0 Mbps | No | $3.410 hourly | $1.696 hourly | $3.891 hourly | $2.260 hourly |
| I2 Eight Extra Large | i2.8xlarge | 244.0 GB | 104 units | 32 vCPUs | 6400.0 GB (8 * 800.0 GB SSD) | 64-bit | 10 Gigabit | N/A | No | $6.820 hourly | $3.392 hourly | $7.782 hourly | $4.521 hourly |
| D2 Extra Large | d2.xlarge | 30.5 GB | 14 units | 4 vCPUs | 6000.0 GB (3 * 2000.0 GB) | 64-bit | Moderate | 750.0 Mbps | No | $0.690 hourly | $0.402 hourly | $0.821 hourly | $0.472 hourly |
| D2 Double Extra Large | d2.2xlarge | 61.0 GB | 28 units | 8 vCPUs | 12000.0 GB (6 * 2000.0 GB) | 64-bit | High | 1000.0 Mbps | No | $1.380 hourly | $0.804 hourly | $1.601 hourly | $0.885 hourly |
| D2 Quadruple Extra Large | d2.4xlarge | 122.0 GB | 56 units | 16 vCPUs | 24000.0 GB (12 * 2000.0 GB) | 64-bit | High | 2000.0 Mbps | No | $2.760 hourly | $1.608 hourly | $3.062 hourly | $1.690 hourly |
| D2 Eight Extra Large | d2.8xlarge | 244.0 GB | 116 units | 36 vCPUs | 48000.0 GB (24 * 2000.0 GB) | 64-bit | 10 Gigabit | 4000.0 Mbps | No | $5.520 hourly | $3.216 hourly | $6.198 hourly | $3.300 hourly |
| HI1. High I/O Quadruple Extra Large | hi1.4xlarge | 60.5 GB | 35 units | 16 vCPUs | 2048.0 GB (2 * 1024.0 GB SSD) | 64-bit | 10 Gigabit | N/A | No | $3.100 hourly | $2.255 hourly | $3.580 hourly | $2.260 hourly |
| High Storage Eight Extra Large | hs1.8xlarge | 117.0 GB | 35 units | 16 vCPUs | 48000.0 GB (24 * 2000.0 GB) | 64-bit | 10 Gigabit | N/A | No | $4.600 hourly | $2.574 hourly | $4.931 hourly | $2.961 hourly |

# ABSTRACTION LEVEL ?

*Snowflake*
↓
*Large* – *No servers*
*Medium* – *to configure*

*So far in this course*

*Input data*
↓ *Output*
*Query* →

| Application |
|---|

Logistic Regression
*or*
*SQL query*

*Spark* →
*RDD* →

| Compute Framework |
|---|

Spark *on a subset of machines*

*Slice up machines*

*VM* →

| Hardware |
|---|

Amazon EC2 –
CloudLab – *Server*
Private Cluster –
…

*Input*
↓ *Output*
*Query*

| Application |
|---|

*Serverless computing?*

| Compute Framework |
|---|

# STATELESS DATA PROCESSING

Intermediate state in Spark/MR was on local disk

local storage is ephemeral so intermediate state needs to be remote!

Storage/State

Compute

resource limits

Redis

| Key Value Store (Low Latency) |

Function, Dependencies DAG

| Function Scheduler |

| Blob Store S3 (High Bandwidth) |

| Container |
| Container |

| Container |
| Container |

# "SERVERLESS" COMPUTING → Provided by Cloud Provider
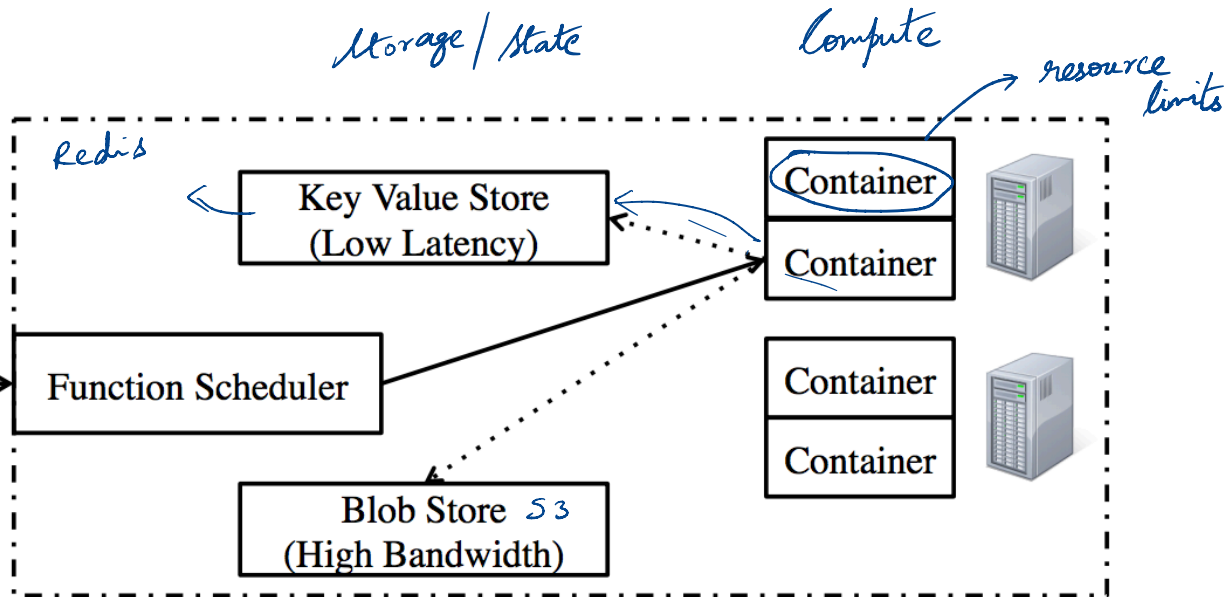
- Submit a function (lambda) to be executed

~~300~~ 900 seconds single-core → Time bound

512 MB in /tmp → Storage ⎫ Bounds

3GB RAM → Memory ⎭

Python, Java, node.js

some processing for each upload | upload

trigger

S3

$00002

cloud database

**Google** Cloud Platform

**CLOUD FUNCTIONS** ALPHA
A serverless platform for building event-based microservices

Microsoft Azure

Azure Functions
Process events with a serverless code architecture

# PYWREN API

python test.py

test.py

Language Integrated !!

```python
import pywren        ←
import numpy as np

def addone(x):        →  use libraries
    return x + 1          like scipy

wrenexec = pywren.default_executor()  ←
xlist = np.arange(10)
futures = wrenexec.map(addone, xlist)  →  map function similar to PySpark

print [f.result() for f in futures]
         ↳ block    similar   to   get   in   Ray API
```

Automatically captures dependencies
and ships them to the cloud
[cloud pickle ~ 2010]

The output is as expected:

```
[1, 2, 3, 4, 5, 6, 7, 8, 9, 10]
```

# PYWREN: HOW IT WORKS

Distributed key
value : get/put

Amazon

S3

Blob Storage

`future = runner.map(fn, data)`

fn    data

Invoke    fn    AWS
Lambda
API

get    get

fetch fn & data

computation

output

Containers

lambda

`future.result()`

poll

variable in
your laptop!

fetch

S3

your laptop

the cloud

# HOW IT WORKS

```
future = runner.map(fn, data)
```

func     data

Serialize func and data

Put on S3

Invoke Lambda

pull job from s3
download anaconda runtime
python to run code
pickle result
stick in S3
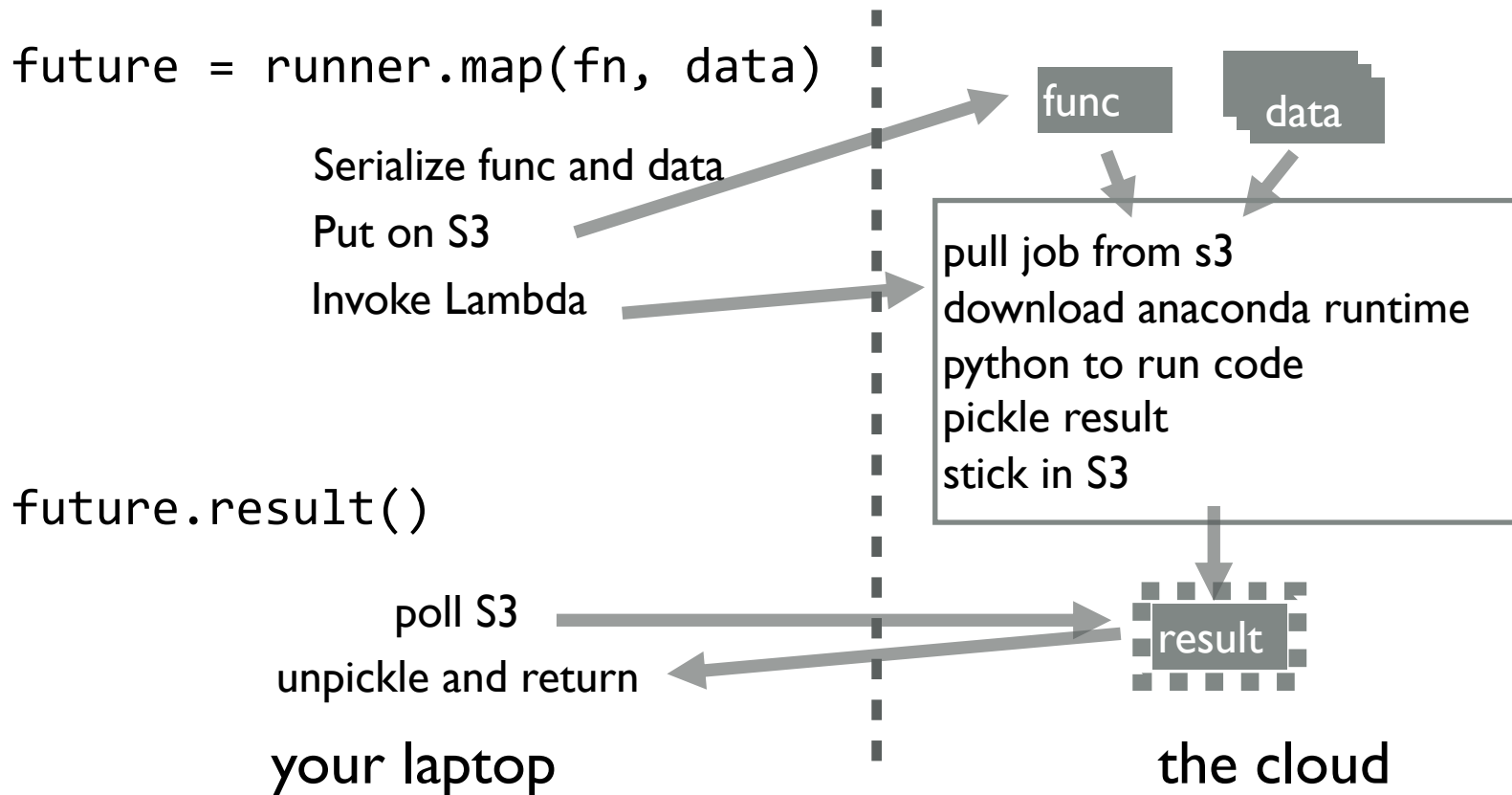
```
future.result()
```

poll S3

unpickle and return

result

your laptop            the cloud

# STATELESS FUNCTIONS: WHY NOW ?

What are the trade-offs ?

→ Need more network I/O

All the data is read over network!

→ But network BW is pretty good! Comparable to local SSD BW!

→ Bottleneck could be S3?

| Storage Medium | Write Speed (MB/s) |
|---|---|
| SSD on c3.8xlarge | 208.73 |
| SSD on i2.8xlarge | 460.36 |
| 4 SSDs on i2.8xlarge | 1768.04 |
| S3 | 501.13 |

# MAP AND REDUCE ?

Shuffle phase in MR is now being done using Redis

Sort benchmark
↳ Same as MapReduce paper

[0 - 100] key1, key2 . . . .
[101 - 200] key3 . . . .

Input Data

Output Data

Redis
Key-value
store - memory

bucket keys
into ranges

small files
not good for blob store like S3

# PARAMETER SERVERS

sparse models
↳ Ad click prediction

read
input

compute gradient

λ

λ

Use lambdas to run "workers"

Parameter server as a service ?

λ

λ

get

update

ML model stored

Redis
or
VMs etc.
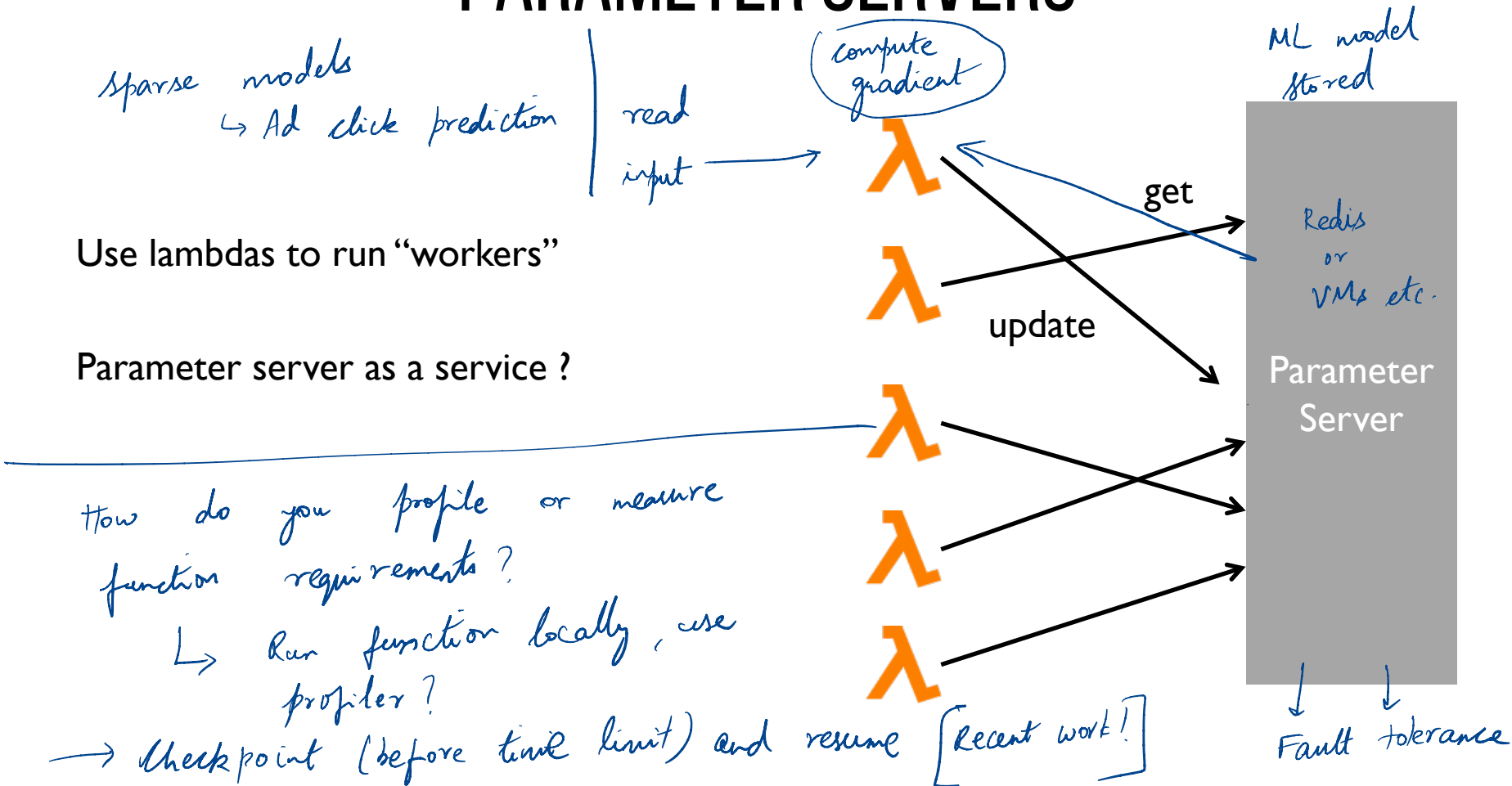
Parameter
Server

How do you profile or measure
function requirements ?
↳ Run function locally, use
profiler ?
→ Check point (before time limit) and resume [Recent work!]

λ

λ

Fault tolerance

# WHEN SHOULD WE USE SERVERLESS ?

| Yes! | Maybe not ? |
|---|---|

Use when we need elasticity

Use when you don't need fine grained comm. across workers

↳ Not all lambdas might be active at the same time!

Not use serverless when you need local state (actors)

Iterative workloads ⌐ might need state from prev. iteration

# SUMMARY

Motivation: Usability of big data analytics

Approach: Language-integrated cloud computing

Features

      - Breakdown computation into stateless functions

      - Schedule on serverless containers
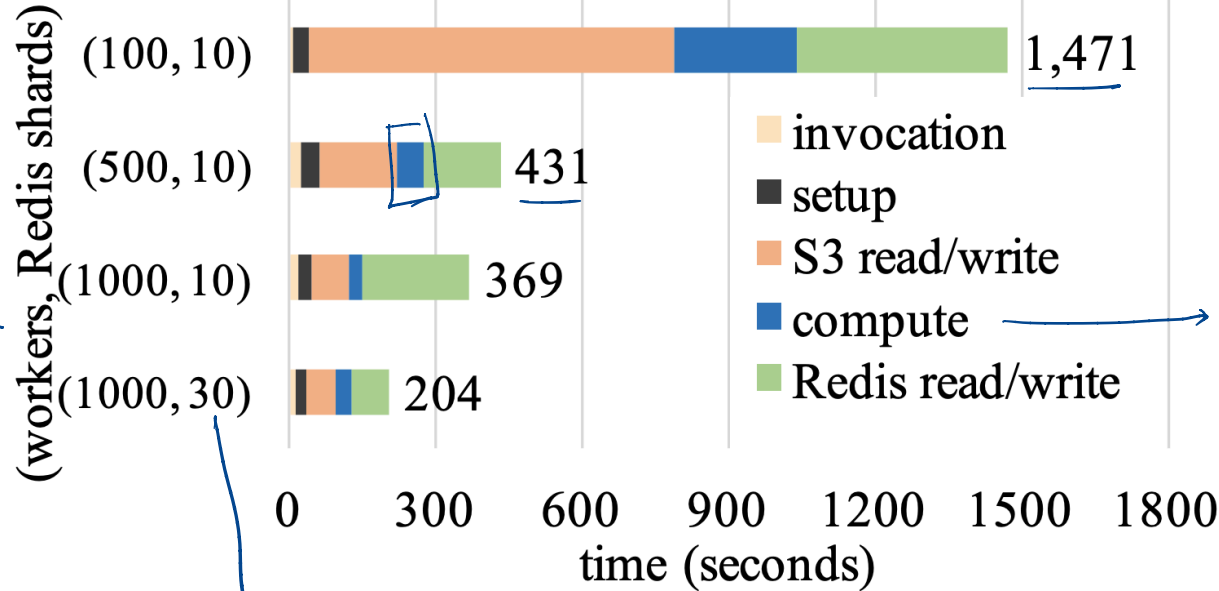
      - Use external storage for state management

Open question on scheduling, overheads

# DISCUSSION

https://forms.gle/PAMDKmwHepmPWDrBA

Scale workers, Redis independently
→ compute, storage

Increasing workers by 5x
!= 5x improvement

Hard to know how to choose num partitions

Compute is very short compared to I/O

more shards reduces time to read/write to Redis

(workers, Redis shards)

(100, 10) — 1,471
(500, 10) — 431
(1000, 10) — 369
(1000, 30) — 204

invocation
setup
S3 read/write
compute
Redis read/write

time (seconds)

0   300   600   900   1200   1500   1800

Consider you are a cloud provider (e.g., AWS) implementing support for serverless. What could be some of the new challenges in scheduling these workloads? How would you go about addressing them?

- Mapping lambda functions ⟶ machines
    How do we do this?

- Locality? Does one lambda talk to some Redis shard?
    Can we infer it?

- When to schedule a new container / when do we reuse?
    ↓
⇐ Need to find opt configuration? Use ML?

- Resource requirements are fixed! 900s, 1 core.
    upto 3 GB

# OPEN QUESTIONS

- Scalable scheduling: Low latency with large number of functions ?

- Debugging: Correlate events across functions ?

- Launch overheads: Fraction of time spent in setup (OpenLambda)

- Resource limits: 15 minute AWS Lambda (Oct 2018)

See you on Thursday!

Cold Starts

    ↳ App side  ⎤ Containers would
                   ⎥ be warm for 5 mins
    ↳ Sched. side  ⎦ ⇒ if you ran
                       one within 5mins

Azure - policy paper
ATC 2020

EBS | Block Service

mount
/ home
write
VM
unmount

EBS

mount   unmount

lambda         lambda

Memory
512 MB to
3 GB