

CSCI 381/780

Cloud Computing

Virtualization 3:

Serverless Computing

Jun Li
Queens College



Why care?

Rapidly growing in popularity

Change the way we write applications and expose new challenges

“The future of AWS”

– Marvin Theimer,
Distinguished Engineer at AWS



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 vCPUs	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

#thecloudistoodamnhard

What type of instances?

How many to spin up?

What base image?

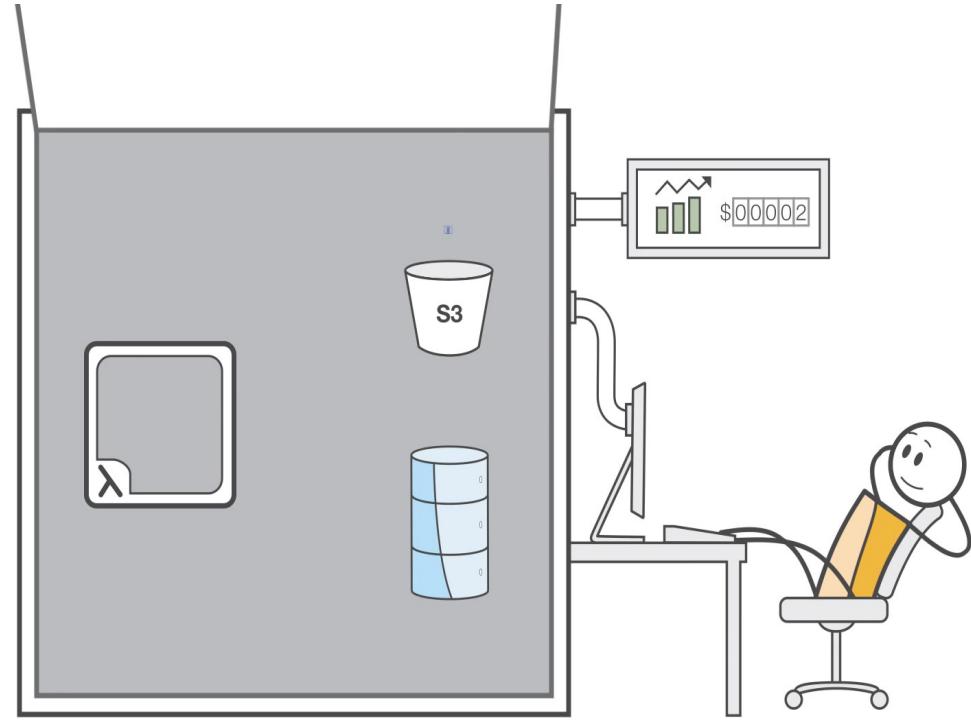
What price spot?

And then wait to start.....

EC2Instances.info Easy Amazon EC2 Instance Comparison

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
Custer Compute Eight Extra Large	c02xlarge	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	
Custer GPU Quadruple Extra Large	c01t2xlarge	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 vCPU	0 GB (EBS only)	32-bit/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 vCPU	0 GB (EBS only)	32-bit/64-bit	Low to Moderate	N/A	Yes	\$0,026 hourly	\$0,018 hourly	\$0,038 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	\$3,072 hourly	
M4.16xlarge	m4.16xlarge	256.0 GB	188 units	64 vCPUs	0 GB (EBS only)	64-bit	20 Gigabit	10000.0 Mbps	Yes	\$3,850 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	32 vCPUs	0 GB (EBS only)	64-bit	10 Gigabit	4000.0 Mbps	Yes	\$1,675 hourly	\$1,242 hourly	\$3,091 hourly	\$2,789 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	30 vCPUs	0 GB (EBS only)	64-bit	10 Gigabit	6000.0 Mbps	No	\$7,200 hourly	\$5,476 hourly	\$8,672 hourly	\$6,648 hourly	
P2 T2xlarge	p2.16xlarge	732.0 GB	186 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.xlarge	15.0 GB	26 units	8 vCPUs	0 GB (EBS only)	64-bit	High	1000.0 Mbps	No	\$0,650 hourly	\$0,474 hourly	\$0,767 hourly	\$0,611 hourly	
G2 Eight Extra Large	g2.2xlarge	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 T2xlarge	x1.16xlarge	976.0 GB	174.5 units	192 vCPUs	192.0 GB (2 * 96.0 GB SSD)	64-bit	10 Gigabit	5000.0 Mbps	No	\$8,669 hourly	\$6,479 hourly	\$9,613 hourly	\$7,523 hourly	
X1 32xlarge	x1.32xlarge	1963.2 GB	349 units	384 vCPUs	384.0 GB (2 * 192.0 GB SSD)	64-bit	20 Gigabit	10000.0 Mbps	No	\$13,338 hourly	\$9,158 hourly	\$16,226 hourly	\$15,046 hourly	
R3 High-Memory Large	r3.large	19.25 GB	6.5 units	2 vCPUs	32.0 GB SSD	64-bit	Moderate	N/A	No	\$0,106 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	\$2,650 hourly	\$1,672 hourly	\$3,500 hourly	\$1,989 hourly	
I2 Extra Large	i2.xlarge	30.5 GB	14 units	4 vCPUs	80.0 GB SSD	64-bit	Moderate	500.0 Mbps	No	\$0,853 hourly	\$0,424 hourly	\$0,973 hourly	\$0,865 hourly	
I2 Double Extra Large	i2.2xlarge	61.0 GB	27 units	8 vCPUs	160.0 GB (2 * 80.0 GB SSD)	64-bit	High	1000.0 Mbps	No	\$1,705 hourly	\$0,848 hourly	\$1,946 hourly	\$1,311 hourly	
I2 Quadruple Extra Large	i2.4xlarge	122.0 GB	53 units	16 vCPUs	320.0 GB (4 * 80.0 GB SSD)	64-bit	High	2000.0 Mbps	No	\$3,110 hourly	\$1,696 hourly	\$3,891 hourly	\$2,260 hourly	
I2 Eight Extra Large	i2.8xlarge	244.0 GB	104 units	32 vCPUs	640.0 GB (8 * 80.0 GB SSD)	64-bit	10 Gigabit	N/A	No	\$6,820 hourly	\$3,392 hourly	\$7,782 hourly	\$4,321 hourly	
I2 Double Extra Large	i2.xlarge	30.5 GB	14 units	4 vCPUs	600.0 GB (2 * 300.0 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.xlarge	61.0 GB	28 units	8 vCPUs	1200.0 GB (2 * 200.0 GB)	64-bit	High	1000.0 Mbps	No	\$1,380 hourly	\$0,804 hourly	\$1,601 hourly	\$0,885 hourly	
D2 Quadruple Extra Large	d2.2xlarge	122.0 GB	56 units	16 vCPUs	2400.0 GB (12 * 200.0 GB)	64-bit	High	2000.0 Mbps	No	\$2,760 hourly	\$1,408 hourly	\$3,092 hourly	\$1,690 hourly	
D2 Eight Extra Large	d2.4xlarge	244.0 GB	116 units	32 vCPUs	4800.0 GB (24 * 200.0 GB)	64-bit	10 Gigabit	4000.0 Mbps	No	\$5,520 hourly	\$3,216 hourly	\$6,198 hourly	\$3,300 hourly	
H1.1 High I/O Quadruple Extra Large	h1.1xlarge	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	\$1,698 hourly	\$3,580 hourly	\$2,200 hourly	
H1.1 High Storage Extra Large	h1.1xlarge	117.0 GB	35 units	16 vCPUs	4800.0 GB (24 * 200.0 GB)	64-bit	10 Gigabit	N/A	No	\$4,600 hourly	\$2,514 hourly	\$4,931 hourly	\$2,661 hourly	
M3 General Purpose Medium	m3.medium	3.75 GB	3 units	1 vCPU	80.0 GB (2 * 40.0 GB SSD)	64-bit	Moderate	N/A	No	\$0,367 hourly	\$0,248 hourly	\$0,130 hourly	\$0,100 hourly	
M3 General Purpose Large	m3.large	7.5 GB	8.5 units	2 vCPUs	32.0 GB SSD	64-bit	Moderate	N/A	No	\$0,153 hourly	\$0,095 hourly	\$0,259 hourly	\$0,199 hourly	
M3 General Purpose Extra Large	m3.xlarge	15.0 GB	13 units	4 vCPUs	80.0 GB (2 * 40.0 GB SSD)	64-bit	High	500.0 Mbps	No	\$2,266 hourly	\$1,050 hourly	\$5,018 hourly	\$3,997 hourly	
M3 General Purpose Double Extra Large	m3.2xlarge	30.0 GB	26 units	8 vCPUs	160.0 GB (2 * 80.0 GB SSD)	64-bit	High	1000.0 Mbps	No	\$5,532 hourly	\$3,080 hourly	\$7,036 hourly	\$5,993 hourly	

What it is?



CLOUD FUNCTIONS ALPHA

A serverless platform for building event-based microservices

Microsoft Azure

Azure Functions

Process events with a serverless code architecture

Classic web stack

RPC →

Application

Server

OS

Hardware

1st Generation: Virtual Machines

RPC →

Application

Server

OS

Hardware

Virtual hardware

1st Generation: Virtual Machines

RPC →

Application Application

Server

Server

OS

OS

Hardware

Virtual hardware

2nd Generation: Containers

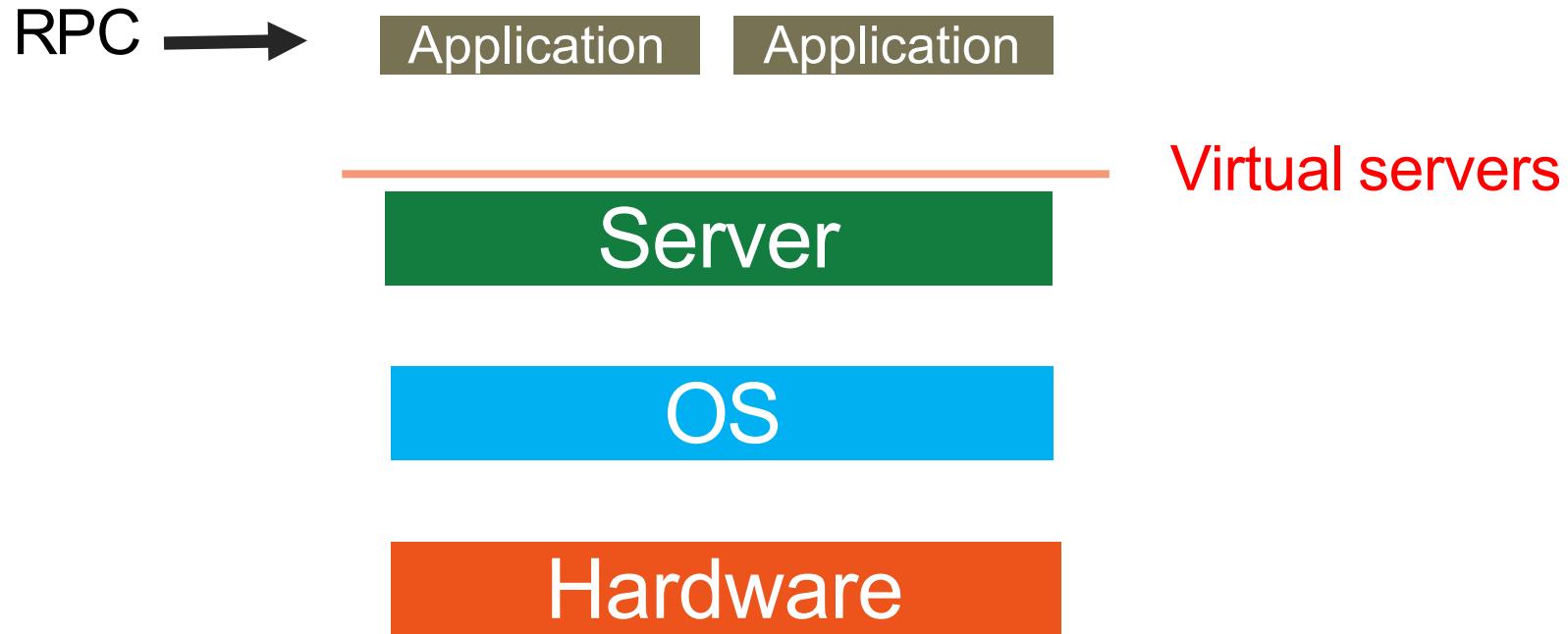
RPC →

Application Application

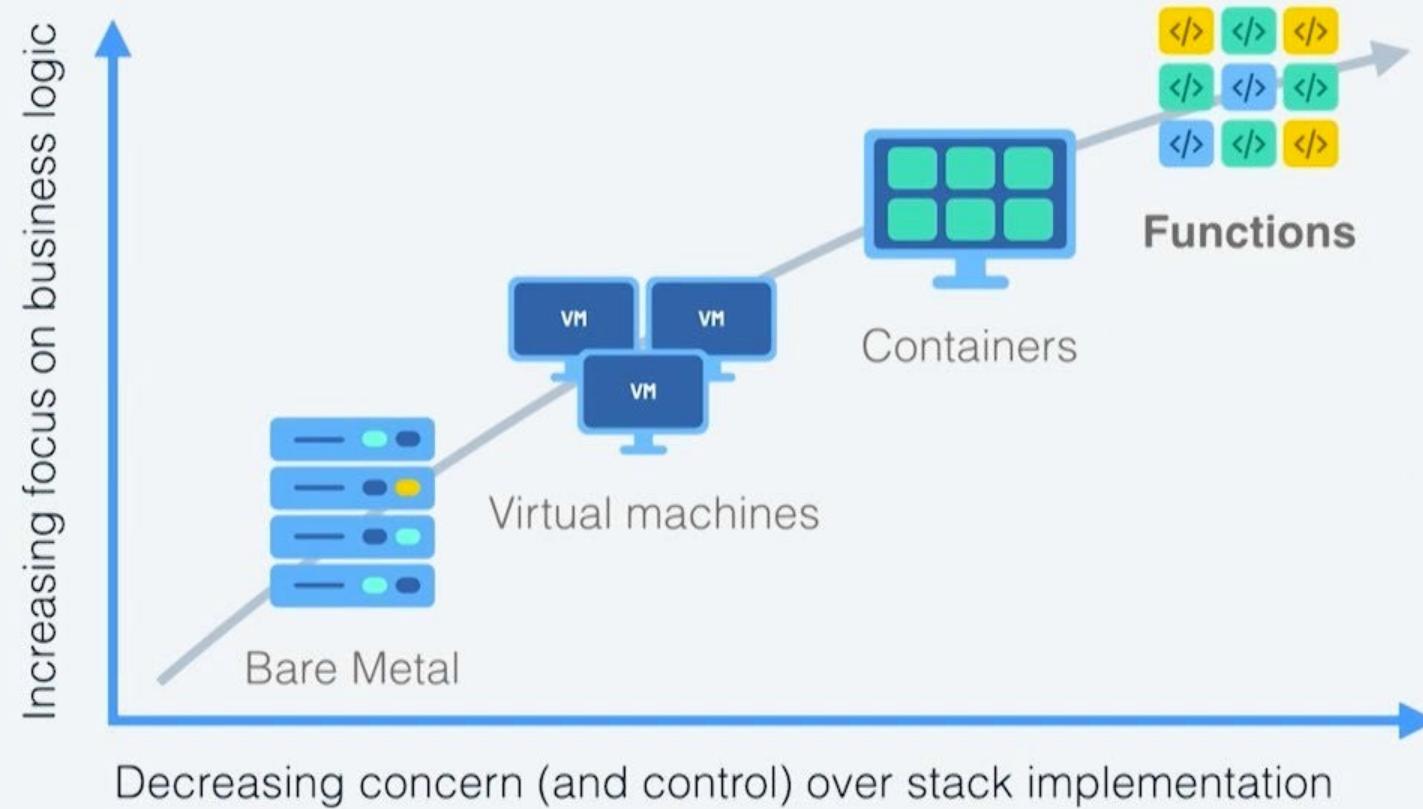


(based on slides at <https://www.usenix.org/node/196323>)

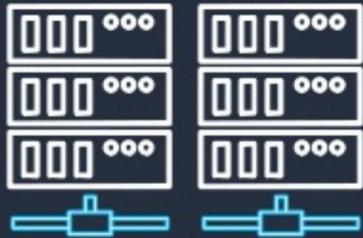
3rd Generation: Lambdas



(based on slides at <https://www.usenix.org/node/196323>)



Serverless means ...



No server or container management



Flexible scaling



High availability



No idle capacity

What is the essence of “Serverless Computing”?

Or, what do people really like about it?

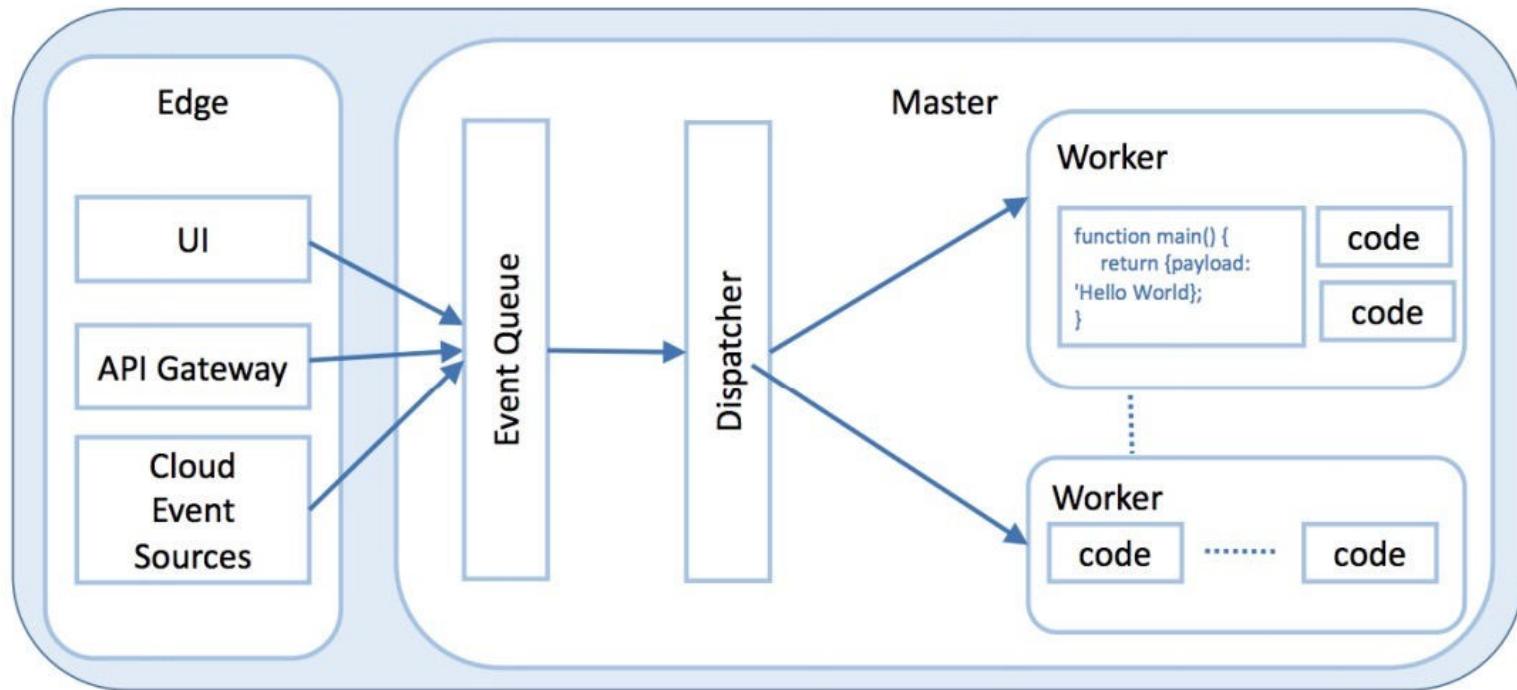
- Management-free
 - No need to handle creation, failure, replication, etc.
- Autoscaling
 - Spin up/down functions quickly based on load
- Only pay for what you use

“ I didn't have to worry about building a platform and the concept of a server, capacity planning and all that “yak shaving” was far from my mind... However, these changes are not really the exciting parts. The killer, the gotcha is the billing by the function... ”

“ This is like manna from heaven for someone trying to build a business. Certainly I have the investment in developing the code but with application being a variable operational cost then I can make a money printing machine which grows with users... ”

What is Today's Serverless Computing Like?

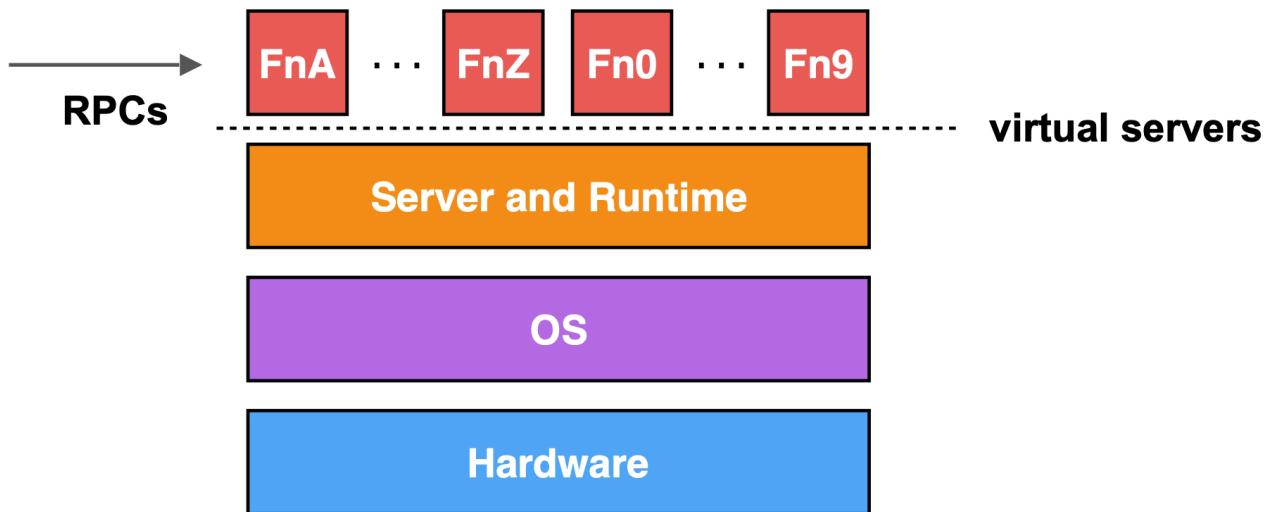
- Largely offered as Function as a Service (FaaS)
 - Cloud users write functions and ship them
 - Cloud provider runs and manages them
- Still runs on servers
- Have attractive features but also many limitations (more later this lecture)



Core capability

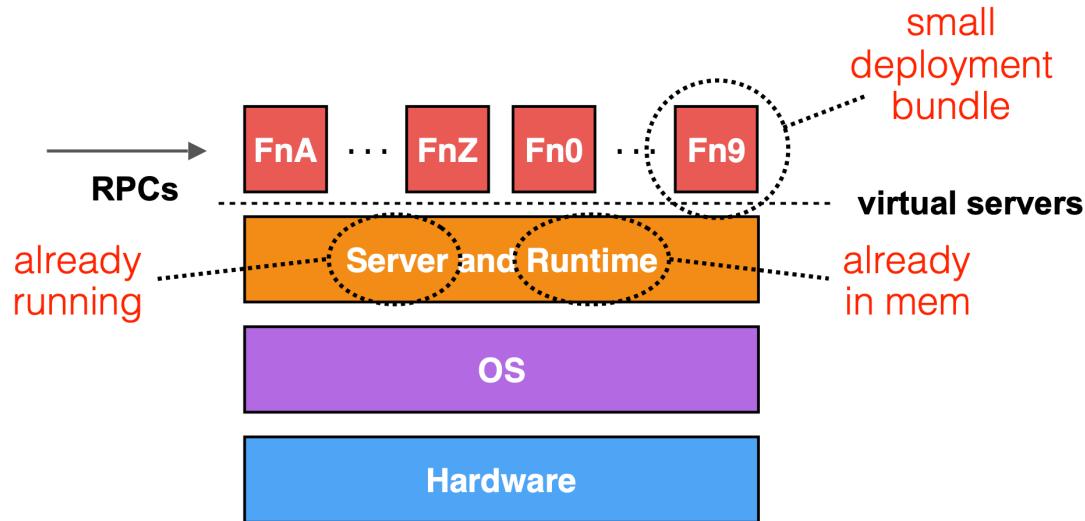
1. Manage a set of user defined functions
2. Take an event sent over HTTP or received from an event source
3. Determine function(s) to which to dispatch the event
4. Find an existing instance of function or create a new one
5. Send the event to the function instance
6. Wait for a response
7. Gather execution logs
8. Make the response available to the user
9. Stop the function when it is no longer needed.

3rd generation: lambdas



decompose application

3rd generation: lambdas



advantages:

- very fast startup
- agile deployment
- share memory

problems:

- not flexible

Serverless Applications

Event source



Changes in
data state



Requests
to
endpoints
Changes in
resource
state



Lambda function



Node.js
Python
Java
C# (.NET Core & Core 2.0)
Go
Ruby
Powershell
BYR – Bring your own Runtime

Services (anything)



AWS Lambda

- An event-driven, serverless computing FaaS platform introduced in 2014
- Functions can be written in Node.js, Python, Java, Go, Ruby, C#, PowerShell
- Each function allowed to take 128MB - 3GB memory and up to 15min
- Max 1000 concurrent functions
- Connected with many other AWS services

Lambda Function Triggering and Billing Model

- Run user handlers in response to events
 - web requests (RPC handlers)
 - database updates (triggers)
 - scheduled events (cron jobs)
- Pay per function invocation
 - No charge when no functions run (no triggering event)
 - Billed by duration of function, configured memory size, and # of functions
 - charge $actual_time * memory_cap$

Share everything

- Share server pool between customers
 - Any worker can execute any handler
 - No spinup time
 - Less switching
- Encourage specific runtime (C#, Node.JS, Python)
 - Minimize network copying
 - Code will be in resident in memory

Multi-node architecture

load balancers



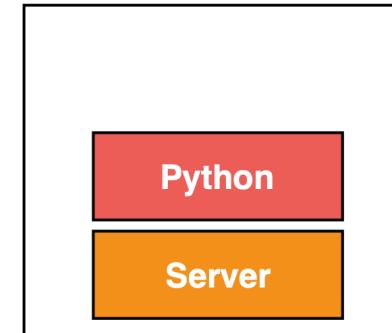
...



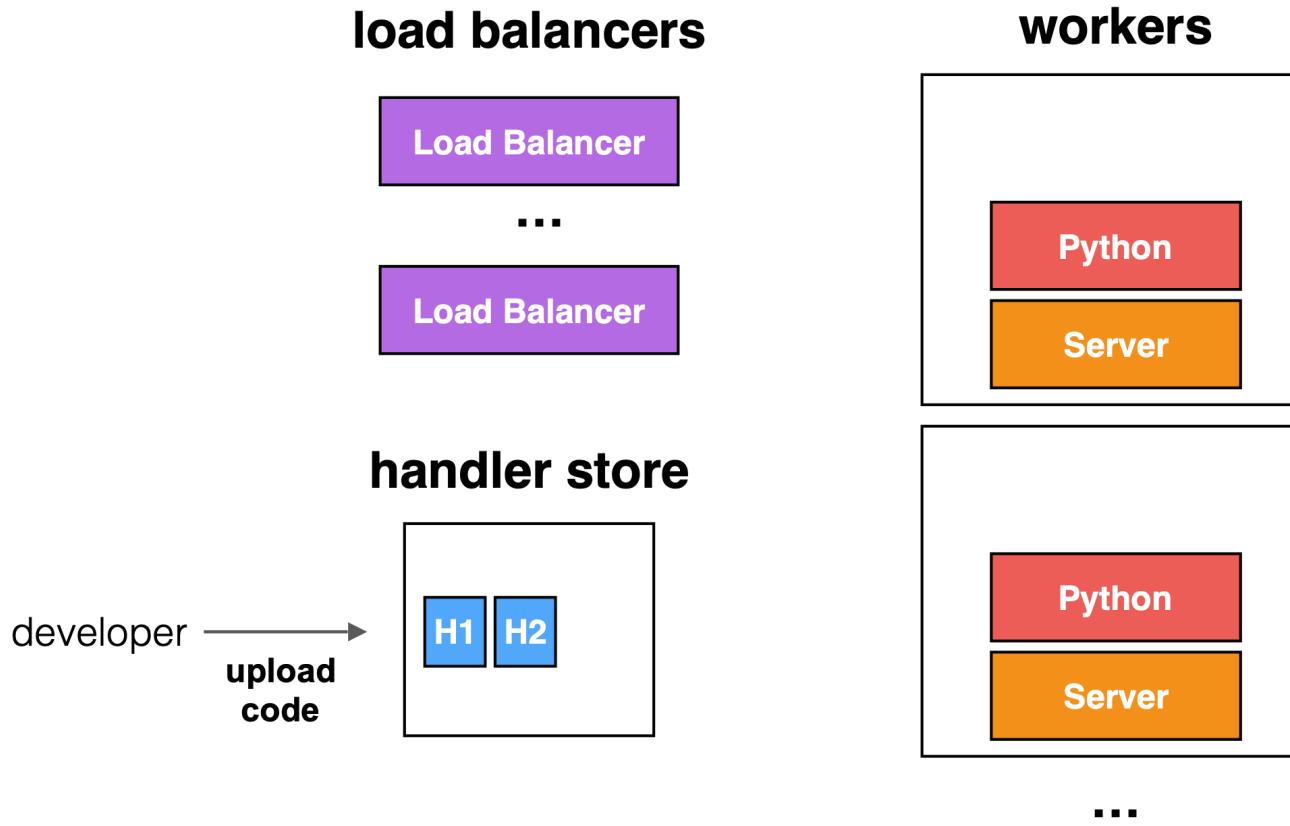
workers



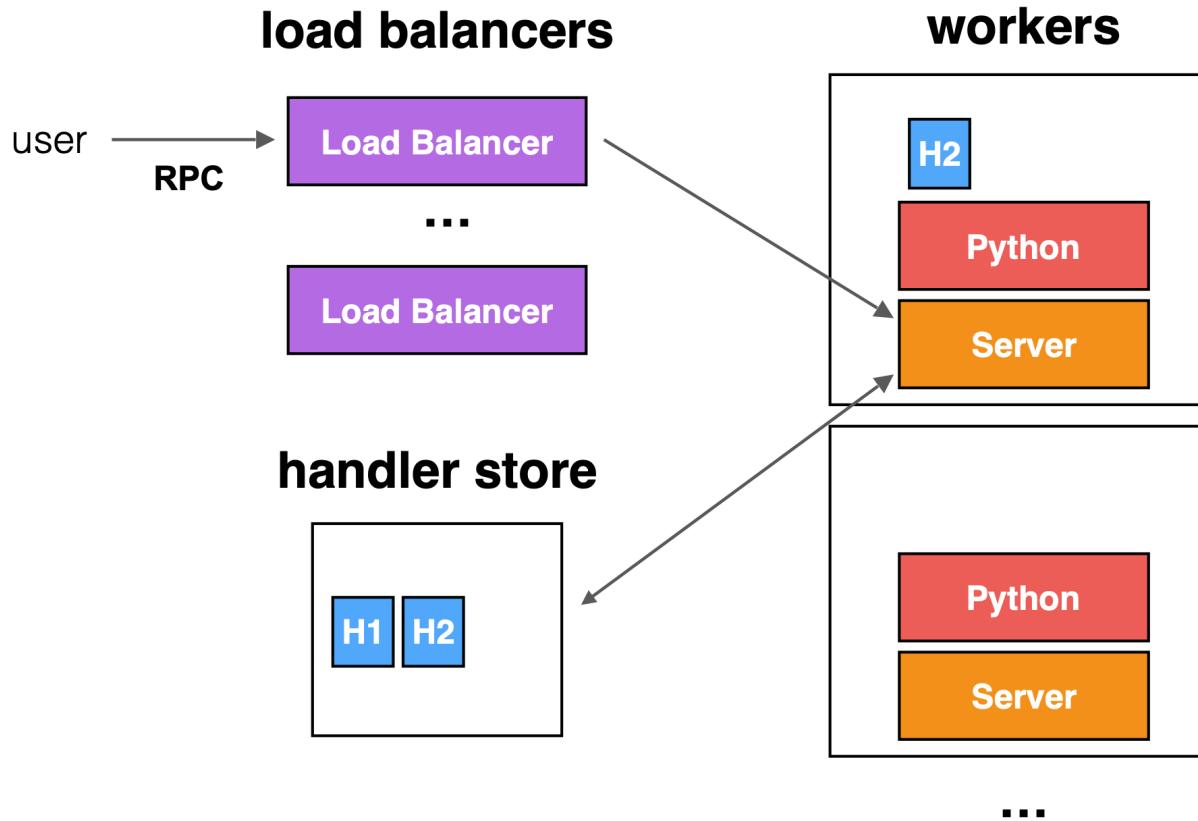
handler store



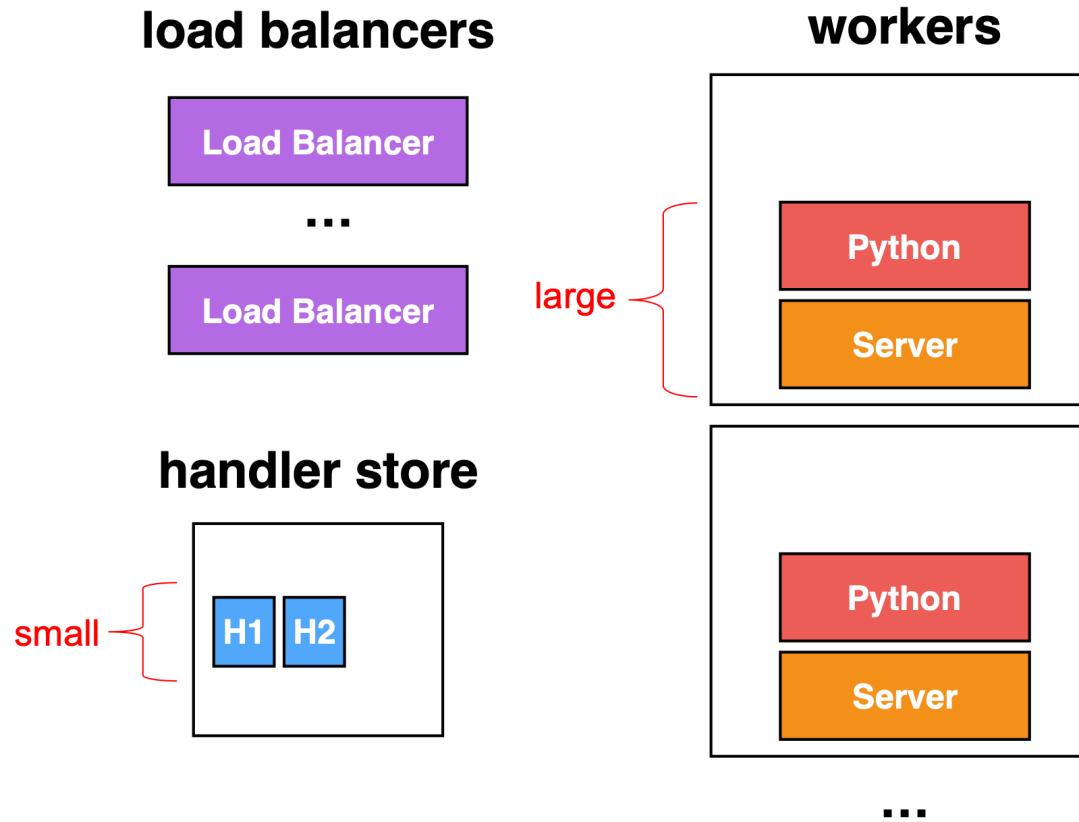
Multi-node architecture



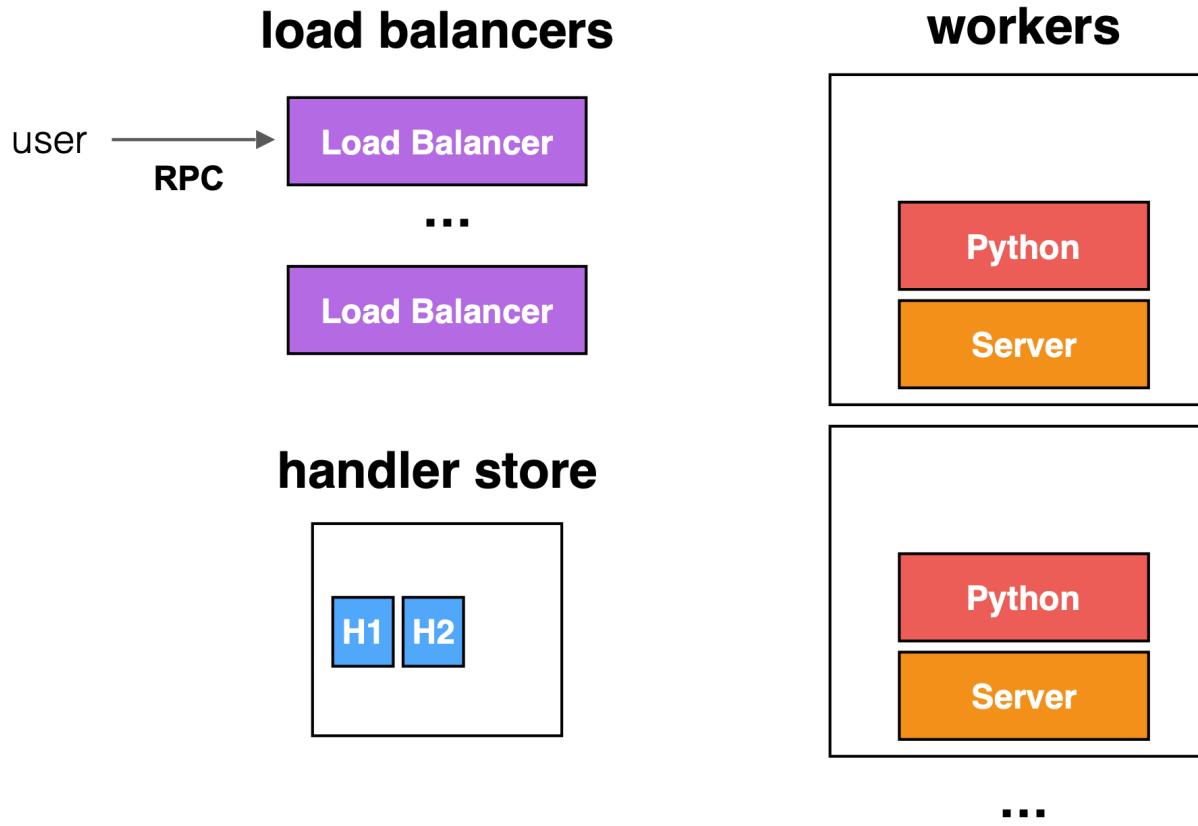
Multi-node architecture



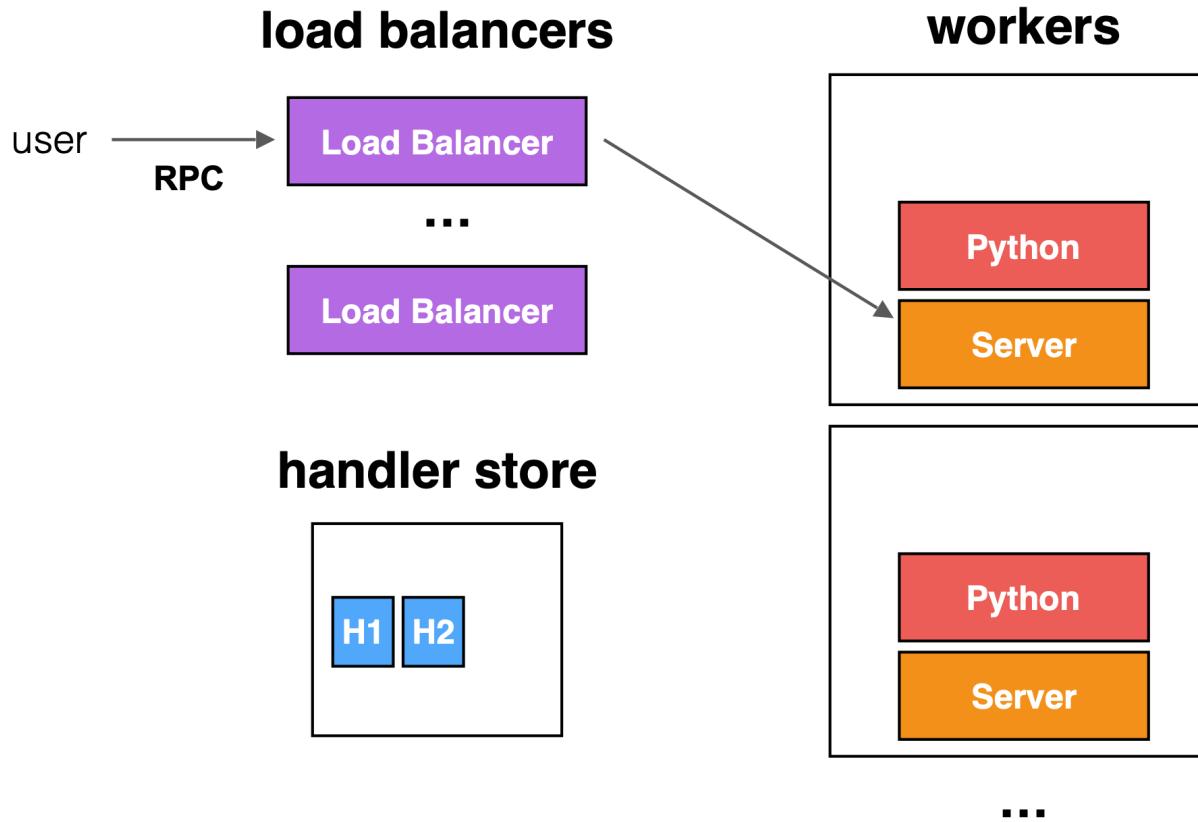
Multi-node architecture



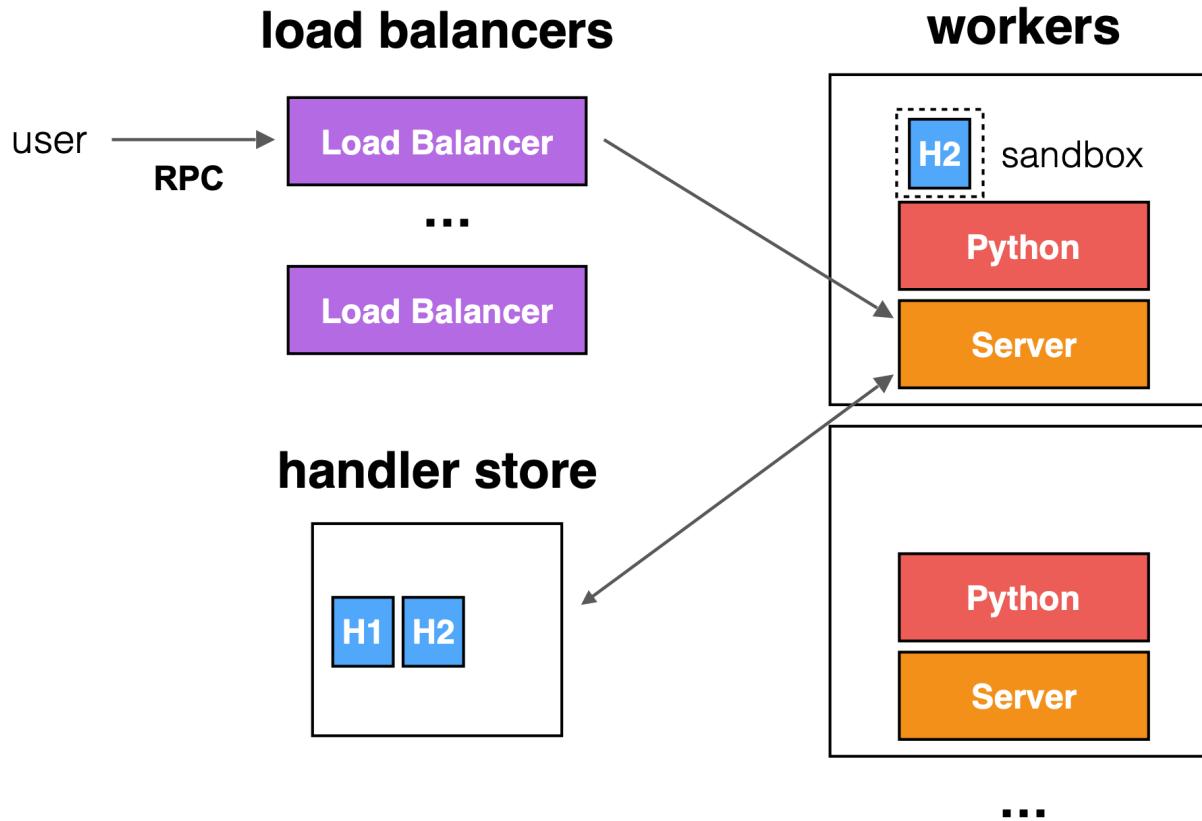
Multi-node architecture



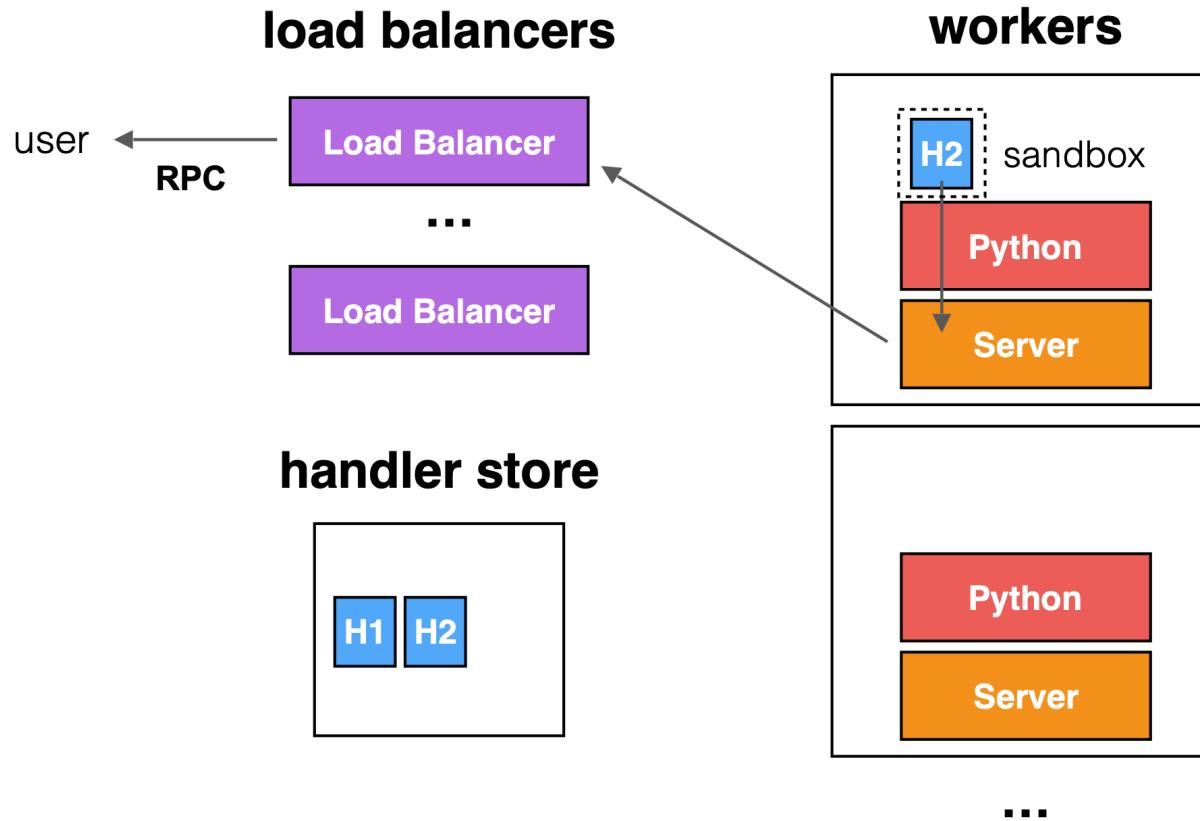
Multi-node architecture



Multi-node architecture



Multi-node architecture

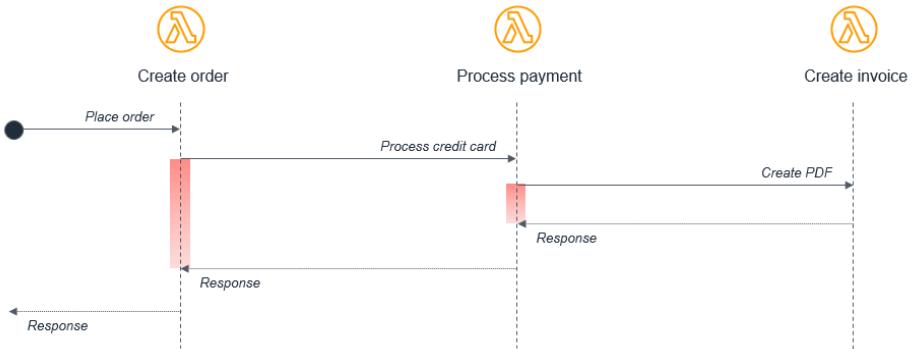


Internal Execution Model

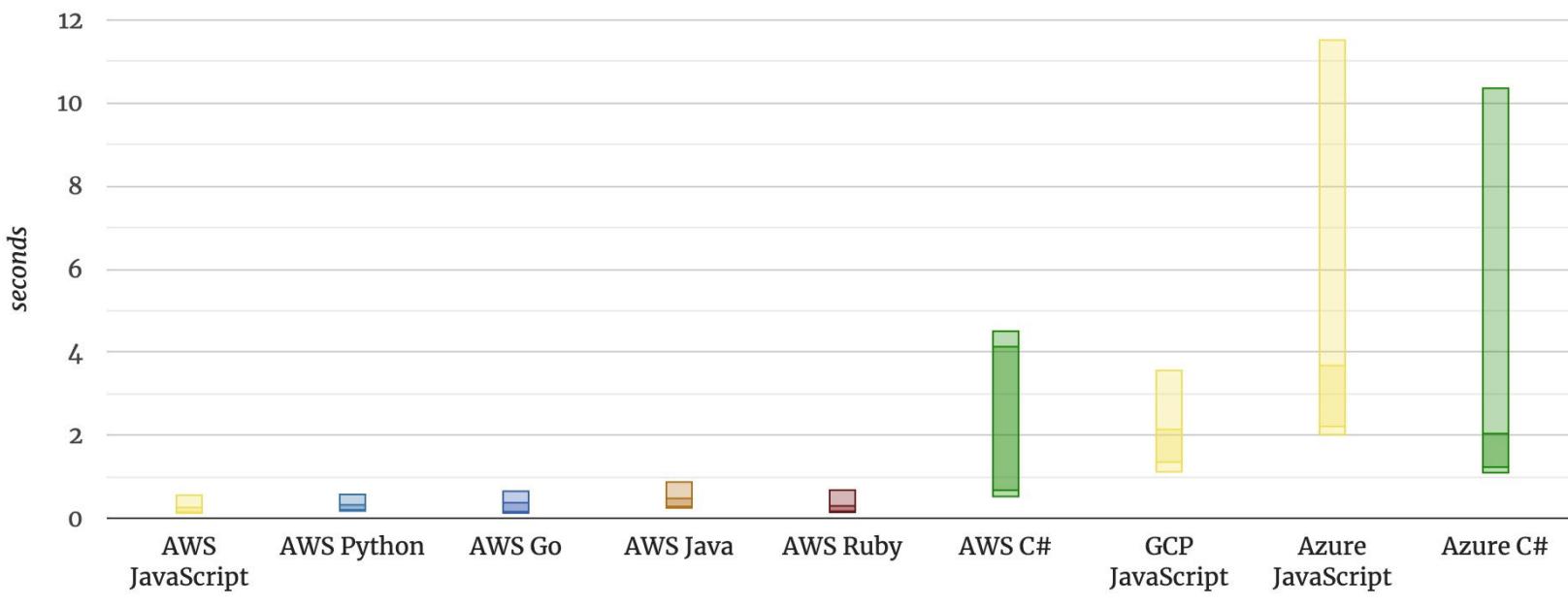
- Developers upload function code to a *handler store* (and associate it with a URL)
- Events trigger functions through RPC (to the URL)
- Load balancers handle RPC requests by starting *handlers* on *workers*
 - Calls to the same function are typically sent to the same worker(s)
- Handlers sandboxed in containers
 - AWS Lambda reuses the same container to execute multiple handlers when possible

Limitations of Today's Serverless Offerings

- Difficult and slow to manage states
 - Have to use (slow) cloud storage!
- No easy or fast way to communicate across functions
 - Have to go through cloud storage or other services
- Functions can only use limited resources
- No control over function placement or locality
 - e.g., starting functions on “cold” machines can be slow
- Billing model does not fit all needs



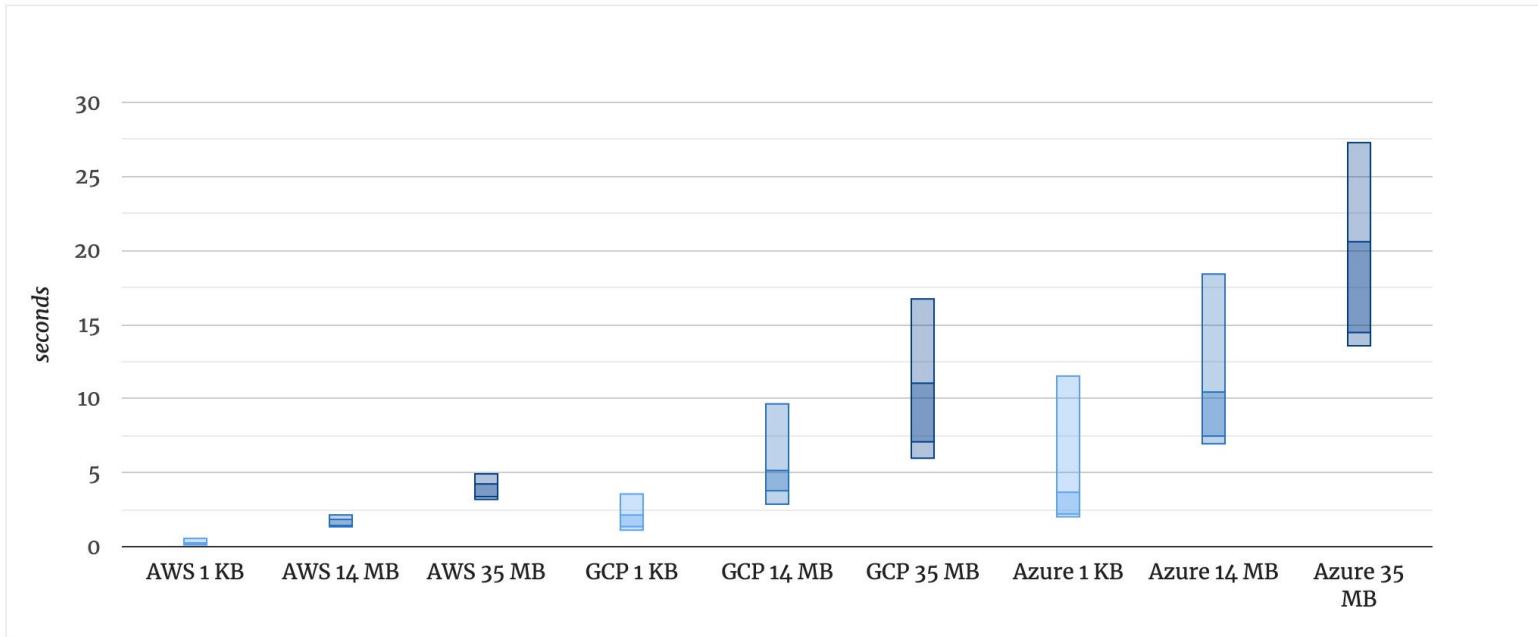
Cold Start Time of Different Languages/Offerings



Typical cold start durations per language

source: <https://mikhail.io/serverless/coldstarts/big3/>

Cold Start Time with Different Package Sizes



source: <https://mikhail.io/serverless/coldstarts/big3/>

Good and Bad Use Cases

- Some good ones:
 - Parallel, independent, stateless tasks
 - Event-triggered, synchronous processing
 - Ephemeral with highly dynamic demand
- Some bad ones:
 - Stateful applications
 - Distributed applications and protocols
 - Applications that demand more resources (esp. memory)

Next: Key-value Store