ECE 590/COMPSI 590 Special Topics: Edge Computing

How Does Edge Help The Cloud?

Monday September 10th, 2018

Last Lecture: Recap

- Higher-end mobile devices
- Cloudlets
 - Current presence
 - Challenges
- · Mobile offloading
- Future directions in mobile offloading

3

Class Outline

- Edge helping cloud
 - Why edge makes sense for the cloud
 - Background: latency and jitter
 - Challenges in supporting low-latency low-jitter solutions with modern cloud architectures
- Telecom and the edge
 - > An infrastructure view of edge computing
 - ➢ 5G and ETSI MEC



Why do Amazon and Microsoft Want to Create Edge Services?





Class Outline Edge helping cloud Why edge makes sense for the cloud Background: latency and jitter Challenges in supporting low-latency low-jitter solutions with modern cloud architectures Telecom and the edge An infrastructure view of edge computing 5G and ETSI MEC

Why do Amazon and Microsoft Want to Create Edge Services?



- Gateways are already already a pervasive reality for IoT deployments
 - Most likely, you will have an IoT gateway, and you will run something on it

8

7







Latency Components

- Latency, in a distributed system:
 - Getting data to and from the execution point
 - + service invocation time
 - + service execution time

Latency with Edge and Cloud

• Cloud:

 \succ Globally **pooled** users \rightarrow central server farm

• Edge:

Duke UNIVERSITY

➤ Local users → local gateway/cloudlet

Latency with Edge and Cloud: Comparison (1/2)

- Cloud communication latency strictly greater than edge latency
 - Speed of light

From:http://ipnetwork.bgtmo.ip.att.n et/pws/network_delay.html

U.S. Netwo	rk I	.ater	ncy																						
Figures are i	n m	s. Th	resho	olds	are	dista	ince	sens	itive																
CITY PAIRS	Atl																								
Austin	27	Aus							<u> </u>	Irro	nt														
Cambridge	28	28 52 Cam Overall																							
Chicago	24	30	- 24	Chi					Ave	ara	10' 10'														
Cleveland	18	35	19	- 7	Cle				3	1 m	90. IS														
Dallas	17	6	50	26	27	Dal			Ŭ																
Denver	37	25	42	19	26	19	Den																		
Detroit	21	39	21	7	4	34	26	Det																	
Houston	18	7	44	30	34	5	24	39	Hou																
Indianapolis				5	8					Ind															
Kansas City	22	15	31	12	18	10	15	18	15		Kan														
Los Angeles	46		67	44	50	28	25	51	33		39	LA													
Madison				5							18		Mad												
Nashville	- 7	22	28	17	11	19	33	14	25		17	44		Nas											
New Orleans	12	13	37	33	31	12	32	29	- 7		21	40		20	NO										
New York	24	49	6	18	13	45	37	17	40		26	61		23	34	NY									
Orlando	11	27	37	33	28	24	43	30	19		33	53		15	14	32	Orl								
Philadelphia	22	43	9	17	10	40	36	15	- 39		25	59		21	32	3	32	Pa							
Phoenix	40	20	64	42	45	20	34	49	23		29	10		35	29	59	43	56	Phx						
San Antonio	26	3	51	32	33	8	25	40	5		17	28		26	11	48	25	43	18	SA					
San Diego	43	27	68		54	26	28	53	29		35	4		47	35		49	63	- 7	25	SD				
San Francisco	51	41	67	45	51	37	28	52	40		43	8		60	48	62	60	61	19	37	13	SF			
St. Louis	18	22	26	- 7	13	17	21	15	23	13	6	45	11	11	29	21	27	19	37	24	44	50	StL		
Seattle	64	54	69	41	50	48	30	48	53		44			57	60	65		62	37	55	31	16	53	Sea	
Washington	18	38	10	20	13	33	38	17	37		22	61		23	29	5	27	4	53	40	59	63	16	60	Was
1																									







Latency Requirements: Often Not Strictly "As Little As Possible"

- Example of going for "as little as possible": highfrequency trading systems
- Not strictly "as little as possible":
 - Human attention
 - Systems bottlenecked by other components
 - ePrivateEye example: 30 FPS camera rate -> no improvement from processing frames faster than 33 ms





- Jitter: deviations from the mean
- Jitter is problematic for voice, gaming, video conferencing, control, augmented reality, ...

Class Outline

- Edge helping cloud
 - > Why edge makes sense for the cloud
 - Background: latency and jitter
 - Challenges in supporting low-latency low-jitter solutions with modern cloud architectures
- Telecom and the edge
 - > An infrastructure view of edge computing
 - ➢ 5G and ETSI MEC

Cloud Latency: Background

- Recognize latency magnitude as an issue
 > E.g., Content Delivery Networks as one solution
- Recognize jitter as an issue
 - > E.g., for multi-player games, VoIP
 - Edge should be able to support applications with tighter latency requirements

Distributed Data Analytics: Stragglers (1/2)

Big data platforms:

- Divide data into small pieces
- Perform calculations on the pieces in parallel
- MapReduce, Dryad, Spark, …
- Task completion latency is set by the time of the slowest task





Latency Variability Sources (1/3)

- Shared Resources
 - CPU cores
 - Processor caches
 - Memory bandwidth
 - Network bandwidth
- In our measurements with AWS t2.micro, we have seen up to 11x increase in latency

From: The Tail at Scale, J. Dean et al, Communications of the ACM, 2013









Specific Measurements of Latency and Latency Variability (3/3)

Game server latency statistics

Interference	Avg. Time	σ	Timeouts
Idle (none)	8.1	10.2	0%
CPU + Disk	6.2	7.9	1.7%
Net (no tc)	N/A	N/A	100%
Net (tc, dedicated)	23.6	29.6	6.7%
Net (tc, sharing)	33.9	16.9	1.7%



From: Empirical Evaluation of Latency-sensitive Application Performance in the Cloud, Barker and Shenoy, MMSys'10, Feb. 2010 31

Duke UNIVERSITY

There are Ways of Improving Cloud Latency Support

- E.g.,
 - For stragglers: speculative, coded, approximate execution
 - For latency caused by shared network or CPU: isolated resources
- But:
 - ➤ All require additional resources
 - New applications need even tighter latencies









Country +	Market value (\$ Bn) 🗢	Revenue +	Profit +
China	213.8	88.8	20.5
USA	200.1	127.3	7.3
USA	137.3	115.7	0.9
UK	135.7	74.4	11.1
Mexico	70.7	60.2	7.1
Spain	67.1	82.3	5.2
Australia	58.4	25.8	3.5
Japan	58.2	127	5.6
Germany	48.8	76.7	-7
Japan	47.2	38.78	3.8
	Country + China USA USA USA Mexico Spain Australia Japan Germany Japan	Country + Market value (\$ Bn) + China 213.8 USA 200.1 USA 137.3 UK 135.7 Mexico 70.7 Spain 67.1 Australia 58.4 Japan 58.2 Germany 47.2	Country + Market value (\$ Bn) + Revenue + China 213.8 88.8 USA 200.1 127.3 USA 137.3 115.7 UK 135.7 74.4 Mexico 70.7 60.2 Spain 67.1 82.3 Japan 58.2 127 Germany 48.8 76.7 Japan 47.2 38.78



Dec. 2017

Verizon peels back curtain on edge computing, deep learning for real-time video analytics

by Mike Dano | Mar 29, 2018 12:59pm

Duke UNIVERSITY

Era

Mobile Offloading: Application View





- The view we have seen so far
- But, there is telecom piping underneath all of it



Mobile Offloading: Infrastructure View

(2/2)

- Cortana Alexa
- Infrastructure:
 - ➢ Pervasive
 - ➤ Expensive
 - Including real estate, laying and maintaining wires, …
 - ➤ Mission-critical

Telecom as an Infrastructure Layer

- Telecom as a utility
 - Commoditization of telecommunication services
 - > "Metered data" services, minutes of voice, number of texts
 - > Hard to differentiate offerings from different companies
- Connectivity services → connected experiences
 - Not exclusive to edge services
 - > ... but very important in edge context

42









Edge Computing is a Part of 5G

- One of the building blocks
- Offers:

- ➤ Lower latency
- Reduced load on core network
- Idea: co-locate edge computing servers with cellular base stations





ETSI MEC: Example Standards

- Study on MEC support for V2x use cases
- UE identity API
- System, host, and platform management
- Bandwidth management API
- UE application interface
- Application lifecycle, rules and requirements management
- Radio Network information API
- Location API
- ...

Duke





Telecom Edge vs. Cloudlet Edge

- Existing pervasive infrastructure
- Minimal possible latency for cellular devices
- Know all about mobility

- Have a concept of location can geo-locate without a GPS
- Know how to handle handoff
 - Computing handoff ≠ cellular hand-off though





