ECE 590/COMPSI 590 Special Topics: Edge Computing

Edge Helping Higher-end Mobile Devices: Mobile Offloading

Wednesday September 5th, 2018

Last Class Recap

- Edge and IoT devices
 - Common IoT architectures
 - Role of the gateway
- Opportunities: edge for responsive IoT applications
 - Hardware
 - Algorithms
 - Edge for system decisions

Duke UNIVERSITY

2









7

Special Case: Camera Installations

- E.g., city, campus security cameras
 - Very common
 - Of major practical importance
 - Often not mobile devices
 - Many video-specific mechanisms
- Time-permitting, will cover later in the course

Duke UNIVERSITY



From: The Design and Implementation of a Wireless Video Surveillance System, Zhang et al, ACM MobiCom'15







Mobile Device vs. a Server

- Isn't a mobile device a desktop in your pocket?
- Server > mobile device
 - ➢ Power constraints → 500 W of power on a high-end GPU, 10 W on a mobile SoC GPU
 - Space constraints



Existing and Possible "Cloudlets"

- On-site computing
- Targeted edge installations
- Resource scavenging

On-Site Computing

- At universities
- ... and other medium and large organizations
 Shrinking but not disappearing
 Usually have low utilization

13

Targeted Installation: Chick-fil-A (1/2)



July 2018

15

Duke







<section-header><list-item><list-item><list-item><list-item><list-item><list-item><list-item><text>

Cloudlet Challenges

- **Mobile** devices \rightarrow supporting mobility
- Cloudlet →does not have the scale of the cloud

Cloudlets Helping Mobile Devices: Challenges: Rapid Service Provisioning

- A scenario: a student comes to Hudson Hall and needs to use our cloudlet
 - Service discovery
 - Provisioning delay
 - Do not have the scale of the cloud: do we prioritize this user over others? Shift workloads with every user?

Cloudlets Helping Mobile Devices: Challenges: Service Handoff

- A scenario: the student moves from Hudson Hall to CIEMAS
 - ➤Do we transfer their workload state?
 - Do we de-provision their Hudson Hall services?

Duke university

Cloudlets Helping Mobile Devices: Challenges

- Platform challenges
 Challenges similar to wireless hand-off
- Workload allocation and scheduling challenges



Our Recent Work: Picking Workloads for Local Execution (1/2)

• Work joint with



- Workloads achieve different utilities when executed in different locations
 - Quality
 - Latency

Duke

Want to maximize the sum of utilities for all workloads

25

Our Recent Work: Picking Workloads for Local Execution (2/2)

- Easy to solve for one capacitated cloudlet
 Pick the workloads with the highest utility gains over the cloud
- Becomes complex with as few as two capacitated cloudlets

26

27



- Current preser
- Challenges
- Mobile offloading
- Future directions in mobile offloading
- Challenges

•



Goals: Reducing Mobile Device Energy Consumption (1/3)

- Need to have:
 - Energy to {transmit data + receive results} < energy to {execute the operation on the mobile device}
- · Design principles:

Duke UNIVERSITY

- > Pick the most compute-intensive parts of the operation
- > Reduce the size of what is transmitted: data and results
- Order-of-magnitude mobile energy savings possible

Example: Face Recognition with MAUI 35 Smartphone only MAUI (Wi-Fi, 10ms RTT) 30 MAUI (Wi-Fi, 25ms RTT) MAUI (Wi-Fi, 50ms RTT) 25 20 20 15 10 MAUI (Wi-Fi, 100ms RTT) MAUI* (3G, 220ms RTT) 5 ONE RUN FACE RECOGNITION From: MAUI: Making Smartphones Last Longer with Code Offload, Cuervo et al., ACM MobiSys'10. 30 Duke UNIVERSITY

Goals: Reducing Mobile Device Energy Consumption (2/3)

- Not minimizing total energy:
 - Combined server + mobile energy spending can be higher than mobile-only energy spending
- System heterogeneity principle:
 - Server energy spending is not as important as mobile device energy spending
 - > Server grade does not factor into energy minimization objective

Goals: Reducing Mobile Device Energy Consumption (3/3)

- Often: transmit partially processed, rather than raw, data
 - Energy to {extract features + transmit extracted features + receive results} < energy to {transmit data + receive results}</p>
 - Energy to {extract features + transmit extracted features + receive results} < energy to {execute the operation on the mobile device}





Mobile Offloading: Need for Scheduling Mechanisms

- Time, energy vary with network connectivity
- Need to make decisions for different conditions
 Different ways of placing different parts of operations
 Offline versus online
 - Joint scheduling of different operations
 - Scheduling that takes into account different local processors and the cloud

Role of the Edge (1/2)

- Short transmission distance helps both transmission energy and latency
 > Better performance of existing offloading scenarios
 - Offloading equations "work out" in more cases
- Potentially, additional privacy

36





Future Directions: "Offload Shaping"

- Adapting operations for offloading
- A form of creative pre-processing
 - Changing application pipelines specifically for offloading
- Some examples from: The Case for Offload Shaping, by Hu et al, ACM HotMobile'15

Duke UNIVERSITY

Offload Shaping: Object Recognition in Video Captures (1/2)

- Object recognition works poorly on blurry frames
 - Can drop blurry frames before transmitting them to the cloud/cloudlet for processing





	Send all	Drop blurry
Bytes transferred	0.51M	0.34M
Glass energy (J)	429(2)	292(3)
Server CPU usage	1.00(0.01)	0.81(0.01)
(normalized)	1.00(0.01)	0.01(0.01)

40

Duke

Offload Shaping: Object Recognition in Video Captures (2/2)

- Results from similar frames are likely to be the same
 B B Fi G
 - Discard frames that are sufficiently similar

		-	-
	No	Drop	Improve-
	shaping	similar	ment
Bytes transferred	$0.51\mathrm{M}$	0.23M	55%
Frames recognized	171(2)	189(1)	11%
Glass power (W)	$1.82_{(0.01)}$	$1.83_{(0.01)}$	-1%
E2E latency (ms)	920(8)	393(2)	57%
Glass energy (J/frame)	1.66(0.01)	0.72(0.01)	57%
Server CPU usage	1.00	0.97	7207
(normalized)	1.00(0.01)	0.21(0.01)	1370

Offload Shaping

- (+) Holistic view of the entire system
 - Fixing inefficiencies that become obvious when we think about the system beginning-to-end
- (-) Solutions likely to be application-specific
 > E.g., blur detection in one of the previous examples



Opportunities: Providing Local Context

- Especially when context is large
- Opportunities for behavior specialization

Side Note: Context Awareness in Applications is Not New

- Traces back to early 1990s
- E.g.:

Duke UNIVERSITY

Active badge location system

OS updates only when a phone is plugged in and is on WiFi

Context Awareness Example: Is a Cell Phone User a Driver or a Passenger?

 Using cell phone's gyroscope, accelerometer data to detect and differentiate specific motions



From: I Am a Smart Phone and I Know My User is Driving, Chu et al, IEEE COMSNETS'14.

46

Large Local Context: 3D Maps of the Environment for AR/VR (1/2)

- Massive amounts of information and processing
 - Useful to not regenerate for all users
 - Useful to not fetch from the cloud



Mesh representing a student dorm room

47

Large Local Context: 3D Maps of the Environment for AR/VR (2/2)



Mesh representing a lab

48

Duke



Opportunities: Thinking Across Multiple Devices and Multiple Applications

New paradigms



- Without the cloudlets, nearby devices have no exposure to each other's actions
 - No single "choke point"

50

Opportunities: Thinking Across Multiple Devices and Multiple Applications

- Same application likely to be invoked on different devices served by one cloudlet
- Invited speaker in October





Reading Material for the Next Class

- Technical:
 - Commoditization of the wireless industry
 - Vodafone perspective on edge computing
- Approaching research (and technology):
 - > Technology and Courage

