#### Latency-aware Virtual Machine Placement for Mobile Edge Computing

#### Wei Wang

#### BUPT Ph.d candidate & UC Davis visiting student Email: <u>weiw@bupt.edu.cn, waywang@ucdavis.edu</u>



Group Meeting, Feb. 10, 2017



- Introduction
- Existing works
- Problem Statement
- Problem Formulation
- Future works



### Introduction

#### Cloud Computing

- Cloud computing is a type of Internet-based computing that provides shared computer processing resources and data to computers and other devices on demand.
- Computer processing resources and data are usually deployed in centralized datacenters, which is far away from end users.
- Drawbacks, long-distance network connection between user and cloud result in long service latency.





### Introduction

### Mobile Edge Computing(MEC)

 Mobile Edge Computing provides an IT service environment and cloud-computing capabilities at the edge of the mobile network, within the Radio Access Network (RAN) and in close proximity to mobile subscribers. The aim is to reduce latency, ensure highly efficient network operation and service delivery, and offer an improved user experience.[1]



[1] Hu, Yun Chao, et al. "Mobile edge computing—A key technology towards 5G." *ETSI White Paper* 11 (2015).



### Introduction

- MEC cloud and overall MEC system
  - Mobile operators are working on Mobile
     Edge Computing (MEC) in which the computing, storage and networking resources are integrated with the base stations.[2]



[2] Lav Gupta and Raj Jain, Mobile Edge Computing – An Important Ingredient of 5G Networks, IEEE Software Defined Networks



### **Existing work**

#### • Mobility-driven service migration

- Problem: mobility of user may increase the distance between user and its VM, and thus increase latency.
- Solution: migrate user's VM across MEC clouds dynamically when user moves.



Wang, Shiqiang, et al. "Dynamic service migration in mobile edge-clouds." *IFIP Networking Conference (IFIP Networking), 2015.* IEEE, 2015.
Bittencourt, Luiz Fernando, et al. "Towards virtual machine migration in fog computing." *P2P, Parallel, Grid, Cloud and Internet Computing (3PGCIC), 2015 10th International Conference on.* IEEE, 2015.



# Existing work (cont.)

- SLA-driven VM Scheduling in Mobile Edge Computing
  - Problem: Service providers' cost of renting VMs at Edge clouds is calculated as \$/time unit. Each type of VM has its maximum capacity to handle request. If the number of requests exceeds its capacity, some requests will go to remote clouds, and thus cause penalty for violating SLA. How to reduce cost while minimizing service penalty?
  - Approach: LYAPUNOV OPTIMIZATION-BASED scheduling algorithm for deploying and releasing VMs dynamically.

Katsalis, Kostas, et al. SLA-driven VM Scheduling in Mobile Edge Computing. 9th International Conference on Cloud Computing, IEEE, 2016



### VM based service for MEC

- Role of VMs in MEC and central clouds
  - In general, user devices play as clients and MEC clouds perform as servers.
  - To accommodate users requests, MEC clouds need to provide not only the computing and storage ability, but also the service specific software and data, which can be packaged in Virtual Machines(VM).

#### • Options for service handling.

- Handled at Local MEC Cloud
- Handled at Other MEC Cloud
- Handled at Centralized Cloud



# Latency of MEC service

- To get network service, users' requests need to go to base station first, and then find its destination VM to get service.
- Define: The MEC cloud, where the first hop base station locates at, is called origination cloud of request.
- These procedures will introduce network and processing latency, which is very critical for Quality of Experience.
- Major parts of overall Latency:
- Transmission Latency : depends on bandwidth;(out of scope)
- Propagation Latency : depends on length of network path;
- Service\_Latency: depends on number\_and-work-load of VMs.
   (Transmission Delay)
   (Propagation Delay)



### VM's influence on Latency

- As the server for user's requests, the location of VMs has significant influence on the distance from user to server, and thus influence the propagation latency.
- The number of VMs decides average load of each VM, and is related with processing latency.



• Service routing map decides which VM is responsible for processing the requests originate from a specific Cloud, and thus has influence on both propagation and processing latency.



### **Problem statement**

### • Problem Description:

- MEC clouds are already deployed at network edge, and each cloud has certain hardware resource to run VMs.
- Network distance (latency) between each MEC cloud pair is known.
- Several kinds of MEC service are already known, and each kind of service has its expected latency threshold.
- Each kind of service corresponds to a specific kind of server VM, which has certain capacity of handling requests.
- Expected request load of each kind of service that originates from each MEC cloud is known.
- How to place VMs at each MEC cloud for each kind of service to meet their latency requirements?



### **Problem Formulation**

#### • Given:

- E: Set of edge clouds
- $L_{e1,e2}$ : Network propagation latency from  $e1 \in E$  to  $e2 \in E$
- $C_e$ : Hardware capacity of edge cloud  $e \in E$
- S: Set of services
- $R_s$ : Computing capacity required to deploy a VM for service  $s \in S$
- $T_s$ : Expected latency requirement of service  $s \in S$
- $u_s$ : Handling capacity at VM for service  $s \in S$
- $\lambda_e^s$ : Request load of service  $s \in S$  that originates from cloud  $e \in E$



#### • Variables:

- n<sup>s</sup><sub>e</sub>: How many VMs need to be deployed for service s ∈ S at edge cloud e ∈ E.
- $x_{src,dst}^s$ : Whether the requests of service  $s \in S$  from  $src \in E$  are processed by edge cloud  $dst \in E$ .

### • Object:

• Minimize required hardware resource (cost) for deploying VMs.

$$\sum_{\mathbf{e}\in E}\sum_{s\in S}\mathbf{n}_{\mathbf{e}}^{s}*\mathbf{R}_{s}$$



- Constraints about hardware capacity:
  - (1) Required resource for all VMs at each MEC cloud  $e \in E$  should not exceed its hardware capacity.

$$\sum_{e \in S} n_e^s * R_s < C_e, \forall e \in E$$



- Constraints about service strategies:
  - (2) Each kind of service s ∈ S needs at least one corresponding VM all overall the MEC clouds.

$$\sum_{e \in E} n_e^s \ge 1, \forall s \in S$$

 (3) The requests of service s ∈ S from one MEC cloud src ∈ E should be processed by one destination cloud dst ∈ E.

$$\sum_{dst\in E} \mathbf{x}_{src,dst}^{s} = 1, \forall s \in S, \forall src \in E$$



- Constraints about service strategies (cont.):
  - (4) There must be VM(s) to handle service s ∈ S at MEC cloud dst ∈ E, if there are requests of s being routed to dst. And vice versa.

$$n_{dst}^{s} - \frac{\sum_{src \in E} x_{src,dst}^{s}}{M} \ge 0, \forall s \in S, \forall dst \in E$$
$$\sum_{src \in E} x_{src,dst}^{s} - \frac{n_{dst}^{s}}{M} \ge 0, \forall s \in S, \forall dst \in E$$



- Constraints about service strategies (cont.):
  - (5) the requests of service s ∈ S originate from MEC cloud dst ∈ E should be processed locally, if there are VM(s) for service s deployed at dst.

$$n_{dst}^{s} - \frac{x_{dst,dst}^{s}}{M} \ge 0, \forall s \in S, \forall dst \in E$$
$$x_{dst,dst}^{s} - \frac{n_{dst}^{s}}{M} \ge 0, \forall s \in S, \forall dst \in E$$



- Constraints about service latency:
  - (6) Requests of service s ∈ S originates from src ∈ E can not be processed at a dst, to where the propagation latency exceeds the threshold latency of service s.

$$(T_s - L_{src,dst}) * x_{src,dst}^s \ge 0, \forall s \in S, \forall src \in E, \forall dst \in E$$

• Processing and queueing Latency?



• Processing and queueing latency at a VM:

- Queueing Theory. M/M/1 System with Poisson Process
- $n_e^s$ : number of VMs for service s at e.
- $\sum_{src \in E} \lambda_{src}^s * x_{src,dst}^s$ : overall load for service s at dst.
- u<sub>s</sub>: departure rate of service s at one VM.
- Assumption: requests to one cloud are distributed evenly to all VMs, the system would be  $n_e^s M/M/1$ .
- Expected average service latency is:

$$\frac{1}{u_s - \frac{\sum_{src \in E} \lambda_{src}^s * x_{src,dst}^s}{n_e^s}}$$



- Constraints about service latency (cont.):
  - (7) The average arrive rate of each service s ∈ S to a VM should not exceed its departure rate, so that the queueing system is stable.

$$\mathbf{u}_{s} * \mathbf{n}_{dst}^{s} - \sum_{src \in \mathbf{E}} \lambda_{src}^{s} * \mathbf{x}_{src,dst}^{s} \ge 0, \forall s \in \mathbf{S}, \forall dst \in \mathbf{E}$$

 (8) The overall latency of the farthest customer should not exceed the required threshold latency.

$$L_{max}(s, dst) + \frac{1}{u_s - \frac{\sum_{e \in E} \lambda_s^e * x_{e,dst}^s}{n_{dst}^s}} \le T_s, \forall s \in S, \forall dst \in E$$





set EDGECLOUDS := SFO LAX NYC MCO PEK HKG NRT; # BOM LCY FRA; set SERVICES := AR MAP GAME FACERECG VEDIO SOCIAL SHOPPING;# BACKUP PAY TRAVEL;

naram:	c1(	oudCanacity:=	naram.	muPer\/M·	-			para	m ne	etLat	:ency	:							
SE0	90	outcupacity		0 02	-				SFO	LAX	NYC	MCO	PEK	HKG	NRT:	=# B(	DM LO	CY F	RA:=
	166		MAP	0.02				SFO	0	9	85	75	185	159	125	#280	151	158	
	100		GAME	0.025				LAX	9	0	71	68	170	157	121	#273	136	153	
MCO	80		EACEPECG	0.02				NYC	85	71	0	34	284	203	157	#203	67	82	
DEV	00		VEDTO	0.03				MCO	75	68	34	0	322	235	164	#238	96	109	
PER	100		VEDIO	0.055				PEK	185	170	284	322	0	127	87	#260	187	390	
HKG	100		SOCIAL	0.04				HKG	159	157	203	235	127	0	49	#279	239	268	
NRT	100	;	SHOPPING	0.035;				NRT	125	121	157	164	87	49	<u>a</u> .	#364	209	273	
#BOM	10	0	#BACKUP	0.04				#BOM	1 280	273	2203	2228	2 260	279	36/	1 0	1/12	135	
#LCY	10	0	#PAY	0.03				#1.01	1 200	1 1 2 7 3	203		: 107	2/9	200	140	142	10	
#FRA	10	0;	<b>#TRAVEL</b>	0.045;				#LCT	101	1 1 2 0		90	200	259	205	142	10	10	
								#FKA	125	5 153	82	103	1 390	268	27:	5 135	10	0;	
param:		vmCapacity:=	param:	required	Latency:=														
AR		3	AR	150	-														
MAP		3	MAP	250	naram lam	bda(tr	<u>۱</u> .												
GAME		5	GAME	180	pur um Ium	SEO		NYC	м		PEK	н	KG	NRT	• _#	BON	1 1	CV	FRA · -
FACERE	CG	3	FACERECG	200	ΔP	a a1	a a2	9 9	3 0	015	a a	3 0	62	a a	" 3 #	a a2	6	15 G	a 015
VEDIO		7	VEDIO	350		0.01	0.02	: 0.0	1 0	010	a a	3 0	02 025	0.0	2 #	0.02	: 0.0	15 (	A 019
SOCIAL		2	SOCIAL	300	CAME	0.02	0.023	0.0	4 U 1 E O	02	0.0	25 0	015	0.0	- 1 - μ	0.01	0.0	17 (	2 021
SHOPPT	NG	4:	SHOPPING	350:	GAME	0.01	0.05	0.0	72 0	010	0.0	20 00	.015	0.0	L #	0.01	0.0	117 6	2.015
#BACKU	P	8	#BACKUP	500	FACERECG	0.05	0.01	0.0	23 0	.018	0.0	26 0	.015	0.0	08 #	0.011	. 0.0	113 6	1.015
#DACKO		1	#DACKOT	250	VEDIO	0.013	0.021	0.0	23 0	.011	0.0	34 0	.016	0.0	3 #	0.015	0.6	126 6	0.031
		+ 2.		250	SOCIAL	0.042	0.025	5 0.0	50 0	.031	0.0	42 0	.017	0.0	21 #	0.03	0.0	1 6	9.012
#IRAVE	L	5;	#IRAVEL	350;	SHOPPING	0.012	0.021	0.0	33 0	.009	0.0	32 0	.044	0.0	3; #	0.029	0.0	)45 (	9.025
					<b>#BACKUP</b>	0.01	3 0.03	8 0.0	04	0.02	1 0.	015	0.02	1 0.0	033	0.035	6.0.	)21 (	0.012
					#PAY	0.05	0.02	25 0.0	033	0.01	1 0.	048	0.02	3 0.0	037	0.036	5 0.0	18 6	029
					#TRAVEL	0.02	0.03	8 0.0	028	0.04	1 0.	032	0.04	5 0.0	055	0.023	0.0	31 (	9.027;



### **Results**

nS
<sup>11</sup> e

•  $X^{S}_{src,dst}$ 

VMINUM [*,*]												
:	AR	FACERECG	GAME	MAP	SHOPPING	SOCIAL	VEDIO	:=				
HKG	2	1	3	0	Θ	4	0					
LAX	2	Θ	Θ	2	2	3	0					
MCO	4	1	4	1	Θ	0	0					
NRT	3	2	Ο	5	Θ	0	5					
NYC	0	3	Ο	2	Θ	0	0					
PEK	3	Θ	3	0	1	0	0					
SF0	1	Θ	3	1	3	0	Θ					

[NYC,*,*]												
:	AR	FACERECG	GAME	MAP	SHOPPING	SOCIAL	VEDIO	:=				
HKG	0	Θ	0	0	0	0	0					
LAX	0	Ο	0	0	0	1	0					
MCO	1	Θ	1	0	0	0	0					
NRT	0	Θ	0	0	0	0	1					
NYC	0	1	0	1	0	0	0					
PEK	0	Θ	0	0	0	0	0					
SF0	0	Θ	Θ	0	1	0	0					



NI.....

#### • Case A: Initial placement

• Place VMs for multiple services on ALL empty MEC clouds.

#### • Case B: Incremental placement

- Place new VMs for one or multiple services on MEC clouds, which already have other VMs.
- Case A and B can be covered by above formulations
- Case C: Dynamic VM management.
  - Heuristics algorithms for service load change.
  - Options: 1)VM clone & migration, 2)VM exchange, 3)Service map optimization.



### Thank you!

Wei Wang

