Optimal Resource Allocation for Analytics and Multimedia Applications in Cloud-to-Things Continuum Platforms

洪華駿

Department of Computer Science, National Tsing Hua University, Taiwan

Agenda

- I. Introduction
- II. Application Deployment Problem
- III. Delay-Sensitive Application Optimization Problem
- IV.Delay-Insensitive Application Optimization Problem
- V. Conclusion & Future Work

Cloud



Cloud-to-Things Continuum Platform



[1] https://www.openfogconsortium.org/

Cloud-to-Things Continuum Framework [1]

	Application Service			
	Application Support			
	Node Management (IB)			
	Virtualization			
	Node Management (OOB)			
	Hardware Security			
	Computing Networks Storage Sensors			
	Hardware Platform Infrastructure			
	Protocol Abstraction			
	Sensors and Actuators			

Intelligent Framework



Three Research Problems

Goal: Serving as many users as possible



Application: Remote Rendering Goal: Reduce latency

Application: Content Delivery Goal: Maximize delivered information

Contributions

Propose an intelligent cloud-to-things continuum framework

Solve application deployment problem

- Deploy applications while considering decompositions
- 1/U approximation factor

Solve delay-sensitive application optimization problem

- Adapt bitrates for game streaming applications
- Optimal solution solved in polynomial time

Solve delay-insensitive application optimization problem

- Provide content delivery plans while considering user behaviors and content representations under challenged networks
- An efficient algorithm outperforms state-of-the art algorithms

Agenda

- I. Introduction
- II. Application Deployment Problem
- III. Delay-Sensitive Application Optimization Problem
- IV.Delay-Insensitive Application Optimization Problem
- V. Conclusion & Future Work

Agenda: Application Deployment Problem

• Problem Statement

- Proposed Algorithm: Ideal Case
- Proposed Algorithm: General Case
- Evaluation & Summary





Application Deployment Problem

Problem: Given a set of requests Q and a set of fog devices V. Each request has three requirements: (i) a splittable application, (ii) a target QoS, and (iii) a specified location. Devices located at various places have U kinds of resources. Our problem is to determine which request $q \in Q$ should be served on which device $v \in V$ without overload the resources. Our goal is to maximize number of served requests.



Challenges:

- Diverse users (requests) \rightarrow diverse resource requirements
- Heterogeneous fog devices \rightarrow application decomposition
- Large problem size \rightarrow polynomial time algorithm

Agenda: Application Deployment Problem

- Problem Statement
- Proposed Algorithm: Ideal Case
- Proposed Algorithm: General Case
- Evaluation & Summary



Application Deployment Problem: Ideal Case

- Splittable applications \bigcirc
 - Arbitrary size of operators
- Transfer link constraints to node constraints
- Grid map \bigcirc
 - Select one of grids
- Reduced device heterogeneity \bigcirc
 - Proportional resource levels

X4 1 core 2 GHz 1 GB RAM

. . .





4 cores 2 GHz 4 GB RAM

. . .



X2

8 cores 2 GHz 8 GB RAM

. . .





✓ Variant of knapsack problem \checkmark NP hard

APproXimation Algorithm (APX)

1. Step 1: Request selection

- Least total required resource first $\min_{q \in \mathbb{Q}} \sum_{u \in \mathbb{U}} F(q, u, \cdots)$
- q: request
- u: resource type
- F: resource consumption model
- 2. Step 2: Device selection
 - Round robin
- 3. Step 3: Request decomposition
 - As large operator as possible



- $O(|Q| \log |Q| + |Q||V||U|)$
- 1/U Approximation Factor
- Q: Request set
- V: Device set
- U: Resource type set

Proof: $\frac{1}{|U|}$ Approximation Factor

- Step1: Design an upper bound solution (OPT')
 - Intuitions are similar to APX
- Step2: Prove results of OPT' >= optimal solution (OPT)
- Step3: Derive approximation factor (APX/OPT')

The algorithm results in optimal solution while |U| = 1

Validating the Approximation Factor

- OPT: formulate as ILP and optimally solve it using CPLEX
- Results of APX are always above theoretical 1/U bound
- APX == OPT while U=1



Limitations: Operators, Links, Location, and Devices → How to make it more generalized?

Agenda: Application Deployment Problem

- Problem Statement
- Proposed Algorithm: Ideal Case
- Proposed Algorithm: General Case
- Evaluation & Summary



Application Deployment Problem: General Case

• Goal: Serve as many requests (users) as possible

	Ideal	Generalized
Operator	Any Size	Predefined
Resource Constraint	Node	Node + Link
Location	Grids	Any Location
Device Heterogeneity	Proportional	Any Device



Problem Formulation

 $\max \sum_{q \in \mathbf{Q}} p_q \qquad \text{Decision variable}$ $st: z_q = \sum \sum \frac{x_{q,i,k}}{|\hat{\mathbf{V}}_{\mathscr{A}_q}|} \quad \forall q \in \mathbf{Q};$ (1b) $i \in \hat{\mathbf{V}}_{\mathscr{A}_{\mathbf{G}}} \ k \in \mathbf{V}_{q,i}$ $z_q - 1 < p_q \leq z_q \ \forall q \in \mathbf{Q};$ (1c) $\sum x_{q,i,k} \le 1 \; \forall q \in \mathbf{Q}, i \in \hat{\mathbf{V}}_{\mathscr{A}_{\mathbf{q}}};$ (1d) $k \in \mathbf{V}_{q,i}$ $\sum \sum \dot{F}(\cdot) x_{q,i,k} \le R_{k,u}$ $\overline{q \in \mathbf{Q}}_{i \in \hat{\mathbf{V}}_{\mathscr{A}_{a}}}$ (1e) $\forall k \in \mathbf{V}_i, u \in \mathbf{U};$ $\sum \sum \sum \sum \sum y_{q,(i,i'),(k,k')}$ $\overline{q \in \mathbf{Q}}_{i \in \widehat{\mathbf{V}}_{\mathscr{A}_{q}}} i' \in \widehat{\mathbf{V}}_{\mathscr{A}_{q}} k \in \overline{\mathbf{V}}_{q,i} k' \in \overline{\mathbf{V}}_{q,i'}$ (1f) $J_{q,i,i'}T_{(k,k'),l}\overline{F}(\cdot) \leq B_l \ \forall l \in \mathbf{E};$ $x_{q,i,k} = \sum_{k' \in \mathbf{V}_{q,i'}} y_{q,(i,i'),(k,k')}$ (1g) $\forall q \in \mathbf{Q}, i \in \mathbf{\hat{V}}_{\mathscr{A}_{q}}, i' \in \mathbf{\hat{V}}_{\mathscr{A}_{q}}, k \in \mathbf{V}_{q,i};$ $x_{q,i',k'} = \sum y_{q,(i,i'),(k,k')}$ (1h) $\forall q \in \mathbf{Q}, i \in \hat{\mathbf{V}}_{\mathscr{A}_{q}}, i' \in \hat{\mathbf{V}}_{\mathscr{A}_{q}} k' \in \mathbf{V}_{q,i};$ $p_q, x_{q,i,k} \in \{0,1\} \ \forall q \in \mathbf{Q}, i \in \hat{\mathbf{V}}_{\mathscr{A}_q}, k \in \mathbf{V}_{q,i}.$ (1i)

^(1a) Objective Throughput maximization Constraints Node Link





Agenda: Application Deployment Problem

- Problem Statement
- Proposed Algorithm: Ideal Case
- Proposed Algorithm: General Case
- Evaluation & Summary



Evaluation Setup

- Requests
 - Poisson arrival / departure rates: 1 min / 10 mins
- Network topology (BRITE [1])
 - Number of fog devices: [10, 25, 50, 75, 100]
 - Location of fog devices and latencies of their edges: based on BRITE
- Resource capacities of fog devices
 - CPU: $[100\% \sim 800\%] \rightarrow$ random
 - RAM: $[1 \text{ GB} \sim 16 \text{ GB}] \rightarrow \text{random}$
- Bandwidth of links:
 - [45 kbps (LoRa), 8 Mbps (WiFi), 25 Mbps (4G)] → Random

Baseline Algorithms

- Designed for fog devices
- Splittable application
- Consider both node and link capacities
- Optimal Data Stream Processing Placement (ODP) algorithm [1]
 - CPLEX
- Fog and Cloud Placement (FCP) algorithm [2]
 - Heuristic
- Linear algorithm
 - Greedily deploys operators to neighbor fog devices
 - Considers all the constraints
- Random algorithm
 - Mimics a platform without centralized server
 - Does not consider any constraint

[2] M. Taneja and A. Davy. 2017. Resource aware placement of IoT application modules in Fog-Cloud Computing Paradigm. In Proc. of IFIP/IEEE Symposium on Integrated Network and Service Management (IM). Lisbon, Portugal.

SSE Serves More Users

• SSE outperforms others in terms of supported users

- 140%, 49%, and 46% compared to Linear, FCP, and ODP
- Random satisfies zero requests after running 14 hours
- Random, FCP, and ODP are too aggressive
 - Overload links by up to 21, 12, and 10 links



Summary: Application Deployment

- Solved ideal and general application deployment problem
 - APX algorithm with $\frac{1}{|U|}$ approximation factor
 - Validated by numerical analysis
 - Satisfy at least 46% more requests compared to state-of-the-art algorithms
- Derived resource consumption models



Our Testbed

Agenda

- I. Introduction
- II. Application Deployment Problem
- III. Delay-Sensitive Application Optimization Problem (2)
- IV.Delay-Insensitive Application Optimization Problem
- V. Conclusion & Future Work

Agenda: Delay-Sensitive Application Optimization Problem

- Problem Statement
- Proposed Algorithm
- Evaluations & Summary





Bitrate Adaptation Problem

Problem: Given a set of gamers G and a shared bottleneck link with available bandwidth B. Our problem is to determine allocating bitrates, which do not overload B to G. Our goal is to maximize gaming experience. The gaming experience model m(g,b) is affected by game types g and allocated bitrates b.



Challenges:

- Various game types and gamers \rightarrow various gaming experience models
- Delay sensitive \rightarrow real-time algorithm

Agenda: Delay-Sensitive Application Optimization Problem

- Problem Statement
- Proposed Algorithm
- Evaluations & Summary



Optimal Bitrate Adaptation Algorithms

• Quality maximization: *EFF_{avg}*

Gamer selection: highest gaming experience improvement first Bitrate adaptation: add one unit (w = 1kbps)

• Quality fairness: *EFF_{mm}*

Gamer selection: worst gaming experience first Bitrate adaptation: add one unit (w = 1kbps)



Optimality of the *EFF_{avg}* Algorithm

- Lemma: EFF_{avg} is optimal if gaming experience model $m^*(g, b)$ is monotonic decreasing
- Proof:
 - EFF_{avg} finds steepest slope at every iteration: $p_1, p_2, p_3 \dots$ $\rightarrow p_i > p_{i+1} \rightarrow No$ alternative allocation
 - ✓ if it is feasible solution (monotonic decreasing) → Optimal
 - Monotonic decreasing: double deviation of $m^*(g, b)$ $\rightarrow m^*(g, b)'' < 0$



 p_n : gaming experience improvement of adding one unit

Optimality of the *EFF_{mm}*Algorithm

- Lemma: EFF_{mm} is optimal if it allocates *w* to gamer *u* who has the worst gaming experience
- Proof:
 - EFF_{mm} allocates w (1kbps) to user u who has the worst gaming experience = l
 - \rightarrow MOS improvement = q
 - Consider alternative user u'
 - \rightarrow Improvement = 0 (the worst gaming experience is still *l*)
 - \rightarrow Need to allocate another *w* to *u* to achieve the same improvement
 - Hence, we cannot find an alternative leading to a better solution $\rightarrow EFF_{mm}$ is optimal

Agenda: Delay-Sensitive Application Optimization Problem

- Problem Statement
- Proposed Algorithm
- Evaluations & Summary



Evaluation Setup

- Gaming experience model: user study
- Available resource = 760 Mbps (PlanetLab)
- Optimal algorithms: CPLEX (an optimization solver)
 - OPT_{avg}
 - *OPT_{mm}*


Gaming Experience Model

- Quadratic model $m(g, f, b) = \alpha_{g,1}f + \alpha_{g,2}b + \alpha_{g,3}f^2 + \alpha_{g,4}b^2 + \alpha_{g,5}fb + \alpha_{g,6}$ Game Frame Bitrate Type Rate • MOS Scores: 1 ~ 7
 - R-square values > 0.98

How to simplify the problem? \rightarrow focus on bitrate allocation



[Hong et al. TCSVT'15, Credits: collaborators @ Academia Sinica]

37

Optimal Frame Rate Selection

- Lemma: m(g, f, b) calculates optimal MOS under certain bitrate if we have $f^*(g, b)$, which gives us optimal frame rate
- Quadratic gaming experience model $m(g, f, b) = \alpha_{g,1}f + \alpha_{g,2}b + \alpha_{g,3}f^2 + \alpha_{g,4}b^2 + \alpha_{g,5}fb + \alpha_{g,6}b^2$
- Partial derivation to f
 → Optimal frame rate: f*(g, b) = (-(α_{g,1} + α_{g,5}b))/(2α_{g,3})
 Optimal MOS:

$$m^{*}(g,b) = \alpha_{g,1}f^{*}(g,b) + \alpha_{g,2}b + \alpha_{g,3}(f^{*}(g,b))^{2} + \alpha_{g,4}b^{2} + \alpha_{g,5}f^{*}(g,b)b + \alpha_{g,6}.$$

 $m(g,f,b) \rightarrow m^*(g,b)$

Our Algorithms are Optimal

- The resulting MOS scores are exactly the same
- Efficient algorithms run much faster than OPT algorithms

# of Gamers	OPT_{avg}	EFF avg	# of Comerc	OPT_{mm}	EFF _{mm}			
	Mean	Mean		Worst	Worst			
100	5.16	5.16	1	5.53	5.53	Same		
200	5.15	5.15	2	5.01	5.01			
400	5.15	5.15	4	4.91	4.91			
800	4.49	4.49	8	4.91	4.91			

[MOS Scores]

# of Gamers	OPT $_{avg}$ EFF		$F_{avg} \parallel \# \text{ of Gamers}$		OPT _{mm}		EFF _{mm}		
	Mean	Max	Mean	Max		Mean	Max	Mean	Max
100	0.02	0.03	0.090	0.097	1	0.02	0.01	0.009	0.009
200	0.15	0.16	0.102	0.103	2	0.36	0.37	0.015	0.016
400	0.95	0.96	0.153	0.157	4	4.47	4.48	0.022	0.023
800	5.01	5.02	0.237	0.248	8	125.8	127.1	0.029	0.029
21X faster			[Running Time]		4]	4K+X faster			

Summary: Delay-Sensitive Applications

- Solved bitrate adaptation problem
 - Two optimal algorithms
 - Outperform baselines by up to 30% and 46%, respectively
- Solved bandwidth estimation
 - ► >80% accuracy
- Solved real-time codec reconfiguration [Hong et al. TCSVT'15, Credits: collaborators @ NCTU]
- Derived gaming experience models [Hong et al. TCSVT'15, Credits: collaborators @ Academia Sinica]



Agenda

- I. Introduction
- II. Application Deployment Problem
- III. Delay-Sensitive Application Optimization Problem
- IV.Delay-Insensitive Application Optimization Problem 3
- V. Conclusion & Future Work

Agenda: Delay-Insensitive Application Optimization Problem

- Problem Statement
- Proposed Algorithm
- Evaluation & Summary



Content Delivery Under Challenged Network



Distribution Planning Problem

Problem: Given a set of contents N and a set of users U. Each content has Lrepresentations. Each user has C contacts (run into fog devices). Each contact has Rresources. Our problem is to determine delivering which representation $l \in L$ of which content $n \in N$ to which user $u \in U$ at which contact $c \in C$ without overloading the resources R. Our goal is to maximize user experience.



Challenges:

- Definition of user experience \rightarrow user study
- Various user behaviors \rightarrow individual distribution plans
- Limited resources and large size of contents \rightarrow multiple representations

Agenda: Delay-Insensitive Application Optimization Problem

- Problem Statement
- Proposed Algorithm
- Evaluation & Summary



Distribution Planning Problem Formulation

Decision $\max \sum_{l=1}^{U} \sum_{l=1}^{N} \sum_{l=1}^{L} \sum_{l=1}^{C_{u'}} \frac{\text{Variable}}{x_{u,n,l,c}} \rho_{nL+l} \psi_{u,n} \quad (1a)$ u=1 n=1 l=1 c=1N = L $st: \sum \sum b_{nL+l} x_{u',n,l,c'} \le R_{u',c'};$ n=1 l=1 $\sum_{c=1}^{C_{u'}} x_{u',n',l',c} \ge \sum_{c=1}^{C_{u'}} x_{u',n',l'',c};$ $\psi_{p_{u',c'},n'} \ge \bar{\psi} x_{u',n',l',c'};$ $C_{u'}$ $\sum x_{u',n',l',c} \le 1;$ c=1

- UE **Objective:** Maximize User Experience
- (1b) Resource Budgets R
- (1c) Layer Dependency
- (1d)Viewing Probability Checker
- (1e)**Duplication Avoidance**

Dynamic Programming

- Termination: $(1b) \sim (1e)$
- Recursion & Memorization:

$$DP(i, R_1, R_2, ..., R_C) = \max\{$$

$$DP(i - 1, R_1 - b[i], R_2, R_3 ...) + UE,$$

$$DP(i - 1, R_1, R_2 - b[i], R_3 ...) + UE,$$

$$DP(i - 1, R_1, R_2, R_3 - b[i] ...) + UE,$$

...

$$DP(i - 1, R_u)\}$$



(1b) Resource Budgets

R

47

- (1c) Layer Dependency
- (1d) Viewing Probability Checker
- (1e) Duplication Avoidance
- Time & Space complexity: $O(NL \Pi R_c)$

DP can optimally solve the problem under delay-tolerable environment

No. Contents	1	2	3	4	5	6	7
Running Time (sec)	0.08	0.11	0.19	0.61	1.8	111.2	457.2
Memory Usage (MB)	21	32	103	269	1918	10197	>64GB

Algorithm: Contact-Driven Round Robin (CDRR)

User Selection Round Robin R **Content Selection** 2. Higher user R experience using less resources first More contents for UE popular users **Contact Selection** Unpopular fog R devices first

O(UNLlog(UNL) + UNLC)



No.	Total User Experience					
Content	DP	CDRR	Perf. Gap			
1	0.43	0.42	3%			
2	0.36	0.34	6%			
3	0.33	0.31	6%			
4	0.31	0.29	6%			
5	0.27	0.25	7%			
6	0.29	0.27	7%			
7	0.29	0.27	7%			

DP vs CDRR Performance Gap: 7%

Agenda: Delay-Insensitive Application Optimization Problem

- Problem Statement
- Proposed Algorithm
- Evaluation & Summary



Evaluation Setup

- User experience: user study
- Trace: a real testbed
- Baseline algorithms
 - CSI [1]: Consider trajectories of users User Behaviors
 - Epidemic [2]: Classic content delivery algorithm in challenged networks Classic Algorithm



CSI: Deliver to users having different trajectories

Epidemic: Broadcast to each other

[1] W. Hsu, D. Dutta, and A. Helmy, "CSI: A paradigm for behavior- oriented profile-cast services in mobile networks," *Ad Hoc Networks*, vol. 10, no. 8, pp. 1586–1602, 2012.
[2] A. Vahdat and D. Becker, "Epidemic routing for partially connected ad hoc networks," Duke University, Tech. Rep., 2000.

User Experience of Different Representations

- Definition of user experience: understanding level
- 182 participants (recruited)
 - 120 participants (filtered) with 587 samples
- Content: Apple Daily News
- 5-layer representations: text, audio, low-, medium-, and high- resolution videos

CROWDSOURCED USER EXPERIENCE SCORES

	Article	Audio	240p Video	360p	480p
Average	55%	68%	71%	74%	77%
No. Samples	112	112	134	104	125
			MP4		

Today, China Post holds a recruitment with 30000~ 43000 NTD monthly salary. It attracts 40775 applications, including 26 PhDs and 2880 masters, which sets a new high. There are 1571 vacancies and 90% of the applyer attenning the interview.



* required field.

Age: _____ * Gender: OFemale OMale *

What's the understanding level? (the higher the better):

○1 ○2 ○3 ○4 ○5 *

Which company holds the recruitment?:

 \bigcirc China Post \bigcirc China Air \bigcirc I do not know *

Are there many masters applying for the job?:

 \bigcirc Yes \bigcirc No \bigcirc I do not know *

How many people attending the interview?:

 \bigcirc >40,000 \bigcirc <40,000 \bigcirc I do not know *

What did China Post hold?:

 \bigcirc Recruitment \bigcirc Exhibition \bigcirc I do not know * submit

Implementation

- Android Application
 - 15 Users
- Fog device: 11 Raspberry PIs
- Location: 3 rural villages + a university
- Content: 46 news from apple daily





CDRR Results in Higher User Experience Using Less Energy

- CDRR Outperforms others:
 - User experience: 1.1 times (CSI) and 2.7 times (Epidemic)
 - Energy efficiency: 1.2 times (CSI) 1.4 times (Epidemic)

Because of detailed individual plans



Summary: Delay-Insensitive Application

- Solved distribution planning problem
- Derived user experience models of different representations
- Evaluated our algorithm with real implementation
 - Outperform state-of-the-art and classic algorithms by 110% and 270% in terms of user experience, respectively



Agenda

- I. Introduction
- II. Application Deployment Problem



- III. Delay-Sensitive Application Optimization Problem
- IV.Delay-Insensitive Application Optimization Problem
- V. Conclusion & Future Work

Conclusion

- Help providers to realize the cloud-to-things continuum platform
 A performance guaranteed deployment algorithm makes providers' life easier to estimate costs/benefits and formulate prices
- Help users to have optimized user experience
 - Game streaming applications with optimal gaming experience
 - Content delivery applications with maximized understanding level





Future Work: Testbed in NTHU

- 8 Smart Street Lamps in NTHU
- Street Lamp
 - Fog devices: Raspberry Pis and IPCs
 - Sensors: camera, air pollution sensors, ...
 - Communications: ethernet, WiFi mesh, LoRa, Zigbee, and Bluetooth
- Analytic applications: object recognition, car plate recognition, ...





LTEON

Killer Applications of Our Platform

- Analytic Applications
 - Complicated deep learning models
 - \rightarrow huge amount of resources and data is required
- IoT Applications [1]
 - Low end devices
 → need longer time to run complicated applications







Ecosystem



For Application Developers

- Challenges
 - High learning curve to implement multi-operator applications → programming model
 - Privacy sensitive data → data management and protection approaches



Models

Developers

Programing Model

- Decomposition
 - Dynamically split applications into any number of smaller operators
- Communication
 - Automatically handle the communications between operators
- Performance Optimization
 - Stream processing
 - Parallel programming



Data Management and Protection Approaches

- Training as a Service (TaaS)
 - Help developers to train their models
- Data analytic APIs
 - Help developers analyze the data and figure out some characteristics
 - Help developers to figure out root cause of low accuracy



Publications

Journal Articles



H. Hong, C. Hsu, T. Tsai, C. Huang, K. Chen, and C. Hsu, "Enabling Adaptive Cloud Gaming in an Open-Source Cloud Gaming Platform," IEEE Transactions on Circuits and Systems for Video Technology (TCSVT), vol. 25, no. 12, December, 2015.



H. Hong, T. El-Ganainy, C. Hsu, K. Harras, and M. Hefeeda, "Disseminating Multi-layer Multimedia Content over Challenged Networks," IEEE Transactions on Multimedia (TMM), vol.20, no.2, February, 2017.

Conference Papers



H. Hong, P. Tsai, A. Cheng, M. Uddin, N. Venkatasubramanian, and C. Hsu, "Supporting Internet-of-Things Analytics in a Fog Computing Platform," in Proc. of IEEE International Conference on Cloud Computing Technology and Science (CloudCom), Hong Kong, December, 2017. **(Best Paper Award)**



H. Hong, P. Tsai, and C. Hsu, "Dynamic Module Deployment in a Fog Computing Platform," in Proc. of IEEE Asia-Pacific Network Operations and Management Symposium (APNOMS), Kanazawa, Japan, October, 2016. **(Best Paper Award)**

Demo Paper



H. Hong, S. Wang, C. Tan, T. El-Ganainy, K. Harras, C. Hsu and M. Hefeeda, "Challenged Content Delivery Network: Eliminating the Digital Divide," in Proc. of ACM Multimedia Demo Paper (MM Demo), Brisbane, Australia, Oct., 2015.

PhD Symposium



Hua-Jun Hong, "From Cloud Computing to Fog Computing: Unleash the Power of Edge and End Devices," in Proc. of IEEE International Conference on Cloud Computing Technology and Science (CloudCom), Hong Kong, December, 2017.

Other Publications

Book Chapter

- C. Hsu, **H. Hong**, T. El-Ganainy, K. Nahrstedt, and N. Venkatasubramanian, "Multimedia Fog Computing: Minions in the Cloud and Crowd," ACM Frontiers of Multimedia Research, Chapter 10, Association for Computing Machinery and Morgan & Claypool, January 2018.

Journal Articles

- **H. Hong**, C. Fan, Y. Lin, and C. Hsu, "Optimizing Cloud-Based Video Crowdsensing," IEEE Internet of Things Journal (JIoT), vol. 3, no. 3, January, 2016.
- **H. Hong**, D. Chen, C. Huang, K. Chen, and C. Hsu, "Placing Virtual Machines to Optimize Cloud Gaming Experience," IEEE Transactions on Cloud Computing (TCC), vol. 3, no. 1, 2014.

Conference and Workshop Papers

- M. Rahman, A. Rahman, A. Afrin, **H. Hong**, P. Tsai, M. Uddin, N. Venkatasubramanian, and C. Hsu, "Adaptive Sensing Using Internet-of-Things with Constrained Communications," in Proc. of ACM Adaptive and Reflective Middleware (ARM), Las Vegas, NV, USA, December, 2017.
- P. Tsai, **H. Hong**, A. Cheng, and C. Hsu, "Distributed Analytics in Fog Computing Platforms Using TensorFlow and Kubernetes," in Proc. of IEEE Asia-Pacific Network Operations and Management Symposium (APNOMS), Seoul, Korea, September, 2017.
- Y. Chen, **H. Hong**, S. Yao, A. Khunvaranont, and C. Hsu, "Gamifying Mobile Applications for Smartphone Augmented Infrastructure Sensing," in Proc. of IEEE Annual Workshop on Network and Systems Support for Games (NetGames), Taipei, Taiwan, June, 2017.
- **H. Hong**, Y. Lin, J. Chuang, and C. Hsu, "Animation Rendering on Multimedia Fog Computing Platforms," in Proc. of IEEE International Conference on Cloud Computing Technology and Science (CloudCom), Luxembourg, December, 2016.
- T. Fan-Chiang, **H. Hong**, C. Hsu, "Segment-of-Interest Driven Live Game Streaming: Saving Bandwidth without Degrading Experience," in Proc. of IEEE Annual Workshop on Network and Systems Support for Games (NetGames), Zagreb, Croatia, December, 2015.
- **H. Hong**, C. Lee, K. Chen, C. Huang, and C. Hsu, "GPU Consolidation for Cloud Games: Are We There Yet?," in Proc. of IEEE Annual Workshop on Network and Systems Support for Games (NetGames), Nagoya, Japan, December, 2014.

Poster and Demo Papers

- Y. Hsieh, **H. Hong**, P. Tsai, Y. Wong, Q. Zhu, M. Uddin, N. Venkatasubramanian, and C. Hsu, "Managed Edge Computing on Internet-of-Things Devices for Smart City Applications," in Proc. of IEEE/IFIP Network Operations and Management Symposium (NOMS), Taipei, Taiwan, April, 2018.
- Q. Zhu, M. Uddin, N. Venkatasubramanian, C. Hsu, and **H. Hong**, "Enhancing Reliability of Community Internet-of-Things Deployments with Mobility," in Proc. of IEEE INFOCOM, Honolulu, HI, USA, April, 2018.
- **H. Hong**, D. Chen, C. Huang, K. Chen, and C. Hsu, "QoS-Aware Virtual Machine Placement for Cloud Games," in Proc. of ACM Annual Workshop on Network and Systems Support for Games (NetGames), Denver, CO, December 2013.



hua.j.hong@gmail.com

Introduction

Fog similar concepts comparisons Timeline (from WSN to clout-to-things continuum) Benefits of Cloud-to-Things Continuum

Cloud, Distributed Cloud, Cyber Foraging, CloudLet, and Fog

	Low Latency	Location Awareness and Mobility Support	Virtualization Support	High Heterogeneity
Cloud	Х	×	\checkmark	×
Distributed Cloud [1]	Δ	Δ	\checkmark	×
Cyber Foraging [2]	\checkmark	Δ	×	Δ
CloudLet and MEC [3,4]	\checkmark	\checkmark	\checkmark	Δ
Fog [5]	\checkmark	\checkmark	\checkmark	\checkmark

[1] P. Endo, A. de Almeida Palhares, N. Pereira, G. Goncalves, D. Sadok, J. Kelner, B. Melander, and J. Mangs, "Resource Allocation for Distributed Cloud: Concepts and Research Challenges," *IEEE Network*, 25(4), 42-46, 2011.

[2] M. Satyanarayanan. Pervasive computing: Vision and challenges. IEEE Personal communications, 8(4):10–17, 2001.

[3] M. Satyanarayanan, V. Bahl, R. Caceres, and N. Davies, "The Case for VM-based Cloudlets in Mobile Computing. *IEEE pervasive Computing*, 8(4), 14-23, 2009.

[4] Mobile-edge computing. https://portal.etsi.org/Portals/0/TBpages/MEC/Docs/Mobile-edge_Computing_-

_Introductory_Technical_White_Paper_V1%2018-09-14.pdf.

[5] F. Bonomi, R. Milito, J. Zhu, and S. Addepalli, "Fog Computing and Its Role in the Internet of Things" *in Proc. of ACM SIGCOMM workshop on Mobile Cloud Computing (MCC)*, Helsinki, Finland, 2012.

Timeline



Benefits of Cloud-to-Things Continuum



Application Deployment Problem

Knapsack problem variations Approximation factor proof Resource consumption model

Knapsack Problem Variations

- Subset Sum Problem (SSP): weight = profit
- Bounded Knapsack Problem (BKP): duplicated items with a boundary
- Unbounded Knapsack Problem (UKP): duplicated items without a boundary
- Multi dimensional Knapsack Problem (d-KP): multiple constraints
- Multiple Knapsack Problem (MKP): multiple knapsacks
- Multiple Knapsack Problem with Identical Capacities (MKP-I): knapsacks have the same capacities
- Multiple-Choice Knapsack Problem (MCKP): items have variations (diff weight & profits), need to pick one variation of each item
- Quadratic Knapsack Problem (QKP): two items result in a weight & profit
- Bin Packing (BP): use least number of knapsack to pack all items
- MKP+SSP -> MSSP
- \circ MKP-I + SSP -> MSSP-I
- $\circ MCKP+d-KP \rightarrow MMKP$

^[1] H. Kellerer, U. Pferschy, and D. Pisinger, "Knapsack Problems," Springer, 2004

^[2] Silvano Martello, "Knapsack Problems: Algorithms and Computer Implementations," Wiley & Sons, 1990

^[3] C. Wilbaut, S. Hanafi, and S. Salhi, "A Survey of Effective Heuristics and Their Application to a Variety of Knapsack Problems. IMA Journal of Management Mathematics, 19(3), 227-244, 2008.



 H. Kellerer, U. Pferschy, and D. Pisinger, "Knapsack Problems," Springer, 2004
 Silvano Martello, "Knapsack Problems: Algorithms and Computer Implementations," Wiley & Sons, 1990
 C. Wilbaut, S. Hanafi, and S. Salhi, "A Survey of Effective Heuristics and Their Application to a Variety of Knapsack Problems. IMA Journal of Management Mathematics, 19(3), 227-244, 2008.
Approximation Factor Proof

- Step1: Design an upper bound solution (OPT')
 - Intuitions are similar to APX
- Step2: Prove results of OPT' >= optimal solution (OPT)
- Step3: Derive approximation factor (APX/OPT')

Step1: Design an Upper Bound Solution (OPT')

- Similar intuitions with APX
- Differences
 - Merge all the fog devices
 - Sum up different type of resources
- Sort requests by total required resources



Step2: Prove that OPT' is Upper Bound

- OPT: \dot{H} satisfied requests $(\dot{r}_1, \dot{r}_2, \cdots \dot{r}_{\dot{H}})$
- OPT': H' requests $(r'_1, r'_2, \cdots, r'_{H'})$, which may/may not be satisfied
 - Results are sorted by total required resources



Step3: Derive Approximation Factor

(1)

(2)

• APX: *H* satisfied requests $(r_1, r_2, \cdots r_H)$





$$\sum_{u \in \mathbf{U}} F(r_h, u, \cdots) \le \sum_{u \in \mathbf{U}} F(r_{h+1}, u, \cdots) \quad \forall 1 \le h \le H$$
(3)

(1) and (2)
$$\frac{H}{H'} \ge \frac{1}{|U|} \frac{B}{A} \implies \frac{H}{H'} \ge \frac{1}{|U|} \implies \frac{H}{\dot{H}} \ge \frac{1}{|U|} \implies \frac{H}{\dot{H}} \ge \frac{1}{|U|}$$
$$B \ge A \qquad H' \ge \dot{H}$$

Remove the Proportional Heterogeneity Assumption



$$\sum_{k \in \mathbf{V}} \sum_{u \in \mathbf{U}} R_{k,u}$$

$$\sum_{k \in \mathbf{V}} \min R_{k,u \in \mathbf{U}}$$

Resource Consumption Models

- Input:
 - Applications
 - Target QoS
 - Hardware specifications
- Output: required resources
 - CPU load
 - RAM usage
 - Network throughput



Testbed

- Fog devices: 5 Raspberry PI3s + 5
 PCs (i7/8GB RAM)
 - Docker + Kubernetes
- Server: a tiny PC (i5/8GB RAM)
 - Kubernetes
- Switch
 - Private networks
- Analytic applications (implemented with TensorFlow)
 - Air quality monitor (QoS: 0.25 4 Hz)
 - Sound classifier (QoS: 6/60 10/60 Hz)
 - Object recognizer (QoS: 5/60 9/60 Hz)





Derived Models on Raspberry PIs

- Sample figures from the object recognizer
 - Constant model for RAM
 - Power models for CPU and network throughput



How about new devices?

Applications	CPU Load	Network Throughput			
Air Quality Monitor	0.99 (0.000023)	0.99 (0.000001)			
Sound Classifier	0.85 (0.023263)	0.84 (0.069494)			
Object Recognizer	0.99 (0.012540)	0.97 (0.004982)			

Bootstrapping Online Regression

- If PC is a new type of fog device adding to our platform
 - Bootstrap from Raspberry PIs' system models
 - 15 Iterations
 - Error: CPU = 2% and RAM = 13% on average



SSE is Resource Efficient and Real-Time Analytics Friendly

- Turns off 6% and 12% fog devices compared to ODP and FCP, respectively
- Use fewer CPU/RAM/bandwidth to satisfy more requests
- Lower latency compared to others



Running Time

No. Devices	10	25	50	75	100	150
SSE	1.07	2.41	5.36	8.88	9.78	13.52
ODP	2.37	22.66	101.93	382.83	>500	>500
FCP	0.07	0.10	0.21	0.27	0.31	0.46

Requests	500		1000		1500		2000		2500	
U	APX	OPT	APX	OPT	APX	OPT	APX	OPT	APX	OPT
1	0.25	1.76	0.71	3.35	1.32	5.25	1.92	7.27	3.46	12.67
2	0.42	2.14	1.06	4.87	2.49	9.7 9	3.69	16.43	5.41	>500
3	1.01	2.64	2.73	7.86	4.44	203.65	6.81	>500	8.74	>500
4	1.94	3.27	4.61	>500	7.57	>500	10.16	>500	12.52	>500
5	2.63	5.33	5.76	>500	8.88	500	11.88	>500	15.21	>500

Any Size of Operators

• A sound classifier application (4-layer NN model)

• 348 operators



Delay-Sensitive Application Problem

Optimality Proof Problem Formulation Bandwidth Estimator Gaming Experience Model Performance Comparing to Baselines Running Time of our Algorithms

Optimality of the *EFF_{avg}* Algorithm

- Lemma: *EFF_{avg}* is optimal if *m**(*g*, *b*) is monotonic decreasing
 Proof:
 - *EFF_{avg}* finds steepest slope at every iteration: p₁, p₂, p₃ ...
 → p_i > p_{i+1} → No alternative allocation
 ✓ if it is feasible solution (monotonic decreasing) → Optimal
 - Monotonic decreasing: double deviation of $m^*(g, b)$ (negative) $\rightarrow m^*(g, b)'' = 2\alpha_{g,4} - \alpha_{g,5}^2/2\alpha_{g,3} < 0$



Our models satisfy this equation

 p_n : quality improvement of adding one unit

Optimality of the *EFF_{mm}*Algorithm

- Lemma: EFF_{mm} is optimal if it allocates w to gamer u who has the worst MOS
- Proof:
 - EFF_{mm} allocates w (1kbps) to user u who has the worst MOS (l) \rightarrow MOS improvement = q
 - Consider alternative *u*′
 - \rightarrow MOS improvement = 0 (the worst MOS is still *l*)

 \rightarrow Need to allocate another *w* to *u* to achieve the same MOS improvement

• Hence, we cannot find an alternative leading to a better solution $\rightarrow EFF_{mm}$ is optimal

Bitrate Adaptation Problem Formulation

Decision $maximize \sum^{U} m_{p_u, \underline{x_u}, \underline{y_u}}$ Variables Objective 1: Gaming Experience Maximization u=1IJ $s.t.: \sum y_u \leq R;$ **Bandwidth Constraint** u=1 $\check{k}_{n'} < y_{n'} < \hat{k}_{n'}, \ \forall \ 1 \le u' \le U;$ $1 \leq x_u \leq F, 1 \leq y_u \leq B, \ \forall \ 1 \leq u \leq U.$ Objective 2:

 $maximize \min_{u=1}^{U} m_{p_u, x_u, y_u}$

Objective 2: Quality Fairness

Bandwidth Estimator

- WBest [1]
 - Estimate bandwidth capacity by sending probing packets



- Our bandwidth estimator
 - Leverage existing video packets as probing packets
 → Avoid additional network overhead

The Bandwidth Estimator is Accurate

- Send W packets and select the median value as estimated bandwidth capacity
- $\circ W = \{100, 200, 300, 400, 500\}$
 - W = 100 \rightarrow accuracy < 50%
 - W = 200 ~ 500 \rightarrow 80% < accuracy < 85%

Gaming Experience User Study

- Gaming experience \rightarrow Mean Opinion Scores (MoS): 1~7
- o 101 Subjects
- 3 type of games = {Batman, FGPX, CoD}
- Bitrate = $\{0.5, 1, 2\}$ Mbps
- Frame rate = $\{10, 20, 30\}$ fps
- Game sessions = $2 \sim 4$ minutes



[Hong et al. TCSVT'15, Credits: collaborators @ Academia Sinica]

Codec Reconfiguration [Hong et al. TCSVT'15, Credits: collaborators @ NCTU]

• Migration from *ffmpeg* to *live555*

- ffmpeg does not have an interface for dynamically reconfiguring video codecs
- live555 offers a more comprehensive RTCP implementation, which allows us to measure and collect the network flow statistics

Our Algorithms Outperform Baselines: Setup

- Optimal algorithms: CPLEX
 - OPT_{avg}
 - *OPT_{mm}*
- Baseline algorithms
 - *Base_{eq}*: equally allocate



- *Base_{norm}*: proportionally (to average MOS scores) allocate
- Bottleneck link capacity R = 760 Mbps (PlanetLab)

Our Algorithms Outperform Baselines

- \circ *EFF_{avg}* results in the best average MOS scores
 - Outperforms baselines by up to 30%
- EFF_{mm} results in the best worst MOS scores
 - Outperforms baselines by up to 46%



Our Algorithms run in Realtime

 \circ 4000 gamers ~ 1 second

~										
# of Gamers	EFI	F_{avg}	\mathbf{EFF}_{mm}							
	Mean	Max	Mean	Max						
500	0.181	0.183	0.179	0.184						
1000 2000	0.296	0.299	0.287	0.290						
	0.523	0.531	0.520	0.533						
4000	1.000	1.104	1.060	1.066						
8000	1.677	1.681	1.654	1.661						

Table 5.5: Running Time in Seconds

Delay-Insensitive Application Problem

Dynamic programming problem CDRR: Practical concerns Block diagram Energy consumption measurement

Dynamic Programming

$$d(i,j) = \max\{d(i-1,j), u_j + d(i-1,j-w_i)\}$$

Weight restriction j (j=1,2, ..., W)

<u>;</u>	\nearrow	1	2	3	4	5	6	7	8	9	10	11
4 :	1	1	1	1	1	1	1	1	1	1	1	1
ς Υ		{1}	{1}	{1}	{1}	{1}	{1}	{1}	{1}	{1}	{1}	{ 1}
, v	2	1	1	1	1	1	1	7	8	8	8	8
<u>i</u>		{1}	{1}	{1}	{1}	{1}	{1}	{2}	{2,1}	{2,1}	{2,1}	{2,1}
ed	3	1	1	1	4	5	5	\bigcirc	8 Targat	8	8	11
oick		{1}	{1}	{1}	{3}	{3,1}	{3,1}	{2}	{2,1}	{2,1}	{2,1}	{3,2}
ns p	4	1	2	3	4	5	6	7	8	9	10	11
lter		{1}	{4 }	{4,1}	{3}	{3,1}	{4,3}	{4,3,1}	{2,1}	{4,2}	{4,2,1}	{3,2}

Complexity: O(IW)

Practical Concerns

- Planed contents are downloaded or plans are not received
 - Determining the downloading order
 - Next outstanding representation of each content (sorted by viewing $\frac{\text{probability}}{\text{content size}}$)
- Video size is too large
 - Segmenting video layers

Block Diagram

- Content Matcher
 - Topia Term Extractor
 - Google Bayesian
- Contact Predictor
 - Frequency-based approach



Prediction Accuracy



DP vs CDRR

No	Ru	inning T	lime (se	c)		No	Total User Experience			
NU.	D	P	CDRR			INU.	рр	CDDD	Dorf Con	
Content	Mean	Max	Mean	Max	Content		DP	CDKK	Peri. Gap	
1	0.08	0.08	0.05	0.05		1	0.43	0.42	3%	
2	0.11	0.11	0.06	0.07		2	0.36	0.34	6%	
3	0.19	0.19	0.63	0.66		3	0.33	0.31	6%	
4	0.6	0.61	0.07	0.07		4	0.31	0.29	6%	
5	1.8	1.8	0.07	0.08		5	0.27	0.25	7%	
6	110.2	111.2	0.08	0.08		6	0.29	0.27	7%	
7	269.4	457.2	0.08	0.09		7	0.29	0.27	7%	

How to Know the Energy Consumptions?

- Power meter: Agilent 66321D
- Read trace:
 - Downloaded news
 - Downloaded plan
 - Uploaded user profiles

• 1-day distribution plan only consumes up to 153 J, which is 0.3%



Fog Testbed and Performance Measurement

Testbed





104

Distributed Analytics Result in Better Performance for Complex Analytics

- Air pollution monitor: too simple to benefit from distributed analytics
- Object recognizer:
 - One more device gives 35.5% and 54.1% improvements



We only consider the object recognizer in the following results

Cut on Equal Complexity Point Results in Better Performance

Cut

Device 2

Device 1

- Setup: run object recognizer on two fog devices (different cut points)
- Cut point 4&5 result in the most no. processed images
- Two fog devices consume the same CPU resources on cut point 4&5
 Equal Complexity Point





Cut Points Affect Network Overhead

Put more computations on device 1 can reduce network overhead



Low Virtualization Overhead

- Setup: with and without Docker
- Overhead from Docker


Low Communication Overhead

- Setup: without Docker and distributed speedup
- Overhead from distributed computing using TensorFlow
 - 10%



109