A Hierarchical Edge Cloud Architecture for Mobile Computing

Liang Tong, Yong Li and **Wei Gao** University of Tennessee – Knoxville



IEEE INFOCOM 2016

Cloud Computing for mobile devices

 Contradiction between limited battery and complex mobile applications



- Mobile Cloud Computing (MCC)
 - Offloading local computations to remote execution
 - Reduced computation delay

Increased communication delay
THEUNIVERSITY OF
TENNESSEE
TENNESSEE
TENNESSEE
TENNESSEE

The limits of Cloud Computing

- Network communication latency of MCC
 - Can be up to 400 ms
 - Performance degrades! Many mobile apps are delay-sensitive

Round trip cities	Max(ms)	Mean(ms)	Min(ms)
Berkeley-Canberra	174.0	174.7	176.0
Berkeley-Troudheim	197.0	197.0	197.0
Pittsburgh-Hong Kong	217.0	223.1	393.0
Pittsburgh-Seatle	83.0	83.9	84.0
Pittsburgh-Dublin	115.0	115.7	116.0



Existing solution

- Small scale cloud servers at the edge
 - Reduce the network latency accessing data center
 - Support user mobility





The limits of Cloudlet

- Cloudlet has limited computing resources
 - A large amount of peak load atency
 - More capacity?



Our solution

Motivation

KNOXVILLE

 Peak loads at different edge cloud servers do not appear at the same time



Our solution

Key idea

KNOXVILLE

- Hierarchical edge cloud architecture
 - Opportunistically aggregate peak loads
 - Improve the resource utilization



Our solution

- Key problems
 - How to efficiently provision edge cloud capacity?
 - How to appropriately place mobile workload at different tiers of servers?
- Our work
 - Formally study the characteristics of the peak load
 - Analyze the efficiency of capacity provisioning
 - Design a workload placement algorithm to further improve the efficiency of program execution



Formal study of the peak load

- System model
 - m tier-1 server and 1 tier-2 server
 - *C*: Computational capacity of the tier-2 server
 - c_i and w_i: computational capacity and workload of the i-th tier-1 server
 - When $w_i > c_i$, a workload of $\eta_i = w_i c_i$ will be offloaded to tier-2. Tier-2 (C) Peak load





IEEE INFOCOM 2016

Formal study of the peak load

- Tier-1 workload model
 - CDF of the peak load Characteristics of workload exceeding C_i • P(n_i ≤ x) = {P(w_i ≤ x + c_i) if x ≥ 0 0 otherwise
- Tier-2 workload model
 - Characteristics of tier-2 workloads

$$\mathbf{P}(\sum_{i=1}^{m} \eta_i) \le x) = \mathbf{P}(\sum_{i=1}^{m-1} \eta_i \le x) \times \mathbf{P}(\eta_m = 0)$$

+
$$\int_{0^+}^{x} \mathbf{P}(\sum_{i=1}^{m-1} \eta_i \le x - t) \, \mathrm{d}\mathbf{P}(\eta_m \le t)$$

Workload of tier-2 server



Formal study of the peak load

Provisioning of edge cloud capacity



- Insights
 - Hierarchical edge cloud has a higher chance to successfully serve the peak loads with the same capacity provisioned.



Objective

- Minimize the total delay of executing all programs
- Our focus
 - Where to place a mobile program
 - How much capacity to each program
- Challenge
 - Computation/communication delay tradeoff
 - delay = computation + communication
 - Higher tiers: less computational delay, but more communication delay





- Problem formulation
 - *m* programs at tier-1, servers in a tree-topology

Computation Communication
delay delay
min
$$f = \sum_{i=1}^{m} \left(\frac{w_i}{\lambda_{i,y_i} c_{\gamma_i}} + (L(\gamma_i) - 1) \right)_{B_{\gamma_i}} s_i \right),$$

s. t. $\sum_{j \in O_j} \frac{\lambda_{i,j}}{\lambda_{i,j}} = 1, j = 1, 2, ..., n$
Capacity allocation
of server *j* to workload *i*

- Nonlinear Mixed Integer Programming
- Challenge: γ_i and $\lambda_{i,j}$ have interdependency



IEEE INFOCOM 2016

Problem transformation

 $\min f(\boldsymbol{\lambda}|\boldsymbol{\gamma} = \boldsymbol{\gamma}^*)$ s. t. $g(\boldsymbol{\lambda}|\boldsymbol{\gamma} = \boldsymbol{\gamma}^*) = 0$ min $f(\boldsymbol{\lambda}, \boldsymbol{\gamma})$ s.t. $g(\boldsymbol{\lambda}, \boldsymbol{\gamma}) = 0$ Non-linear Mixed Convex optimization with variable λ Integer Programming How to determine optimal workload placement γ ? Integer Programming



IEEE INFOCOM 2016

Solution: Simulated Annealing (SA)

- Basic idea
 - Local optima avoidance: accepting a new state which has a worse value with an acceptance probability
- Settings
 - State: workload placement vector γ
 - Value $f(\mathbf{\gamma})$: optimal value of corresponding convex optimization problem
 - Acceptance probability

$$P = exp(-\frac{\Delta d}{T})$$
 Annealing temperature,
decreases in each iteration

Convergence



- **Solution: Simulated Annealing**
 - Value (total delay)



System experimentation

- Comparisons
 - Flat edge cloud
- Evaluation metric
 - Average completion time: indicates computational capacity
- Experiment settings
 - Workload rate
 - Provisioned capacity



Evaluation setup

- Evaluation with a computing-intensive application
 - SIFS of images
- Edge cloud topology
 - Flat edge cloud: two tier-1 servers
 - Capacity is equally provisioned to each server
 - Hierarchical edge cloud: two tier-1 and one tier-2 server
 - Capacity is provisioned to the tier-2 server and tier-1 servers
- Experiments
 - 5 minutes with different size of images



Offloading performance

- Maximum capacity: 4 concurrent threads
 - More capacity provisioned, more improvement





Offloading performance

Maximum capacity: 4 concurrent threads

Only limited improvement at low workload





Simulation experimentation

Comparisons

Four edge clouds with different topologies and capacity provisioning



- Evaluation metric
 - Average delay: includes both computation and communication delay



Simulation setup

- Evaluation with real trace from Wikipedia
 - Randomly select one segment
- Computational capacity provisioning
 - 40 GHz to be provisioned to each topology
- Network setup
 - Two edge cloud servers are connected via 100 Mbps Ethernet
- Experiments
 - 1000 user requests during each simulation



Effect of computation amounts

- Workload placement algorithm is used
 - Data size: normal distribution with an average of 5 MB



Up to 40% delay deduction



IEEE INFOCOM 2016

Effect SA cooling parameter

- Performance when cooling parameter varies
 - Insights: there exists a tradeoff between performance and overhead of workload placement





Summary

- Offloading computations to remote cloud could hurt the performance of mobile apps
 - Long network communication latency
- Cloudlet could not always reduce response time for mobile apps
 - Limited computing resources
- Hierarchical edge cloud improve the efficiency of resource utilization
 - Opportunistically aggregate peak loads



Thank you!

- Questions?
- The paper and slides are also available at: http://web.eecs.utk.edu/~weigao

