



COLLEGE OF ENGINEERING
Electrical & Computer
Engineering

QoE and Power Efficiency Tradeoff for Fog Computing

Yong Xiao and Marwan Krunz

Research Assistant Professor

NSF BWAC Center Manager

Department of Electrical and Computer Engineering

University of Arizona



THE UNIVERSITY
OF ARIZONA

Outline

- Introduction
- QoE and Power Efficiency Tradeoff
 - Fog Computing without Cooperation
 - Cooperative Fog Computing
- An ADMM-based Distributed Optimization Algorithm
 - Introduction of ADMM
 - ADMM via Variable Splitting
- Conclusion and Future work

Cloud Computing Challenges

- Global data center IP traffic will grow **3-fold** from 2015 to 2020, reaching **15.3 zettabytes** by the end of 2020

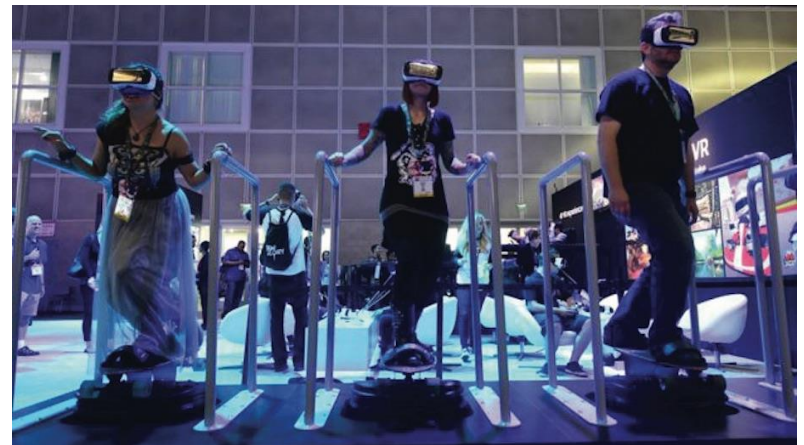


Latency, Latency, Latency!!!

Big drops in sales and traffic have been found when pages took longer to load

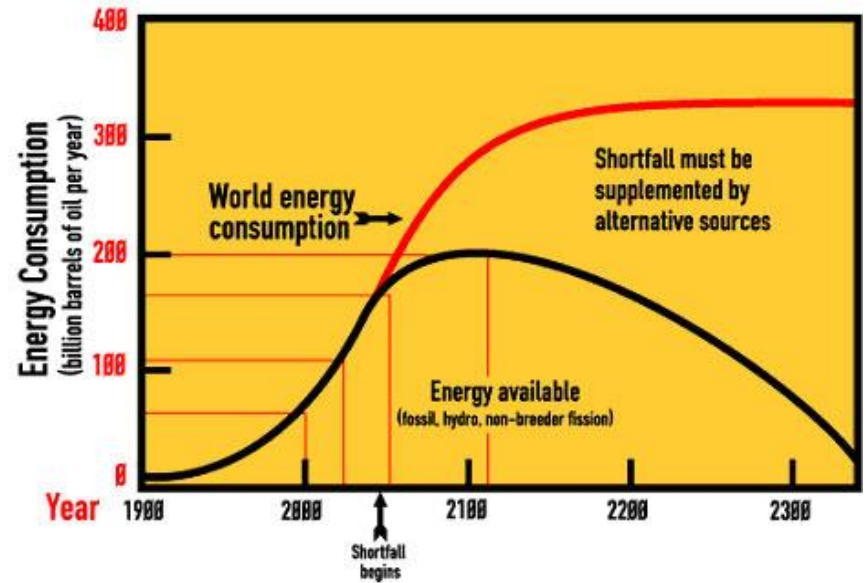
- 0.5s delay will cause a 20% drop in Google's traffic
- 0.1s delay can cause a drop in 1% of Amazon's sales

Many future applications become more sensitive to latency.



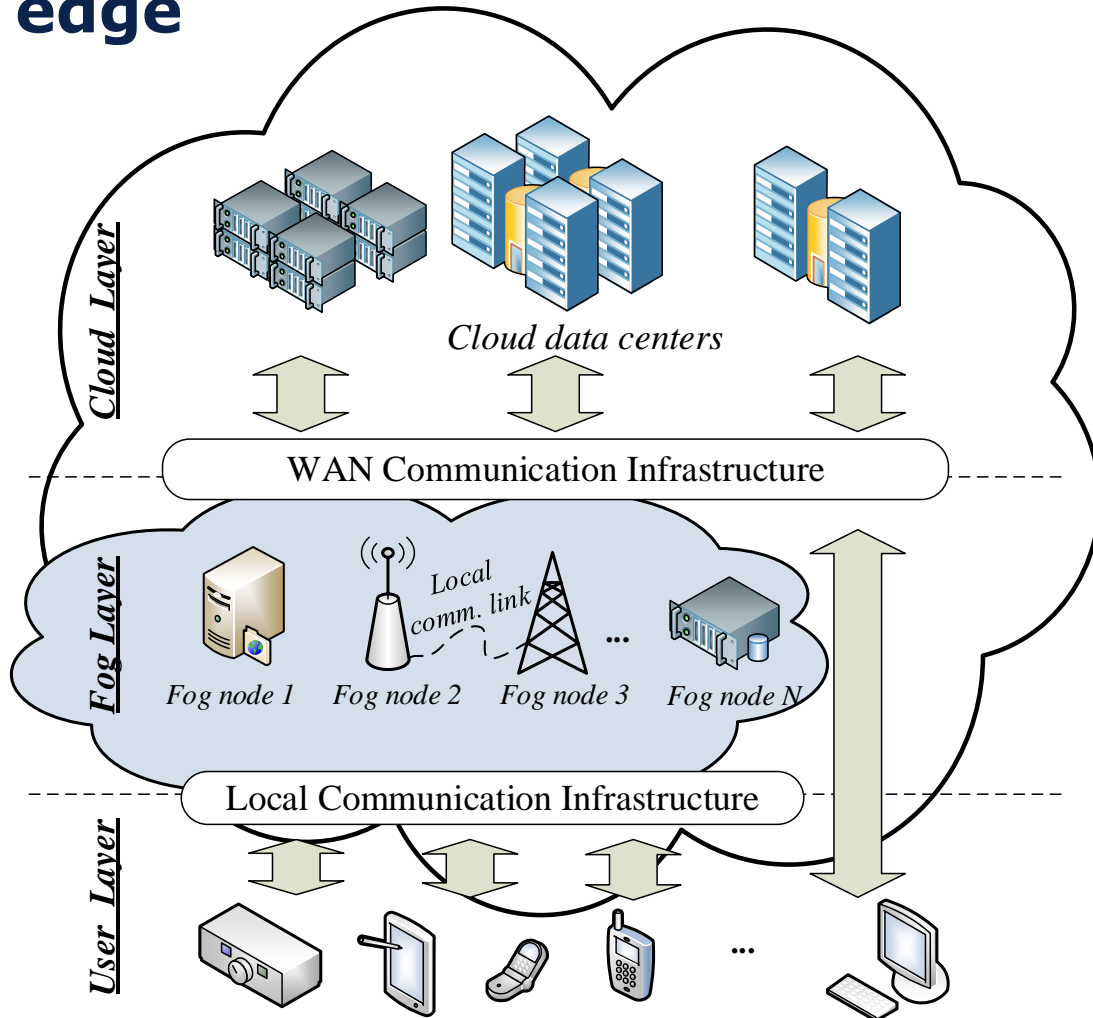
Energy, Energy, Energy!!!

- By the year 2040, world energy consumption would exceed the available energy produced from existing sources



Fog Computing Architecture

Digitization drives data and infrastructure to the edge



Data centers usually located in remote area

Fog nodes are deployed closer to the users

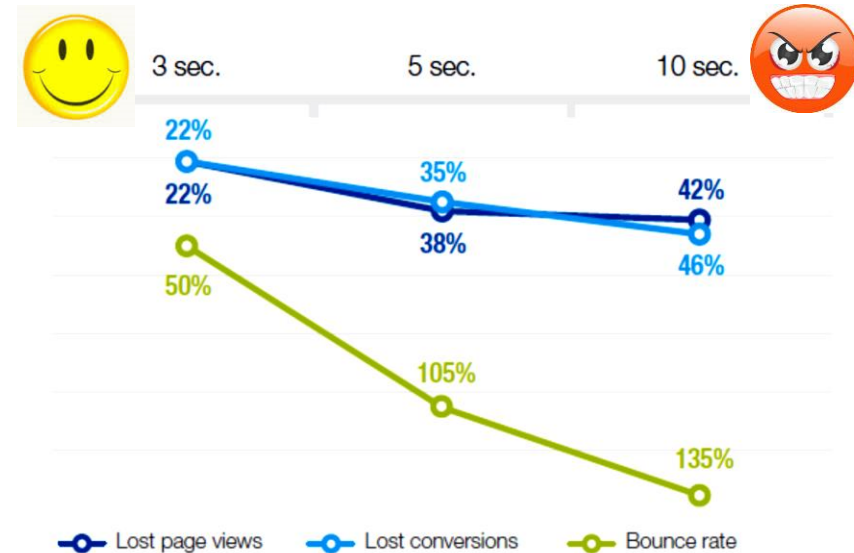
Users desire high QoE services

Key Contributions

- Characterize the fundamental **tradeoff** between **QoE** and **Power Efficiency** for fog computing
- Propose **offload forwarding strategy** for cooperative fog computing
- Propose a new distributed ADMM via **variable splitting approach** to optimize the cooperative fog computing networks

QoE for Fog Computing

- We focus on the QoE of users measured by the average service response-time influenced by
 - ❖ Round-trip workload transmission time:
 - ✓ Non-cooperative fog computing
 - ✓ Cooperative fog nodes
 - ❖ Queueing delay.



Response-time Analysis

- No Offloading:

$$R_j^{W1} = \tau_j^u + \tau^c$$

Upper bound

Workload tx time between fog nodes and cloud

Workload tx time between users and fog nodes

- Full Offloading:

$$R_j^{W2}(\alpha_j) = \tau_j^u + \frac{1}{\mu_j - \lambda_j}$$

Queueing delay

- Partial Offloading:

$$R_j^{W3}(\alpha_j) = \tau_j^u + \alpha_j \left(\frac{1}{\mu_j - \alpha_j \lambda_j} \right) + (1 - \alpha_j) \tau^c$$

Portion of offloaded workload

Maximizing QoE

- Response-time minimization problem:
 - ❖ For non-cooperative fog computing:
each fog node j

$$\begin{aligned} \min_{0 \leq \alpha_j \leq 1} R_j(\alpha_j) \\ \text{s.t. } \eta_j(\alpha_j) \leq \bar{\eta}_j. \end{aligned}$$

Portion of offloaded workload

Power efficiency constraint

Power Efficiency

- We define *power efficiency* as the power consumption per unit of offloaded workload by the fog layer:
 - ❖ Total power consumption for each fog node j :

$$w_j = e_j (w_j^S + w_j^D \alpha_j \lambda_j)$$

Static power consumption/leakage power

Dynamic power consumption

Power usage effectiveness (PUE)

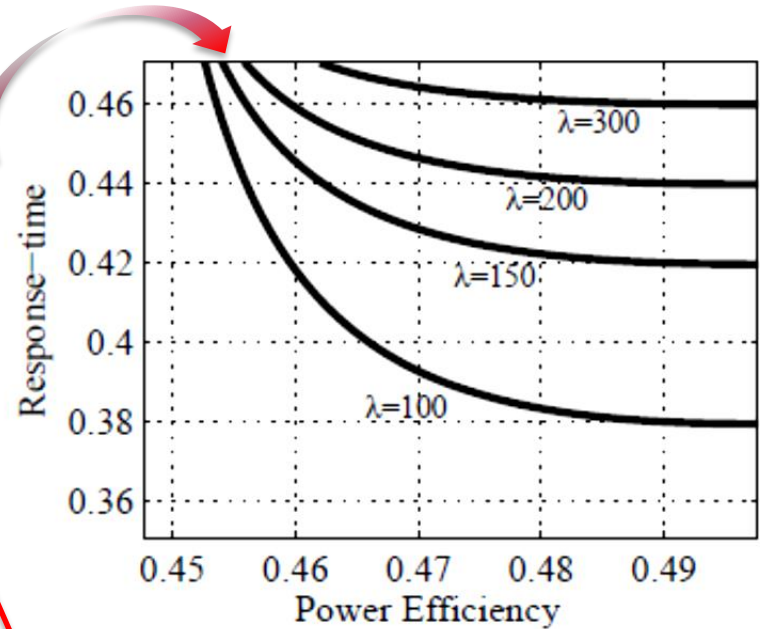
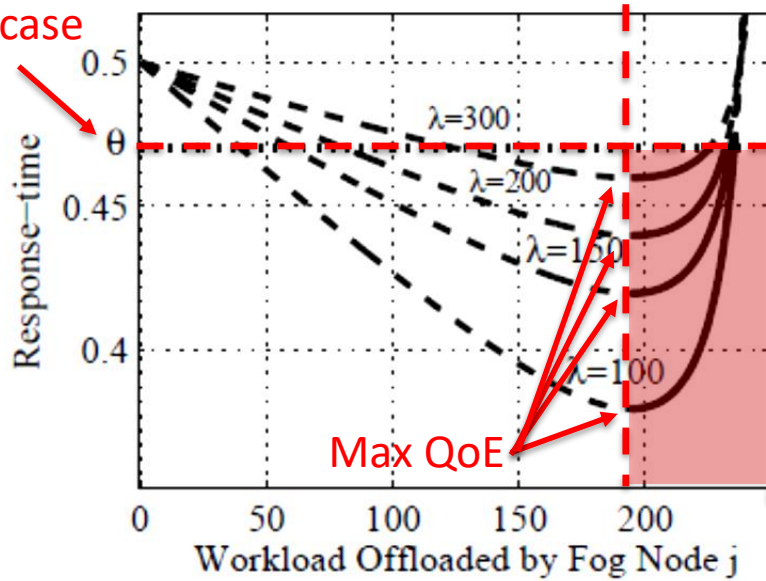
- ❖ Power efficiency:

$$\eta_j(\alpha_j) = \frac{w_j}{\alpha_j \lambda_j} = e_j \left(\frac{w_j^S}{\alpha_j \lambda_j} + w_j^D \right)$$

Workload offloaded by fog node j

QoE and Power Efficiency Tradeoff

Guaranteed
worst-case
QoE



Optimal Tradeoff Region

Cooperative Fog Computing

- Performance of cooperative fog computing is closely related to the cooperation strategy.
- We propose offload forwarding strategy:
 - ❖ Each fog node forwards part of its offloaded workload to others to further improve users' QoE.
 - ❖ Fog nodes can then be divided into
 - ✓ Requesters: require help from others.
 - ✓ Servers: can help processing workload for others.

Response-time Analysis

- Cooperative fog computing with offload forwarding
 - ❖ Fog node j forwards the offloaded workload to a set of neighboring fog nodes \mathcal{C}_j

$$R_j^{C^3}(\xi_j, \varphi_{j\bullet}) = \tau_j^u + \frac{1}{\sum_{i \in \mathcal{F}} \lambda_i} \left[\varphi_{jj} \left(\frac{1}{\mu_j - \varphi_{jj}} \right) + \sum_{i \in \mathcal{C}_j} \varphi_{ji} \left(\tau_{ji} + \frac{1}{\mu_i - \sum_{k \in \mathcal{F}} \varphi_{ki}} \right) \right] + \varphi_{ic} \tau^c,$$

Partition of workload to be forwarded from fog node j to fog node i

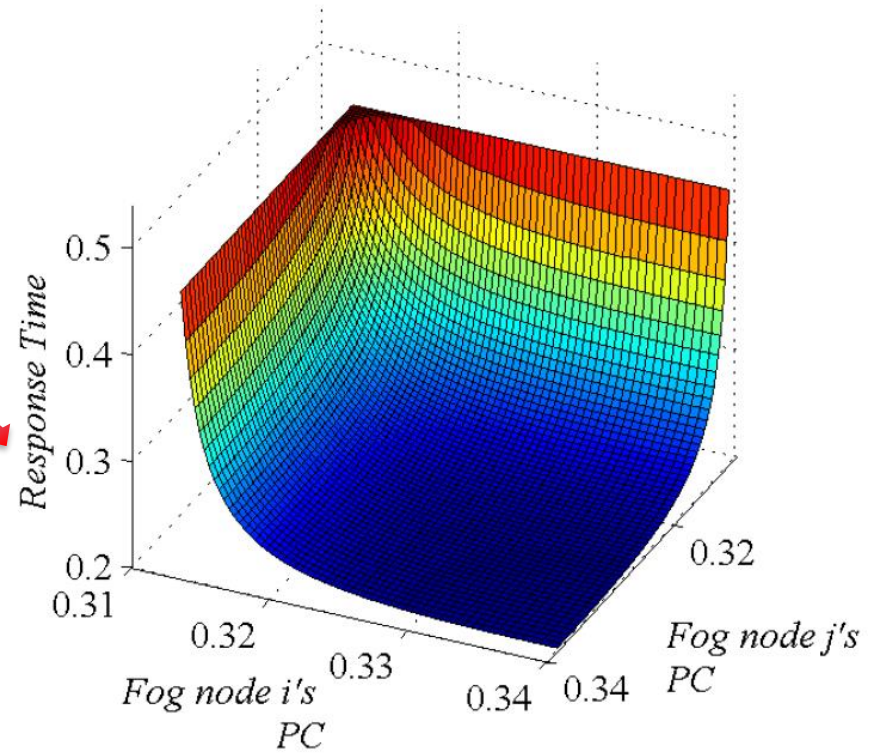
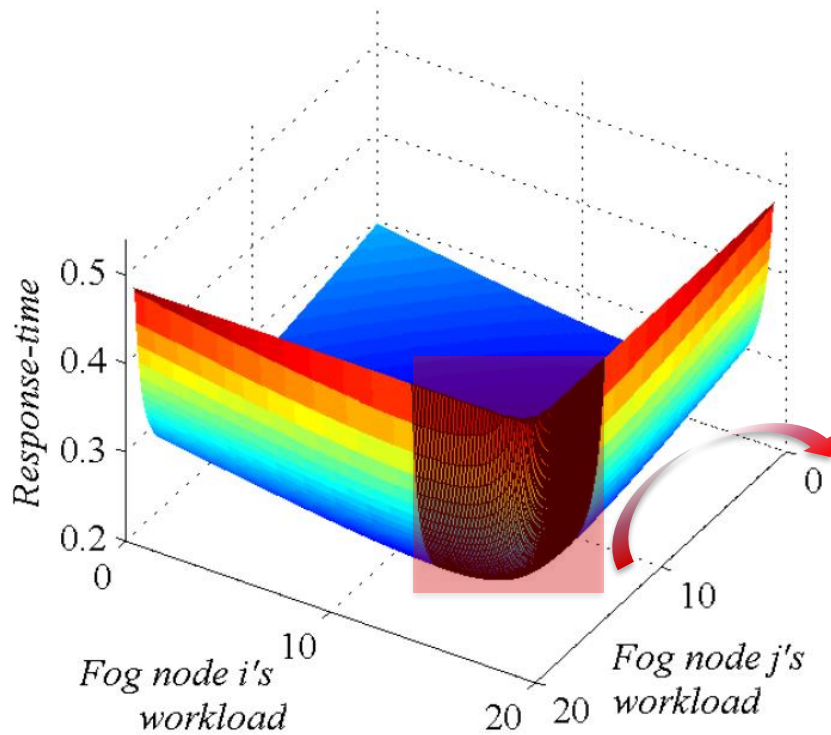
Maximizing QoE

- Response-time minimization problem

$$\begin{aligned} & \min_{\varphi_{1\bullet}, \dots, \varphi_{N\bullet}} \sum_{j=1}^N R_j^{C3}(\xi_j, \varphi_{j\bullet}) \\ & \text{s.t.} \quad \sum_{k \in \mathcal{C}_j} \varphi_{jk} + \varphi_{jj} + \varphi_{jc} = \lambda_j, \\ & \sum_{k \in \mathcal{F}} \varphi_{kj} \leq \min\{\mu_j, \chi_j\}, 0 \leq \varphi_{kj} \leq \lambda_k, \forall k, j \in \mathcal{F} \end{aligned}$$

The maximum amount of workload offloaded by fog node j under the power efficiency constraint $\eta_j(\alpha_j) \leq \bar{\eta}_j$.

QoE and Power Efficiency Tradeoff



Outline

- Introduction
- QoE and Power Efficiency Tradeoff
 - Fog Computing without Cooperation
 - Cooperative Fog Computing
- An ADMM-based Distributed Optimization Algorithm
 - Introduction of ADMM
 - ADMM via variable splitting
- Conclusion and Future work

Why Apply ADMM to Optimize Fog Computing

- ADMM approach is suitable to optimize fog computing networks:
 - ❖ Objective function (Users' QoE) is convex;
 - ❖ Distributed optimization for fog nodes;
 - ❖ With equality constraints:
offloaded + unprocessed workload = workload arrival rate;

Standard ADMM Approach

Optimization Problem

$$\begin{array}{ll} \text{minimize} & f(x) + g(z) \\ \text{subject to} & Ax + Bz = c \end{array}$$



ADMM Solution

$$\begin{array}{l} x^{k+1} := \underset{x}{\operatorname{argmin}} L_{\rho}(x, z^k, y^k) \\ z^{k+1} := \underset{z}{\operatorname{argmin}} L_{\rho}(x^{k+1}, z, y^k) \\ y^{k+1} := y^k + \rho(Ax^{k+1} + Bz^{k+1} - c) \end{array}$$

Problems for Applying ADMM to Fog Computing

- Standard ADMM cannot be directly applied because:
 - 1) Inequality constraints: forwarded workload \leq workload arrival rate;
 - 2) From two blocks to multiple blocks;
 - 3) No communication among fog nodes;
- Objective:
 - ❖ Extending standard ADMM to solve the optimal tradeoff problem

Problem for standard ADMM

$$\begin{array}{ll} \text{minimize} & f(x) + g(z) \\ \text{subject to} & Ax + Bz = c \end{array}$$



Our Problem

$$\begin{array}{ll} \max_{\varphi_{1\bullet}, \dots, \varphi_{N\bullet}} & \sum_{j=1}^N R_j^{C3}(\xi_j, \varphi_{j\bullet}) \\ \text{s.t.} & \sum_{k \in \mathcal{C}_j} \varphi_{jk} + \varphi_{jj} + \varphi_{jc} = \lambda_j, \\ & \sum_{k \in \mathcal{F}} \varphi_{kj} \leq \min\{\mu_j, \chi_j\}, 0 \leq \varphi_{kj} \leq \lambda_k, \end{array}$$

Proposed Distributed Optimization Framework

- A distributed ADMM via variable splitting approach:
 - 1) Introduce indicator functions and auxiliary variables to remove the inequality constraint

$$\begin{aligned} \min_{\varphi_{\bullet 1}, \dots, \varphi_{\bullet N}, \psi} & \sum_{i \in \mathcal{F}} (R_i^{C3}(\xi_i, \varphi_{i\bullet}) + \mathbf{I}_{\mathcal{G}_i}(\varphi_{i\bullet})) \\ & + \mathbf{I}_{\mathcal{G}_c}(\psi) \\ \text{s.t.} & \varphi_{i\bullet} - \psi_i = 0, \forall i \in \mathcal{F}. \end{aligned}$$

- 2) Convert the original problem with multiple random variables into the form with two blocks via variable splitting;

$$\begin{aligned} \varphi^{t+1} &= \arg \min_{\varphi} \mathcal{L}_{\rho}(\varphi_{\bullet 1}, \varphi_{\bullet 2}, \dots, \varphi_{\bullet N}, \psi^t, \Lambda^t) \\ \psi^{t+1} &= \arg \min_{\psi} \frac{\rho}{2} \|\varphi^{t+1} - \psi^t + \frac{1}{\rho} \Lambda^t\| + \mathbf{I}_{\mathcal{G}_c}(\psi) \\ \Lambda^{t+1} &= \Lambda^t - \rho(\varphi^{t+1} - \psi^{t+1}) \end{aligned}$$

Distributed Algorithm

Algorithm 1: Distributed Optimization for Workload Forwarding

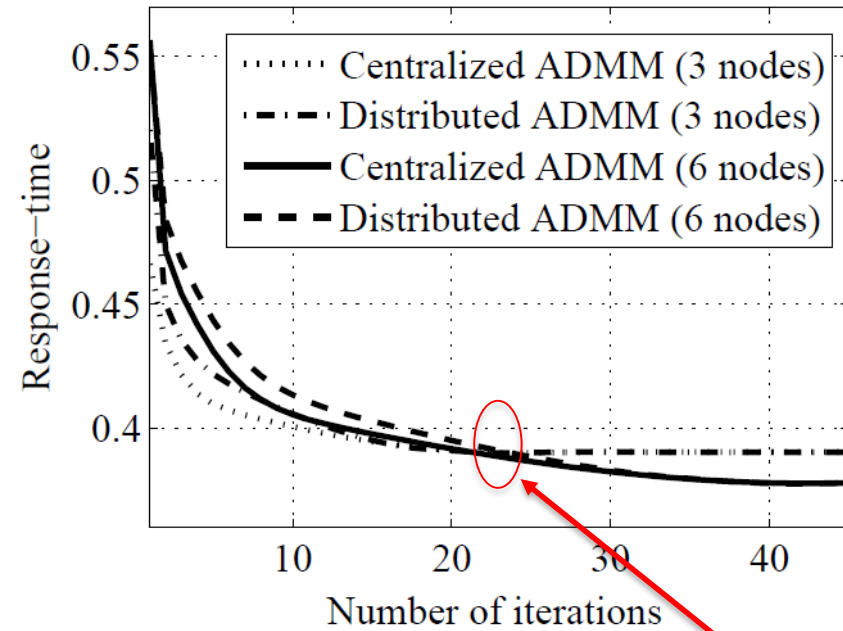
Initialization: Each fog node i chooses an initial service vector $\varphi_{\bullet i}^0$ and WFC chooses an initial dual variable Λ^0 .

WHILE $t=0, 1, \dots$

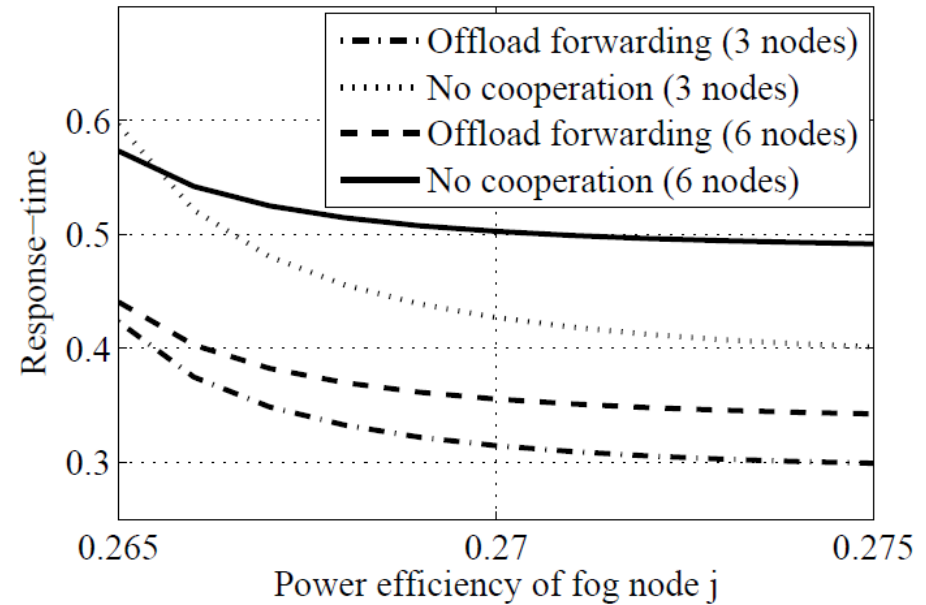
- i) *Fog node updating:* Each fog node i calculates $\varphi_{\bullet i}^{t+1}$ by solving (*) and then sends the resulting $\varphi_{\bullet i}^{t+1}$ and λ_k to the cloud
- ii) *WFC Updating:* cloud calculates ψ^{t+1} by solving ψ -updating problem in (18).
- iii) *Dual Variable Updating* cloud updates dual variables $\Lambda^{t+1} = \Lambda^k - \rho (\varphi^{t+1} - \psi^{t+1})$ and sends φ_i^{t+1} and Λ_i^{t+1} to fog node i .

ENDWHILE

Simulation results (I)

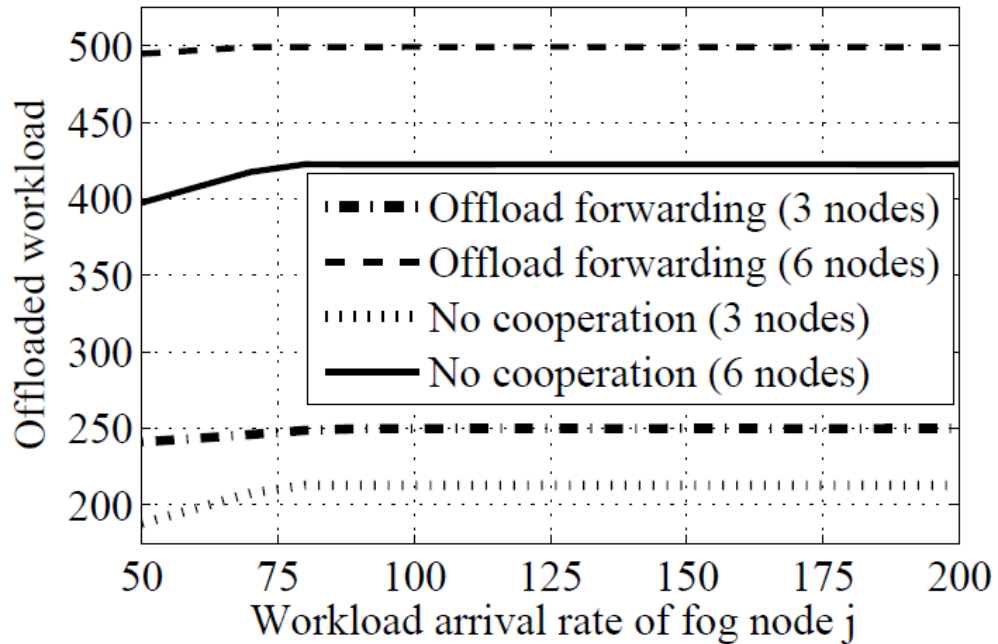
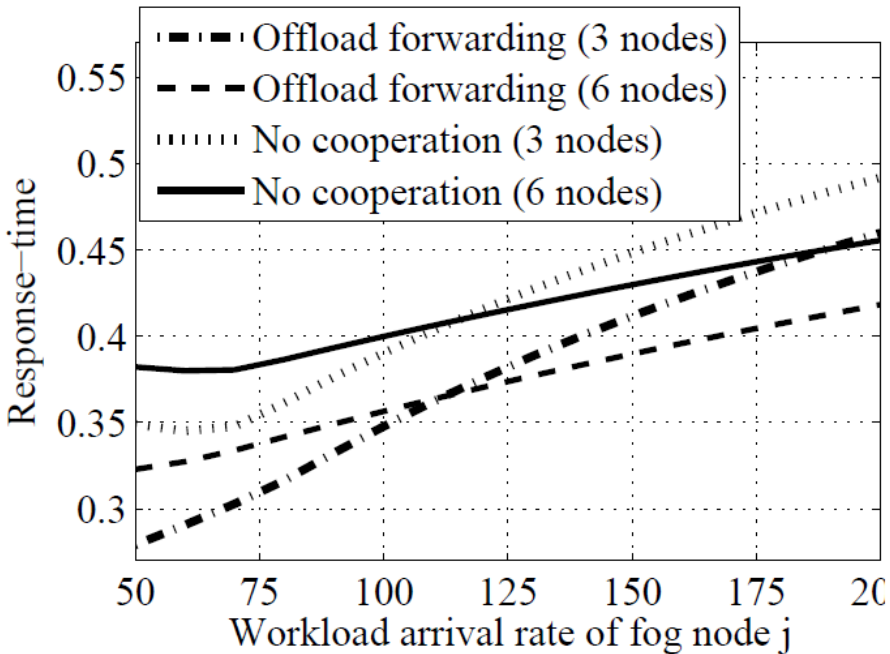


Converge in only 22 iterations



Observation: the number of fog nodes does not affect the convergence speed.

Simulation results (II)



Conclusion

- Characterize the fundamental tradeoff between QoE and Power Efficiency for fog computing
- Propose offload forwarding strategy for cooperative fog computing
- Propose a new distributed ADMM via variable splitting algorithm
- Future work:
 - Extending into stochastic environment
 - Study the QoE and power efficiency tradeoff in more complex fog computing networks, e.g., with other cooperation strategies

THANK
YOU!

E-mail: xyong_2005@yahoo.com

Homepage: <https://sites.google.com/site/xyong2007/>