

Towards Transcoding as a Service in Multimedia Cloud: Energy-Efficient Job Dispatching Algorithm

Weiwen Zhang, Yonggang Wen, *Member, IEEE*, Jianfei Cai, *Senior Member, IEEE*,
and Dapeng Oliver Wu, *Fellow, IEEE*

Abstract—In this paper, we investigate energy-efficient job dispatching algorithm for transcoding as a service (TaaS) in a multimedia cloud. We aim to minimize the energy consumption of service engines in the cloud while achieving low delay for TaaS. We formulate the job dispatching problem as a constrained optimization problem under the framework of Lyapunov optimization. Using the *drift-plus-penalty* function, we propose an online algorithm that dispatches the transcoding jobs to service engines, with an objective to Reduce Energy consumption while achieving the QUEUE STability (REQUEST). We first characterize the fundamental tradeoff between energy consumption and queue delay for the REQUEST algorithm numerically, and obtain its performance bound theoretically. Second, we study the robustness of the REQUEST algorithm, with numerical results indicating that the REQUEST algorithm is robust to the inaccuracy of estimating the transcoding time. Third, we compare the performance of the REQUEST algorithm with the other two algorithms, i.e., Round Robin and Random Rate algorithms. By simulation and real trace data, we show that by appropriately choosing the control variable, the REQUEST algorithm outperforms Round Robin and Random Rate algorithms, with smaller time average energy consumption and time average queue length. The proposed REQUEST algorithm can be applied in cloud-assisted multimedia transcoding service.

Index Terms—Energy Efficiency, Transcoding as a Service, Job Dispatching.

I. INTRODUCTION

WITH the popularity of mobile devices, users have an increasing demand of online video consumption on devices. According to Cisco VNI report [1], the global Internet video traffic will contribute 69 percent for all Internet traffic in 2017, up from 57 percent in 2012. This trend of video consumption, however, may be hampered by the limited bandwidth and inherent nature of stochastic wireless channels (e.g., multi-path fading and shadowing effects), which can degrade user's experience while watching videos.

Transcoding technology [2] is introduced to adapt the videos according to the available bandwidth or different users' requirements. Basically, a content provider can transcode the same video into multiple rates or multiple formats for users' need. In addition, the resolution size of a video can be reduced such that users can view the video smoothly over the network.

Copyright (c) 2013 IEEE. Personal use of this material is permitted. However, permission to use this material for any other purposes must be obtained from the IEEE by sending a request to pubs-permissions@ieee.org.

Weiwen Zhang, Yonggang Wen and Jianfei Cai are with School of Computer Engineering, Nanyang Technological University, Nanyang Avenue, Singapore 639798. Email: {wzhang9, ygwen, asjfc} @ntu.edu.sg.

Dapeng Oliver Wu is with Department of Electrical & Computer Engineering, University of Florida, FL 32611-6130. Email: wu@ece.ufl.edu.

However, such a transcoding process is computation-intensive for the content provider. It is a challenge for the content provider to maintain the low delay for transcoding when many requests arrive. Therefore, a large-scale platform should be designed to support the transcoding process.

Cloud computing [3], due to its elasticity of resource allocation, offers a natural way to process a very large number of transcoding jobs. A large number of servers in the cloud can perform transcoding jobs on behalf of the content provider. In this case, the content provider can benefit from the cloud for video consumption from users. This has become an opportunity to deliver transcoding as a service (TaaS) [4], [5].

A generic cloud-assisted transcoding system is illustrated in Figure 1(a). Particularly, users request a content with specific requirement (e.g., bit rate and resolution size), which is determined by the physical capability of the devices and the available bandwidth. If a particular content is available at the content provider, the content can be rendered immediately. Otherwise, the content provider will send a transcoding job to the cloud in order to cater for the requirement of users. In the cloud, there is a dispatcher at the front end and a large number of service engines at the back end. The arriving transcoding job is routed by the dispatcher and completed by one service engine in the cloud.

In this paper, we consider how to dispatch transcoding jobs to a set of available service engines in order to save the energy consumption. To support the transcoding process, a significant portion of energy will be consumed on service engines in the cloud. Hence, we aim to minimize the energy consumption while maintaining low delay for TaaS, by intelligently dispatching transcoding jobs to service engines in the cloud. The dispatching algorithm should be aware of the CPU speed and queue backlog of service engines. Intuitively, if the dispatcher routes many transcoding jobs to the service engine with slow CPU speed, it can reduce the energy consumption; however, it would make the queue arbitrarily long and incur long delay. Therefore, we consider the energy-delay tradeoff in designing dispatching algorithm for TaaS.

We formulate the job dispatching as an optimization problem under the framework of Lyapunov optimization [6]. We model the service engines as a set of parallel queues. Based on *drift-plus-penalty* function, we propose an online algorithm that dispatches the transcoding jobs to the service engines in order to Reduce Energy consumption while achieving the QUEUE STability (REQUEST).

Our contributions in this research are multi-fold:

- Adopting the framework of Lyapunov optimization, we

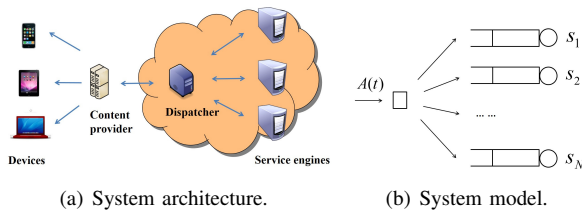


Fig. 1. Overview of cloud-assisted multimedia transcoding platform. The content provider can send transcoding jobs to the cloud. A dispatcher at the front end of the cloud receives transcoding jobs and dispatches them to a set of service engines at the back end for transcoding.

propose the control algorithm REQUEST to dispatch transcoding jobs. We characterize the energy-delay trade-off of the REQUEST algorithm numerically and derive the performance bounds theoretically.

- We study the robustness of the REQUEST algorithm. Numerical results show that, given the inaccuracy of estimating the transcoding time, the error of the time average energy consumption and queue backlog is small. Therefore, the REQUEST algorithm is robust to inaccuracy of the transcoding time estimation.
- We compare the performance of the REQUEST algorithm with Round Robin and Random Rate algorithms using simulation and real trace data. The results show that by appropriately choosing the control variable, the REQUEST algorithm outperforms the other two algorithms, with smaller time average energy consumption while achieving queue stability.

The rest of this paper is organized as follows. In Section II, the review of related work is presented. In Section III, we present the models of distributed transcoding in the cloud. In Section IV, we formulate the optimization problem by Lyapunov optimization and propose an online algorithm, with the performance and robustness analysis. Numerical characterization is given in Section V. In Section VI, we compare the performance of REQUEST algorithm with Round Robin and Random Rate algorithms. Section VII summarizes this paper and provides future directions.

II. RELATED WORK

Prior works [7], [8] have investigated transcoding in distributed systems. [7] scheduled tasks for a cluster-based web server to process, in order to minimize the total processing time by predicting the processing time per individual task. [8] estimated the transcoding time and imported an estimation model for load distribution among distributed servers. Those two works did not investigate the robustness of the scheduling algorithms for the case that the estimation model is not accurate. In this paper, our proposed algorithm is robust to the inaccuracy of the estimated time.

Another line of research [4], [5], [9], [10], [11], [12], [13], [14], [15] leverages cloud computing to enhance the performance of transcoding. [4], [9] utilized a Hadoop-based cloud for transcoding media content, which can greatly improve encoding times. [5] proposed Cloud Transcoder to bridge the gap between videos and mobile devices, reducing the transcoding burden on mobile devices. [10] provided a simulation for cloud

transcoding system with cache capability, and explored the proper cache sizes and the number of computers to operate effectively in the cloud. [11] provided load sharing algorithm in a transcoding cluster. [12] presented a scalable distributed media transcoding system that can reduce the transcoding time. [13] used queue waiting time of transcoding servers to make an admission control for video streams and job dispatching for video transcoding to prevent jitters. [14] considered the cost-efficient virtual machine provision for video transcoding. [15] provided mechanisms for allocation and deallocation of virtual machines to video transcoding servers.

In addition, [16], [17] attempted to minimize the transmission energy on the mobile devices. However, neither of them considered the energy consumption of servers for transcoding. This paper aims to minimize the energy consumption in the cloud for TaaS while maintaining the queue stability. The tradeoff between energy consumption and queue stability is characterized under the framework of Lyapunov optimization.

III. SYSTEM MODELS AND PROBLEM FORMULATION

A. Arrival Model

We consider a discrete time slot model. The length of a time slot is τ . We assume that τ is small such that there is at most one transcoding job arriving to the dispatcher for each time slot. We denote p as the probability of one arrival to the dispatcher for each time slot and $1-p$ if there are no arrivals.

We assume the transcoding time needed for an arriving job at each time slot is associated with the CPU speed of the service engine. Suppose that we have N service engines for transcoding. Each service engine can operate in different CPU speed s_i , where $i = 1, 2, \dots, N$. Without loss of generality, we assume $s_1 \leq s_2 \leq \dots \leq s_N$. The service engine in faster CPU speed can have less completion time for transcoding. We denote $A(t)$ as the transcoding time needed for the arrival at time slot t by a baseline server, which has a CPU speed S . We assume that by statistical learning, the dispatcher can estimate the transcoding time of each job, $A(t)$. Then, if the transcoding job is dispatched to the i th service engine, the transcoding time at the i th service engine is $A_i(t) = \frac{SA(t)}{s_i}$. The transcoding time of the same arrival can be different for service engines, and thus the dispatcher needs to decide which service engine should process the arriving job.

We denote $u(t)$ as the decision variable. Intuitively, if there are no arrivals, the dispatcher does not need to make the decision. Otherwise, the dispatcher decides the routing of the arriving job. The decision variable is chosen from the set $\Phi = \{0, 1, 2, \dots, N\}$, given by,

$$u(t) = \begin{cases} 0, & \text{if no arrival occurs} \\ i, & \text{if dispatched to service engine } i, \end{cases} \quad (1)$$

where $i = 1, 2, \dots, N$.

B. Queueing Model

We model the service engines as a set of queues, as shown in Figure 1(b). To characterize the dynamics of these queues, we define queue length $\mathbf{Q}(t)$ as the unfinished transcoding

time of jobs in each service engine at time slot t , i.e., $\mathbf{Q}(t) = \{Q_1(t), Q_2(t), \dots, Q_N(t)\}$. The queue of the i th service engine evolves according to

$$Q_i(t+1) = \max[Q_i(t) - \tau, 0] + A_i(t)\mathbf{1}_{\{u(t)=i\}}, \quad (2)$$

where $A_i(t)$ is the transcoding time of an arrival at time slot t for the i th service engine, and $\mathbf{1}$ is an indicator function that is 1 if $u(t) = i$ and 0 otherwise. If $u(t) = i$, the arrival is dispatched to the i th service engine and the queue length is increased by $A_i(t)$; otherwise, no arrival occurs to the i th service engine. We can observe the queue length of service engines for each time slot by Eq. (2).

To guarantee the delay of transcoding jobs, we require all the queues to be stable¹, defined as

$$\bar{Q} = \lim_{T \rightarrow \infty} \sup \frac{1}{T} \sum_{t=0}^{T-1} \sum_{i=1}^N \mathbb{E}\{Q_i(t)\} < \infty, \quad (3)$$

where \bar{Q} is the long term time average queue length and the expectation is taken over the randomness of $A(t)$. To address the responsiveness of real-time transcoding, we will show that the short time delay is achieved by the increase of energy consumption on the service engines.

C. Energy Consumption Model

We consider each service engine as a physical machine². Particularly, we only consider the computation energy consumption in the service engine, which is a dominant term for the energy consumption in the distributed servers [18]. As such, we ignore other sources of energy consumption in the service engine, e.g., memory and network.

We assume that each service engine operates in a constant CPU speed when processing transcoding jobs. Its resulted energy consumption is assumed to be a function of CPU speed. If the dispatcher dispatches the transcoding job to the i th service engine at time slot t , the energy consumption on the i th service engine is $A_i(t)\kappa s_i^\alpha$, where $A_i(t)$ is the transcoding time for the i th service engine and κs_i^α is the power that is a convex function of CPU speed [19], [20]. Normally, α is set to be 3 [19]. Also, without loss of generality, we set the constant parameter $\kappa = 1$. If the job is dispatched to the service engine with fast CPU speed, it will result in high energy consumption. If the service engine has no transcoding jobs to process, the service engine can be set to sleep mode, resulting in very small energy consumption that can be negligible.

In addition, we ignore the resulted energy and time due to a transition from sleep mode to running mode of a service engine. For the computation-intensive transcoding jobs, the computation energy and time are the first order component for the total energy and time, while the energy and time due to the transition overhead from the sleep mode to running mode

¹According to Little's theorem, the average queue length is proportional to average delay. In this paper, we aim to satisfy the queue stability.

²For the virtual machine, its energy consumption model can be more complicated. One can adjust the energy consumption model and then adopt our mathematical framework. It is our future work to consider virtualized services enabled by virtual machine.

are the second order component. In other words, the computation energy and time are the dominant terms. Although the overhead is also critical, its effect could be ignored for the decision of job dispatching. Thus, the energy consumption for completing the job is

$$E_i(t) = A_i(t)\kappa s_i^\alpha \mathbf{1}_{\{u(t)=i\}}, \quad (4)$$

where $\mathbf{1}_{\{u(t)=i\}}$ is the indicator function that denotes 1 if $u(t) = i$ and 0 otherwise. In this paper, we consider the long term time average energy consumption, given by

$$\bar{E} = \lim_{T \rightarrow \infty} \sup \frac{1}{T} \sum_{t=0}^{T-1} \sum_{i=1}^N \mathbb{E}\{E_i(t)\}, \quad (5)$$

where the expectation is taken over the randomness of $A(t)$.

D. Problem Formulation

Intuitively, if the dispatcher routes the transcoding job to the service engine with the least queue backlog, it can reduce the delay for the transcoding job; however, it would incur large energy consumption if many transcoding jobs are dispatched to service engines with fast CPU speed. If the dispatcher routes many transcoding jobs to the service engine with the slow CPU speed, it can reduce the energy consumption; however, it would make the queue arbitrarily long and incur long delay. Therefore, we consider the tradeoff between energy consumption and time delay.

In this paper, we aim to minimize the long term time average energy consumption subject to the constraint that time average queue length should not go to infinity. Mathematically, the constrained optimization problem is written as

$$\min_{\{u(t)\}} \bar{E}, \quad (6)$$

$$\text{s.t. } \bar{Q} < \infty, \quad (7)$$

$$u(t) \in \Phi, \quad (8)$$

where $u(t)$ is the decision variable, Eq. (7) denotes the queue stability and Eq. (8) denotes the feasibility constraint.

IV. ONLINE DISPATCHING ALGORITHM

In this section, we adopt the Lyapunov optimization framework to solve the optimization problem (6) and design an online dispatching algorithm.

A. Algorithm Design under i.i.d. Transcoding Time Model

We first assume transcoding time of an arriving job by the baseline server $A(t)$ is i.i.d. for every time slot. Then, $A_i(t)$ is also i.i.d. for the i th service engine. We will discuss how the obtained results can be extended to the non-i.i.d. model in Section IV-D.

We define the quadratic Lyapunov function

$$L(\mathbf{Q}(t)) = \frac{1}{2} \sum_{i=1}^N Q_i(t)^2. \quad (9)$$

Then, we define the one-slot Lyapunov drift as $\Delta(\mathbf{Q}(t)) = \mathbb{E}\{L(\mathbf{Q}(t+1)) - L(\mathbf{Q}(t)) | \mathbf{Q}(t)\}$. Specifically,

$$\begin{aligned} & L(\mathbf{Q}(t+1)) - L(\mathbf{Q}(t)) \\ &= \frac{1}{2} \sum_{i=1}^N Q_i(t+1)^2 - \frac{1}{2} \sum_{i=1}^N Q_i(t)^2 \\ &= \frac{1}{2} \sum_{i=1}^N \{\max[Q_i(t) - \tau, 0] + A_i(t) \mathbf{1}_{\{u(t)=i\}}\}^2 \\ &\quad - \frac{1}{2} \sum_{i=1}^N Q_i(t)^2. \end{aligned}$$

Suppose all $A_i(t)$ are upper bounded by A_{max} for all i and t . Then using the fact $(\max[x - y, 0] + z)^2 \leq x^2 + y^2 + z^2 + 2x(z - y)$ for $\forall x, y, z \geq 0$, we have

$$\begin{aligned} & L(\mathbf{Q}(t+1)) - L(\mathbf{Q}(t)) \\ &\leq \frac{1}{2} [A_{max}^2 + N\tau^2] + \sum_{i=1}^N Q_i(t) A_i(t) \mathbf{1}_{\{u(t)=i\}} \\ &\quad - \tau \sum_{i=1}^N Q_i(t). \end{aligned}$$

Hence,

$$\begin{aligned} \Delta(\mathbf{Q}(t)) &\leq B - \mathbb{E}\left\{\sum_{i=1}^N \tau Q_i(t) | \mathbf{Q}(t)\right\} \\ &\quad + \mathbb{E}\left\{\sum_{i=1}^N Q_i(t) A_i(t) \mathbf{1}_{\{u(t)=i\}} | \mathbf{Q}(t)\right\}, \end{aligned} \quad (10)$$

where B is a finite constant satisfying $B = \frac{1}{2}(A_{max}^2 + N\tau^2)$. The minimization of the right hand side of Eq. (10) will guarantee the queue stability [6].

However, achieving the queue stability does not necessarily lead to the minimum energy consumption on service engines. Thus, we consider the *drift-plus-penalty* function for the dispatching algorithm, which is a weighted sum of drift and penalty, i.e., $\Delta(\mathbf{Q}(t)) + V\mathbb{E}\{E(t) | \mathbf{Q}(t)\}$, where $V \geq 0$ and $E(t) = \sum_{i=1}^N E_i(t)$. We can have the bound of the *drift-plus-penalty* function $\Delta(\mathbf{Q}(t)) + V\mathbb{E}\{E(t) | \mathbf{Q}(t)\}$

$$\begin{aligned} &\leq B - \mathbb{E}\left\{\sum_{i=1}^N \tau Q_i(t) | \mathbf{Q}(t)\right\} \\ &\quad + \mathbb{E}\left\{\sum_{i=1}^N Q_i(t) A_i(t) \mathbf{1}_{\{u(t)=i\}} | \mathbf{Q}(t)\right\} \\ &\quad + V\mathbb{E}\left\{\sum_{i=1}^N A_i(t) \kappa_i s_i^\alpha \mathbf{1}_{\{u(t)=i\}} | \mathbf{Q}(t)\right\}. \end{aligned} \quad (11)$$

We define $F(\mathbf{Q}(t))$ as the bound of the *drift-plus-penalty* function and design the algorithm as in Algorithm 1 by minimizing $F(\mathbf{Q}(t))$. Note that if $V = 0$, only drift function is considered and we choose the queue with the minimum $Q_i(t) A_i(t) = Q_i(t) \frac{A_i(t) S_i}{s_i}$. Since $A_i(t)$ and S_i are unrelated to the i th service engine, the minimum $\frac{Q_i(t)}{s_i}$ determines the job dispatching. Thus, the algorithm is reduced to the policy of shorter queue and faster server.

Algorithm 1 REQUEST Algorithm

Input: $\mathbf{Q}(t)$

Output: $u(t)$

- 1: At the beginning of each time slot t , observe the queue length $\mathbf{Q}(t)$.
- 2: Determine $u(t)$ that minimizes the bound of *drift-plus-penalty* function

$$\sum_{i=1}^N Q_i(t) A_i(t) \mathbf{1}_{\{u(t)=i\}} + V \sum_{i=1}^N A_i(t) \kappa_i s_i^\alpha \mathbf{1}_{\{u(t)=i\}} \quad (12)$$

i.e.,

$$u(t) = \arg \min_i \{A_i(t)(Q_i(t) + V\kappa_i s_i^\alpha)\}. \quad (13)$$

- 3: Update the queue $\mathbf{Q}(t)$ according to (2).
-

B. Performance Analysis

Theorem 1: We assume that the arrival of transcoding jobs is strictly within the capacity region. Capacity region³ is defined as the set Λ of all non-negative $a_i(t) = A_i(t) \mathbf{1}_{\{u(t)=i\}}$, for which

$$\mathbb{E}\{a_i(t)\} < \tau, \quad \forall i. \quad (14)$$

We also assume that $\mathbb{E}\{L(\mathbf{Q}(0))\} < \infty$. Then for any control variable $V > 0$, the online algorithm can stabilize the system, with a resulted time average energy consumption and queue backlog satisfying the following inequalities:

$$\bar{E} \leq E^* + \frac{B}{V}, \quad (15)$$

$$\bar{Q} \leq \frac{B + VE^*}{\varepsilon}, \quad (16)$$

where ε is a constant and E^* is a theoretical lower bound on the time average energy consumption.

Proof: See Appendix A. ■

Theorem 1 shows that by choosing the control variable V , one can achieve a time average energy consumption \bar{E} arbitrarily close to E^* . But this achievement is at the cost of a long delay as the queue length grows linearly with V . Such a tradeoff is important for the cloud operator.

C. Robustness Analysis

In the previous subsections, we assume that we can observe the queue length accurately at the beginning of each time slot before making the decision of dispatching transcoding jobs. However, this observation may not be accurate. In this subsection, we study the robustness of the REQUEST algorithm under the inaccurate queue length information.

Theorem 2: Suppose the estimated queue length $\hat{Q}_i(t)$ satisfies $|\hat{Q}_i(t) - Q_i(t)| \leq q_e$, where $q_e \geq 0$. And we use $\hat{Q}_i(t)$ in place of $Q_i(t)$ for the dispatching algorithm. Consider the bound of the *drift-plus-penalty* function F . Let F^+ be optimal

³This means that the expected transcoding time experienced by a service engine for one time slot should be less than the length of one slot. This can also be interpreted as the incoming expected workload to each service engine should not exceed the workload that service engine can process, if we multiply s_i on both sides of Eq. (14).

if we use $\hat{Q}_i(t)$ and F^* be optimal if we use $Q_i(t)$. Then, by use of $\hat{Q}_i(t)$ we can have a C -additive approximation with some finite constant C as follows,

$$F^+ \leq F^* + C, \quad (17)$$

$$\bar{Q} \leq \frac{B + VE^* + C}{\varepsilon}, \quad (18)$$

$$\bar{E} \leq E^* + \frac{B + C}{V}. \quad (19)$$

Proof: See Appendix B. ■

Theorem 2 shows that we can still minimize the energy consumption and provide queue stability with the inaccurate queue length information by choosing a large V .

D. Extension to non-i.i.d. Arrival

Theorem 1 and 2 are derived under the assumption that $A(t)$ is i.i.d. We can extend those results to the case when $A(t)$ is Markovian. The REQUEST algorithm can still achieve the $[O(1/V), O(V)]$ energy-delay tradeoff,

$$\bar{E} \leq E^* + O\left(\frac{1}{V}\right), \quad (20)$$

$$\bar{Q} \leq O(V), \quad (21)$$

where E^* is a theoretical lower bound on the time average energy consumption.

This can be proved by using multi-slot drift analysis. More details can be found in [6].

V. NUMERICAL CHARACTERIZATION OF REQUEST ALGORITHM

In this section, we first build a statistical model to estimate the transcoding time. Following that, we characterize the trade-off between the energy consumption and the queue backlog for the REQUEST algorithm. Then, we study the robustness of the REQUEST algorithm given that the estimated transcoding time is not accurate.

A. Statistical Model of Transcoding Time

We can model the transcoding time as a function of file size of a video, given by Eq. (22),

$$A = LX, \quad (22)$$

where L is the file size and X is a random variable that denotes transcoding time for a unit of file size. Specifically, X reflects the complexity of the transcoding process, which is determined by the conversion of resolution size, bit rate and frame rate, etc. In this paper, we only consider the conversion of resolution size; other transcoding parameters remain our future work.

To model the transcoding time, we measure the elapsed time of the video transcoding that converts a set of video files into different resolution cases. We consider the application scenario of transcoding flv files into mp4 files in six commonly used resolution cases for the output⁴, i.e., 320x240, 427x240,

⁴Our model is not restricted to these six resolution cases, but can be extended for any resolution cases.

TABLE I
PARAMETER SETTINGS

Arrival probability	$p = 0.8$
Number of service engines	$N = 10$
Normalized CPU speed of a baseline server	$S = 3.2$
Normalized CPU speed of service engines	$\{s_i\} = 2.0 : 0.1 : 2.9$
Parameters of energy model	$\kappa = 1, \alpha = 3$

480x360, 640x360, 640x480 and 854x480. The original flv files are in 1920x1080, with equal duration time but different file size (ranging from 0.1 to 5MB, with the mean 1.87MB). We find that the random variable X can be modeled by a Gamma distribution, with the probability density function as

$$p_X(x) = \frac{1}{b\Gamma(a)} \left(\frac{x}{b}\right)^{a-1} e^{-\frac{x}{b}}, \quad \text{for } x > 0, \quad (23)$$

depending on two parameters (the shape a and the scale b). The CDF fitting by Gamma distribution are shown in Figure 2. It indicates that the transcoding time can be well modeled with a Gamma distribution.

B. Settings for Numerical Characterization

Table I summarizes the parameter settings. We set $p = 0.8$. The CPU speed of a baseline server is $S = 3.2$ GHz. We assume there are 10 service engines, the CPU speeds of which are ranging from 2.0GHz to 2.9GHz incremented by 0.1GHz. We also set $\kappa = 1$ and $\alpha = 3$ for the energy model. Each service engine is assumed to have an empty queue at the first time slot.

C. Energy-Delay Tradeoff of REQUEST Algorithm

We first characterize the energy-delay tradeoff of the REQUEST algorithm for the scenario in which the requirement of transcoding jobs is fixed to a specific resolution size for a period of time. In this case, the transcoding time can be assumed to be i.i.d.. We plot the tradeoff between energy consumption and queue length for each video in Figure 3. For each resolution case, with the increase of V , the time average energy consumption decreases and converges to the optimal value. However, with the increase of V , the time average queue length grows linearly. Hence, the variable V controls the energy-delay tradeoff of the REQUEST algorithm. These results in Figure 3 are consistent with Theorem 1.

We then characterize the energy-delay tradeoff of the REQUEST algorithm for the scenario in which the requirement of transcoding jobs is changed among the resolution cases. Particularly, we assume we have the knowledge of each resolution case being requested. The requested resolution size is determined by the screen size of end devices and the available bandwidth [21]. Each resolution size is assumed to be requested by users with the probability given in Table II. More realistic models can be adopted and our proposed algorithm is still valid. We plot the tradeoff between energy consumption and queue length for this scenario in Figure 4. It is shown that we can have similar energy-delay tradeoff analysis (i.e., $O(1/V)$ and $O(V)$): the time average energy consumption is reduced by the expense of increasing the time average queue length for the non-i.i.d. transcoding time.

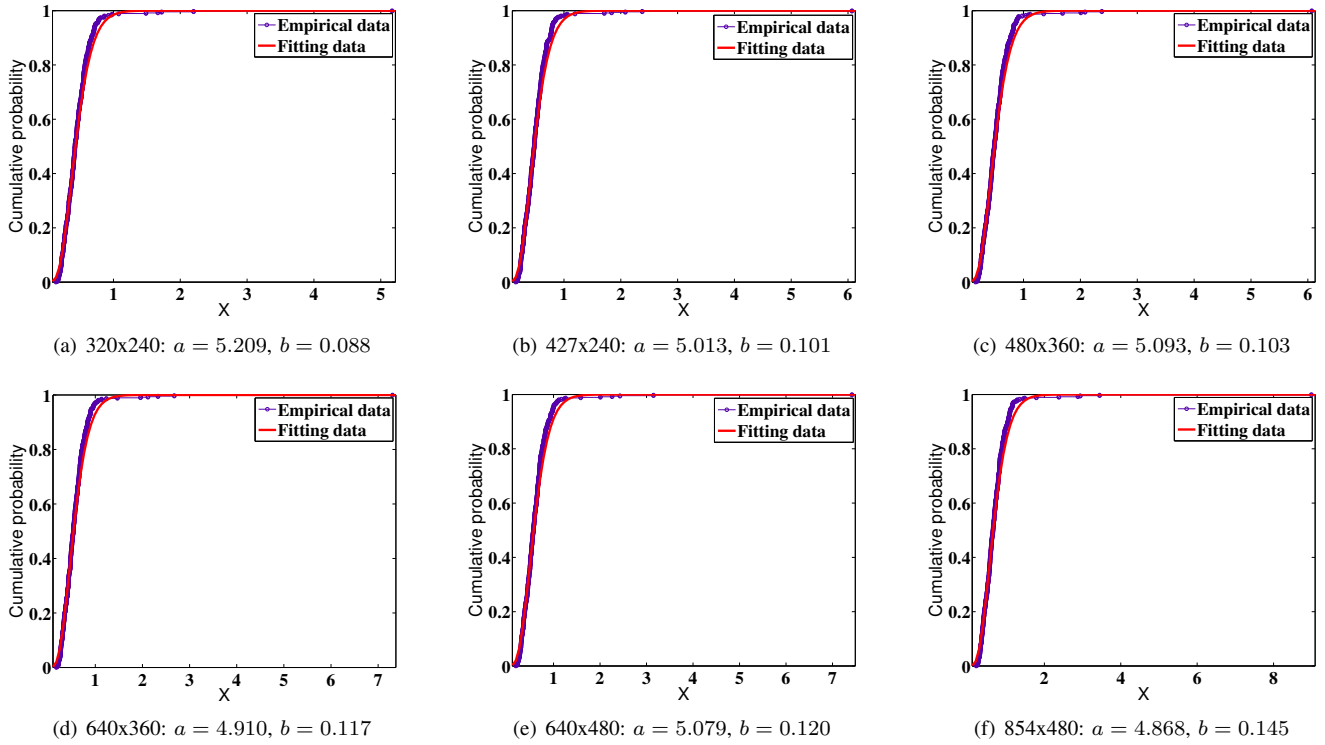


Fig. 2. Gamma distribution CDF fitting for transcoding time in six resolution cases.

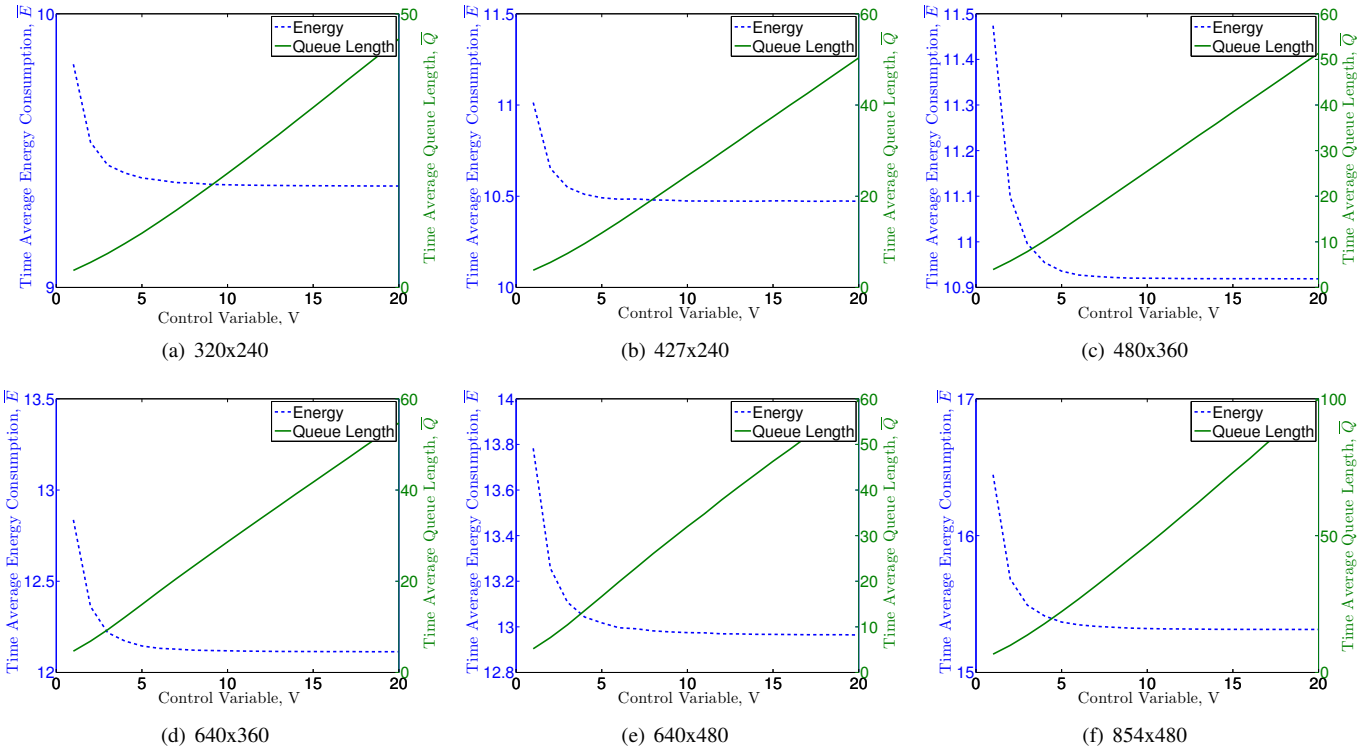


Fig. 3. Time average energy consumption and queue length under different V values for various resolution cases (i.i.d. transcoding time). $V = [1 : 1 : 20]$. $T = 100000$. $\tau = 0.5s$.

TABLE II
PROBABILITY OF REQUESTS FOR DIFFERENT RESOLUTION CASES

resolution case	probability
320x240	$p_1 = 0.1$
427x240	$p_2 = 0.1$
480x360	$p_3 = 0.15$
640x360	$p_4 = 0.15$
640x480	$p_5 = 0.25$
854x480	$p_6 = 0.25$

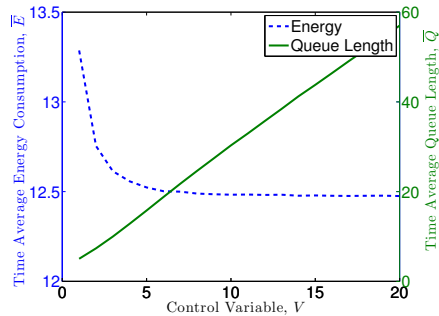


Fig. 4. Time average energy consumption and time average queue length under different V values for non-i.i.d. transcoding time. $V = [1 : 1 : 20]$. $T = 100000$. $\tau = 0.5s$.

D. Robustness of REQUEST Algorithm

Although we can estimate the time of transcoding jobs, this estimation may not be accurate and the observation of the queue length is thus also inaccurate. Therefore, it is necessary to study the robustness of the algorithm under the inaccurate queue length information.

To study the robustness of the REQUEST algorithm, we add a random estimation error that is uniformly distributed, with the range of $\pm 50\%$ for the transcoding time. We consider the relative error between the value of time average energy consumption (or time average queue length) using inaccurate queue length and the value using accurate queue length. The relative error of time average energy consumption (or time average queue length) is defined as the ratio between the difference due to the inaccuracy and the value using accurate queue length. We plot the errors of time average energy consumption and queue length in Figure 5 for varying V . It shows that both the errors are small. Therefore, the REQUEST algorithm is robust to the transcoding time estimation.

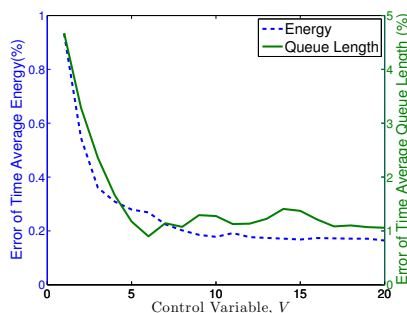


Fig. 5. Error of time average energy consumption and queue length under different V values. $V = [1 : 1 : 20]$. $T = 100000$. $\tau = 0.5s$.

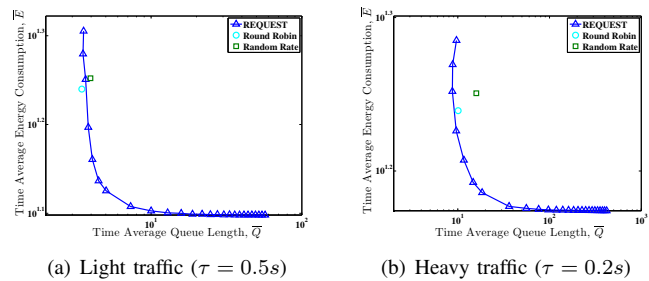


Fig. 6. Energy-delay tradeoff for dispatching algorithms. $T = 100000$. $V = [0, 0.07, 0.14, 0.3, 0.5, 0.75, 1 : 1 : 20]$.

VI. PERFORMANCE COMPARISON OF DISPATCHING ALGORITHMS

In this section, we compare the performance of dispatching algorithms, including Round Robin, Random Rate and REQUEST, under simulated traffic and real trace data.

The Round Robin and Random Rate algorithms are illustrated as follows:

- Round Robin: transcoding jobs are scheduled in a cyclical fashion among N service engines.
- Random Rate: transcoding jobs are dispatched to the i th service engine with the probability $\frac{s_i}{\sum_{i=1}^N s_i}$, which is proportional to the CPU speed of service engines.

Round Robin and Random Rate algorithms are similar, in the sense that they attempt to make load balance among the service engines. However, these two algorithms are static and unaware of the arrivals, which limits their performance for achieving small energy consumption.

A. Simulated Traffic

In this subsection, we compare the performance of dispatching algorithms under simulated traffic.

We first plot the fundamental tradeoff between time average energy consumption and time average queue length for these three dispatching algorithms in Figure 6. In our simulation, we set $\tau = 0.5s$ to model the light traffic and $\tau = 0.2s$ to model the heavy traffic. In Figure 6(a), under the light traffic, the tradeoff of Round Robin and Random Rate algorithms are close with the boundary of the REQUEST algorithm. Note that the axes are plot with logarithmic scale. In Figure 6(b), under the heavy traffic, the tradeoff of Round Robin and Random Rate algorithms are out of the boundary of the REQUEST algorithm; the REQUEST algorithm is more effective to minimize the time average energy consumption and the queue length. It is also shown that by choosing different control variable V , the REQUEST algorithm is adaptive to balance the tradeoff between time average energy consumption and time average queue length.

We then compare these three algorithms under the light traffic ($\tau = 0.5s$). We plot the time average energy consumption, time average queue length and the file size in each time slot that reflects the traffic in Figure 7 from top to bottom, respectively. Particularly, for the REQUEST algorithm, we set $V = 0, 1, 5$, respectively. It is shown that the REQUEST algorithm ($V = 0$) has close time average queue length

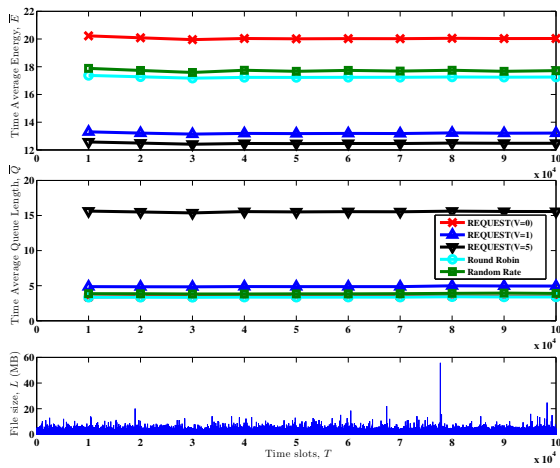


Fig. 7. Performance comparison under light traffic. $\tau = 0.5s$.

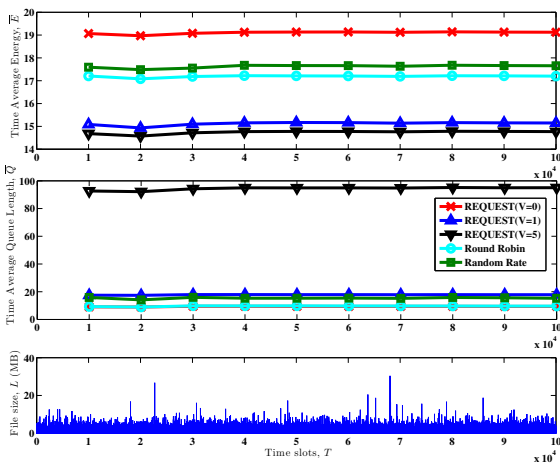


Fig. 8. Performance comparison under heavy traffic. $\tau = 0.2s$.

with Round Robin and Random Rate algorithms but has the highest time average energy consumption. For the REQUEST algorithm ($V = 1$) and the REQUEST algorithm ($V = 5$), they can have small time average energy consumption, which is about 30% smaller than Round Robin and Random Rate algorithms. Although the REQUEST algorithm ($V = 1$) has a slightly larger time average queue length than Round Robin and Random Rate algorithms, it can still maintain the queue length within a margin.

We also compare these three algorithms under the heavy traffic ($\tau = 0.2s$). We plot the time average energy consumption, time average queue length and the file size in each time slot in Figure 8 from top to bottom, respectively. It is shown that the REQUEST algorithm ($V = 0$) has the highest time average energy consumption. The REQUEST algorithm ($V = 5$) can have the smallest time average energy consumption, but its time average queue length is high due to the heavy traffic in this setting. The REQUEST algorithm ($V = 1$) has a slightly larger time average energy consumption than the REQUEST algorithm ($V = 5$) but achieves a much smaller time average queue length. The Round Robin algorithm and Random Rate algorithm perform similarly due to the proportional dispatching strategy.

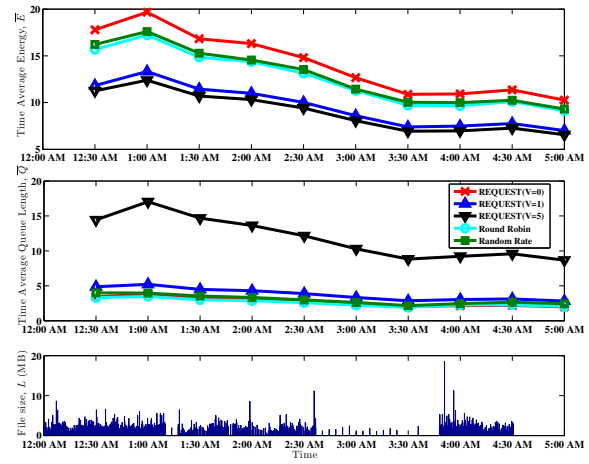


Fig. 9. Performance comparison under real trace (12:00AM-5:00AM).

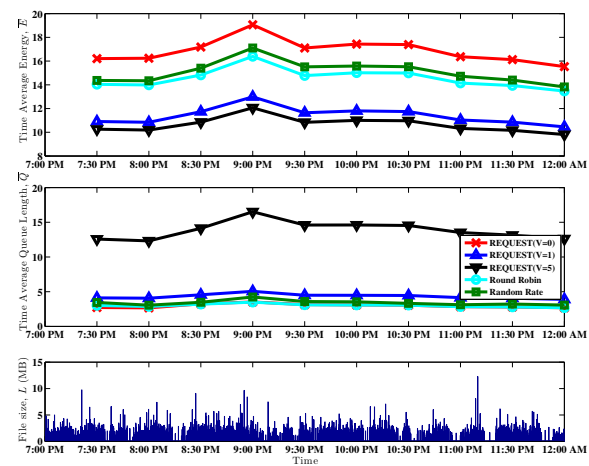


Fig. 10. Performance comparison under real trace (7:00PM-12:00AM).

B. Real Trace Data

In this subsection, we evaluate and compare the performance of the dispatching algorithms by using a real trace data.

The trace data captures the video requests to a CDN node in China. We consider two periods of time, i.e., 7:00PM-12:00AM and 12:00AM-5:00AM in a day. The data traffic of 12:00AM-5:00AM is lighter than that of 7:00PM-12:00AM, as shown in the bottom subfigures of Figure 9 and 10. For every 30 minutes, we plot the time average energy consumption and time average queue length for these two periods in Figure 9 and 10, respectively. For the period of 7:00PM-12:00AM, since its traffic is heavier, it results in larger time average energy consumption and time average queue length for all dispatching algorithms. For the period of 12:00AM-5:00AM, when there are small amount of video requests (e.g., from 3:00AM to 4:00AM), both the time average energy consumption and time average queue length decrease for all dispatching algorithms.

Using this trace data, we compare the performance of dispatching algorithms with the analysis as follows. The REQUEST algorithm ($V = 0$) can achieve short time average queue length at the high cost of time average energy con-

sumption. The REQUEST algorithm ($V > 0$) is more energy-efficient than Round Robin and Random Rate algorithms. It can also manage the tradeoff between energy consumption and queue length in dispatching transcoding jobs. Therefore, combined with the previous simulation results, we can have the insight that the cloud operator can tune the control variable V of the REQUEST algorithm such that it can outperform Round Robin and Random Rate algorithms.

VII. CONCLUSION

We investigated dispatching algorithms on how to route transcoding jobs in the multimedia cloud. To minimize the energy consumption by cloud service engines, we formulated the job dispatching policy as an optimization problem under the framework of Lyapunov optimization. We characterized the energy-delay tradeoff and the robustness of the REQUEST algorithm. The simulation results showed that the REQUEST algorithm is more energy-efficient than Round Robin and Random Rate algorithms. The insight is that the cloud operator can dynamically tune the control variable of the REQUEST algorithm in order to reduce the energy consumption while maintaining the queue stability.

In the future, we will build up a more general transcoding time model by considering the bit rate adaptation. In addition, we will take virtual machine into consideration for virtualized services. Finally, we will evaluate the performance of the proposed algorithm in the real multimedia platform.

APPENDIX A PROOF OF THEOREM 1

Proof: Since we assume that the arrival process is within the capacity region, there exists at least one stationary randomized control policy that can stabilize the queue [22], with $\mathbb{E}\{E(t)\} = E^*$ and $\tau \geq \mathbb{E}\{a_i(t)\} + \varepsilon, \forall i$, where $\varepsilon > 0$. Since $A(t)$ is i.i.d., so is $A_i(t)$. In addition, since $A_i(t)$ is independent of the current queue backlog $\mathbf{Q}(t)$, we have $\mathbb{E}\{A_i(t)|\mathbf{Q}(t)\} = \mathbb{E}\{A_i(t)\}$. Thus, the *drift-plus-penalty* function satisfies

$$\begin{aligned} & \Delta(\mathbf{Q}(t)) + V\mathbb{E}\{E(t)|\mathbf{Q}(t)\} \\ & \leq B - \mathbb{E}\left\{\sum_{i=1}^N \tau Q_i(t) | \mathbf{Q}(t)\right\} \\ & \quad + \mathbb{E}\left\{\sum_{i=1}^N Q_i(t) A_i(t) \mathbf{1}_{\{u(t)=i\}} | \mathbf{Q}(t)\right\} + VE^* \\ & \leq B - \sum_{i=1}^N Q_i(t) \varepsilon + VE^*. \end{aligned}$$

Taking a conditional expectation over $\mathbf{Q}(t)$ for this *drift-plus-penalty* function and using iterative expectation law, we can have

$$\begin{aligned} & \mathbb{E}\{L(\mathbf{Q}(t+1)) - L(\mathbf{Q}(t))\} + V\mathbb{E}\{E(t)\} \\ & \leq B - \sum_{i=1}^N \varepsilon \mathbb{E}\{Q_i(t)\} + VE^*. \end{aligned}$$

Then, summing over all time slots $t \in \{0, 1, \dots, T-1\}$, and dividing by T , we obtain

$$\begin{aligned} & \frac{\mathbb{E}\{L(\mathbf{Q}(T)) - L(\mathbf{Q}(0))\}}{T} + \frac{V}{T} \sum_{t=0}^{T-1} \mathbb{E}\{E(t)\} \\ & \leq B - \frac{1}{T} \sum_{t=0}^{T-1} \sum_{i=1}^N \varepsilon \mathbb{E}\{Q_i(t)\} + VE^*. \end{aligned}$$

In this case, we have

$$\frac{1}{T} \sum_{t=0}^{T-1} \sum_{i=1}^N \varepsilon \mathbb{E}\{Q_i(t)\} \leq B + VE^* + \frac{\mathbb{E}\{L(\mathbf{Q}(0))\}}{T}. \quad (24)$$

Then,

$$\frac{1}{T} \sum_{t=0}^{T-1} \sum_{i=1}^N \mathbb{E}\{Q_i(t)\} \leq \frac{B + VE^* + \frac{\mathbb{E}\{L(\mathbf{Q}(0))\}}{T}}{\varepsilon}. \quad (25)$$

Taking a lim sup as $T \rightarrow \infty$, we have

$$\bar{Q} \leq \frac{B + VE^*}{\varepsilon}. \quad (26)$$

Similarly, taking the same rationale, we have

$$\bar{E} \leq E^* + \frac{B}{V}. \quad (27)$$

APPENDIX B PROOF OF THEOREM 2

Proof: Using $\hat{Q}_i(t)$, we still try to minimize the bound of the *drift-plus-penalty* function,

$$\begin{aligned} F(\hat{Q}_i(t)) & = B - \mathbb{E}\left\{\sum_{i=1}^N \tau \hat{Q}_i(t) | \mathbf{Q}(t)\right\} \\ & \quad + \mathbb{E}\left\{\sum_{i=1}^N \hat{Q}_i(t) A_i(t) \mathbf{1}_{\{u(t)=i\}} | \mathbf{Q}(t)\right\} \\ & \quad + V\mathbb{E}\left\{\sum_{i=1}^N \kappa A_i(t) s_i^\alpha \mathbf{1}_{\{u(t)=i\}} | \mathbf{Q}(t)\right\}. \end{aligned} \quad (28)$$

Denote $q_i(t) = \hat{Q}_i(t) - Q_i(t)$. Plugging $\hat{Q}_i(t)$ into Eq. (28), we have

$$\begin{aligned} F(\hat{Q}_i(t)) & = B - \mathbb{E}\left\{\sum_{i=1}^N \tau Q_i(t) | \mathbf{Q}(t)\right\} \\ & \quad + \mathbb{E}\left\{\sum_{i=1}^N Q_i(t) A_i(t) \mathbf{1}_{\{u(t)=i\}} | \mathbf{Q}(t)\right\} \\ & \quad + V\mathbb{E}\left\{\sum_{i=1}^N \kappa A_i(t) s_i^\alpha \mathbf{1}_{\{u(t)=i\}} | \mathbf{Q}(t)\right\} \\ & \quad - \mathbb{E}\left\{\sum_{i=1}^N \tau q_i(t) | \mathbf{Q}(t)\right\} \\ & \quad + \mathbb{E}\left\{\sum_{i=1}^N q_i(t) A_i(t) \mathbf{1}_{\{u(t)=i\}} | \mathbf{Q}(t)\right\}. \end{aligned} \quad (29)$$

We denote the minimum value of $F(\hat{Q}_i(t))$ as F^+ , and the minimum value of $F(Q_i(t))$ as F^* . Then we have

$$F^+ \leq F^* - \mathbb{E}\left\{\sum_{i=1}^N \tau q_i(t) | \mathbf{Q}(t)\right\} + \mathbb{E}\left\{\sum_{i=1}^N q_i(t) A_i(t) \mathbf{1}_{\{u(t)=i\}} | \mathbf{Q}(t)\right\} \quad (30)$$

Since $|q_i(t)| \leq q_e$ and $|A_i(t) \mathbf{1}_{\{u(t)=i\}}| \leq A_{max}$, we have $F^+ \leq F^* + q_e(N\tau + A_{max})$. Let $C = q_e(N\tau + A_{max})$, then we obtain Eq. (17). This indicates that Eq. (11) still holds if $\mathbf{Q}(t)$ is replaced by $\hat{\mathbf{Q}}(t)$ and B is replaced by $\hat{B} = B + C$.

Then, we prove Eq. (18) and (19) as follows. Using the approach in Appendix A, we have

$$\frac{\mathbb{E}\{L(\mathbf{Q}(T)) - L(\mathbf{Q}(0))\}}{T} + \frac{V}{T} \sum_{t=0}^{T-1} \mathbb{E}\{E(t)\} \leq B + C - \frac{1}{T} \sum_{t=0}^{T-1} \sum_{i=1}^N \varepsilon \mathbb{E}\{Q_i(t)\} + VE^*.$$

Therefore,

$$\frac{1}{T} \sum_{t=0}^{T-1} \sum_{i=1}^N \mathbb{E}\{Q_i(t)\} \leq \frac{B + C + VE^* + \frac{\mathbb{E}\{L(\mathbf{Q}(0))\}}{T}}{\varepsilon}.$$

Taking a lim sup as $T \rightarrow \infty$, we have Eq. (18). Similarly, taking the same rationale, we have Eq. (19).

ACKNOWLEDGMENT

The authors would like to thank the reviewers for their insightful comments. W.W. Zhang and Y.G. Wen were supported by Start-Up Grant from NTU, MOE Tier-1 Grant (RG 31/11) from Singapore MOE and EIRP02 Grant from Singapore EMA. Y. G. Wen and J.F. Cai also acknowledge the support from by the Singapore National Research Foundation under its IDM Futures Funding Initiative and administered by the Interactive & Digital Media Programme Office, Media Development Authority. D.P. Wu acknowledges support from NSF under grant ECCS-1002214 and NSFC under grant No. 61228101.

REFERENCES

- [1] *Cisco Visual Networking Index: Forecast and Methodology, 2012 - 2017*, Cisco, 2013.
- [2] A. Vetro and C. W. Chen, "Rate-reduction transcoding design for wireless video streaming," in *2002 International Conference on Image Processing*, vol. 1, 2002, pp. 29–32.
- [3] M. Armbrust, A. Fox, R. Griffith, A. Joseph, R. Katz, A. Konwinski, G. Lee, D. Patterson, A. Rabkin, I. Stoica, and M. Zaharia, "A view of cloud computing," *Communication of the ACM*, vol. 53, no. 4, pp. 50–58, 2010.
- [4] A. Garcia, H. Kalva, and B. Furth, "A study of transcoding on cloud environments for video content delivery," in *Proceedings of the 2010 ACM multimedia workshop on Mobile cloud media computing*. ACM, 2010, pp. 13–18.

- [5] Z. Li, Y. Huang, G. Liu, F. Wang, Z.-L. Zhang, and Y. Dai, "Cloud transcoder: Bridging the format and resolution gap between internet videos and mobile devices," in *Proceedings of the 22nd international workshop on Network and Operating System Support for Digital Audio and Video*. ACM, 2012, pp. 33–38.
- [6] M. J. Neely, "Stochastic network optimization with application to communication and queueing systems," *Synthesis Lectures on Communication Networks*, vol. 3, no. 1, pp. 1–211, 2010.
- [7] J. Guo and L. N. Bhuyan, "Load balancing in a cluster-based web server for multimedia applications," *IEEE Transactions on Parallel and Distributed Systems*, vol. 17, no. 11, pp. 1321–1334, 2006.
- [8] D. Seo, J. Kim, and I. Jung, "Load distribution algorithm based on transcoding time estimation for distributed transcoding servers," in *2010 International Conference on Information Science and Applications (ICISA)*. IEEE, 2010, pp. 1–8.
- [9] A. Garcia and H. Kalva, "Cloud transcoding for mobile video content delivery," in *2011 IEEE International Conference on Consumer Electronics (ICCE)*. IEEE, 2011, pp. 379–380.
- [10] S. Ko, S. Park, and H. Han, "Design analysis for real-time video transcoding on cloud systems," in *Proceedings of the 28th Annual ACM Symposium on Applied Computing*. ACM, 2013, pp. 1610–1615.
- [11] J. Guo and L. Bhuyan, "Load sharing in a transcoding cluster," in *Distributed Computing-IWDC 2003*. Springer, 2003, pp. 330–339.
- [12] H. Sanson, L. Loyola, and D. Pereira, "Scalable distributed architecture for media transcoding," in *Algorithms and Architectures for Parallel Processing*. Springer, 2012, pp. 288–302.
- [13] A. Ashraf, F. Jokhio, T. Deneke, S. Lafond, I. Porres, and J. Lilius, "Stream-based admission control and scheduling for video transcoding in cloud computing," in *2013 13th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing (CCGrid)*. IEEE, 2013, pp. 482–489.
- [14] A. Ashraf, "Cost-efficient virtual machine provisioning for multi-tier web applications and video transcoding," in *2013 13th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing (CCGrid)*. IEEE, 2013, pp. 66–69.
- [15] F. Jokhio, A. Ashraf, S. Lafond, I. Porres, and J. Lilius, "Prediction-based dynamic resource allocation for video transcoding in cloud computing," in *2013 21st Euromicro International Conference on Parallel, Distributed and Network-Based Processing (PDP)*. IEEE, 2013, pp. 254–261.
- [16] Y. Cui, X. Ma, J. Liu, and Y. Bao, "Energy-efficient on-demand streaming in mobile cellular networks," *IEEE COMSOC MMTC E-LETTER*, vol. 8, no. 5, 2013.
- [17] Y. Cui, S. Xiao, X. Wang, M. Li, H. Wang, and Z. Lai, "Performance-aware energy optimization on mobile devices in cellular network," in *IEEE INFOCOM (accepted)*, 2014.
- [18] P. Bohrer, E. N. Elnozahy, T. Keller, M. Kistler, C. Lefurgy, C. McDowell, and R. Rajamony, "The case for power management in web servers," in *Power aware computing*. Springer, 2002, pp. 261–289.
- [19] Y. Chen, A. Das, W. Qin, A. Sivasubramaniam, Q. Wang, and N. Gautam, "Managing server energy and operational costs in hosting centers," in *ACM SIGMETRICS Performance Evaluation Review*, vol. 33, no. 1. ACM, 2005, pp. 303–314.
- [20] R. Stanojevic and R. Shorten, "Distributed dynamic speed scaling," in *INFOCOM, 2010 Proceedings IEEE*. IEEE, 2010, pp. 1–5.
- [21] *Cisco webex network bandwidth*, Cisco, June 2013.
- [22] M. J. Neely, "Energy optimal control for time-varying wireless networks," *IEEE Transactions on Information Theory*, vol. 52, no. 7, pp. 2915–2934, 2006.

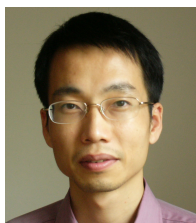


Weiwen Zhang received his Bachelor's degree in software engineering and Master's degree in computer science from South China University of Technology (SCUT) in 2008 and 2011, respectively. He is currently a PhD student in the School of Computer Engineering at Nanyang Technological University (NTU) in Singapore. His research interests include cloud computing and mobile computing.



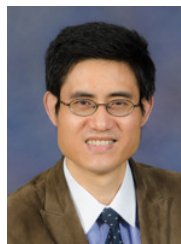
Yonggang Wen (S99-M08) received his Ph.D. degree in electrical engineering and computer science from Massachusetts Institute of Technology (MIT), Cambridge, USA in 2008. He is currently an Assistant Professor with School of Computer Engineering at Nanyang Technological University, Singapore. is an assistant professor with school of computer engineering at Nanyang Technological University, Singapore. Previously he has worked in Cisco to lead product development in content delivery network, which had a revenue impact of 3 Billion US dollars

globally. Dr. Wen has published over 90 papers in top journals and prestigious conferences. His latest work in multi-screen cloud social TV has been featured by global media (more than 1600 news articles from over 29 countries) and recognized with ASEAN ICT Award 2013 (Gold Medal) and IEEE Globecom 2013 Best Paper Award. He serves on editorial boards for IEEE Transactions on Multimedia, IEEE Access Journal and Elsevier Ad Hoc Networks. His research interests include cloud computing, green data center, big data analytics, multimedia network and mobile computing.



Jianfei Cai (S'98-M'02-SM'07) received his PhD degree from the University of Missouri-Columbia. He is currently an Associate Professor and has served as the Head of Visual & Interactive Computing Division and the Head of Computer Communication Division at the School of Computer Engineering, Nanyang Technological University, Singapore. His major research interests include visual computing and multimedia networking. He has published more than 100 technical papers in international conferences and journals. He has been actively

participating in program committees of various conferences. He has served as the leading Technical Program Chair for IEEE International Conference on Multimedia & Expo (ICME) 2012 and the leading General Chair for Pacific-rim Conference on Multimedia (PCM) 2012. He was an invited speaker for the first IEEE Signal Processing Society Summer School on 3D and high definition high contrast video process systems in 2011. He is currently an Associate Editor for IEEE Trans on Image Processing (T-IP) and has served as an Associate Editor for IEEE Trans on Circuits and Systems for Video Technology (T-CSVT), and Guest Editor for IEEE Trans on Multimedia (TMM), ELSEVIER Journal of Visual Communication and Image Representation (JVCI), etc. He is also a senior member of IEEE.



Dapeng Oliver Wu (S'98-M'04-SM06-F'13) received B.E. in Electrical Engineering from Huazhong University of Science and Technology, Wuhan, China, in 1990, M.E. in Electrical Engineering from Beijing University of Posts and Telecommunications, Beijing, China, in 1997, and Ph.D. in Electrical and Computer Engineering from Carnegie Mellon University, Pittsburgh, PA, in 2003.

He is a professor at the Department of Electrical and Computer Engineering, University of Florida, Gainesville, FL. His research interests are in the

areas of networking, communications, signal processing, computer vision, and machine learning. He received University of Florida Research Foundation Professorship Award in 2009, AFOSR Young Investigator Program (YIP) Award in 2009, ONR Young Investigator Program (YIP) Award in 2008, NSF CAREER award in 2007, the IEEE Circuits and Systems for Video Technology (CSVT) Transactions Best Paper Award for Year 2001, and the Best Paper Awards in IEEE GLOBECOM 2011 and International Conference on Quality of Service in Heterogeneous Wired/Wireless Networks (QShine) 2006.

Currently, he serves as an Associate Editor for IEEE Transactions on Circuits and Systems for Video Technology, Journal of Visual Communication and Image Representation, and International Journal of Ad Hoc and Ubiquitous Computing. He is the founder of IEEE Transactions on Network Science and Engineering. He was the founding Editor-in-Chief of Journal of Advances in Multimedia between 2006 and 2008, and an Associate Editor for IEEE Transactions on Wireless Communications and IEEE Transactions on Vehicular Technology between 2004 and 2007. He is also a guest-editor for IEEE Journal on Selected Areas in Communications (JSAC), Special Issue on Cross-layer Optimized Wireless Multimedia Communications. He has served as Technical Program Committee (TPC) Chair for IEEE INFOCOM 2012, and TPC chair for IEEE International Conference on Communications (ICC 2008), Signal Processing for Communications Symposium, and as a member of executive committee and/or technical program committee of over 80 conferences. He has served as Chair for the Award Committee, and Chair of Mobile and wireless multimedia Interest Group (MobiG), Technical Committee on Multimedia Communications, IEEE Communications Society. He was a member of Multimedia Signal Processing Technical Committee, IEEE Signal Processing Society from Jan. 1, 2009 to Dec. 31, 2012. He is an IEEE Fellow.