

Understanding Watching Patterns of Live TV Programs on Mobile Devices: A Content Centric Perspective

Yuheng Li¹ and Qianchuan Zhao¹

¹Center for Intelligent and Networked Systems,
Department of Automation and TNList
Tsinghua University, Beijing, China, 100084
[e-mail: liyuheng07@mails.tsinghua.edu.cn, zhaoqc@tsinghua.edu.cn]

*Received December 10, 2014; revised March 30, 2014; revised March 28, 2015; accepted July 3, 2015;
published September 30, 2015*

Abstract

With the rapid development of smart devices and mobile Internet, the video application plays an increasingly important role on mobile devices. Understanding user behavior patterns is critical for optimized operation of mobile live streaming systems. On the other hand, volume based billing models on cloud services make it easier for video service providers to scale their services as well as to reduce the waste from oversized service capacities. In this paper, the watching behaviors of a commercial mobile live streaming system are studied in a content-centric manner. Our analysis captures the intrinsic correlation existing between popularity and watching intensity of programs due to the synchronized watching behaviors with program schedule. The watching pattern is further used to estimate traffic volume generated by the program, which is useful on data volume capacity reservation and billing strategy selection in cloud services. The traffic range of programs is estimated based on a naïve popularity prediction. In cross validation, the traffic ranges of around 94% of programs are successfully estimated. In high popularity programs (>20000 viewers), the overestimated traffic is less than 15% of real happened traffic when using upper bound to estimate program traffic.

Keywords: Mobile video streaming, network measurement, human factor, traffic estimation

This work was supported by NSFC (Nos. 61425027, 61074034, 61021063, 61174072, 61174105), the National Science and Technology Support Program (2013BAG18B00), the National International Collaboration Project (B06002) and the TNList Cross-Discipline Foundation. This work was also supported in part by Innovation Joint Research Center for Cyber-Physical-Society System. We would thank Beijing Shandong Technology for the operational data of Dopool TV as well as providing computational resources.

1. Introduction

Audience behavior pattern is one of the fundamental problems in mass media, which is the basis of program production, scheduling and advertising. Recent years, much attention is paid on network video streaming services from both practical and academic perspectives [1,2]. With the proliferation of smart devices, mobile video is booming and changed the landscape of mass media. More people tend to watch TV on mobile devices, especially young people [3]. From the technical perspective, mobile video already represented half of mobile Internet data at the beginning of 2012, and will grow at a CAGR of 69% between 2013 and 2018 according to Cisco's forecast [4]. The traffic of video service comes from watching behaviors of its audience, and makes up large proportion of cost of the service. Characterizing user watching behaviors has a great importance on system simulation [5–7], traffic planning [8–10], and optimization of system design [11–13] of video services.

Acquiring accurate large scale behavioral data of traditional TV relies on phone survey or people-meter logs from a small group of audience that it is impossible to collect data in a large scale and may lead to inaccurate estimations. In network based media systems, it is possible to collect usage data from the whole audience population for the first time. Generally network based video services falls into two categories: live streaming and video-on-demand (VoD). Live streaming service, including IPTV and on-line live streaming, is similar to the traditional TV broadcasting. Video content on live streaming is only available when it is broadcast on the channel. Viewers have to join the channel when the program is on the channel to watch it. On the other hand, VoD provide the service for pre-recorded contents to users that the content is available any time. Measurements are adopted on various types of services, including online live streaming systems [13–15], IPTV systems using set-top boxes [5,16–18], as well as VoD systems [19–23]. Most recent studies about video streaming on mobile devices focus on VoD systems [22,24–26], whereas live streaming on mobile devices are with much less studies [27]. Since this paper falls into the category of live streaming, we will concentrate on live streaming in the rest of the paper.

Channel popularity measured by number of accesses and length of each access (also known as session length) are two important aspects of watching patterns in live streaming system. Most existing measurement studies concerned both aspects, while only a few studied the relationship between them. Longer overall session length in popular channels is reported in PPLive [14] by comparing the CDF of session length distribution of a popular channels with an unpopular one. Cha et al. [16] observed higher surfing probability in unpopular channels of an IPTV system. Based on calculated correlation coefficient of channel popularity and session length, [15] reported weak correlation in online streaming, while little to no correlation is reported in a recent study of commercial IPTV [17].

Intuitively, it is the program on a channel rather than the channel itself attracting users to watch. [5] reported spikes of user joining/departure events on IPTV channels caused by the beginning/end of programs, suggesting user behaviors are highly correlated to programs. Watching pattern of users should be studied at program level. However, in existing studies on IPTV [5,16,17], on-line live streaming [14,15], and our previous measurement study [27], watching patterns are studied on the level of channels instead of programs.

Mobile device provides a good opportunity to understand user engagement in videos. On large screens (PC/IPTV), logged watching behaviors on one device might be a shared experience of several viewers. In contrast to large screens, mobile watching is a more personal

experience that the device can directly correspond to its owner. Viewers have to keep the playback in foreground and it is unlikely for them to leave their mobile devices unattended as they do on large screens [16]. The playback time is more likely to be the actual watching time of the viewer.

On the other hand, video streaming is a traffic intensive service. Traditionally, the video service provider has to rent large enough bandwidth upfront to fulfill the demands of the audience at peak hours, while most of the rented bandwidth is idle as there are not many viewers most of the time. What is more, some live services only broadcast specific events. Cloud services provide more flexible options for content distribution that video service providers do not need to oversize their infrastructures any more. The cloud-centric system design for media services receives more and more attention in recent years [28,29]. From the system operation perspective, there are various billing models of public cloud services. In the traffic volume based billing models adopted by main stream cloud services [30], the cloud service will auto scale the bandwidth on edge servers to fulfill requests from users, and it is no need for content providers to actively allocate bandwidth capacity in traffic volume based models. Amazon CloudFront provides on-demand and reserved capacity options for transferred traffic volume [30]. The leading cloud service provider in China, Ali Cloud [31], provides various options to its customers, the usage can be charged by the peak bandwidth in the granularity of day or the traffic volume at the granularity of hour. Generally, discounts are designed for reserved capacity in traffic volume based billing models. Thus, the estimation of traffic generated by video contents is important for video service providers to choose proper cloud service billing options to save distribution costs. Existing traffic prediction methods of video services focus on the dynamics of bandwidth requirements in the near future using time series analysis methods [8], which require large amount of historical data of bandwidth dynamics.

In this paper, we focus on the watching pattern of live TV programs on mobile devices with a content-centric manner to study the total watching time and traffic volume generated in programs rather than user behaviors on channel level. The watching pattern is further utilized to estimate traffic volume generated by the program, which is useful on data volume capacity reservation as well as billing model selection in cloud services. The analysis of watching pattern is based on realistic data collected from a large scale commercial mobile live TV system which broadcasts satellite TV channels in China on mobile devices. We find that the total watching time of long entertainment genre programs which contributes major proportion at traffic peaks distributed in a linear band with increasing program audience size, where the slope of the band is the maximum watching intensity among all the programs. A strong correlation is observed between program popularity and watching intensity. There are both programs with low and high watching intensity among unpopular programs, whereas all the popular programs are with high watching intensity. It turns out the content-aware analysis is able to capture the correlation between popularity and watching intensity in programs.

We propose a program traffic range estimation framework based on watching patterns of mobile TV, which is suitable for traffic volume reservation on cloud services using traffic volume based billing models. The impact of popularity error on traffic range is analyzed, and cross validation is adopted to verify results of popularity prediction and traffic range estimation. It successfully forecasts the traffic range of over 90% of programs. The average over-provisioned traffic is lower than 15% of program traffic in popular programs when the estimated upper traffic bound is used to reserve service capacity.

The remainder of this paper is organized as follows. Section 2 gives an overview of the system and the dataset. Program watching patterns are presented and analyzed in Section 3. In

Section 4, watching patterns of programs are exploited to forecast program traffic range. Error analysis and cross validation are presented as well. Finally, we discuss implications from our findings and conclude the paper in section 5.

2. Preliminaries

2.1 System Background

We study a mobile video streaming system “Dopool TV”¹, which provides free on-line TV live streaming of more than 350 channels covering most satellite TV channels in China on mobile devices. To some extent, it is a substitution of conventional TV on mobile devices, since many channels on it are broadcast simultaneously on TV. It adopts HTTP Live Streaming to transfer video contents to various types of mobile devices running iOS, Android system. Live video signals are encoded in H.264/AAC format with resolution of 320×240 at 256 kbps for video and 64kbps for audio on encoding servers. Channels are with single bitrate that bitrate adaptation is not applied in the system. The system uses a C/S architecture, and CDN is used to distribute video contents to viewers. Requests from clients are dispatched to different edge servers with DNS-based request routing.

The player software offers a channel catalogue, on which EPG (Electronic Program Guide) information is shown along with channel names. Users may choose from the catalogue or search for channels to play. It does not provide functionalities of pause and seeking in live streaming.

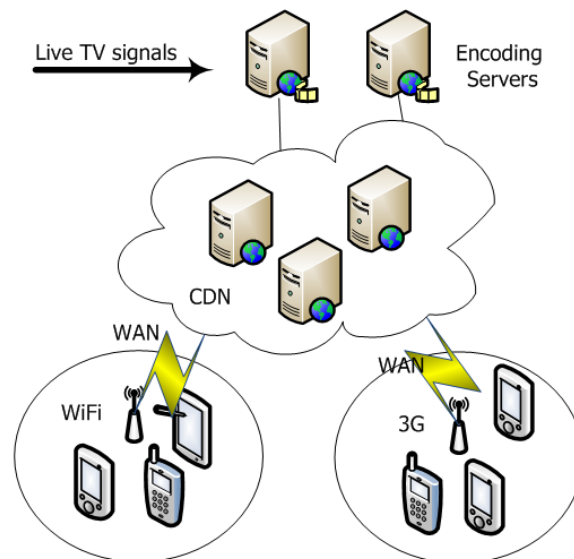


Fig. 1. System architecture of Dopool TV

2.2 Log Collection

To study user behaviors at program level, we integrate the watching trace of users with EPG information of TV programs.

The watching trace is reported by Dopool TV players at the start and end of playing. Playback sessions are distinguished by the watcher, the channel, starting time, and the duration

¹ http://www.dopool.com/?page_id=986

of the session. The watcher is defined by a universal unique ID that is assigned to each copy of the player and remains the same on the device. The watching length does not count the buffering time of the video. For privacy concern, the ID of the watcher is a random globally unique identifier, and the player will not collect any information related to the owner of the device. Clients also report other information such as network connection type, IP addresses in the trace.

The system does not keep the EPG data of channels. The EPG dataset is a program table of channels collected from TVMAO², a website providing schedule and detailed introduction of TV programs in China. Each EPG entry includes: channel, start time, end time, title, and genre of the program. The three item tuple (channel, start time, end time) defines a TV program, whose content is described by its title and genre. Genres provide a good granularity to put similar programs together as well as to distinguish different types of programs.

We integrate watching session data and EPG data by correlating starting time of sessions with program periods. If a playback session covers the period on the corresponding channel of a program, then the playback is considered as a playback session of the program. We further aggregate the length of playback sessions of the same watcher by programs to the time spent on viewing (TSV). TSV is the total time of the viewer spent on the channel within the program period. The watcher may switch network connection between 3G or WiFi, or respond to messages received by the device, resulting in multiple sessions in the program. We focus on the TSV in programs because it reflects the interests of the watcher to the program content, and insensitive to interruptions caused by the device, network and in-program commercial breaks.

2.3 Data Overview

The watching time and EPG data covers a period from November 24th to December 31st, 2012. Among the watching time data, we pick all the long entertainment shows with length ranging from 105 to 120 minutes on the top 15 channels that form the dataset for analysis. The statistics of the dataset are listed in [Table 1](#).

We focus on the top 15 channels because of the availability of matched EPG data. The top 15 channels are satellite public TV channels accounting for around 70% of active devices, covering most popular main stream shows in China as well as some unpopular ones. Entertainment shows are chosen because entertainment is the main motivation of users' engagement in mobile TV [32] that remarkable proportion of the mobile TV traffic is from the entertainment shows genre. The content pattern of these shows are similar, while other genres, take TV series for example, may have many sub-genres like Sci-Fi, Actions, Crime, etc. Moreover, the content of entertainment shows is generally independent to previous episodes. Thus, correlations among selected programs are not considered. The program length is limited to the range between 105 to 120 minutes, which is about two hours. This length range covers most popular main stream entertainment shows in China, which are typically broadcast at the evenings of weekends. Among all the entertainment programs, these 100 programs contribute 59.1% of accesses and 69.5% of total watching time. In the weekly traffic peak happened in weekend evenings in the top channels, traffic from these programs occupies a proportion of 71.9%. Understanding the watching pattern of these shows will be helpful to understand the peak traffic. Because users cannot control playback progress in live streaming, a relative long program length is good for collecting enough samples as well.

² <http://www.tvmao.com>

Since factors such as program length, content pattern, and correlation among programs are eliminated, the long entertainment shows provide good samples to study the impact of content popularity on its total watching time. It is also meaningful from the traffic perspective in that the samples play an important role in peak traffic. In the rest part of the paper, we use the aggregated watching time records (TSV) from sessions in long entertainment shows as the dataset to study user watching patterns.

Table 1. Overview of watching data used in the paper

Num. Programs	Len. Of Programs (min)	Num. Devices	Num. Sessions	Num. TSV records	Total watching Times (s)	Ratio of weekly peak traffic
100	105 - 120	5.06×10^5	1.29×10^6	7.45×10^5	6.18×10^8	71.9%

3. Watching patterns of programs

We focus our analysis on 100 representative long entertainment shows from popular channels on Dopool mobile TV dataset. These programs are all from the most populous channels that the sample size could support our analysis. Entertainment programs with length around 2 hours are chosen to avoid the impact of type and length of programs. We believe it is representative for online live event broadcasting to mobile devices.

Audience watching pattern is studied by viewer at program level. The popularity of a program is defined as the number of viewers engaged in the program, denoted by x . The time spent watching (TSV) of a user for a program is the aggregated duration of watching time during the period of the program. It represents the total length of watched program content and total network traffic generated by the watching behavior, as bitrate adaptation is not applied in the system. In contrast to single playback session, it is insensitive to commercial breaks. The watching intensity of a program, denoted by w , is defined as the per-capita TSV of audience who watched it. The total watching time, denoted by y , is the sum of all the time spent on the program of all the viewers. It can be converted to the total network traffic generated by watching behaviors given the video bitrate.

3.1 Watching time of viewers at program level

We rank programs by their popularity, and then plot their popularity against their ranks in [Fig. 2](#) in log-scales. Program popularity exhibits some power law properties with a dropped tail, where $x \sim i^{-0.88}$ shown as the red dashed line in the figure. Audience attention is concentrate on a small number of programs, while the audience size in most programs is small.

The total watching time is plotted against the popularity of programs as scatter plot in [Fig. 3](#). Each point represents a program. It is obvious $y_i = w_i x_i$, where i is the identifier of program and the watching intensity w_i is the slope to the origin.

The minimum and maximum of w_i in our data are $w_{\min} = 179.89$ seconds and $w_{\max} = 1258.43$ seconds. These two observed extreme values are used as minimum and maximum baselines of watching intensity of such type of programs. The guide lines L1 and L2 are $y = w_{\max} x$ and $y = w_{\min} x$, representing the maximum and minimum possible watching time of program with popularity x according the two baselines of watching intensity.

However, the total watching time of programs are distributed in a band below L1 rather than in the whole area bounded by L1 and L2. The Pearson correlation coefficient between x_i and y_i is $r = 0.9876$, suggesting linear relationship can be used to describe the data.

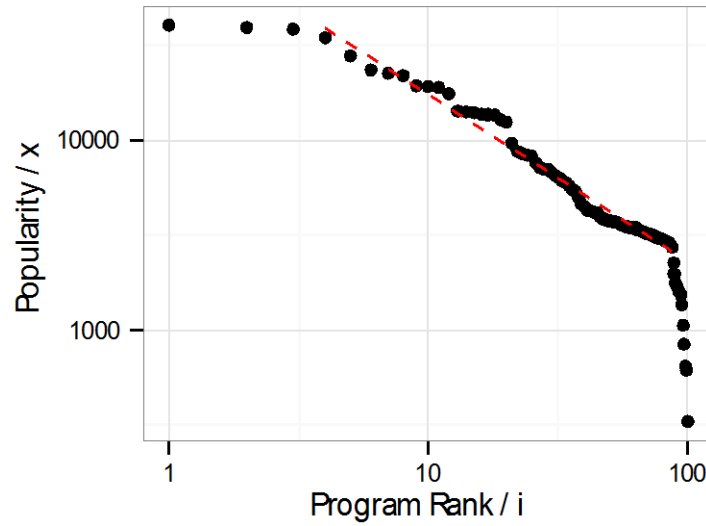


Fig. 2. Program popularity vs. program rank

We use linear upper and lower bounds to describe the band area in Fig. 3. The upper bound and lower bound of the total watching time of a program with x viewers are denoted by $Y_U(x)$ and $Y_L(x)$ respectively. As our data exhibits linear relationship between y and x , a linear form $Y_U(x) = \alpha_0 + \alpha_1 x$ is used for the upper bound, where parameters are decided by the optimal solution of the LP problem in Eq. (1). The constrain $Y_U(0) = 0$ comes from the fact that no watching time is generated if there are no viewers.

$$\begin{aligned}
 &\underset{\alpha_0, \alpha_1}{\text{minimize}} && \sum_i Y_U(x_i) - y_i \\
 &\text{subject to} && Y_U(x_i) \geq y_i \\
 &&& Y_U(0) = 0
 \end{aligned} \tag{1}$$

The optimal solution $\alpha_0^* = 0, \alpha_1^* = w_{\max}$ gives the upper bound in Eq. (2).

$$Y_U(x) = w_{\max} x \tag{2}$$

The lower bound $Y_L(x)$ is given by a piecewise linear function having the form $Y_L(x) = \max(w_{\min} x, \beta_0 + \beta_1 x), x \geq 0$. The first term $w_{\min} x$ is controlled by minimum average watching time of this type of program, as line L2 in Fig. 3. The latter term describes the lower bound of the band when x becomes large, as the line L3. Parameter β_0 and β_1 are determined by the optimal problem Eq. (3). The constrain $\beta_1 \leq \alpha_1^*$ guarantees $Y_L(x) \leq Y_U(x)$ for all $x > 0$.

$$\begin{aligned}
& \underset{\beta_0, \beta_1}{\text{minimize}} && \sum_i y_i - Y_L(x_i) \\
& \text{subject to} && Y_L(x_i) \leq y_i \\
& && \beta_1 \leq \alpha_1^* \\
& && Y_L(0) = 0.
\end{aligned} \tag{3}$$

The optimal solution of the problem gives $\beta_0^* = -7.48 \times 10^6$, and $\beta_1^* = w_{\max}$. As $\alpha_1^* = \beta_1^*$, L1 is parallel to L3, and the gap is constant $y_0 = -\beta_0^*$. The value of y_0 is actually decided by the active constrain of Eq. (3), which corresponds to the program on the lower bound. L2 and L3 intersect at $x_c = y_0 / (w_{\max} - w_{\min}) = 6933.9$, thus $Y_L(x)$ can be rewritten as Eq. (4).

$$Y_L(x) = \begin{cases} w_{\min} x & 0 \leq x \leq x_c \\ w_{\max} x - y_0 & x > x_c \end{cases} \tag{4}$$

The upper and lower bounds of total watching time are plotted in Fig. 3 as solid and dashed lines in different colors.

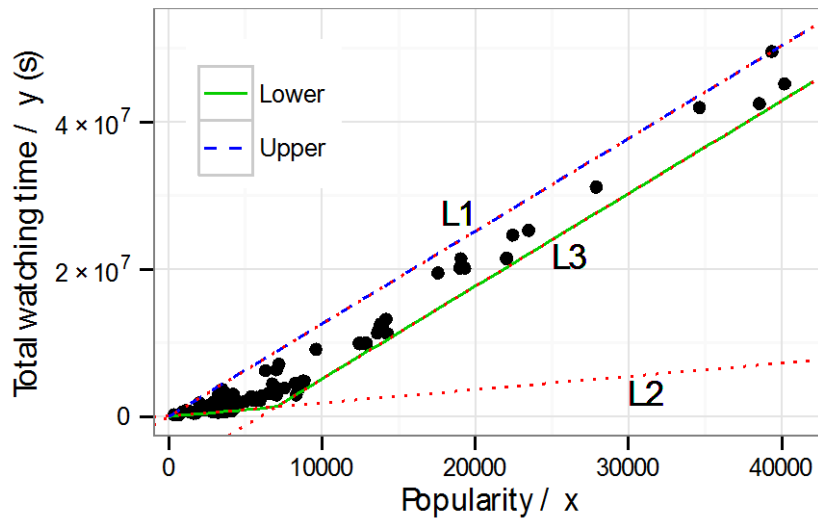


Fig. 3. Total watching time against popularity of entertainment programs. Each point represents a program. Guide lines L1, L2, and L3 are plotted in red dotted lines.

Total watching time is the cumulated TSV of its audience that is related to the size of its audience (popularity) and the watching intensity representing the average TSV of audience. Intuitively, higher watching intensity of a program means that its audience tends to stay longer in it. The watching intensity bounds of a program at a given popularity x can be derived from the bounds of watching time in Eq. (2) and Eq. (4), which have forms in Eq. (5) and Eq. (6), respectively.

$$W_U(x) = \frac{Y_U(x)}{x} = w_{\max} \quad x > 0 \quad (5)$$

$$W_L(x) = \frac{Y_L(x)}{x} = \begin{cases} w_{\min} & 0 < x \leq x_c \\ w_{\max} - \frac{x_c}{x}(w_{\max} - w_{\min}) & x > x_c \end{cases} \quad (6)$$

The lower and upper bounds are depicted in Fig. 4 as solid and dashed lines, respectively. When $x > x_c$, the second term of $W_L(x)$ represents the gap between the two bounds of watching intensity. It decreases reciprocally with increasing popularity x . As a consequence, the average watching time of hot programs are near w_{\max} .

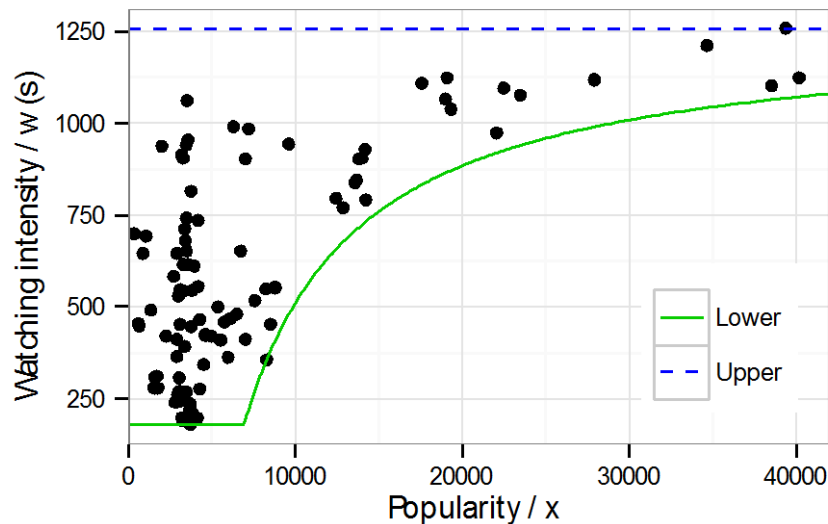


Fig. 4. Watching intensity of programs

3.2 The impact of programs on the watching pattern

To investigate the impact of programs on the watching pattern, we further examine departures of viewers in programs. The departure time of a viewer is the end time of her last playback session of the program. Obviously, a viewer tends to stay in a program, if she is attracted by the program content. Synchronized departures are observed near the end of programs. In programs with more than 10000 viewers, on average there are 25.1% of viewers depart in the last 10 minutes of the program, and these viewers contribute on average 51.8% of the total watching time of the program. In programs with lower popularity (popularity < 10000), the mean ratio of audience departs at the end of the program is 16.6%, and on average they contribute 33.3% of program watching time. The fact that considerable proportion of audience tend to watch to the end of the program reflects that watching behaviors is synchronized to video contents, especially in popular programs. The synchronized departures suggest that the program is the motivation of watching for users. Larger proportion of audience watching to the end of the program in popular programs also explains the trend of large watching intensity in popular programs.

The Pearson correlation coefficient between popularity and watching intensity of programs $r = 0.685$ suggests a strong correlation between the two factors. On the other hand, the Pearson correlation coefficient between the number of playback sessions and average playback session length of channels 0.225 indicates a weak correlation at channel level, which is comparable to the correlation at channel level of online live streaming ($r = 0.2$) [15]. The strong correlation at program level suggests that the intrinsic watching pattern of mobile TV caused by content can be better captured at program level.

Based on the observation of synchronized departures and the comparison of correlation between popularity and watching intensity, the analysis at program level is able to capture the intrinsic pattern of watching behaviors of mobile TV caused by the high correlation of user behaviors and the video content.

3.3 Discussions

w_{\max} , w_{\min} , and x_c are characteristic variables on describing the mobile TV watching pattern and its traffic. The maximum watching intensity w_{\max} among all the 2 hour long programs is around 21 minutes. It characterizes the duration of video consumption on small screens, and determines the traffic needed for an increased audience size in our traffic model as well.

The linear band distributed of program watching time gives the relation between the total watching time of a program and its popularity. The biased distribution of program popularity together with the watching pattern suggests that the demand of traffic varies in a wide range and only a few programs have intensive demand of resources. If the system capacity is planned using fixed upfront capacity (as used by Dooool TV), a large proportion of resources will be idle most of the time. Cloud services can be adopted to scale the system as well as save distribution costs.

If the program popularity can be estimated before broadcasting, the watching pattern can be used to estimate the total watching time and traffic volume generated by viewers. Because the width of total watching time range is constant, it may give relative good estimations in popular programs. It is useful for reserve traffic capacity in cloud services. We will detail the estimation of program traffic in the next section.

The watching pattern at program level does not exist anymore in the analysis on channels without considering the program schedule. The observed watching pattern in programs can be attributed to the impact of content because of highly synchronized behaviors with program schedule.

4. Program Traffic Estimation

As network traffic comes from watching behaviors of the audience, the watching pattern of mobile TV can be exploited to forecast traffic volume in programs. Cloud services offer more flexible options for video services to scale. The estimated traffic is useful for video service provider to reserve cloud data volume to save content distribution costs or supporting billing model selection for cloud services.

According to our previous results, the total watching time of programs are distributed near the upper bound which increases linearly as the audience size of the program, and the gap between the upper and lower bounds of total program watching time is constant. If program popularity could be predicted, program traffic range can be estimated accordingly.

In this section, the total watching time is used to represent the generated traffic for coherence. It can be easily converted to data volume on network using video bit rate. As the bit

rate of channels of Dopool TV is 320 kbps (video + audio), each second of video correspond to 40 kilobytes of transferred data.

4.1 Prediction of Popularity

There are many ways to estimate audience rating of TV programs. Traditional methods are based on historical data and various measurements of television [33,34], while new methods estimate audience rating based on trends of Internet search engines and activities on social networks [35,36]. Results of audience rating (i.e. program popularity) estimation methods can be used as the input of our traffic estimation framework. As we do not have data of historical audience rating or social network trends, we proposed a naive on-line method based on two observations of aggregated number of viewers after the program begins. It is still meaningful to estimate the traffic volume after observation moments only by the naive popularity estimation. The result of naive popularity estimation can be used to tune the estimation of other program popularity prediction methods as well.

As lengths of programs differ slightly, the time of programs is normalized as a number from 0 to 1. The moment 0 and 1 correspond to the beginning and the end of the program, respectively. The counting process $N(t)$ represents the cumulated number of viewer who accessed the program until moment t . If a viewer starts playback multiple times during the program, only the first time will be counted in $N(t)$. It is obvious $N(0)=0$ and $N(1)=x$ from the definition, where x is the program popularity. Fig. 5 gives the normalized cumulative audience size of three representative programs. Since program popularity varies in a wide range, the number of cumulative audience size is normalized by $N(1)$. Different programs are distinguished by colors and markers. The red dotted reference line has the slope of 1, which corresponds to the constant arrival rate of viewers. The arrival process above the reference line indicates there is a higher arrival rate at the beginning of the program, and the arrival rate decreases during the program. On contrary, arrival processes below the guideline suggest increasing arrival rate during the program. The relative smooth shape of the cumulative audience size suggests that the arrival rate of new viewers in programs changes progressively rather than abruptly.

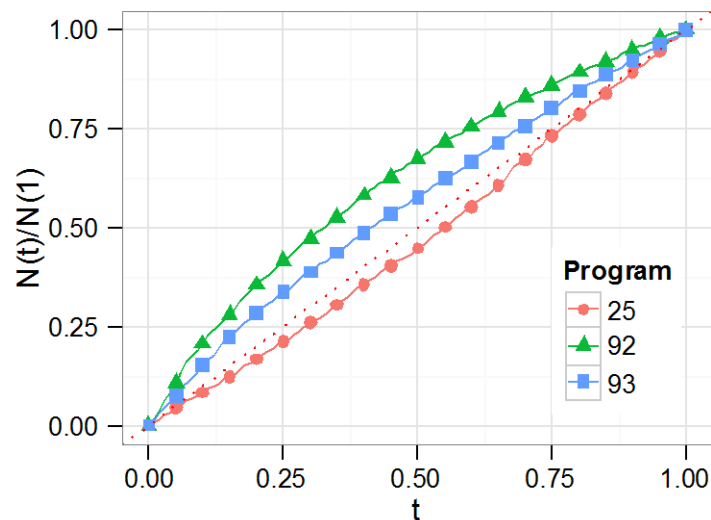


Fig. 5. Normalized cumulative audience size of representative programs

For simplicity, we use two observations of cumulative number of program viewers to estimate the popularity of the program. Let $s_k = N(t_k)$ ($k=1,2$) be the observation at t_k , where the moments t_1, t_2 satisfy $0 < t_1 < t_2 < 1$. A linear model is proposed to describe the relationship between total viewers x and the two observations s_1 and s_2 , as shown in Eq. (7). It is able to capture the absolute arrival rate of viewers as well as the first derivative of arrival rate in the program. Linear regression is used to estimate parameter γ_0, γ_1 , and γ_2 .

$$x = \gamma_0 + \gamma_1 s_1 + \gamma_2 s_2 \quad (7)$$

Observation moments t_1 and t_2 may influence the accuracy of popularity estimation, as the arrival process of viewers is random. Fig. 6 gives the relationship between the two observation points with square root of the variance of random error in the regression (as RMSE in the figure). Different colors and markers are used to represent different values of t_2 . It can be seen that generally the level of RMSE is decided by the values of t_2 . More accurate estimation could be obtained if t_2 is larger. The residual error is small when t_1 is near 0.12 for $t_2 \geq 0.2$, then the RMSE fluctuates when t_1 is near 0.16. The RMSE rises when t_1 approaches t_2 . As the watching behaviors happened within t_2 is the sample used to estimate the program audience population, the observation period cannot be too long. We restrict the observation period within 20% of programs, namely set $t_2 = 0.2$. t_1 is set to 0.12 to acquire relative accurate estimation based on the observations in Fig. 6.

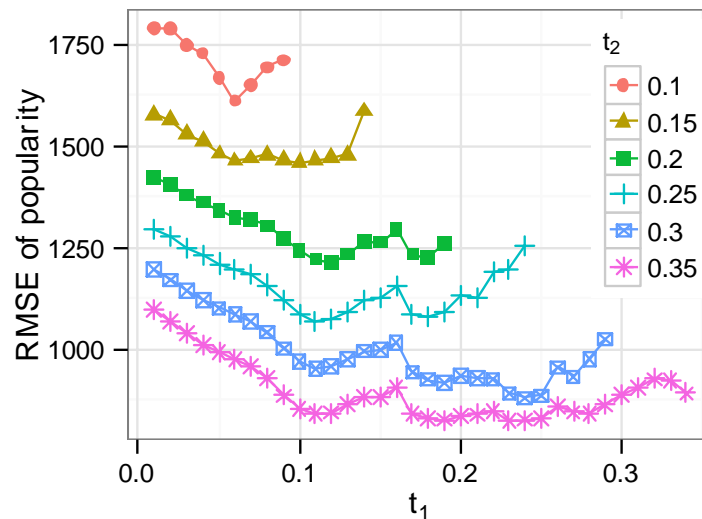


Fig. 6. Impacts of observation positions on the regression results

Table 2 gives the regression results when observing at $t_1 = 0.12$ and $t_2 = 0.2$. The effects of the three parameters are all significant. The adjusted R^2 and the significance level of the F-statistic suggest the model is valid to describe the relationship between the popularity and the two observations of cumulative program audience sizes.

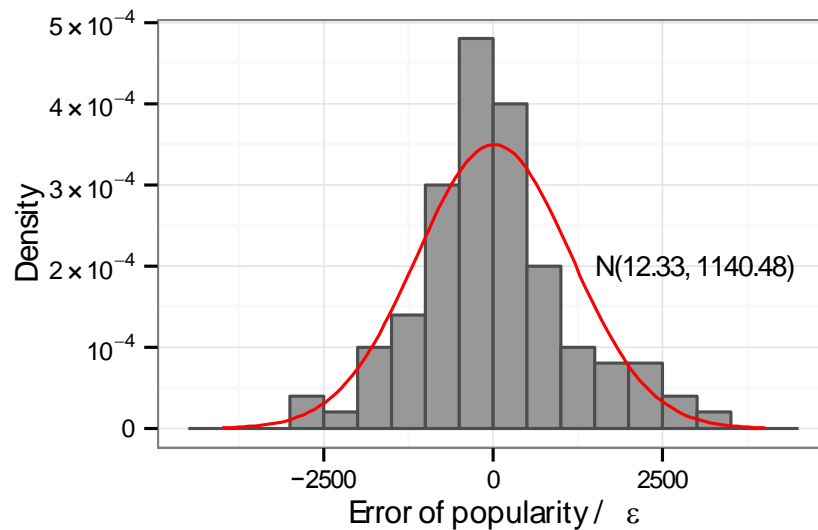
Table 2. Linear regression results of program popularity

Term	Estimate	Std. Error	t-statistic	Significance
γ_0	1153.1449	140.7471	8.193	***
γ_1	-11.3972	1.1017	-10.345	***
γ_2	10.3327	0.7019	14.720	***

Number of samples: 100
Residual standard error: 1099 (df = 97)
Adjusted R^2 : 0.983
F-statistic: 2860 (df = 2; 97) ***

*** $p < 0.001$, ** $p < 0.01$, * $p < 0.05$

Leave-one-out cross validation is adopted to validate the linear model for popularity prediction. Each time a program is chosen as test sample. The parameters of Eq. (7) are estimated from the training set excluding the test sample. Estimation error is calculated as $\hat{q} = \hat{x}_i - x_i$, and its distribution is shown in Fig. 7. It turns out $\hat{q} \sim N(12.33, 1140.48)$ with $p = 0.4794$ given by KS test. The density function of the fitted distribution is plotted as the red curve in the figure.

**Fig. 7.** Error of program popularity estimation

4.2 Traffic Estimation based on Predicted Program Popularity

The lower and upper bounds of program traffic are calculated by Eq. (2) and Eq. (4) based on predicted popularity \hat{x} . The error of popularity prediction will affect the estimation of traffic bounds, as the bounds are calculated based on the estimated popularity. We care more about the traffic estimation when there are large number of viewers in the program, thus we concentrate the analysis on the programs with $x > x_c$. With the error of popularity estimation $\hat{o} = \hat{x} - x$, the estimations of traffic bounds have the form in Eq. (8) and Eq. (9).

$$\hat{y}_L = Y_L(\hat{x}) = w_{\max}(x + \delta) - y_0 \quad (8)$$

$$\hat{y}_U = Y_U(\hat{x}) = w_{\max}(x + \delta) \quad (9)$$

We introduce $u = Y_U(x) - y$ to represent the distance of the program's traffic to the upper bound, where $0 \leq u \leq y_0$ in observed program traffic data. The distribution of u is denoted by $f_u(t)$, where the empirical density is given as the histogram in **Fig. 8**. The traffic of most programs are distributed in the middle of the two bounds that errors on traffic bounds introduced by δ do not necessarily lead to failed estimation of the traffic range.

The error of popularity estimation δ is assumed to follow a normal distribution $N(0, \sigma)$, with probability density function is $f_\delta(t) = \phi(t / \sigma)$. We say the traffic of a program is overestimated if the real traffic is below the estimated lower bound ($y < \hat{y}_L$), whereas the traffic is under estimated if the real traffic is above the estimated upper bound ($y > \hat{y}_U$). From Eq. (8) and the definition of u , program traffic is overestimated if $\delta > (y_0 - u) / w_{\max} = \xi$, with probability $p_L = P(y < \hat{y}_L) = \int_\xi^\infty f_\delta(t) dt = 1 - \Phi(\xi / \sigma)$. Similarly, from Eq. (9), program traffic is underestimated if $\delta < -u / w_{\max} = \eta$, with probability $p_U = P(y > \hat{y}_U) = \int_{-\infty}^\eta f_\delta(t) dt = \Phi(\eta / \sigma)$. The expectations of the two probabilities depend on the distribution of u among programs, as in Eq. (10) and Eq. (11).

$$E[p_L] = \int_0^{y_0} [1 - \Phi(\xi / \sigma)] f_u(t) dt \quad (10)$$

$$E[p_U] = \int_0^{y_0} \Phi(\eta / \sigma) f_u(t) dt \quad (11)$$

Using the empirical distribution of u , the expected over/under-estimation probabilities under different σ is shown in **Fig. 9**. The standard deviation of the error of our popularity prediction $\sigma = 1140.48$ is plotted as red dotted vertical line in the figure, where expected overestimation probability $E[p_L] = 9.48\%$ and underestimation probability $E[p_U] = 2.99\%$. The expected overestimation probability increases with σ much faster than underestimation probability.

If we want to guarantee that the expected out-of-bound probabilities do not surpass a given target δ , the estimation of traffic bounds need to be extended. We extend the estimated lower and upper traffic bounds by l_L and l_U respectively, and have the extended traffic bounds in Eq. (12) and Eq. (13).

$$\tilde{y}_L = Y_L(\hat{x}) - l_L = w_{\max}(x + \delta) - y_0 - l_L \quad (12)$$

$$\tilde{y}_U = Y_U(\hat{x}) + l_U = w_{\max}(x + \delta) + l_U \quad (13)$$

The expected out-of-bound probabilities of programs in Eq. (14) and Eq. (15) can be derived with similar procedures as in Eq. (10) and Eq. (11).

$$E[y < \tilde{y}_L] = \int_0^{y_0} \left[1 - \Phi\left(\frac{y_0 - u + l_L}{w_{\max} \sigma}\right) \right] f_u(t) dt \tag{14}$$

$$E[y > \tilde{y}_U] = \int_0^{y_0} \Phi\left(-\frac{u + l_U}{w_{\max} \sigma}\right) f_u(t) dt \tag{15}$$

The value of l_L and l_U corresponding to $E[y < \tilde{y}_L] = \delta$ and $E[y > \tilde{y}_U] = \delta$ can be calculated numerically using the empirical distribution of u with different σ . The numeric results of l_L^* and l_U^* are plotted in Fig. 10 (a) and (b) and represented as multiples of y_0 , where red dotted line represent $\sigma = 1140.48$ of the naïve popularity prediction. Negative values of l_L^* and l_U^* mean the out-of-bound probability is already less than the given target δ at that level of σ . Guaranteed out-of-bound probability is acquired in exchange of expanded traffic range, which will lead to more over-provisioned data in service capacity reservation. A relative accurate popularity estimation is important to forecast program traffic range.

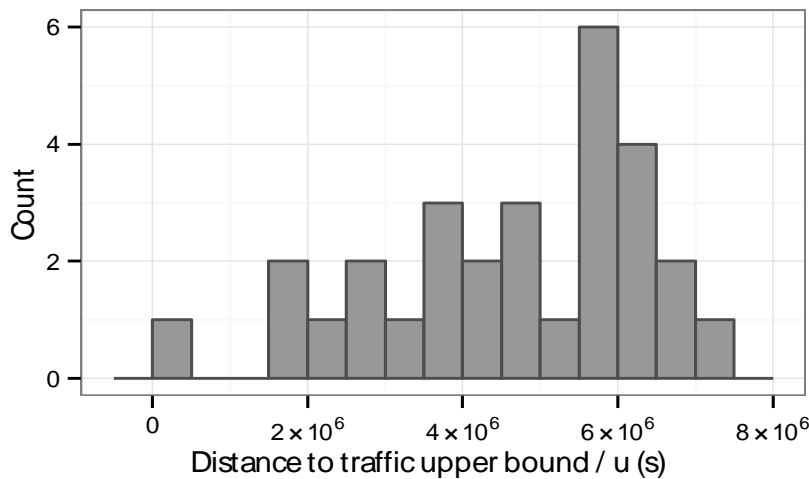


Fig. 8. Empirical distribution of u for popular programs with $x > x_c$

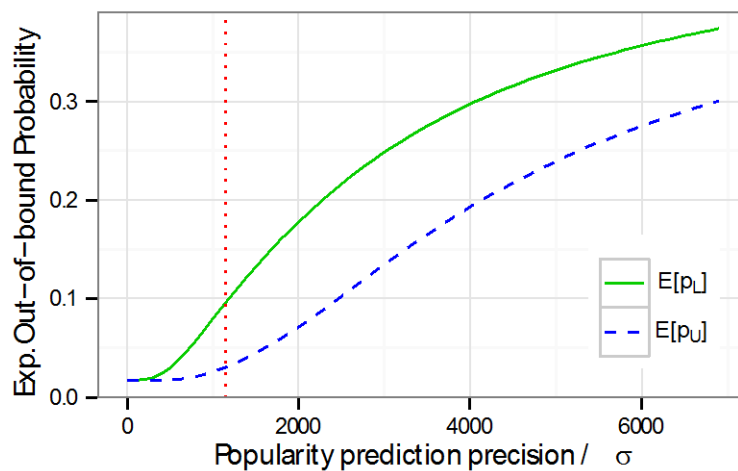


Fig. 9. Expected out-of-bound probability under different accuracy of popularity estimation.

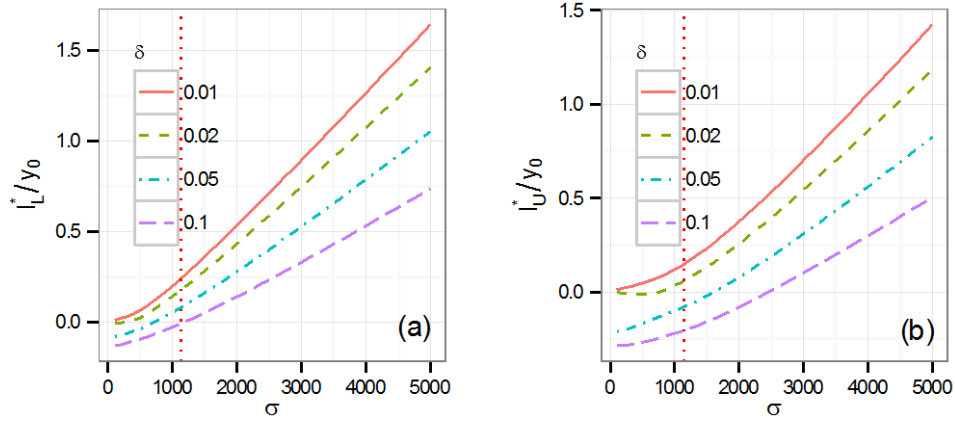


Fig. 10. Extension of bounds required to achieve the expected out-of-bound probability

4.3 Cross Validation of Traffic Estimation

Leave-one-out cross validation is adopted to test the estimation of program popularity and traffic range. The validation procedure goes through all the programs as follows. Each time a program is selected as the test sample, and other programs form the training set. Model parameters, including parameters in Eq. (7), (1) and (3), are derived from the training set. Then the parameters are used to predict the popularity and calculate traffic bounds of the test sample. The real value of the test sample is used to validate the estimation results. If the test sample satisfies $Y_L(\hat{x}_i) \leq y_i \leq Y_U(\hat{x}_i)$, the real value of traffic situated between the estimated bounds. Otherwise, the real value is situated out of estimated bounds, and the model failed to estimate the traffic range of the program.

In the cross validation, the real value of traffic of 6 (6% of the total 100) programs are not situated in the estimated bounds, as listed in **Table 3**. d is defined as the distance of the traffic to the nearest bound for out-of-bound programs, as in Eq. (16). Positive values represent distances above upper bound, while negative ones represent those below the lower bound. The ratio of d_i to the gap between the two bounds at \hat{x} are calculated. The first two programs in the table are with popularity $x > x_c$, where the gap size is y_0 . The distances to the bounds from the upper and lower in the two programs correspond to 1.55% and 28.22% of the gap between the two bounds. The relative high distance to the lower bound of program is due to the over estimation of the program popularity. However, the value out-of-bound distance only correspond to 4.67% of the traffic of the program, which is acceptable.

$$d_i = \begin{cases} y_i - Y_U(\hat{x}_i) & y_i > Y_U(\hat{x}_i) \\ y_i - Y_L(\hat{x}_i) & y_i < Y_L(\hat{x}_i) \end{cases} \quad (16)$$

If the upper bound is used to provision service resources in the cross validation, the histogram of over-estimated traffic of programs are shown in **Fig. 11**. Two programs out of the $[0, y_0]$ range correspond to program 2 and 30 in **Table 3**. The average over provisioned traffic is 3.19×10^6 seconds (solid line in **Fig. 11**), which is about 42.6% of the gap size y_0 , or around

15% of the traffic of a program with 20000 viewers. The fraction of over provisioned traffic to program traffic will be even lower in popular programs with larger audience size.

Table 3. Programs with estimated traffic out of bounds

i	x_i	\hat{q}	y_i	d_i	% of gap
2	40147	3364.322	45167140.4	-2110345.342	-28.22%
30	22443	-2987.6105	24598899.3	115660.576	1.55%
38	3715	-116.3816	668302	-22811.718	-0.59%
51	331	1001.2665	231665.5	-7999.744	-0.56%
59	647	998.0026	290885.7	-5038.478	-0.28%
78	614	970.1464	279502.3	-5474.322	-0.32%

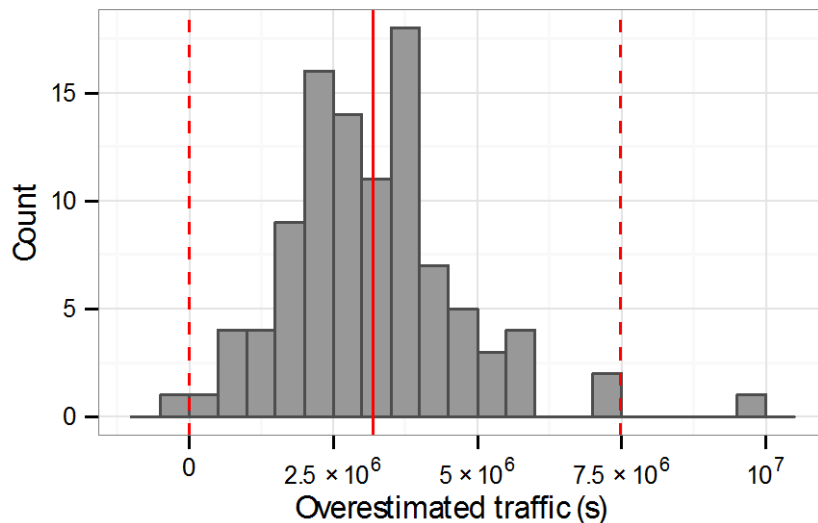


Fig. 11. Overestimated traffic using upper bound. Red dashed lines corresponds to 0 and the gap y_0 . Solid line corresponds to the mean value.

5. Conclusion

In this data-driven study of a large commercial mobile TV system, the watching pattern of viewers is studied at program level from a viewer-content perspective. Compared to content-unaware manner, our approach can better capture watching pattern of mobile TV. We find that the total watching time of programs distributed in a band linearly increasing with program popularity. A relative strong correlation exists between popularity and watching intensity of programs. The lower bound of watching intensity approaches the constant upper bound as the popularity of program increases. That is, there are both “unpopular but long-watched” and “unpopular and short-watched” programs, whereas only “popular and long-watched” programs exist.

We exploited the traffic pattern on estimating watching time / traffic consumption of programs on mobile TV based on predictions of program popularity. A naïve program popularity prediction method is provided based on our data. Other audience rating prediction

method also can be used in our traffic estimation framework. The impact of error of program popularity on traffic estimation was analyzed, and the amount of extension on traffic bounds to achieve given estimation accuracy was discussed as well. Leave-one-out cross validation is used to test traffic estimation of mobile TV. It successfully forecasts the traffic range of 94% of programs. If the upper bound of traffic range is used to reserve traffic capacity for programs, the average over-provisioned traffic is lower than 15% of program traffic in popular programs with more than 20000 viewers.

The total watching time of programs correspond to their network traffic volume, which is important in capacity allocation or choosing billing model in cloud services. The estimated total watching time of a video content can also be used support advertisement pricing, which is the important revenue source for video services. The watching intensity of programs is the per capita attention time of viewers spent on the program. Its characterization is useful for customized and balanced advertisement exposures to viewers. There are other factors that may affect the watching pattern, such as the mobility of user, video bit rate, which will be considered in the future to give more insights of mobile video streaming. High watching intensity of popular programs also suggests opportunities of peer-assisted system design, which may greatly reduce service capacity requirements.

References

- [1] W. Hui, C. Lin, and Y. Yang, "MediaCloud: A New Paradigm of Multimedia Computing.," *KSII Transactions on Internet and Information Systems (TIIS)*, vol. 6, pp. 1153–1170, 2012. [Article \(CrossRef Link\)](#)
- [2] J. Lee, J. Hwang, N. Choi, and C. Yoo, "SVC-based Adaptive Video Streaming over Content-Centric Networking.," *KSII Transactions on Internet and Information Systems (TIIS)*, vol. 7, pp. 2430–2447, 2013. [Article \(CrossRef Link\)](#)
- [3] B. Stelter, "Youths are watching, but less often on TV.," *The New York Times*, February 8, 2012. <http://www.nytimes.com/2012/02/09/business/media/young-people-are-watching-but-less-often-on-tv.html>
- [4] Cisco, "Cisco Visual Networking Index: Global Mobile Data Traffic Forecast Update, 2013–2018.," *White Paper*, Feb. 2014.
- [5] T. Qiu, Z. Ge, S. Lee, J. Wang, J. Xu, and Q. Zhao, "Modeling User Activities in a Large IPTV System.," in *Proc. of the 9th ACM SIGCOMM Conference on Internet Measurement Conference*, Chicago, Illinois, USA: ACM, pp. 430–441, 2009. [Article \(CrossRef Link\)](#)
- [6] W. Tang, Y. Fu, L. Cherkasova, and A. Vahdat, "Modeling and generating realistic streaming media server workloads.," *Computer Networks*, vol. 51, pp. 336–356, 2007. [Article \(CrossRef Link\)](#)
- [7] E. Veloso, V. Almeida, W. Meira, A. Bestavros, and S. Jin, "A hierarchical characterization of a live streaming media workload.," *IEEE/ACM Transactions on Networking*, vol. 14, pp. 133–146, Feb. 2006. [Article \(CrossRef Link\)](#)
- [8] D. Niu, Z. Liu, B. Li, and S. Zhao, "Demand forecast and performance prediction in peer-assisted on-demand streaming systems.," in *Proc. of IEEE INFOCOM 2011*, pp. 421–425, 2011. [Article \(CrossRef Link\)](#)
- [9] D. Niu, H. Xu, B. Li, and S. Zhao, "Quality-assured cloud bandwidth auto-scaling for video-on-demand applications.," in *Proc. of IEEE INFOCOM 2012*, pp. 460–468, 2012. [Article \(CrossRef Link\)](#)
- [10] C. Wu, B. Li, and S. Zhao, "On Dynamic Server Provisioning in Multichannel P2P Live Streaming.," *IEEE/ACM Transactions on Networking*, vol. 19, pp. 1317–1330, Oct. 2011. [Article \(CrossRef Link\)](#)

- [11] K.-W. Hwang, D. Applegate, A. Archer, V. Gopalakrishnan, S. Lee, V. Misra, K.K. Ramakrishnan, and D.F. Swayne, "Leveraging Video Viewing Patterns for Optimal Content Placement," *NETWORKING 2012*, R. Bestak, L. Kencl, L. Li, J. Widmer, and H. Yin, eds., Springer Berlin Heidelberg, pp. 44–58, 2012. [Article \(CrossRef Link\)](#)
- [12] Z. Liu, C. Wu, B. Li, and S. Zhao, "Why Are Peers Less Stable in Unpopular P2P Streaming Channels?," *NETWORKING 2009*, L. Fratta, H. Schulzrinne, Y. Takahashi, and O. Spaniol, eds., Springer Berlin Heidelberg, pp. 274–286, 2009. [Article \(CrossRef Link\)](#)
- [13] C. Wu, B. Li, and S. Zhao, "Diagnosing Network-wide P2P Live Streaming Inefficiencies," *ACM Transactions on Multimedia Computing, Communications and Applications*, vol. 8, pp. 13:1–13:19, Feb. 2012. [Article \(CrossRef Link\)](#)
- [14] X. Hei, C. Liang, J. Liang, Y. Liu, and K.W. Ross, "A Measurement Study of a Large-Scale P2P IPTV System," *IEEE Transactions on Multimedia*, vol. 9, pp. 1672–1687, Dec. 2007. [Article \(CrossRef Link\)](#)
- [15] K. Sripanidkulchai, B. Maggs, and H. Zhang, "An Analysis of Live Streaming Workloads on the Internet," in *Proc. of the 4th ACM SIGCOMM Conference on Internet Measurement*, Taormina, Sicily, Italy: ACM, pp. 41–54, 2004. [Article \(CrossRef Link\)](#)
- [16] M. Cha, P. Rodriguez, J. Crowcroft, S. Moon, and X. Amatriain, "Watching Television over an IP Network," in *Proc. of the 8th ACM SIGCOMM Conference on Internet Measurement*, Vouliagmeni, Greece: ACM, pp. 71–84, 2008. [Article \(CrossRef Link\)](#)
- [17] N. Liu, H. Cui, S.-H.G. Chan, Z. Chen, and Y. Zhuang, "Dissecting User Behaviors for a Simultaneous Live and VoD IPTV System," *ACM Transactions on Multimedia Computing, Communications and Applications*, vol. 10, pp. 23:1–23:16, Apr. 2014. [Article \(CrossRef Link\)](#)
- [18] G. Yu, T. Westholm, M. Kihl, I. Sedano, A. Aurelius, C. Lagerstedt, and P. Odling, "Analysis and characterization of IPTV user behavior," *IEEE International Symposium on Broadband Multimedia Systems and Broadcasting 2009*, pp. 1–6, 2009. [Article \(CrossRef Link\)](#)
- [19] M. Cha, H. Kwak, P. Rodriguez, Y.-Y. Ahn, and S. Moon, "I tube, you tube, everybody tubes: analyzing the world's largest user generated content video system," in *Proc. of the 7th ACM SIGCOMM conference on Internet measurement*, New York, NY, USA: ACM, pp. 1–4, 2007. [Article \(CrossRef Link\)](#)
- [20] Y. Chen, B. Zhang, Y. Liu, and W. Zhu, "Measurement and Modeling of Video Watching Time in a Large-Scale Internet Video-on-Demand System," *IEEE Transactions on Multimedia*, vol. 15, pp. 2087–2098, 2013. [Article \(CrossRef Link\)](#)
- [21] P. Gill, M. Arlitt, Z. Li, and A. Mahanti, "Youtube traffic characterization: a view from the edge," in *Proc. of the 7th ACM SIGCOMM conference on Internet measurement*, New York, NY, USA: ACM, pp. 8–15, 2007. [Article \(CrossRef Link\)](#)
- [22] Z. Li, J. Lin, M.-I. Akodjenou, G. Xie, M.A. Kaafar, Y. Jin, and G. Peng, "Watching Videos from Everywhere: A Study of the PPTV Mobile VoD System," in *Proc. of the 2012 ACM Conference on Internet Measurement Conference*, Boston, Massachusetts, USA: ACM, pp. 185–198, 2012. [Article \(CrossRef Link\)](#)
- [23] H. Yu, D. Zheng, B.Y. Zhao, and W. Zheng, "Understanding user behavior in large-scale video-on-demand systems," *ACM SIGOPS Operating Systems Review*, vol. 40, pp. 333–344, 2006. [Article \(CrossRef Link\)](#)
- [24] I. Ketykó, K. De Moor, T. De Pessemier, A.J. Verdejo, K. Vanhecke, W. Joseph, L. Martens, and L. De Marez, "QoE Measurement of Mobile YouTube Video Streaming," in *Proc. of the 3rd Workshop on Mobile Video Delivery*, Firenze, Italy: ACM, pp. 27–32, 2010. [Article \(CrossRef Link\)](#)
- [25] Y. Liu, F. Li, L. Guo, B. Shen, and S. Chen, "A server's perspective of Internet streaming delivery to mobile devices," in *Proc. of IEEE INFOCOM 2012*, 2012, pp. 1332–1340. [Article \(CrossRef Link\)](#)
- [26] Y. Xiao, R.S. Kalyanaraman, and A. Yla-Jaaski, "Energy Consumption of Mobile YouTube: Quantitative Measurement and Analysis," in *Proc. of the Second International Conference on Next Generation Mobile Applications, Services and Technologies*, pp. 61–69, 2008. [Article \(CrossRef Link\)](#)

- [27] Y. Li, Y. Zhang, and R. Yuan, "Measurement and Analysis of a Large Scale Commercial Mobile Internet TV System," in *Proc. of the 2011 ACM SIGCOMM Conference on Internet Measurement Conference*, Berlin, Germany: ACM, pp. 209–224, 2011. [Article \(CrossRef Link\)](#)
- [28] Y. Jin, Y. Wen, H. Hu, and M.-J. Montpetit, "Reducing Operational Costs in Cloud Social TV: An Opportunity for Cloud Cloning," *IEEE Transactions on Multimedia*, vol. 16, pp. 1739–1751, Oct. 2014. [Article \(CrossRef Link\)](#)
- [29] Y. Wen, X. Zhu, J.J.P.C. Rodrigues, and C.W. Chen, "Cloud Mobile Media: Reflections and Outlook," *IEEE Transactions on Multimedia*, vol. 16, pp. 885–902, Jun. 2014. [Article \(CrossRef Link\)](#)
- [30] Amazon, "Amazon CloudFront Pricing," URL: <https://aws.amazon.com/cloudfront/pricing/>.
- [31] Alibaba Group, "Ali Cloud CDN Pricing (in Chinese)," URL: http://help.aliyun.com/knowledge_detail.htm?knowledgeId=5975217.
- [32] E. Kaasinen, M. Kulju, T. Kivinen, and V. Oksman, "User Acceptance of Mobile TV Services," in *Proc. of the 11th International Conference on Human-Computer Interaction with Mobile Devices and Services*, Bonn, Germany: ACM, pp. 34:1–34:10, 2009. [Article \(CrossRef Link\)](#)
- [33] P. Danaher and T. Dagger, "Using a nested logit model to forecast television ratings," *International Journal of Forecasting*, vol. 28, pp. 607–622, 2012. [Article \(CrossRef Link\)](#)
- [34] P.J. Danaher, T.S. Dagger, and M.S. Smith, "Forecasting television ratings," *International Journal of Forecasting*, vol. 27, pp. 1215–1240, 2011. [Article \(CrossRef Link\)](#)
- [35] Y.-H. Cheng, C.-M. Wu, T. Ku, and G.-D. Chen, "A Predicting Model of TV Audience Rating Based on the Facebook," *Social Computing (SocialCom), 2013 International Conference on*, pp. 1034–1037, 2013. [Article \(CrossRef Link\)](#)
- [36] Y.-Y. Huang, Y.-A. Yen, T.-W. Ku, S.-D. Lin, W.-T. Hsieh, and T. Ku, "A Weight-Sharing Gaussian Process Model Using Web-Based Information for Audience Rating Prediction," *Technologies and Applications of Artificial Intelligence*, S.-M. Cheng and M.-Y. Day, eds., Springer International Publishing, pp. 198–208, 2014. [Article \(CrossRef Link\)](#)



Yuheng Li received his B.E. degree in automatic control in 2007 from Tsinghua University, Beijing, China. He is currently a Ph.D candidate of the Center for Intelligent and Networked Systems (CFINS), Department of Automation, Tsinghua University. His research interests include measurement and modeling of online multimedia systems.



Qianchuan Zhao received the B.E. degree in automatic control, the B.S. degree in applied mathematics, and the M.S. and Ph.D. degrees in control theory and its applications from Tsinghua University, Beijing, China, in 1992 and 1996, respectively. He is currently a Professor and Associate Director of the Center for Intelligent and Networked Systems (CFINS), Department of Automation, Tsinghua University. He was a Visiting Scholar at Carnegie Mellon University, Pittsburgh, PA, and Harvard University, Cambridge, MA, in 2000 and 2002, respectively. He was a Visiting Professor at Cornell University, Ithaca, NY, in 2006. He has published more than 80 research papers in peer-reviewed journals and conferences. His current research focuses on modeling, control and optimization of complex networked systems with applications in smart buildings, smart grid and manufacturing automation. Dr. Zhao received the China National Nature Science Award for the project "Optimization Theory and Optimization for Discrete Event Dynamic System" in 2009. He is an Associate Editor for *Journal of Optimization Theory and Applications*, was an Associate Editor for *IEEE TRANSACTIONS ON AUTOMATION SCIENCE AND ENGINEERING* and is an Associate Editor for *IEEE TRANSACTIONS ON CONTROL OF NETWORK SYSTEMS*. He serves as a Chair of the Technical Committee on Smart Buildings of IEEE RAS.