

# A Server’s Perspective of Internet Streaming Delivery to Mobile Devices

Yao Liu<sup>1</sup> Fei Li<sup>1</sup> Lei Guo<sup>2</sup> Bo Shen<sup>3</sup> Songqing Chen<sup>1</sup>

<sup>1</sup>Dept. of Computer Science

<sup>2</sup>Dept. of CSE

<sup>3</sup>Vuclip

George Mason University

Ohio State University

Milpitas, CA, USA

{yliud, lifei, sqchen}@cs.gmu.edu lguo@cse.ohio-state.edu bshen@vuclip.com

**Abstract**—Receiving Internet streaming services on various mobile devices is getting more and more popular. To understand and better support Internet streaming delivery to mobile devices, a number of studies have been conducted. However, existing studies have mainly focused on the client side resource consumption and streaming quality. So far, little is known about the server side, which is the key for providing successful mobile streaming services.

In this work, we set to investigate the Internet mobile streaming service at the server side. For this purpose, we have collected a one-month server log (with 212 TB delivered video traffic) from a top Internet mobile streaming service provider serving worldwide mobile users. Through trace analysis, we find that (1) a major challenge for providing Internet mobile streaming services is rooted from the mobile device hardware and software heterogeneity. In this workload, we find over 2800 different hardware models with about 100 different screen resolutions running 14 different mobile OS and 3 audio codecs and 4 video codecs. (2) To deal with the device heterogeneity, transcoding is used to customize the video to the appropriate versions at runtime for different devices. A video clip could be transcoded into more than 40 different versions in order to serve requests from different devices. (3) Compared to videos in traditional Internet streaming, mobile streaming videos are typically of much smaller size (a median of 1.68 MBytes) and shorter duration (a median of 2.7 minutes). Furthermore, the daily mobile user accesses are more skewed following a Zipf-like distribution but users’ interests also quickly shift. Considering the huge demand of CPU cycles for online transcoding, we further examine server-side caching in order to reduce CPU cycle demand. We show that a policy considering different versions of a video altogether outperforms other intuitive ones when the cache size is limited.

## I. INTRODUCTION

Recently, mobile devices are getting increasing popularity. For example, according to International Data Corporation, the total number of smartphones sold worldwide in 2010 is 302.6 million [1], which is a 74.4% increase from the previous year (173.5 million). By September 2010, over 58.7 million people in the US owned smartphones [2].

Besides general web surfing on the Internet, these days more and more accesses from mobile devices are directed to all kinds of Internet streaming services. For example, YouTube [3] is the among the earliest to provide streaming services to mobile devices such as iPhone. Today both iOS and Android have native support for YouTube. Other popular streaming service providers, including Netflix [4] and Hulu [5], also provide streaming services to subscribed mobile users via APPs built in various mobile operating systems. Placeshifting

services like Orb [6] and AirVideo [7] allow mobile users to access media content stored on their home computers. Qik [8] allows users to upload from mobile devices and then broadcast the video content to their friends. Different from the above services, Vuclip [9] lets users search and play all kinds of Internet videos on their mobile devices regardless of their mobile device types.

To understand the key challenges of Internet mobile streaming and the difference from traditional Internet streaming, a number of studies have been performed. As today the majority of Internet mobile streaming services are delivered in a client-server architecture, many studies have focused on the resource consumption and streaming quality received on the mobile device. For example, Xiao et al. [10] studied energy consumption when watching YouTube on mobile devices. Huang et al. [11] investigated fetching policies of different mobile video players, and Finamore et al. [12] examined the potential causes for inferior streaming quality of mobile YouTube accesses.

However, these studies mainly focus on the client side by examining specific devices [10], [11] or via local experiments [12]. As the key to the current Internet mobile streaming services, the server side plays a critical role in the entire streaming delivery process. Unfortunately, so far, little is known about the server side, possibly due to the limited availability of data from the server side.

To provide in-depth understanding of the current Internet mobile streaming services, in this study, we set to investigate the server side in streaming delivery to mobile devices. For this purpose, we have collected a 30-day server log from a top Internet mobile streaming service provider. In 30 days, about 105 million video sessions were served with about 212 TB video traffic delivered. Through our analysis, we have a number of findings. While the details are presented later in the paper, some highlights are as follows:

- A unique challenge for Internet streaming delivery to mobile devices is rooted from the fact that mobile devices are very heterogeneous. In this workload, we find over 2800 different hardware models with 92 different screen resolutions running 14 different mobile OSes, using 3 audio codecs and 4 video codecs. This greatly challenges the traditional Internet streaming delivery infrastructure where the bottleneck often lies in the limited bandwidth.
- To deal with the device heterogeneity, runtime transcoding is used to customize a video to the appropriate

versions on the fly for different devices. A video clip could be transcoded into more than 40 different versions in order to serve requests from different devices.

- Compared to videos in traditional Internet streaming, mobile streaming video clips are typically of much smaller size (with a median of 1.68 MBytes) and the video duration is shorter as well (with a median of 2.7 minutes). Furthermore, the daily mobile user accesses are more skewed following a Zipf-like distribution but users’ interests also shift quickly, resulting in a stretched-exponential distribution in the long term.

To reduce the huge CPU cycles demanded for transcoding on the fly, we further explore caching at the server side by trading off storage for CPU cycles. Our study shows that a policy that considers different versions of a video altogether outperforms other intuitive ones (e.g., a file based one) when the cache size is limited. As far as we know, we are the first to provide a server-side analysis on a Vuclip-like Internet mobile streaming service. Our findings provide new insights and lay some foundations to improve the current Internet mobile streaming delivery.

The rest of the paper is organized as follows. We describe some background and the workload overview in section II and study the device hardware and software heterogeneity in section III. We examine various mobile video properties in section IV and further explore the trade-off between the storage and the CPU at the server side in section V. Some related work is described in section VI and we make concluding remarks in section VII.

## II. BACKGROUND AND WORKLOAD OVERVIEW

To investigate how current Internet streaming services are delivered to mobile devices, we have collected a 30-day server log from one of the largest Internet mobile streaming service providers, Vuclip [9]. Vuclip provides mobile users with the search-and-delivery services. It allows users to search for and watch any videos on any video-enabled mobile phones and devices.

Different from many existing services that only provide streaming services to specific mobile devices, Vuclip can serve any type of mobile devices that are capable of streaming playback. Vuclip allows any mobile user to search for interested video available on the Internet, and transcodes them on-demand and on-the-fly based on the type of the mobile device. To serve different types of mobile devices, Vuclip employs on-demand transcoding at the server side. Transcoding is a process to convert the requested video clip to the appropriate codecs, format, and size at runtime upon a request so that the video can be properly rendered and played on the requesting mobile device. Vuclip transcodes a video into different versions by choosing the best audio/video codecs, frame size, frame rate, and quality level combination for the mobile device. According to our analysis, each video was accessed in more than 2 versions on average (as shown in Table I), and the most popular video was accessed in 41 different transcoded versions (as shown later in Figure 7).

To deliver video content, Vuclip uses the traditional client/server (C/S) architecture. The video file is delivered via pseudo streaming over HTTP. That is, when the requested content is available on the server, the client would issue an HTTP GET request to download the content. A video may be downloaded via several HTTP GET requests with different partial ranges specified (i.e., range requests). To differentiate video requests from HTTP requests, we define a *request* as a single HTTP transfer between the client and the server, and a *session* as the set of requests that are involved in downloading an entire video clip.

TABLE I  
SUMMARY OF WORKLOAD

Workload Length	30 Days
# of Sessions	105,389,370
# of Requests	192,255,173
# of Requests from Mobile Devices	181,556,344
# of Unique Videos Accessed	4,052,740
AVG. # of Formats Per Video	2.31
Total Traffic Volume	212 TB

The data we collected is from Nov. 1st to Nov. 30th, 2010. In this one-month server log, there are about 105 million sessions watching more than 4 million different videos. There are a total of about 192 million HTTP requests. The total traffic delivered from the server in these 30 days is about 212 TB. Table I gives a summary of this workload. Note that among all these requests, some are from desktop/laptop computers instead of mobile devices. In order to focus on the requests from mobile users, we differentiate them in the server log through the User-Agent strings specified in each HTTP request. By analyzing the User-Agent, we find there is a total of 150,072 unique User-Agent strings. Among them, 84,281 (56%) represent mobile devices. However, examining the received requests, we find most of them come from User-Agent strings representing mobile users: more than 94% (181 million out of 192 million) requests are from mobile devices.

With the exclusion of desktop/laptop traffic, Figure 1 gives an overview of the server side traffic in 30 days. Note that the left *y*-axis represents the total number of requests per day, while the right *y*-axis represents the total traffic volume per day. During this one-month period, despite a small decrease in the middle, the number of the requests and the delivered traffic amount kept increasing, indicating the popularity of Vuclip.

Figure 2 shows the hourly mobile streaming access patterns in a day. The figure indicates that hourly accesses peak around 17:00 GMT. Furthermore, the total number of requests and the traffic volume served during peak hours almost double these in non-peak hours. Figure 3 further depicts the hourly pattern from Nov. 8th to Nov. 15th (a week). The figure shows clear peak and off-peak hourly patterns for each day. The figure shows some drop after Nov. 12th. It is likely due to the fact that Nov. 13th was a Saturday and Nov. 14th was a Sunday. We can observe the increase of accesses again on Monday.

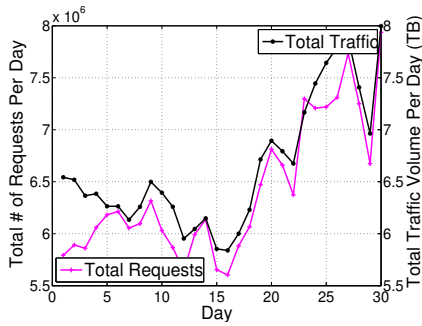


Fig. 1. Daily Accesses in Nov. 2010

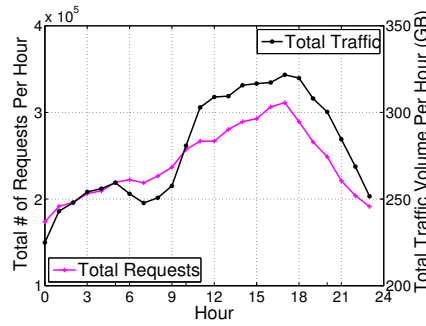


Fig. 2. Hourly Accesses On Nov. 1st 2010

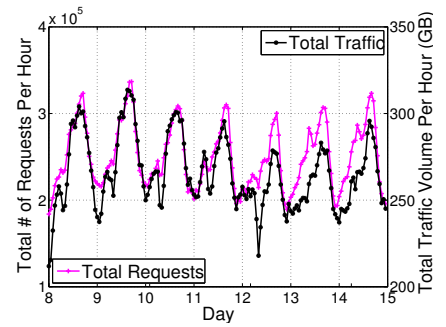


Fig. 3. Weekly Accesses from Nov. 8 to Nov. 15 2010

### III. CHARACTERIZATION OF MOBILE DEVICE HETEROGENEITY

#### A. Mobile System Heterogeneity

To provide Internet streaming services to all kinds of mobile devices like Vuclip, a unique challenge is the heterogeneity among mobile devices. Different from the pre-coding approach that was taken by many other service providers to serve specific types of mobile devices, transcoding has to be used to customize the video into a proper format for the requesting mobile device. Although transcoding is very flexible and desirable to serve heterogeneous mobile devices, transcoding demands huge CPU cycles on the fly.

To get a realistic picture of the mobile device heterogeneity, we retrieve detailed device information from WURFL [13] based on the User-Agents information we have extracted from the server log. Among the 84,281 User-Agents that represent mobile devices, we are able to get the brand and model information from more than 74,708 (88.64%) distinct User-Agent strings. The rest only have browser information.

TABLE II  
SYSTEM HETEROGENEITY OF MOBILE DEVICES

Models	2864
Resolution	92
Mobile OSes	14

As shown in Table II, accesses to Vuclip in these 30 days came from 2864 different device models. These devices have different screen sizes that can support video playback with different resolution rates. Delving into this, we find that these devices have 92 different resolutions (width and height combinations), ranging from  $84 \times 48$  to  $1600 \times 1200$ . Figure 4 shows the most popular resolutions, including  $320 \times 240$ ,  $480 \times 360$ , and  $480 \times 320$ . They also run on 14 different mobile operating systems.

#### B. Audio/Video Codec Heterogeneity

To play video on a mobile device, both audio and video codecs are required. On different devices, the supported codecs may be different as well. Such heterogeneity would further increase the load for the server if the server conducts transcoding

for the mobile device. Note that if such transcoding is done at the client side, it would lead to excessive battery power consumption.

To examine the codec heterogeneity, we further look into the supported audio/video codecs on these 2864 hardware models. We find that typically there are 3 audio codecs being used, namely AAC, AMR, and WMA, and there are 4 video codecs being used, namely H.263, H.264, MPEG-4, WMV. Figures 5 and 6 show the popularity of these codecs. As shown in these figures, AMR is the most popular audio codec, as more than 59% devices support it, and H.263 and MPEG-4 are the most popular video codecs.

TABLE III  
VIDEO CODECS

Type	Video	Audio	#Videos	Sessions
ASF	WMV	WMA	293,025	2,031,161
3GP	H.264	AAC	692,004	12,636,639
3GP	H.263	AMR	2,805,494	46,790,565
3GP	MPEG-4	AMR	138,022	1,213,319
3GP	MPEG-4	AAC	1,762,132	36,552,760
3GP in Total			3,746,548	97,193,283

With 3 audio codecs and 4 video codecs, we expect a total of 12 combinations of different audio/video codecs. In practice, however, not all these combination of the audio and video codec are used. In the workload, we only find 5 combinations. Table III shows the 5 video+audio encoding schemes used. While more than 4 million (4,052,740) unique videos were accessed, we are able to extract about 3.7 million (3,789,229) that are accessed as videos. The rest were only accessed as audio. Table III shows that H.263+AMR and MPEG4+AAC are the most popular encoding schemes, accounting for 84% of total viewing sessions. This is not surprising as H.263 and MPEG4 are the most widely supported video codecs on the 2864 models of mobile devices.

In addition to different codecs, video files are also encoded into two different formats, i.e. two types of containers, 3GP and ASF. 3GP is the 3GPP file format, which is a multimedia container format defined by the Third Generation Partnership Project (3GPP) for 3G UMTS multimedia services. 3GP is often used on 3G mobile phones. On the other hand, ASF

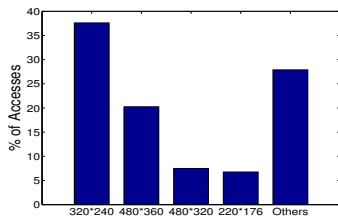


Fig. 4. Most Popular Resolutions

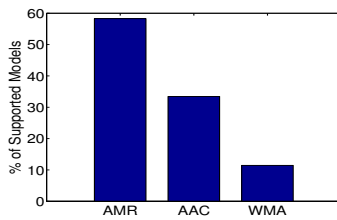


Fig. 5. Audio Codecs

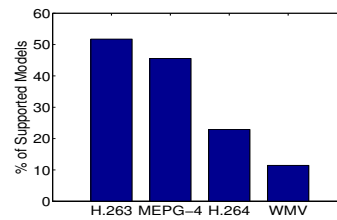


Fig. 6. Video Codecs

(Advanced Systems/Streaming Format) belongs to Microsoft Media framework and it is a proprietary digital audio/digital video container format. Apparently, 3GP is much more widely used in practice than ASF for mobile videos.

TABLE IV  
VIDEO RESOLUTION AND ENCODING RATE

Quality	Frame Width	Encoding Rate (Kbps)
Low	176	51 - 55
Low	320, 360	71 - 187
High	176	81 - 147
High	320, 360	172 - 335
WiFi	320, 360	358 - 423

Besides the above hardware and software heterogeneity, mobile devices may have different network speed, due to various reasons, such as accessing through cellular network or WiFi. To support different mobile Internet access speed, Vuclip also transcodes video clips into 3 different quality levels: Low Quality, High Quality and WiFi Quality. Table IV shows the corresponding range of object encoding rate for different quality levels. Consider the variety of resolutions, videos are also customized into 3 different frame widths: 176, 320, and 360. As we can observe from the table that a larger resolution (width) video does not necessarily come with a high encoding rate. On the other hand, a video with a high encoding rate typically comes with a larger resolution.

TABLE V  
VIDEO QUALITY

Quality	# of Videos	# of Sessions
Low	1,694,108	26,365,900
High	3,323,211	72,392,729
WiFi	91,761	465,815

Table V further shows the number of videos that are of Low, High, and WiFi Quality as well as the number of their requested sessions. The videos in High Quality are mostly requested: 73% viewing sessions are for videos encoded with High Quality. Consider that Vuclip transcodes the video content on-demand, it is not surprising that 87% of video contents have at least one version encoded with High Quality. WiFi Quality, however, is the least requested quality level. This is likely due to the relatively slow mobile accessing speed and tiered data plan billing model today.

Since Vuclip transcodes the original video to accommodate mobile devices with different codecs, frame width, and quality level, for the ease of presentation, we use *versions* to refer to different transcoded video files for each video in the rest of the paper. On the other hand, we use *videos* to refer to a set of video clips that correspond to the same content. Figure 7 shows the CDF of number of versions each video has. As shown in the figure, in this workload, about 59% videos have only one version, and about 3% videos are accessed in 10 or more versions. The largest version number is 41.

#### IV. CHARACTERIZATION OF MOBILE STREAMING VIDEOS

The previous section has shown that mobile device heterogeneity is a great challenge to the service provider. With such a level of heterogeneity, what kind of video clips are being served is of our great interest. In this section, we further analyze the mobile video clips that we have collected from the server log in order to reveal the commons with and differences from the traditional Internet streaming content.

##### A. Video Playback Duration and File Size

Figure 8 depicts the distribution of video playback duration in seconds. In this figure, videos are sorted in decreasing order of the playback duration, and the  $y$ -axis is in log scale. As shown in the figure, video clips accessed by mobile users are mostly short in terms of playback duration: more than 97% videos are less than 10 minutes long, and the median playback duration is 162 seconds (less than 3 minutes). Compared to the longer duration of traditional Internet streaming video clips, such a shorter duration makes it more feasible for mobile devices because video streaming consumes a lot of limited resources on mobile devices, including the network for data receiving, the CPU for decoding, and the display for rendering. Such resource consumption can drain the limited battery power supply at a very high rate.

Correspondingly, Figure 9 shows the file size (bytes) distribution. Again, we sort the video files (versions) based on their sizes in decreasing order. As shown in Figure 9, the video file distribution is similar to that of the duration as shown in Figure 8. Note that, here in this figure, each video may have been accessed in several versions in different formats and file sizes. As we can see, most video files accessed by mobile devices are smaller than 8 MBytes, with a mean file size of 2.78 MBytes and the median file size 1.68 MBytes. This shows that videos accessed by mobile devices are mostly small in terms of bytes. This can reduce the total network

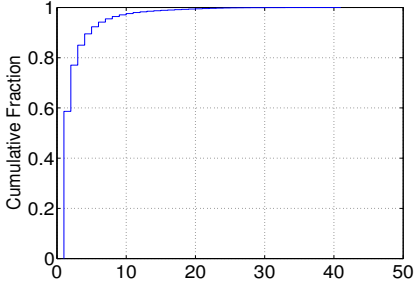


Fig. 7. # of Versions Per Video (CDF)

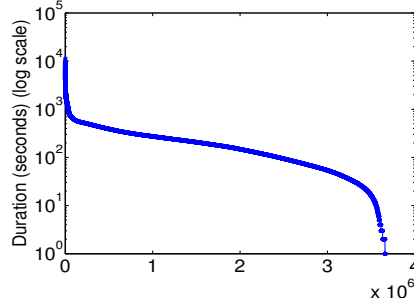


Fig. 8. Video Playback Duration Distribution

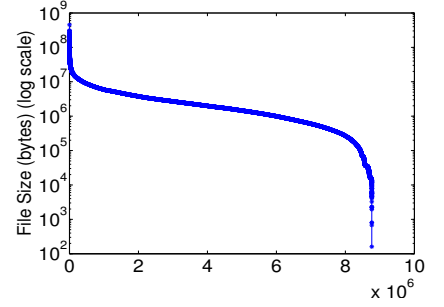


Fig. 9. Video File Size Distribution

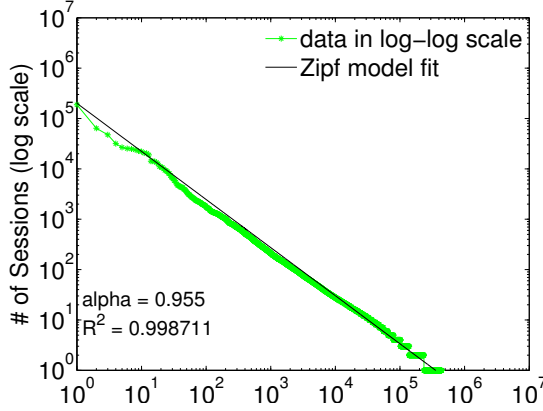


Fig. 10. Daily Video Popularity Distribution

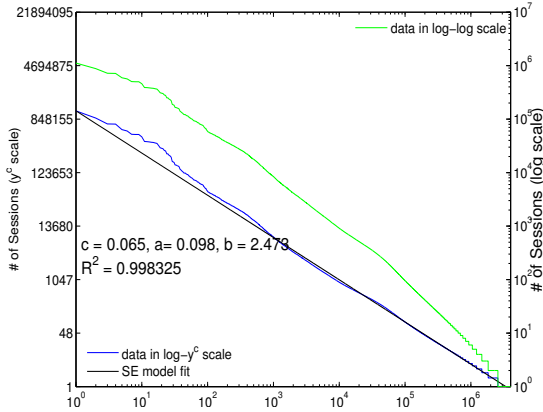


Fig. 11. Monthly Video Popularity Distribution

## B. Popularity of Mobile Videos

Figure 10 shows the popularity pattern of videos accessed site-wide on Nov. 1st. In this figure, the  $x$ -axis represents videos ranked by the number of requested sessions in decreasing order, plotted in log-scale, while the  $y$ -axis represents the number of viewing sessions of this video, also plotted in log scale. This figure shows that, in log-log scale, the popularity distribution of videos accessed can be well fitted with a Zipf-like distribution

$$y_i \propto \frac{1}{i^\alpha},$$

where  $i$  is the popularity rank of the video,  $y_i$  is the number of requested sessions for the video, and  $\alpha$  is the skewness parameter. Moreover, we find  $\alpha = 0.955$  fits our data very well with the goodness of fit value  $R^2$  very close to 1, indicating the popularity distribution is not only Zipf-like, but also very close to the Zipf's law where  $\alpha = 1$ . Similar patterns have been found for other days in the workload.

The Zipf-like distribution is known to be efficient in modeling web traffic, and is the premise for efficient web caching. Specifically,  $\alpha$  is an indicator of request concentration, and proxy caching can be more efficient with a larger  $\alpha$  value. For example, it was reported in [18] that  $\alpha$  varies between 0.64 and 0.83 for web traffic, while it tends to be smaller for media traffic (for example, work [17] reports 0.56 for YouTube traffic). Different from previous measurement studies where data were collected at edge networks, the mobile video accesses are highly concentrated at the server side. Such discrepancy is reasonable as collecting traffic at edge networks can only reflect the local users' accesses, while the server logs can provide a complete view of the video popularity. Furthermore, this also means caching at the server side is more effective than caching at the edge/client side, if caching at the server side is needed. Note for content delivery, caching at the server side is typically not for reducing network traffic as caching at the client side.

While Figure 10 shows short-term (one day) popularity distribution, Figure 11 shows the corresponding distribution in long-term, spanning over the entire 30 days of our workload. In this figure, the left  $y$ -axis is in powered scale while the right  $y$ -axis is in log scale. The  $x$ -axis is in log scale as well. As shown in the figure, the monthly video popularity deviates from a straight line in log-log scale, meaning not a Zipf-like

transmission for downloading the video file. Note the network interface card could consume 30% to 40% of the total battery power consumed during a streaming session to a mobile device [14], [15].

Moreover, compared to the size of the traditional Internet video files [16], [17], the size distribution we find in this server log is much smaller. This provides a great opportunity for reducing the transcoding cost as we discuss later in section V.

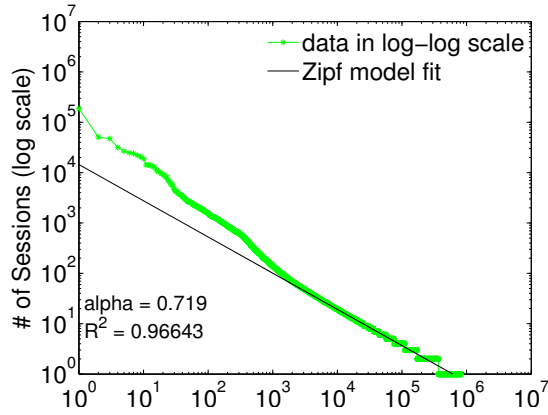


Fig. 12. Daily Version Popularity Distribution

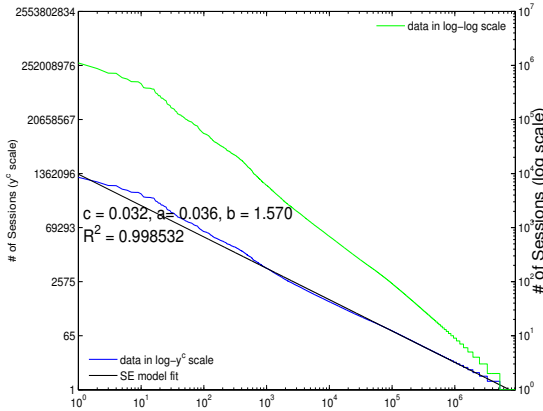


Fig. 13. Monthly Version Popularity Distribution

distribution. Instead, it can roughly be fit with a stretched exponential (SE) distribution as shown by the left  $y$ -axis in powered (by a constant  $c$ ) scale [19]. With SE distribution, the rank distribution function can be expressed as

$$y_i^c = -a \log i + b.$$

An SE distribution is fit by several parameters as shown in the Figure. For example, parameter  $c$  is also called the *stretch factor*, which characterizes the median file size of workload [19]. It was reported that for media workloads with a median file size  $< 5$  MBytes, the stretch factor is  $\leq 0.2$ . Our analysis confirms this with a  $c$  of 0.065 and a median file size of 1.68 MBytes. Parameter  $a$  in an SE distribution increases with the duration of workload as well as the ratio of media request rate to new content birth rate, and it causes the distribution to deviate from a straight line in log-log scale.

The stretched exponential distribution has been used to characterize many natural and economic phenomena, as well as the access patterns of Internet media traffic [19]. It was shown that under an SE distribution, media caching is much less efficient than under a Zipf distribution. This poses a new challenge if long-term caching is needed on the server side.

TABLE VI  
SUMMARY OF VIDEO ACCESSES

Average Daily New Videos	92 K
Average Daily Accesses Videos	502 K
Percentage of Daily New Videos	18%
Average Daily Accessed Videos Difference	292 K
Percentage of Average Daily Difference	58%
Average Daily Requested Sessions	3512 K
Total Accessed Videos in 30 Days	4052 K
Percentage of New Videos in 30 Days	68%

### C. Popularity of Different Video Versions

We have shown in Figure 10 that the daily video popularity follows a Zipf-like distribution. However, as discussed before, each video may be accessed by very diverse mobile devices, resulting in multiple transcoded versions. We thus further examine the popularity distribution of all versions accessed on Nov 1st where each version is counted as a distinct object. Figure 12 shows that when different versions are considered as different objects, the popularity cannot be well-fitted with the Zipf distribution. On the one hand, due to the increased number of video versions (2.31 versions per video on average), the skewness factor  $\alpha$  decreases from around 0.95 to 0.7. On the other hand, as shown in the figure, accesses to the Top-1000 versions are much larger than what Zipf predicts, indicating significant deviation from Zipf-like distribution.

Figure 13 further depicts the monthly version popularity. This figure shows that although the daily version popularity cannot fit well with either Zipf or SE distributions, the monthly version popularity pattern can be well fitted with an SE distribution. The goodness value of the fit is 0.998532, very close to 1.

The version-based popularity patterns we have examined above indicate that different from video based popularity, accesses to different video versions are more distributive. This means poorer caching performance if caching is needed for these versions.

### D. Popularity Evolution

We have shown that mobile users' daily accesses for mobile videos are highly concentrated, but monthly accesses patterns are flatter. In this subsection, we further examine how video popularity changes over time, aiming to shed light on such popularity changes. We first study the commons between accesses in consecutive days, and then, more specifically, we consider the temporal locality characteristics of these accesses.

Table VI summarizes the daily accesses and the corresponding videos that were requested in the system. According to our analysis, 292K out of 502K (58%) unique video clips accessed daily are not accessed in the previous day on average. Among these 292K video clips, about 92K are new video clips. This indicates that about 18% video clips accessed every day are new. The rest 200K (40%) are unpopular ones that were in the system, but were not frequently accessed. In total, new video clips account for about 68% of total unique video clips

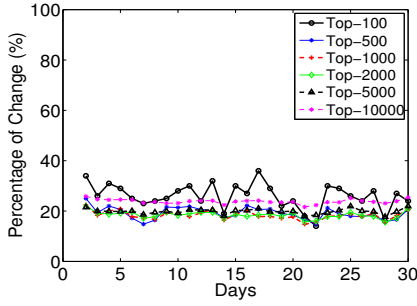


Fig. 14. Percentage of Change in Top Videos

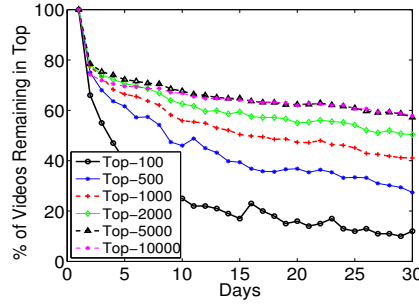


Fig. 15. % of Videos Remaining in Top

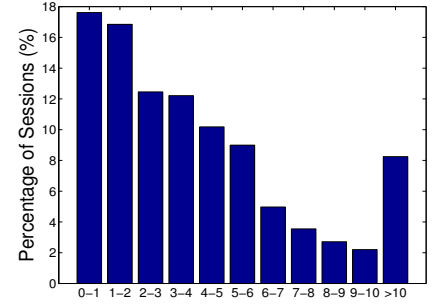


Fig. 16. % of Sessions for Videos of Different Length (minute)

accessed during 30 days. Since new video clips are generated at a high rate of 18%, this confirms the implication of SE-distribution that a monthly static caching scheme may not be so efficient as a more frequently updated one.

Unlike traditional video on-demand streaming systems, Vuclip has a larger repository as well as a faster new content generation rate. We next examine if temporal locality is helpful in predicting what will be popular in the future in such a highly dynamic system.

First, we examine the temporal locality at the server side. We have shown in Table VI that about 18% new video clips are added into the video repository daily. We further examine if top accessed videos change at a similar rate. Figure 14 shows the percentages of change in Top-100, Top-500, Top-1000, Top-2000, Top-5000 and Top-10000 accessed videos every day in 30 days. The percentages of change in Top-500, Top-1000, Top-2000 and Top-5000 requested video clips are all about 20%, which is similar to the birth rate of new video clips everyday. The Top-100 videos change at a higher rate, fluctuating between 20% and 35%. Previous studies on a video-on-demand (VoD) system report the daily rate of change in Top-100 videos is less than 15% [16]. Compared to the traditional VoD system, the result here indicates that mobile users tend to shift their interest faster. The change rate of Top-10000 videos is relatively stable at about 25%, higher than that of Top-500, Top-1000, Top-2000, and Top-5000 videos. This is likely due to the fact that the Top-10000 list also includes videos that are less popular, e.g., new videos that were only active for a short period of time (one-day), and the 40% old videos that were not accessed in the previous day.

Figure 15 keeps track of the top videos on the first day (Nov. 1st) of our workload, and examines how many of them would remain on the top list for the next 29 days til Nov. 30th. We notice that during the first 5 days, a great percentage of top video clips are purged out of the top list, indicating quick user interest shift. However, such a change becomes relatively stable beyond the first 5 days, indicating videos that were popular during the past  $N$  days are likely to be popular in the future. Nevertheless, despite such a quick change of user interests in a short period, video popularity is still relatively stable in the long term. For example, although only about 10% video clips remain in Top-100 after 29 days, this percentage

increases to 30% for Top-500, 40% for Top-1000, 50% for Top-2000 and about 60% for Top-5000 and Top-10000, which indicates that with a large cache size, the cached video clips can be flushed less frequently.

### E. Correlation Between Popularity and Video Length

We have shown in section IV-A that the mobile video files are often short and the video popularity in a short-term has a Zipf-like distribution. Thus one may wonder whether shorter videos get more accesses. To examine this, we group video files into 1 minute interval based on their lengths. Figure 16 shows the total number of requested sessions decreases as the video length increases, and videos shorter than 5 minutes account for about 70% total requests. We further calculate the correlation between video popularity and video length, and examine if shorter videos tend to be more popular. Our results show that the correlation coefficient is 0.006, which indicates the correlation is weak. We have also conducted similar tests based on versions, and the results are similar. This indicates the large percentage of requests for short videos are due to the large population of short video contents instead of user preferences.

## V. TRADE-OFF BETWEEN CPU AND STORAGE

As aforementioned, in order to conduct on-demand transcoding to serve all kinds of heterogeneous mobile devices, Vuclip has to employ a lot of computer clusters. This challenges service providers both technically and economically with the growing popularity of Internet mobile streaming services and it is thus very desirable to reduce the huge demand of CPU cycles for such transcoding.

In the previous section, we have shown that the mobile users' accesses are more concentrated (skewed) than those in the traditional Internet streaming services. Thus, caching on the server side, sometimes called reverse caching, could be explored to temporarily cache some transcoded objects so that on-the-fly transcoding would not be necessary if the same type of mobile devices access the same video. Such full-object caching is possible for mobile videos because mobile video objects are typically smaller with a median of 1.68 MBytes as we showed before. Note that different from the traditional caching objectives such as web proxy caching for reducing

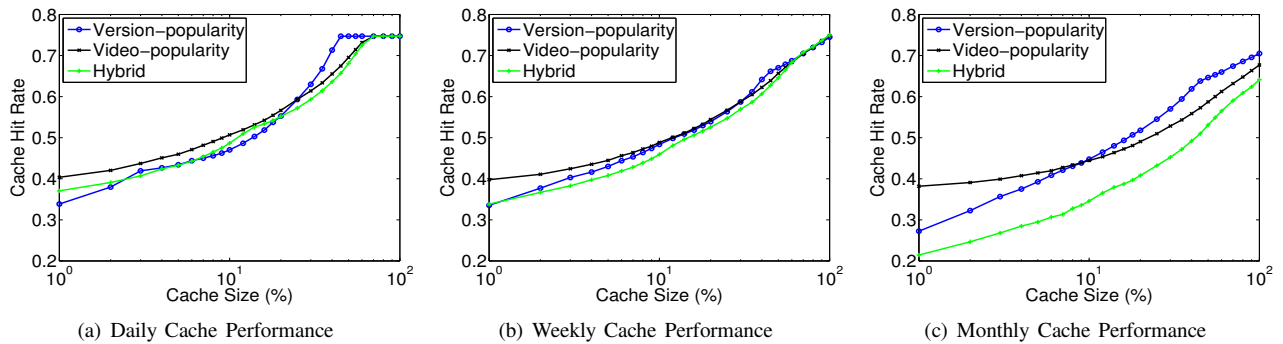


Fig. 17. Cache Performance over Different Time

the network traffic, caching at the server side here is to reduce the CPU cycles demanded for transcoding. That is, a trade-off between the storage and the CPU.

On the other hand, we have also shown that that mobile users' interests shift quickly. This provides some hints for cache replacement. Regardless whether the cache is implemented via disk and/or memory, the cache replacement policy is the key to the cache performance. Typical cache replacement strategies (such as popularity-based policies) may work, however, the complexity added by Vuclip-like services comes from the fact that a video often has multiple transcoded versions. Intuitively, these different versions could be considered as separate objects in the cache. However, if we consider video popularity, different versions of a same video are internally related. Therefore, we next explore different replacement strategies via simulations for Vuclip-like systems.

A simple strategy is to ignore the internal relationship of different versions of a video, and consider each version as a distinct object. Under this assumption, the existing web proxy cache replacement policies can be adopted. Since we are dealing with video objects, we thus first consider a *version-popularity* based replacement policy, in which a utility function is defined as the ratio of the version access number to the storage size occupied by that version. The version with the least utility is the victim to be purged from the cache.

On the contrary, if we consider that different versions of a video are related because along the diminishing popularity of a video, all of its versions may get fewer and fewer accesses, we can also consider a policy in which all different versions of a video are bundled together as one object in the cache. Taking the popularity of this object as the sum of popularity of all its versions, we can design a *video-popularity* based replacement policy, in which the entire object (with all versions) is replaced once it is identified as the victim based on the utility function defined as the ratio of the total access number to this video to the total size of occupied storage.

With the above two strategies, naturally, one may wonder if a *hybrid* policy could perform better. That is, neither considering the version independently as the version-popularity based policy does, nor considering all versions of a video as one object as the video-popularity based policy does. In a hybrid strategy, a utility function can be defined for each video as in

the *video-popularity* based replacement policy, but the victim is the least popular version of the least popular video.

To study the effectiveness of these different strategies, we conduct trace-driven simulations using the collected workload to compare their performance. With the accesses of Nov. 1st, Figure 17(a) shows that when the cache size is smaller than 27% of the total size of accessed objects of that day, a video-popularity based replacement policy can achieve the highest cache hit rate, and save roughly about 55% CPU cycles, while a version-popularity based strategy performs the worst. This is likely due to the highly concentrated video access pattern as shown in Figure 10 compared to less concentrated version access pattern shown in Figure 12. The hybrid strategy performs consistently worse than the video-popularity based one, but still a little better than the version-popularity based policy. These results are consistent with our analysis on the video and version popularity. When the cache size increases, the version-popularity based policy has more flexibility in choosing the best version to cache, and thus achieves the best cache hit ratio among the three.

Figure 17(b) further shows the results when one-week trace from Nov. 1st to Nov. 8th was simulated. The cache size percentage used in this simulation is based on total accessed objects of Nov. 1st. We find that the cache performance over one week can reach as high as that of a day. However, the cache performance over a month as shown in Figure 17(c) is worse. This is because with an SE distribution for monthly video accesses, caching is much less efficient than with a Zipf distribution for daily accesses.

## VI. RELATED WORK

The Internet has witnessed the sharp increase of Internet video traffic in the recent years with all kinds of Internet streaming systems, such as VoD and Internet P2P-based streaming systems. Lots of research has been conducted to study these Internet streaming systems. For example, Yu et al. [16] examined server logs of a traditional VoD system with a total of over 6700 unique videos, and analyzed user access patterns, session length, and video popularity. Yin et al. presented the access logs of a live VoD system [20], which shows different user and content properties compared to [16]. Krishnappa et al. collected Hulu traffic at a campus edge



network, and examined the potential benefits of performing caching and prefetching at edge networks [21]. For P2P-based streaming systems, Wu et al. investigated P2P streaming topologies in UUSEE [22]. Huang et al. conducted a large scale measurement to study the PPLive-based on-demand streaming [23].

Along the increasing popularity of user-generated content (UGC), studies have also been conducted to characterize UGC videos. For example, Cha et al. studied user behaviors and video popularity of YouTube, and compared them with non-UGC content from Netflix [24]. Work [17] examined the traffic characteristics of YouTube at a campus edge network.

With the rapid increase of Internet-capable mobile devices in recent years, mobile Internet video services and accesses are surging. A few studies have been conducted to investigate the performance of mobile streaming applications, mostly focusing on the resource utilization for receiving streaming data on mobile devices. For example, Xiao et al. studied the power consumption of mobile YouTube [10]. Finamore et al. collected traffic from several edge locations and studied the potential reasons for the inferior streaming experience of mobile YouTube users [12]. Previously, we have also conducted measurements to study the resource utilization of different streaming approaches to mobile devices [15]. Furthermore, in order to save battery power, we had designed and implemented BlueStreaming, a system that can leverage low-power of Bluetooth to help P2P streaming to mobile devices [25].

However, to the best of our knowledge, no prior work has examined Internet mobile streaming services from the server side, which is the key to provide Internet mobile streaming services. In this study, we have investigated the commons and differences of mobile Internet streaming services with/from the traditional Internet streaming services. Our study reveals a critical challenge in Internet mobile streaming services is the hardware and software heterogeneity of mobile devices. Our analysis also shows different access patterns of mobile videos from traditional Internet streaming videos. Furthermore, we have also shown that caching at the server side with a proper replacement policy can significantly reduce the resource consumption for Vuclip-like Internet streaming systems in dealing with heterogeneity.

## VII. CONCLUSION

The wide adoption of mobile devices in practice has made pervasive Internet streaming possible. While a number of studies have been conducted to examine the streaming services from the client's perspective, in this work, we have studied the Internet mobile streaming services from the server side via one-month server log collected from one of the largest Internet mobile streaming service providers. Through detailed analysis, we have shown the great hardware and software heterogeneity of mobile devices, different characteristics of mobile videos, and different user access patterns from those in traditional Internet streaming services. As a great challenge that Vuclip-like system faces is the huge demand of CPU resources for online transcoding to deal with heterogeneity, we show that

caching at the server side with a proper replacement policy can effectively trade-off limited storage size for great savings on CPU cycles. These results provide some basic guidelines for building and optimizing future Internet mobile streaming systems.

## VIII. ACKNOWLEDGEMENT

We appreciate constructive comments from anonymous referees. The work is partially supported by NSF under grants CNS-0746649, CNS-1117300, CCF-0915681, CCF-1146578, and AFOSR under grant FA9550-09-1-0071.

## REFERENCES

- [1] "IDC - Press Release - prUS22689111," <http://www.idc.com/about/viewpressrelease.jsp?containerId=prUS22689111>.
- [2] "September 2010 U.S. Mobile Subscriber Market Share," [http://www.comscore.com/Press\\_Events/Press\\_Releases/2010/11/comScore\\_Reports\\_September\\_2010\\_U.S.\\_Mobile\\_Subscriber\\_Market\\_Share](http://www.comscore.com/Press_Events/Press_Releases/2010/11/comScore_Reports_September_2010_U.S._Mobile_Subscriber_Market_Share).
- [3] "Mobile YouTube," <http://m.youtube.com>.
- [4] "Netflix," <http://www.netflix.com>.
- [5] "Hulu," <http://www.hulu.com>.
- [6] "OrbLive," <http://itunes.apple.com/app/id290195003>.
- [7] "Air Video," <http://itunes.apple.com/app/id306550020>.
- [8] "Qik," <http://www.qik.com>.
- [9] "Vuclip," <http://m.vuclip.com/>.
- [10] Y. Xiao, R. Kalyanaraman, and A. Yla-Jaaski, "Energy Consumption of Mobile YouTube: Quantitative Measurement and Analysis," in *Proc. of NGMAST*, 2008.
- [11] J. Huang, Q. Xu, B. Tiwana, Z. M. Mao, M. Zhang, and P. Bahl, "Anatomizing Application Performance Differences on Smartphones," in *Proc. of MobiSys*, 2010.
- [12] A. Finamore, M. Mellia, M. Munafo, R. Torres, and S. G. Rao, "YouTube Everywhere: Impact of Device and Infrastructure Synergies on User Experience," in *Proc. of ACM IMC*, 2011.
- [13] "WURFL," <http://wurfl.sourceforge.net/>.
- [14] J. Adams and G. Muntean, "Adaptive-Buffer Power Save Mechanism for Mobile Multimedia Streaming," in *Proc. of IEEE ICC*, 2007.
- [15] Y. Liu, L. Guo, F. Li, and S. Chen, "An Empirical Evaluation of Battery Power Consumption for Streaming Data Transmission to Mobile Devices," in *Proc. of ACM Multimedia*, 2011.
- [16] H. Yu, D. Zheng, B. Y. Zhao, and W. Zheng, "Understanding User Behavior in Large-Scale Video-on-Demand Systems," in *Proc. of EuroSys*, 2006.
- [17] P. Gill, M. Arlitt, Z. Li, and A. Manhanti, "YouTube Traffic Characterization: A View From the Edge," in *Proc. of ACM IMC*, 2007.
- [18] L. Breslau, P. Cao, L. Fan, G. Phillips, and S. Shenker, "Web Caching and Zipf-like Distributions: Evidence and Implications," in *Proc. of IEEE INFOCOM*, 1999.
- [19] L. Guo, E. Tan, S. Chen, Z. Xiao, and X. Zhang, "The Stretched Exponential Distribution of Internet Media Access Patterns," in *Proc. of PODC*, 2008.
- [20] H. Yin, X. Liu, F. Qiu, N. Xia, C. Lin, H. Zhang, V. Sekar, and G. Min, "Inside the Birds Nest: Measurements of Large-Scale Live VoD from the 2008 Olympics," in *Proc. of ACM IMC*, 2009.
- [21] D. K. Krishnappa, S. Khemmarat, L. Gao, and M. Zink, "On the Feasibility of Prefetching and Caching for Online TV Services: A Measurement Study on Hulu," in *Proc. of PAM*, 2011.
- [22] C. Wu, B. Li, and S. Zhao, "Exploring large-scale peer-to-peer live streaming," in *IEEE Transactions on Parallel and Distributed Systems*, June 2008.
- [23] Y. Huang, T. Z.J Fu, D.-M. Chiu, J. C.S. Lui, and C. Huang, "Challenges, design and analysis of a large-scale p2p-vod system," in *Proc. of ACM SIGCOMM*, Seattle, WA, USA, August 2008.
- [24] M. Cha, H. Kwak, P. Rodriguez, Y.-Y. Ahn, and S. Moon, "I Tube, You Tube, Everybody Tubes: Analyzing The Worlds Largest User Generated Content Video System," in *Proc. of ACM IMC*, 2007.
- [25] Y. Liu, F. Li, L. Guo, Y. Guo, and S. Chen, "BlueStreaming: Towards Power-Efficient Internet P2P Streaming to Mobile Devices," in *Proc. of ACM Multimedia*, 2011.