

Understanding Internet Video Sharing Site Workload: A View from Data Center Design

Xiaozhu Kang
Columbia University
xk2001@columbia.edu

Haifeng Chen
NEC Laboratories America
haifeng@nec-labs.com

Hui Zhang
NEC Laboratories America
huizhang@nec-labs.com

Xiaoqiao Meng
NEC Laboratories America
xqmeng@nec-labs.com

Guofei Jiang
NEC Laboratories America
gfj@nec-labs.com

Kenji Yoshihira
NEC Laboratories America
kenji@nec-labs.com

ABSTRACT

In this paper we measured and analyzed the workload on Yahoo! Video, the 2nd largest U.S. video sharing site, to understand its nature and the impact on online video data center design. We discovered interesting statistical properties on both static and temporal dimensions of the workload including file duration and popularity distributions, arrival rate dynamics and predictability, and workload stationarity and burstiness. Complemented with queueing-theoretic techniques, we further extended our understanding on the measurement data with a virtual design on the workload and capacity management components of a data center assuming the same workload as measured, which reveals key results regarding the impact of Service Level Agreements (SLAs) and workload scheduling schemes on the design and operations of such large-scale video distribution systems.

Categories and Subject Descriptors

C.2.0 [Computer-Communication Networks]: General

General Terms

Measurement, Design, Performance

Keywords

Measurement, Data Center Design, Online Video, Workload Management, Capacity Planning, SLA, Queueing Model

1. INTRODUCTION

Internet Video sharing web sites such as YouTube [1] have attracted millions of users in a dazzling speed during the past few years. Massive workload accompanies those web sites along with their business success. In order to understand the nature of such unprecedented massive workload and the impact on online video data center design, we analyze Yahoo! Video, the 2nd largest U.S. video sharing site in this paper. The main contribution of our work is an extensive trace-driven analysis of Yahoo! Video workload dynamics.

We crawled all 16 categories on the Yahoo! Video site for 46 days (from July 17 to August 31 2007), and the data was collected every 30 minutes. This measurement rate was chosen as a tradeoff between analysis requirement and resource

constraint. Due to the massive scale of Yahoo! Video site, we limited the data collection to the first 10 pages of each category. Since each page contains 10 video objects, each time the measurement collects dynamic workload information for 1600 video files in total. Throughout the whole collection period, we recorded 9,986 unique videos and a total of 32,064,496 video views. This can be translated into a daily video request rate of 697064, and gave approximately 5.54% coverage on the total Yahoo! Video workload in July 2007, based on [2].

2. WORKLOAD STATISTICS

2.1 Static Properties

Video Duration: We recorded 9,986 unique videos in total, and the video durations range from 2 to 7518 seconds. Among them, 76.3% is less than 5 minutes, 91.82% is less than 10 minutes, and 97.66% is less than 25 minutes. The mean video duration is 283.46 seconds, and the median duration is 159 seconds.

File Popularity: File popularity is defined as the distribution of stream requests on video files during a measurement interval. We performed goodness-of-fit test with several distribution models, and found Zipf with an exponential cutoff fits best and well on the file popularity at four time scales - 30 minutes, 1 hour, 1 day, and 1 week.

2.2 Temporal Properties

Job Size Stationarity: Job size distribution is defined as the distribution of stream requests on video durations during a measurement interval. We use *histogram intersection distance* [4] to measure the change between two job size distributions, and calculated the pair-wise histogram intersection distance of two adjacent data points during the measurement. Figure 1 shows the CDFs of histogram intersection distance distribution for 3 time scales. We can see that within 30-minute and one-hour scale, the histogram distance is very small for most of the time. For example, 90% of the time it is no more than 0.15. But from day to day, the difference of request size distributions is obvious. This indicates that short-term dynamic provisioning only needs to focus on request arrival rate dynamics, while capacity planning at daily or longer basis has to take into account both arrival rate and job size dynamics.

Arrival Rate Predictability: We calculate the autocorrelation coefficient of the arrival rates at the time scale of 30 minutes, and from Figure 2 we can see that the workload is highly correlated in short term. We also use Fourier

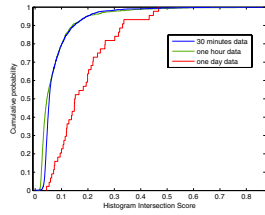


Figure 1: Histogram intersection distance distribution

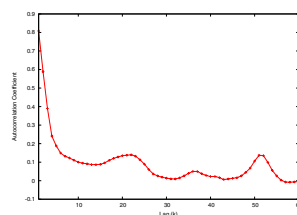


Figure 2: Workload autocorrelation coefficient

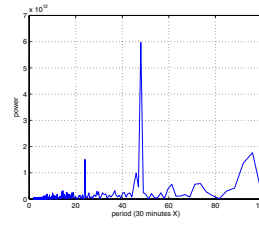


Figure 3: Workload periodicity

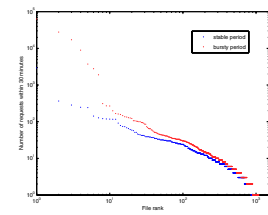


Figure 4: Comparison of a bursty interval and its preceding interval

analysis to discover the possible periodicity in the workload dynamics after removing a few workload spikes. As shown in Figure 3, the maximum value on the figure indicates that the period is one day. With the strong periodicity components, well-known statistical prediction approaches can be applied to further improve the accuracy in capacity planning.

Burstiness: While we can not predict unexpected spikes in the workload, it is necessary to learn the nature of the burstiness and find out an efficient way to handle it once a bursty event happens. The comparison of the request (popularity) distribution during one spike interval and that in the preceding interval is shown in Figure 4. We can see that the workload can be seen as two parts: a base workload similar to the workload in the previous normal period, and an extra workload that is due to several very popular files.

3. WORKLOAD AND CAPACITY MANAGEMENT: A VIRTUAL DESIGN

3.1 Methodology

System model: We model a single video server as a group of virtual servers with First Come First Served (FCFS) queues. The virtual server number corresponds to the physical server's capacity, which is defined as the maximum number of concurrent streams delivered by the server without losing a quality of stream. In the analysis, the number 300 is chosen for the capacity of a video server based on the empirical results in [3]. In this way, we model the video service center as a queueing system with multiple FCFS servers.

Workload Scheduling Schemes: We choose two well-known scheduling schemes to study: random dispatching, which doesn't make use of any information of the servers and just sends each incoming job to one of s server uniformly with probability $1/s$; Least workload Left (LWL) scheme, which tries to achieve load balancing among servers by making use of the per-server workload information and assigns the job to the server with the least workload left at the arrival instant.

Service Level Agreements: We consider two QoS metrics for SLAs: the stream quality for an accepted connection, and the waiting time of a video request in the queue before accepted for streaming. Assume enough network bandwidth, then QoS on stream quality within the data center side can be guaranteed through admission control based on server capacity. For the waiting time W , we consider the bound on the tail of the waiting time distribution (called W_{tail}), defined as $P[W > x] < y$ with $x > 0$, $y < 1$. For example, SLA could be that 90% of the requests experience no more than 5 seconds delays, i.e., $x = 5$ and $y = 90\%$.

3.2 Results

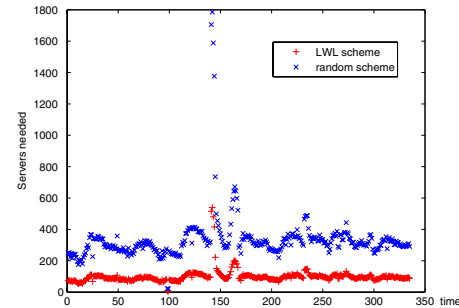


Figure 5: Server Demands with different scheduling schemes

Taking the one-week measurement data from Aug 13th to Aug 20th, we numerically calculated the server demands of random and LWL dispatching schemes with Poisson arrivals (based on results in [5]) and set the SLA requirement as W_{tail} : $\Pr[W > \text{mean service time}] < 0.3$. Figure 5 shows the results; on average 69.9% of servers could be saved with LWL scheme as compared to random dispatching scheme.

4. CONCLUSIONS

In this paper, we present the measurement study of a large Internet video sharing site - Yahoo! Video. With a clear goal to facilitate the data center design, this paper gives a comprehensive workload characterization and proposes a set of guidelines for workload and capacity management in a large-scale video distribution system. The immediate work for next step is expanding the measurement to cover larger portion of the workload on Yahoo! Video.

5. REFERENCES

- [1] Youtube. <http://www.youtube.com>.
- [2] ComScore Video Metrix report: U.S. Viewers Watched an Average of 3 Hours of Online Video in July, 2007.
- [3] L. Cherkasova and L. Staley. Measuring the Capacity of a Streaming Media Server in a Utility Data Center Environment. In *MULTIMEDIA '02*, New York, NY, USA, 2002. ACM.
- [4] M. J. Swain and D. H. Ballard. Color Indexing. *Int. J. Comput. Vision*, 7(1):11–32, 1991.
- [5] W. Whitt. Partitioning Customers into Service Groups. *Management Science*, 45(11):1579–1592, Nov 1999.