



Long-term Streaming Media Server Workload Analysis and Modeling

Wenting Tang, Yun Fu¹, Ludmila Cherkasova, Amin Vahdat¹
Internet Systems and Storage Laboratory
HP Laboratories Palo Alto
HPL-2003-23
January 29th, 2003*

E-mail: {wenting.tang, lucy.cherkasova@hp.com, {fu, vahdat}@cs.duke.edu

enterprise
media
servers,
workload
analysis,
long-term
trends
analysis,
performance
modeling,
synthetic
workload
generator

Currently, Internet hosting centers and content distribution networks leverage statistical multiplexing to meet the performance requirements of a number of competing hosted network services. Developing efficient resource allocation mechanisms for such services requires an understanding of both the short-term and long-term behavior of client access patterns to these competing services. At the same time, streaming media services are becoming increasingly popular, presenting new challenges for designers of shared hosting services. These new challenges result from fundamentally new characteristics of streaming media relative to traditional web objects, principally different client access patterns and significantly larger computational and bandwidth overhead associated with a streaming request. To understand the characteristics of these new workloads we use two long-term traces of streaming media services to develop MediSyn, a publicly available streaming media workload generator. In summary, this paper makes the following contributions: i) we model the *long-term* behavior of network services capturing the process of file introduction and changing file popularity, ii) we present a novel *generalized* Zipf-like distribution that captures recently-observed popularity of both web objects and streaming media not captured by existing Zipf-like distributions, and iii) we capture a number of characteristics unique to streaming media services, including file duration, encoding bit rate, session duration and non-stationary popularity of media accesses.

* Internal Accession Date Only

¹ Dept of Computer Science, Duke University, Durham, NC 27708

© Copyright Hewlett-Packard Company 2003

Approved for External Publication

Long-term Streaming Media Server Workload Analysis and Modeling *

Wenting Tang[†], Yun Fu[‡], Ludmila Cherkasova[†], and Amin Vahdat[‡]

[†]Hewlett-Packard Laboratories
1501 Page Mill Road, Palo Alto, CA 94303
email: {wenting.tang, lucy.cherkasova}@hp.com

[‡]Dept of Computer Science, Duke University
Durham, NC 27708
email: {fu,vahdat}@cs.duke.edu

Abstract

Currently, Internet hosting centers and content distribution networks leverage statistical multiplexing to meet the performance requirements of a number of competing hosted network services. Developing efficient resource allocation mechanisms for such services requires an understanding of both the short-term and long-term behavior of client access patterns to these competing services. At the same time, streaming media services are becoming increasingly popular, presenting new challenges for designers of shared hosting services. These new challenges result from fundamentally new characteristics of streaming media relative to traditional web objects, principally different client access patterns and significantly larger computational and bandwidth overhead associated with a streaming request. To understand the characteristics of these new workloads we use two long-term traces of streaming media services to develop MediSyn, a publicly available streaming media workload generator. In summary, this paper makes the following contributions: i) we model the *long-term* behavior of network services capturing the process of file introduction and changing file popularity, ii) we present a novel *generalized* Zipf-like distribution that captures recently-observed popularity of both web objects and streaming media not captured by existing Zipf-like distributions, and iii) we capture a number of characteristics unique to streaming media services, including file duration, encoding bit rate, session duration and non-stationary popularity of media accesses.

1 Introduction

Two recent trends in network services motivate this work, a move toward shared service hosting centers and the growing popularity of streaming media. Traditionally, service providers over-provision their sites to address highly bursty client access patterns. These access patterns can vary by an order of magnitude on an average day [10] and by three orders of magnitude in the case of flash crowds. In fact, services are often most valuable exactly when the unexpected takes place. Consider the example of a news service when an important event takes place; load on CNN reportedly doubled every seven minutes shortly after 9 AM on September 11, 2001 [7].

Thus, we are pursuing a vision where large-scale hosting infrastructures simultaneously provide “resource-on-demand” capabilities to competing Internet services [20, 8]. The idea is that the system can use statistical multiplexing and efficient resource allocation to dynamically satisfy the requirements of services subject to highly bursty access patterns. For instance, surplus resources resulting from “troughs” in accesses to one service may be

*This work is partially completed while Yun was a summer intern at HP Labs during 2002. A. Vahdat and Y. Fu are supported in part by the National Science Foundation (EIA-9972879). A. Vahdat is also supported by an NSF CAREER award (CCR-9984328).

reallocated to satisfy the requirements of a second service experiencing a peak. Further, Service Level Agreements (SLAs) may specify that, under resource constraints, one service should preferentially receive resources over other services.

A second emerging trend is the growing popularity of streaming media services. Streaming media takes the form of video and audio clips from news, sports, entertainment, and educational sites. Streaming media is also gaining momentum in enterprise intranets for training purposes and company broadcasts. These workloads differ from traditional web workloads in many respects, presenting a number of challenges to system designers and media service providers [13, 18]. For instance, transmitting media files requires more computing power, bandwidth and storage and is more sensitive to network jitter than web objects. Further, media access lasts for a much longer period of time and allows for user interaction (pause, fast forward, rewind, etc.).

The long-term goal of our work is to study resource provisioning and resource allocation at the confluence of the above two trends: network service hosting infrastructures for next-generation streaming workloads. A key obstacle to carrying out such a study is the lack of understanding of changing client access patterns over a long period of time. For both hosting centers and content distribution networks (CDNs), we require such an understanding to determine, for example, how to place objects at individual sites (potentially spread across the network) and how to allocate resources to individual streams and to individual clients.

Thus, we use long-term traces from two streaming media services to construct an open-source media workload generator called MediSyn. For MediSyn, we develop a number of novel models to capture a broad range of characteristics for network services. We also demonstrate how these models generalize to capture the characteristics of traditional web services. Overall, this paper makes the following contributions:

- A primary contribution of our work is its focus on the *long-term* behavior of network services. Among the features of our synthetic generator is the ability to reflect the dynamics and evolution of content at media sites and the change of access rate to this content over time. Existing workload generators assume that there is a set of active objects fixed at the beginning of the “trace”. Similarly, existing techniques assume that object popularity remains the same over the entire duration of the experiment. While these are reasonable assumptions for experiments designed to last for minutes, we are interested in long-term provisioning and resource allocation, as well as the resource allocation for simultaneous competing services (consider a CDN simultaneously hosting hundreds of individual services).
- It was observed [2, 12] that the popularity distribution in media workloads collected over significant period of time (more than 6 months) does not follow a Zipf-like distribution. We present a novel *generalized* Zipf-like distribution to capture the popularity distribution in such workloads. The traditional Zipf-like distribution is a special case of the proposed *generalized* Zipf-like distribution.
- We designed a set of new models to capture a number of characteristics critical to streaming media services, including file duration, file access prefix duration, non-stationary file popularity, new file introduction process and diurnal access patterns.

The rest of this paper is organized as follows. Section 2 outlines the workload properties that MediSyn attempts to capture and the workloads used to develop the models for MediSyn. We present the workload generation process adopted by MediSyn in Section 3. Section 4 introduces the models used in MediSyn and discusses their specifics. We review previous related work in Section 5. Finally, we conclude with a summary and future work in Section 6.

2 Media Workload Properties

Accurate workload characterization is critical for successful generation of realistic workloads. A synthetic media workload generator can produce traces with targeted, controllable parameters and desired distributions for performance experiments studying effective streaming media delivery architectures and strategies. For such experiments, the generated workload must not only mimic the highly dynamic resource-utilization patterns found on today’s media systems but also provide flexible means to generate more intensive, bursty and diverse workloads for future media systems. Challenges to designing a useful analytical workload generator include:

- identifying essential properties of workloads targeted by synthetic workload generators, and those that most affect the behavior of hosting centers;
- designing appropriate mathematical models that closely reproduce the identified workload properties from real traces.

In this section, we highlight the main properties of streaming media workloads modeled in MediSyn, and explain why we believe these properties are important. In this analysis, we attempt to summarize the results of earlier studies [13, 18, 12, 3, 21] characterizing streaming media workloads. Throughout this paper, we use two representative streaming media server logs, collected over a period of years, to demonstrate the chosen properties and to validate our mathematical models introduced to reflect these properties. The streaming media server logs represent two different media services: *HP Corporate Media Solutions Server (HPC)* and *HPLabs Media Server (HPL)*. We define a *session* as a client access to a particular file¹. Table 1 briefly summarizes the workloads.

	HPC	HPL
Log duration	29 months	21 months
Number of files	2,999	412
Number of sessions	666,074	14,489

Table 1: Summary for two media logs used to develop property models in MediSyn.

We partition media workload properties in two groups: *static* and *temporal* properties.

- **Static properties** provide the characteristics of the underlying media files, reflect the aggregate, quantitative properties of client accesses (independent of the access time), and present the properties of individual file accesses. Static properties include:
 - *file duration*,
 - *file encoding bit rate*,
 - *file access popularity*,
 - *file access prefix*.
- **Temporal properties** reflect the dynamics and evolution of accesses to media content over time, and determine the ordering and the timing of session arrivals. The temporal properties of media workloads include:
 - *new file introduction process*,
 - *file life span*,
 - *diurnal access pattern*.

Sections 2.1 and 2.2 describe these properties in detail.

¹A session may consist of multiple client’s requests reflecting history of the client’s interactivity. We plan to extend MediSyn to cover interactivity patterns in our future work.

2.1 Static Properties

- **File Duration.** Different from web, media files have two dimensions: *duration* of the file measured in *time* and the corresponding media file *size* measured in *bytes*, a direct product of file duration and the corresponding encoding bit rate. Earlier media workload studies report different distributions characterizing file durations. However, the duration distribution is largely determined by the nature of the media content stored at a site, which could be a short (a few minutes) media clip or a long (one to two hours) movie. Thus, any single distribution, e.g, the heavy-tail distribution used to capture web object size, may fail to characterize media file duration. For example, our two logs represent two significantly different file duration distributions as shown in Table 2.

Duration	≤ 10 minutes	10-30 minutes	≥ 30 minutes
HPC Log	42%	23%	34%
HPL Log	14%	7%	79%

Table 2: File duration distributions of the HPC and HPL logs, respectively.

The HPC trace has well represented groups of short, medium, and long videos, while the HPL trace is skewed toward longer videos.

File duration influences the duration of playback (download) and how long system resources should be committed to a particular stream. This information is potentially necessary to perform admission control. While a web object retrieval typically takes less than a few hundred milliseconds, most media files require minutes to hours to transmit. Thus, it is important to represent the file duration distribution accurately.

- **File Encoding Bit Rate.** The encoding bit rate of a media file determines the bandwidth and system resources required to deliver the file. Typically, there is a set of popular encoding bit rates offered by commercial encoding software [19] and guided by the bandwidth from different groups of Internet audience.
- **File Popularity.** File popularity is defined as the number of accesses to a file within a certain period of time. The popularity of a file over the entire trace is its global popularity. Recent studies observe highly uneven file popularity both in web and streaming media workloads [13, 12, 4], implying that most accesses to a site are concentrated on a relatively small set of files. Both media workloads used in our study exhibit strong access locality: 14%(30%) of the files accessed on the server account for 90% of the media sessions for the HPC(HPL) trace. The uneven file popularity is one of the major sources of temporal locality.

Traditionally, a Zipf-like distribution is used to capture the global file popularity for web workloads. However, several studies [2, 12] analyzing the properties of media workloads collected over significant periods of time (more than 6 months) report that the popularity distribution of these sites does not follow a Zipf-like distribution. The corresponding global file popularity distribution curve (on a log-log scale) exhibits a “circular”-shape [2, 12]. One explanation of this phenomenon is that a long-term trace can collect enough files with similar popularities. This property is tightly related to the change of rate in file accesses (or the change of file popularity rank) defined by a file life span discussed below. Considering that the life span of most of these files is shorter than the trace duration, there must be many files with equivalent volume of accesses in the trace, which cannot be captured by a Zipf-like distribution very well. Since the design and evaluation of caching and content distribution strategies require a proper understanding of file popularity, accurately reflecting these skewed file popularity distributions is very important.

- **File Access Prefix.** Prior studies [12, 3] show that many clients do not finish the playback of a full video/audio clip. Typically, this reflects the browsing nature of client accesses, client time constraints, or QoS-related issues. Most incomplete sessions (i.e. terminated by clients before the video is finished entirely) access only the initial segments of media files. In the HPC(HPL) trace, only 29%(12.6%) of the accesses finish the playback of the files. 50%(60%) of the accesses in the HPC(HPL) trace last less than 2 minutes. This high percentage of incomplete accesses as well as a high number of sessions accessing only the initial part of the file create a very special resource usage model, which is widely considered for streaming media cache design [23].

2.2 Temporal Properties

- **Causes of Temporal Locality in Media Workloads Collected Over Long Period of Time.**

Temporal reference locality, which is universally observed in web and media workloads [13, 12, 4], is the primary factor that affects session arrival ordering. Temporal locality states that recently accessed objects are likely to be accessed in the near future in the access stream. Two factors can cause the temporal locality in the access stream: *skewed popularity distribution* and *temporal correlation* [11, 16]. Since popular files have a higher probability to be accessed within the access stream, a file’s popularity contributes to its temporal locality. If we randomly permute the access stream, the temporal locality caused by skewed popularity is still preserved under reordering. However, temporal locality caused by temporal correlation cannot be preserved under random permutation. In MediSyn, to generate a stream of session arrivals exhibiting proper temporal locality, we need to clearly understand and distinguish the causes of temporal locality.

To check the existence of possible temporal correlation among sessions for the same files in our traces, we compare reference distances [4] between the original HPC trace and a randomly permuted HPC trace. Let s_1, s_2, \dots, s_n be the stream of accesses representing the media sessions of the entire HPC trace. Let $s_{i_1}, s_{i_2}, \dots, s_{i_m}$ be the sequence of sessions to the same file f_i . The *reference distances* of $s_{i_2}, s_{i_3}, \dots, s_{i_m}$ are defined as $i_2 - i_1, i_3 - i_2, \dots, i_m - i_{m-1}$. We calculate the reference distances for all the files and their sessions over the entire HPC trace. Then we apply similar procedure to the randomly permuted HPC trace.

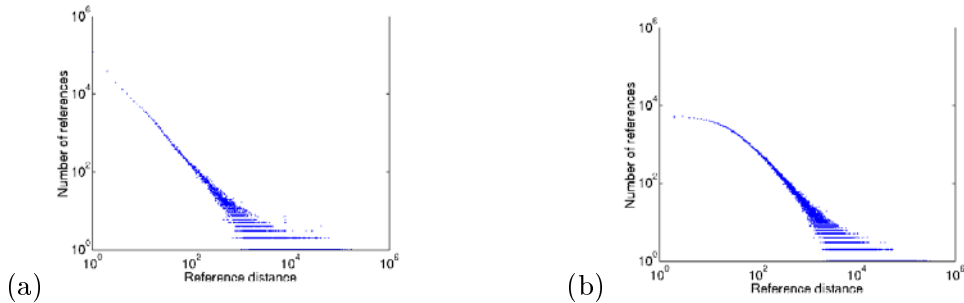


Figure 1: (a) Reference distances calculated over the entire trace period in the original HPC trace. (b) Reference distances calculated over the entire trace period in the permuted HPC trace.

Figure 1 (a) shows the histogram of the reference distances in the original HPC trace on a log-log scale. The X-axis shows the reference distances, and the Y-axis shows the number of sessions in the original HPC trace for the corresponding reference distance. Figure 1 (b) shows the histogram of the reference

distances in the randomly permuted HPC trace on a log-log scale. It can be observed that two curves are not the same. In Figure 1 (b), there are less references with small distances. For example, there are 122,765 references with distance 1 in Figure 1 (a). But there are only 5,468 references with distance 1 in Figure 1 (b). Since in Figure 1 (b), reference distances are only determined by the file popularity, we assume the reason that Figures 1 (a) and (b) are not the same is that there is temporal correlation among sessions for the same files over the entire trace.

To verify whether our media traces exhibit short term temporal correlation, we calculate reference distances within every day and sum the number of references with the same distance over all days for the HPC trace. Then, we permute accesses within every day and calculate reference distances again for the permuted trace. The results are shown in Figures 2 (a) and (b). We observe that there is almost no difference between the original trace and the permuted trace on reference distances. This implies that there is no temporal correlation for sessions within a single day and that temporal locality of sessions within a single day is purely determined by the file popularity distribution within that day.

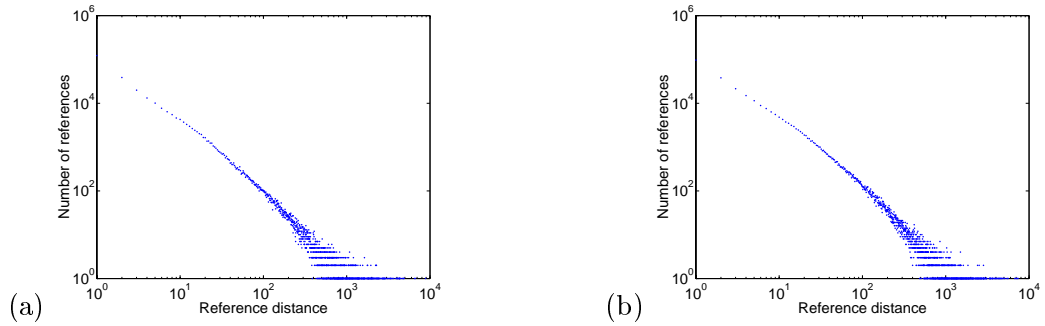


Figure 2: (a) Reference distances calculated within each day in the original HPC trace. (b) Reference distances calculated within each day in the permuted HPC trace.

The analysis above implies that temporal correlation in media workloads collected over long period of time is due to long-term temporal correlation exhibited on a daily time scale, and that there is no temporal correlation for sessions within a single day. These observations motivate our choice of temporal properties in MediSyn. The temporal properties described below intend to reflect the dynamics and evolution of accesses to media content over time and to define the proper temporal locality and long-term temporal correlation found in media workloads.

- **New File Introduction.** One recent study [12] observes that accesses to new files constitute most of the accesses in any given month for enterprise media servers. We envision that this access pattern is even more pronounced for media news and sports sites. While for educational media sites the rate of new file introduction and accesses to them might be different, we aim to design a generic parameterized model capable of capturing the specifics of new content introduction for different media workloads. Among the design goals of our synthetic generator is the ability to reflect the evolution of media content provided by different media sites over a long period of time (months). Since we design MediSyn to support detailed resource allocation studies, we must account for the dynamic introduction of new content and its relative popularity.
- **File Life Span.** A new property called *life span* has recently been proposed in [12] to measure the change in access rate of newly introduced files. Life spans reflect the timeliness of accesses to the introduced files.

We observe that accesses to a media file are not uniformly distributed over the entire trace period. Instead, most accesses for a file occur shortly after the file is introduced, with access frequency gradually decreasing over time. For example, for the HPC(HPL) log, 52%(51%) of the accesses occur during the first week after file introduction, while only 16%(10%) of the accesses occur during the second week, etc. Hence, the file access frequency (file popularity) changes over time. In other words, file popularity is *non-stationary* over the trace period. This phenomenon implies that session arrivals might be very bursty when new files are introduced at a site. This makes media services an ideal candidate for shared hosting infrastructures and requires that special care be taken for efficient resource allocation in the hosting infrastructures.

Life span distribution defines the file popularity changes over a daily time scale. Our analysis of the HPC and HPL traces shows that life spans of streaming media files can be classified into two categories: *news-like* and *regular* life spans according to how the popularity changes over time. A file with a news-like life span has a fast changing rate of popularity during its lifetime. A file with regular life span has much more stable popularity during its lifetime.

The long term temporal correlation in media traces is caused by a non-stationary file popularity over considered period of time. Since the rate of change in file popularity over time can be captured by the file life span combined with the new file introduction process, we use this property to model the temporal correlation in the synthetic media workloads.

- **Diurnal Access Patterns.** Earlier studies observed the diurnal access patterns for streaming media workloads [13, 18, 2, 3, 17]. The diurnal access pattern defines how the number of accesses to a site varies during a given period of time, e.g., a day. Diurnal access patterns might significantly vary for different media sites. For mixed media workloads utilizing a shared infrastructure, the diurnal access patterns have to be taken into account when designing the optimal support for efficient resource allocation. Additionally, the diurnal access pattern is important for capturing the burstiness of resource consumption within a given time period. The diurnal access patterns are defined using the second time scale in our synthetic workload generator, e.g. within a day.

In summary, the temporal properties we introduced not only explain the temporal locality observed in streaming media workloads, but also have fundamental implications on resource sharing and allocation within hosting infrastructures and content distribution networks. The temporal locality exhibited in media workloads is not an independent property: it is a combined *effect* of the skewed popularity distribution, new file introduction process and life span distributions of media files. Thus, we do not attempt to model temporal locality as a stand-alone property in MediSyn: it is a combined result of multiple models aiming to reflect the temporal properties of media workloads discussed above.

3 Workload Generation Process in MediSyn

MediSyn’s goal is to generate a synthetic trace representing a sequence of file accesses to media service. This process consists of two steps: *file property generation* and *file access generation*.

- **File property generation**

The first step is to generate values for all properties introduced in Section 2 for each media file. Static properties define file set parameters (duration and encoding bit rate) and the principal access patterns

(popularity and prefix). Temporal properties define the ordering and timing of media sessions. MediSyn defines these properties using a parameterized set of distributions in the input configuration file.

If a property can be described by a value such as global file popularity, duration, encoding bit rate, MediSyn first generates a sequence of values according to the given distribution, and then selects a value for each file. If a property is modeled as a distribution, the choice of the distribution and parameter(s) of the distribution are generated for each file.

Thus, a file is the basic unit to which the property values are propagated at the first step of workload generation. At the end of the first step, the set of corresponding static and temporal properties shown in Table 3 is generated for each file. Section 4 describes each property generation and correlations among the properties in detail.

File id	Duration	Bit rate (Kbps)	Popularity	File Introduction Time (sec)	Life span	Life span parameters	...
1	3600	112	20000	100	Pareto	1.0	...
2	200	350	14300	50	Lognormal	2.0,10.0	...
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
n	600	28.8	1	10000	Lognormal	1.0, 1.0	...

Table 3: Properties generated for each file.

- **File Access Generation**

Taking the assigned file popularities as the basis, MediSyn independently generates the arrival of media sessions to each file using: i) the initial file introduction time, ii) the life span of the file, and iii) the diurnal access pattern of the file. Each file access includes the following three fields:

- *timestamp* indicating the session arrival time,
- *file id* specifying the target file accessed during the media session,
- *file access prefix* describing the duration of the media session.

Once a sequence of media sessions is generated for each file, all the media sessions are sorted according to a global time and merged together to generate the synthetic trace.

4 Main Models of Workload Generation in MediSyn

This section describes the models used in MediSyn to capture static and temporal properties of streaming media workloads.

4.1 Duration

Prior studies [12, 3] observed that media files might be classified into a set of groups according to their durations. Different workloads can be grouped based on the content of media files hosted by a streaming service. For example, music sites may have file durations from 3 to 5 minutes, while movie sites may have file durations from one and half to two hours. While a particular workload might be captured by a certain statistical distribution, the same distribution may fail to capture another workload. In our case, although the file duration distribution of the HPC trace can be modeled by a heavy-tail distribution such as a *Weibull* distribution [15], the same distribution fails to capture the file duration distribution of the HPL trace.

As shown in Figure 3 (a) and Figure 4 (a), the file durations in our traces are concentrated around a set of hot points. These hot points are usually some common durations, semantically meaningful to a particular type of media content. Based on this observation, we classify these hot points into a set of groups and use a set of *normal* distributions to model the grouped file duration distribution as shown in Figure 3 (a) and Figure 4 (a). Here, each group is modeled by a normal distribution with the mean (μ) of each distribution defined by the hot point of that group. The standard deviation (σ) of each normal distribution determines the concentration of the durations within that group.

Note that we do not use segmented PDFs (Probability Density Functions) to model the duration distribution. We assume a hot point can affect the entire duration scope rather than just a segment. Thus, we use an aggregated distribution, whose PDF sums the PDFs of all normal distributions proportionally. To proportionally sum all duration groups, we associate each group with a ratio determined by the number of files in the group compared with the total number of files in the trace. So the normal distribution PDF of each group is normalized against the ratio of that group. If only a fraction of a normal distribution for a group is used, normalization is performed on the adopted fraction of the distribution. For example, since the mean of the first group in Table 4 is 0, only half of the normal distribution is used. Tables 4 and 5 present the mean (μ), the standard deviation (σ) and the ratio of each normal distribution for the HPC and HPL trace respectively. They show that the HPC and HPL traces have different hot points.

In MediSyn, users can specify a set of duration groups with different μ , σ and ratios based on the nature of the media workload they want to generate. For each duration group, MediSyn generates a sequence of durations according to the ratio and the normal distribution of the group. We use the rejection method [15] to generate the duration sequence according to the parameterized normal distribution. Figure 5 (a) and Figure 5 (b) show the durations generated by MediSyn to simulate the HPC workload based on the parameters presented in Table 4.

Group	1	2	3	4
μ	0	2000	3300	5821
σ	600	400	600	1223
ratio	63%	10%	18%	9%

Table 4: Parameters of the normal distributions for the HPC trace.

Group	1	2	3	4	5
μ	117	2900	4200	5160	6300
σ	1200	240	360	180	1000
ratio	19%	26%	30%	10%	15%

Table 5: Parameters of the normal distributions for the HPL trace.

4.2 Encoding Bit Rate

Since most of commercial media servers are designed to stream media files encoded at some constant bit rates, the current version of MediSyn is designed to only generate a set of constant bit rates for the underlying fileset. MediSyn models encoding bit rates by a discrete distribution, where the value of each bit rate and the ratio of the bit rate occupied in the fileset can be specified. Based on the discrete distribution provided by users, MediSyn generates a sequence of bit rates for the fileset and matches the bit rate with the file duration randomly, since we observe that there is no correlation between them in our traces (the correlation coefficient is 0.0144).

4.3 Popularity

Earlier studies [13, 18] found that media file popularity can often be captured by a *Zipf-like* distribution. A Zipf-like distribution states that the access frequency of the i -th most popular file is proportional to $1/i^\alpha$. If the

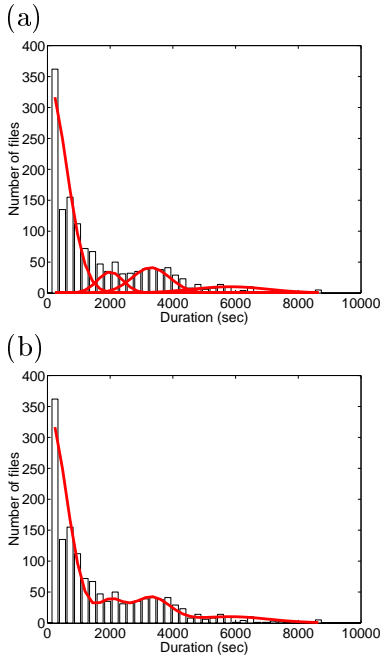


Figure 3: PDF of the HPC duration distribution. (a) 4 normal distributions to capture the 4 peaks. (b) The aggregate distribution of the 4 normal distributions.

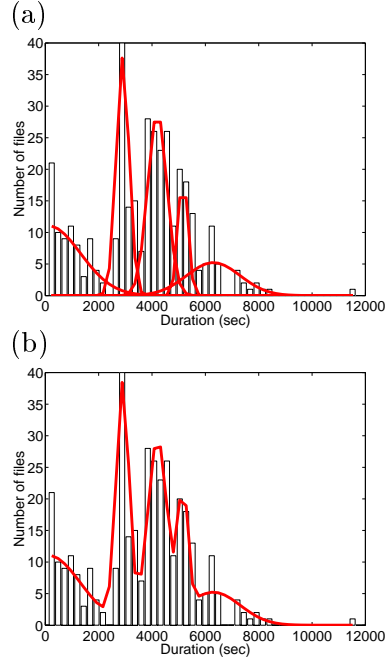


Figure 4: PDF of the HPL duration distribution. (a) 5 normal distributions to capture the 5 peaks. (b) The aggregate distribution of the 5 normal distributions.

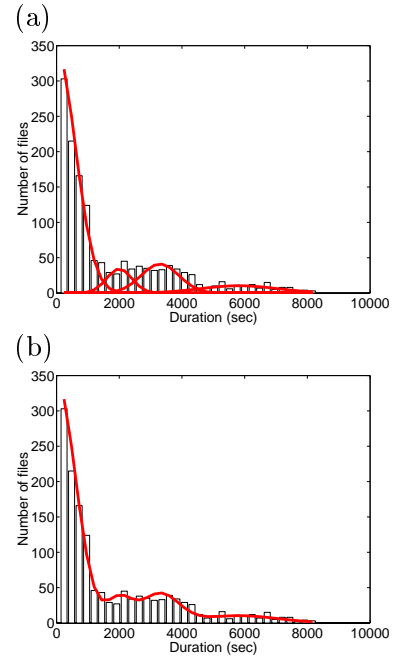


Figure 5: PDF of the MediSyn duration distribution. (a) 4 normal distributions to capture the 4 peaks. (b) The aggregate distribution of the 4 normal distributions.

frequencies of files and the corresponding popularity ranks are plotted on a log-log scale, a Zipf-like distribution can be fitted by a straight line. A larger α implies more sessions are concentrated on the most popular files. Some synthetic workload generators [6, 17] also adopt a Zipf-like distribution in generating file popularity.

However, recent studies [2, 12, 3, 5, 9] observed that for some web and streaming media workloads, a Zipf-like distribution does not accurately capture the file popularity distribution. The popularity distribution of these workloads shows a circular curve on a log-log scale.

For reference, Figure 6 shows the file popularity distributions of the HPC and the HPL traces over the entire trace periods on a log-log scale. They are more like circular curves similar to those distributions noticed in previous web studies [5, 9] and media workload studies [2, 12].

If we use a straight line (a Zipf-like distribution) to fit the circular curve and generate session frequencies based on the value of α obtained by curve fitting, the generated frequencies must be skewed from the original session frequencies. Breslau et al [9] calculated α by excluding the top 100 files. For our traces, not only the beginning but also the end of the curves cannot be fitted by straight lines. Moreover, since the most popular files are especially important for synthetic streaming media workloads, we cannot ignore the first 100 files.

Thus, an important contribution of this work is that a *generalized Zipf-like* distribution is proposed as a unified method to capture file popularity distributions of both Zipf-like and circular-curve shapes. A generalized Zipf-like distribution can be fitted by a straight line on a log-log scale after a Zipf k -transformation. The Zipf k -transformation for file popularity is defined as follows: assume x is a file rank, y is the corresponding access frequency for the file, k_x and k_y are scale parameters, the k -transformation transforms the original x and y as follows:

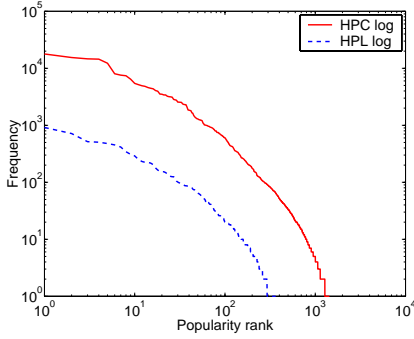


Figure 6: The original popularity distributions of the HPC trace and the HPL trace on a log-log scale.

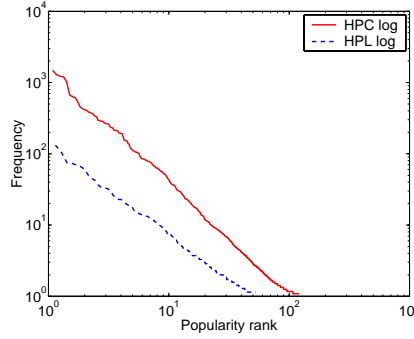


Figure 7: The popularity distributions after Zipf k -transformation on a log-log scale.

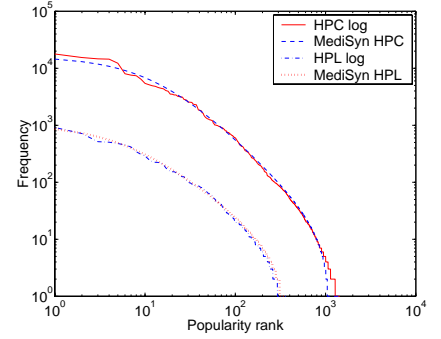


Figure 8: Comparison between the popularity distribution generated by MediSyn and the original traces on a log-log scale.

$$x_k = \frac{x + k_x - 1}{k_x} \quad (1)$$

$$y_k = \frac{y + k_y - 1}{k_y} \quad (2)$$

After the k -transformation, x_k and y_k denote the transformed file rank and frequency respectively. Figure 7 shows the relationship between x_k and y_k of the HPC and HPL traces on a log-log scale respectively. We observe that they are perfectly straight lines. The α value of the Zipf k -transformation is derived through linear regression [14]. Table 6 shows some critical parameters related to the curve fitting of the two workloads. As shown in the table, R^2 is 0.995 for both the HPC and HPL traces, indicating that straight lines fit both the distributions very well.

Trace	α	R^2	k_x	k_y	Maximum frequency	Number of files
HPC	1.561	0.995	12	12	17831	1434
HPL	1.23	0.995	7	7	961	364

Table 6: Parameters of Zipf k -transformation of the HPC and the HPL traces.

The reason that the original traces do not show perfectly straight lines at the heads of the curves is that there is little differentiation in the frequencies of the most popular files (with smaller x). It can be attributed to the fact that a long-term trace can collect enough files with similar popularities over time, and thus these files can be considered as a group (equivalence class), where a group rank will be a better reflection of the file popularities. Intuitively, the effect of the k -transformation is that the popularity follows a Zipf-like distribution if we check every group of k_x files. We divide x by k_x to scale the file ranks so that the $((i - 1)k_x + 1)$ -th rank becomes the i -th rank and reflect now the corresponding file group rank. Other file ranks are transformed to float numbers evenly distributed among integral ranks. So we actually move all points on the log-log scale along the x -axis to the left and squeeze the points to a more straight line. Similarly, the reason that the traces do not show perfectly straight lines at the tails of the curves is that there is not enough differentiation in the number of files with the lowest frequencies. So we divide y by k_y to squeeze those points along the y axis on the log-log scale. The value of k_y is not necessarily the same as k_x . However, MediSyn uses the same value for k_x and k_y based on our observations for both the HPC and HPL traces, which we simply refer to as k . The scale parameter k of our k -transformation is similar to the scale parameter k of a *generalized* Pareto distribution [1].

Another explanation for the k -transformation is that the original frequency sequence cannot be fitted by a Zipf-like distribution starting from rank 1, but it can be fitted into part of a Zipf-like distribution starting from rank k_x . To describe this file rank starting from k_x by a Zipf-like distribution, we have to divide its original rank by k_x . Similar explanation can be applied for k_y . Clearly, Zipf-like distributions are a special case of generalized Zipf distributions when $k = 1$.

To generate a sequence of frequencies following a generalized Zipf-like distribution, users of MediSyn specify the maximum frequency M for the most popular file, the number of files N , the scale parameter k and the Zipf-like distribution parameter α . MediSyn computes the frequency of the most popular x -th file ($x \in [1, N]$) using the following formula

$$\left(\frac{M_k}{\left(\frac{x+k-1}{k}\right)^\alpha} - 1 \right) k + 1, \quad (3)$$

where $M_k = \frac{M-1}{k} + 1$. Figure 8 compares the frequencies generated by MediSyn with the original frequencies in our traces.

To determine whether there is a correlation between file duration and file popularity, we compute the correlation coefficient between file popularity and file duration for both of our workloads. Table 7 shows these results. We use both the file frequency and the file rank as the popularity metric to compute the correlation coefficient. Overall, we observe no strong correlation between file popularity and file duration. So file duration and file popularity are randomly matched in MediSyn.

Workload	HPC frequency	HPL frequency	HPC rank	HPL rank
Correlation Coefficient	-0.03	0.05	-0.20	-0.002

Table 7: Correlation coefficient between file popularity and file duration.

We also check for possible correlation between popularity and encoding bit rate. Once again, there is no correlation between them, so MediSyn matches popularity and encoding bit rate randomly.

4.4 Prefix

One major characteristics of streaming workloads is that a significant amount of clients do not finish playing an entire media file. We refer to the duration between the start of a media session and the time when the session is terminated by the client as the *prefix duration* of the session, or simply the *prefix*. Figure 9 and Figure 10 show the PDF for the prefixes of two typical example files in the HPC trace. The “spikes” in the figures correspond to successfully completed media sessions for the files, while the other prefixes in the figures are incomplete sessions. We observe that there is a strong correlation between the file duration and the prefix distribution:

- *Complete sessions.* The fraction of complete sessions of a file highly depends on the file duration. A short file tends to have more complete sessions. For example, the file durations in Figure 9 and Figure 10 are 723 seconds and 4133 seconds respectively. The file in Figure 9 has more complete sessions than that in Figure 10. We use r_c to denote the ratio of complete sessions for a file compared with the total number of sessions for the file. Figure 11 shows the relationship between file duration and r_c . We can observe that the r_c of each file highly depends on the file duration.
- *Incomplete sessions.* The prefix distribution of incomplete sessions of a file depends on the file duration. Figure 9 reflects that the prefixes of incomplete sessions for a short-duration media file can be captured by

an exponential distribution. While for a long-duration file as shown in Figure 10, the prefixes of incomplete sessions follows a heavy-tail distribution, which cannot be captured by an exponential distribution.

Thus, the overall prefix distribution of a media workload highly depends on each file’s prefix distribution, which in turn depends on the duration of the file. There is not a straightforward solution to directly capture the overall prefix distribution for the entire workload. To generate each file’s prefix distribution, we first generate r_c for the file, then model the distribution of incomplete sessions for the file based on the assigned r_c .

To generate r_c for a file, we need to determine the relationship between r_c and the file duration as shown in Figure 11. We observe that the contour of the dotted area in Figure 11 follows a Zipf-like distribution. To obtain this curve, we segment the durations of all files into 1-minute bins. The maximum r_c value of each bin (denoted as r_c^{max}) constitutes the contour of the dotted area in Figure 11. Figure 12 shows the relationship between r_c^{max} of each bin and the duration (in minutes) for the corresponding bin on a log-log scale. Because the maximum value of r_c^{max} is 0.74, the curve is flat in the beginning. The other points can be fitted with a straight line. Thus, we can use a Zipf-like distribution to capture the distribution of r_c^{max} for all bins. We also observe that the r_c values for other files within a bin, are uniformly distributed between 0 and the r_c^{max} value of the bin. So, to generate r_c for each file, MediSyn first classifies files into 1-minute bins based on their durations. Then, MediSyn generates r_c^{max} for each bin. For files in each bin, their r_c values are chosen according to a uniform distribution between 0 and r_c^{max} of the bin. Through this process, the r_c of each file can be determined.

After the r_c value of each file is determined, the distribution of incomplete sessions needs to be determined. As mentioned above, depending on the file duration, it could be captured by an exponential distribution or some heavy-tail distribution. Additionally, both Figure 9 and Figure 10 show a similar shape in the beginning of the distributions within a certain range of duration. We observe that for more than 90% of the media files in the HPC trace, the distributions of prefixes within the first 5 minutes can be fitted by exponential distributions. These results confirm similar findings for an educational workload studied by Almeida et al [3]. Given the fact that prefixes within a certain duration range (e.g., the first 5 minutes) occupy a high percentage of total incomplete sessions, we introduce a cut-off point and use the following method to model the prefix distribution of a given media file:

- If a media file duration *is less than the cut-off point*, its incomplete prefixes are modeled by an *exponential* distribution.
- If a media file duration is *longer than the cut-off point*, the distribution of incomplete prefixes is modeled by the *concatenation of two distributions*. The distribution of incomplete prefixes less than the cut-off point is modeled by an *exponential* distribution. The distribution of the remaining incomplete prefixes longer than the cut-off point is approximated by a *uniform* distribution.

In the HPC trace, the cut-off point is 5 minutes. Users of MediSyn can specify their own cut-off point. We use the following denotations:

- r_e defines the ratio of incomplete sessions whose prefixes are within the cut-off point compared with the total number of sessions of the file,
- r_u defines the ratio of incomplete sessions whose prefixes are longer than the cut-off point compared with the total number of sessions of the file. If a file duration is less than the cut-off point, r_u is 0.

Given r_c for a file, MediSyn needs to generate the values of r_e and r_u for the file. The strategy is to generate r_e for the file and to set $r_u = 1 - r_c - r_e$. Figure 13 shows the relationship between r_c and r_e for all files in the HPC trace. We can see that for a given r_c , the value of r_e is bounded on both the lower and upper sides. If we

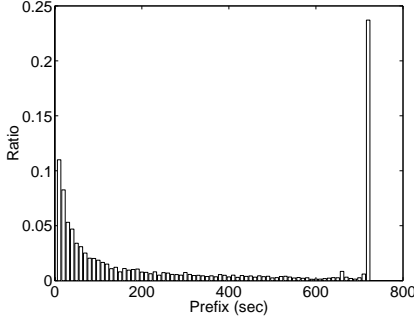


Figure 9: A typical access prefix distribution of short duration media file.

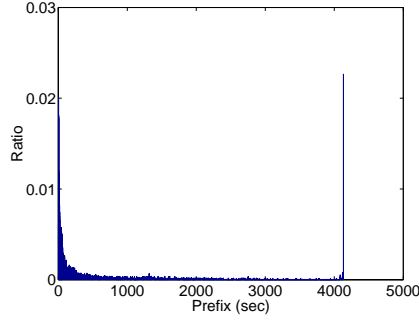


Figure 10: A typical access prefix distribution of long duration media file.

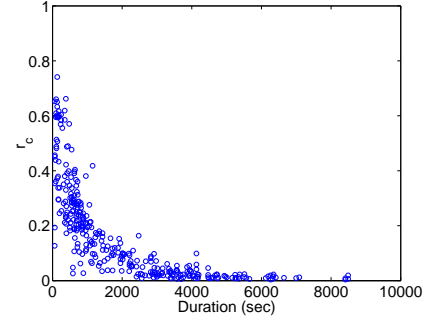


Figure 11: Fraction of complete sessions (r_c 's) versus the corresponding media file durations.

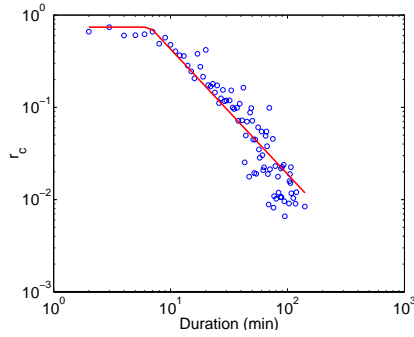


Figure 12: r_c^{max} of all bins versus the corresponding bin duration on a log-log scale.

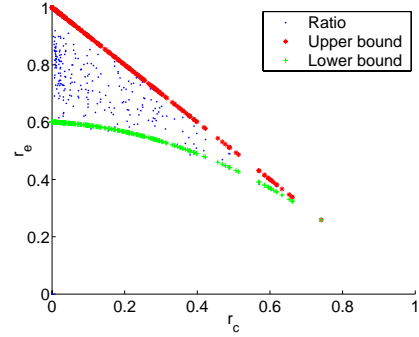


Figure 13: r_e versus r_c .

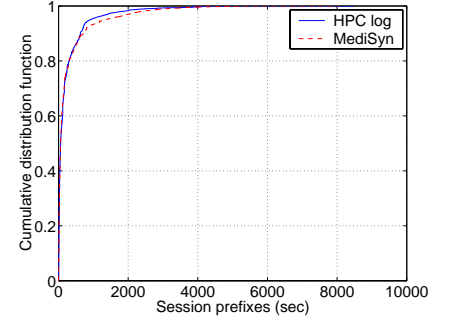


Figure 14: CDFs of the prefixes generated by MediSyn and the HPC trace.

denote the maximum of all r_c values as R_c^{max} (0.74 for the HPC trace), the upper bound r_e^{upper} and the lower bound r_e^{lower} of a given r_c can be computed by Equation 4 and Equation 5 respectively.

$$r_e^{upper} = 1 - r_c \quad (4)$$

$$r_e^{lower} = r_e^{upper} * (0.6 + 0.4 * r_c / R_c^{max}) \quad (5)$$

Figure 13 plots these calculated upper and lower bounds. The upper bound is caused by the limitation that the sum of r_c , r_e , and r_u is 1. The lower bound changes from 60% of the upper bound to the same as the upper bound while r_c increases from 0 to R_c^{max} . So during workload generation, for a given r_c , MediSyn generates the value of r_e according to a uniform distribution between the corresponding upper bound r_e^{upper} and the lower bound r_e^{lower} .

After generating r_c , r_e and r_u for each media file, MediSyn still needs to generate the mean (μ) of the exponential distribution for the incomplete prefixes. Since each file has its own exponential distribution for prefixes, we have to derive the distribution for the μ 's of all media files. Analysis of session prefixes in the HPC trace shows that μ 's of all media files follow a *normal* distribution.

MediSyn generates a sequence of prefixes according to the generated ratios of r_c , r_u , r_e and μ for each file. Then, these prefixes are randomly matched with all the sessions of the file. In this way, MediSyn can generate all session prefixes for every file. Figure 14 shows the cumulative distribution function (CDF) of the prefixes in the HPC trace compared with the CDF of the prefixes generated by MediSyn.

4.5 New File Introduction Process

The process of new file introduction mimics how files are introduced at a media site and attempts to answer the following questions:

- What is the new file arrival process on a daily time scale?
- What is the new file arrival process within an introduction day?

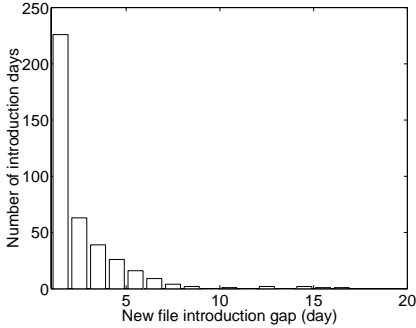


Figure 15: New file introduction gaps measured in days for the HPC trace.

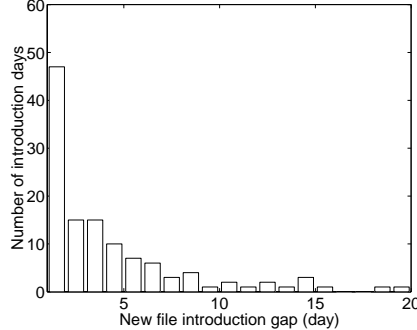


Figure 16: New file introduction gaps measured in days for the HPL trace.

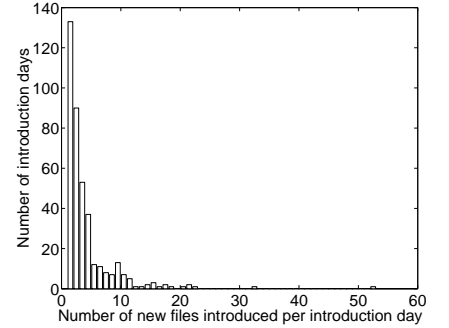


Figure 17: The number of new files introduced per introduction day.

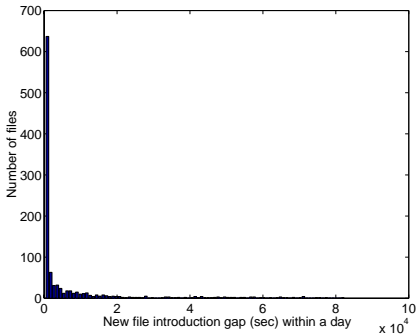


Figure 18: New file introduction time gaps within an introduction day.

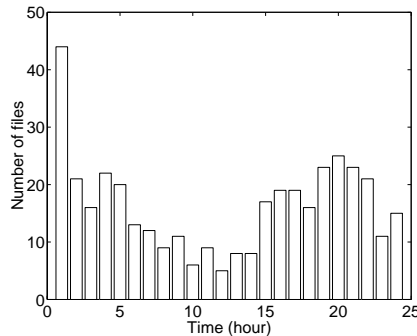


Figure 19: The start times of new file introduction within introduction days.

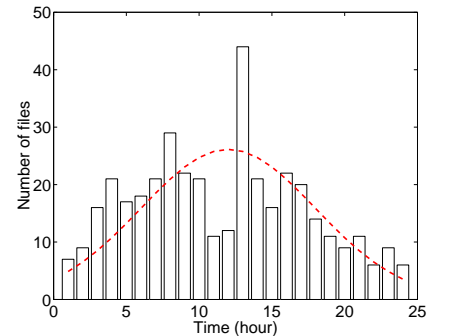


Figure 20: The rotated start times of new file introduction within introduction days.

To model new file arrival on a daily level, we capture the time gap measured in days between two introduction days and the number of new files introduced in each introduction day. Figure 15 shows the distribution of new file introduction gaps measured in days for the HPC trace. The distribution depicted in Figure 15 can be captured by a *Pareto* distribution with $\alpha = 2.0164$. Figure 16 shows the introduction gap distribution for the HPL trace, which can be captured by an *exponential* distribution with $\mu = 4.2705$.

MediSyn can generate new file introduction time gaps according to one of three possible distributions: 1) a *Pareto* distribution, 2) an *exponential* distribution, 3) a *fixed interval*. If users specify a Pareto distribution for the new file introduction process, files tend to be introduced into the system clustered over time. If the introduction process is specified by an exponential distribution, the new file arrival process is a Poisson arrival process, which means the interarrival times are independent. The fixed interval is used to model some artificial introduction process with regular patterns.

Since there may be multiple new files introduced in a given day, we must also model the number of files introduced per introduction day. Figure 17 shows the distribution for the number of files introduced in a given day for the HPC trace. The distribution can be fitted by a *Pareto* distribution with $\alpha = 1.1323$.

After determining the number of files introduced in a given day, MediSyn needs to model the new file arrival process within that day. We model this process by capturing the gap between two file arrivals. Figure 18 shows the time gaps for new files introduced within a day. Since the distribution is too sparse on time scale of seconds, we measure the time gaps at multiples of 900 seconds (15 minutes). The distribution can be captured by a *Pareto* distribution with $\alpha = 1.0073$.

Due to the properties of Pareto distribution, if we only model the time gap between two file arrivals and start to introduce new files from the beginning of a day, then most of the files will be introduced in the beginning of every day. So we also capture the start times of new file introduction process within every introduction day. Figure 19 shows this distribution. Since it looks like a rotated *normal* distribution with the peak at 0, we rotate the the distribution by 12 hours as shown in Figure 20. This is a *normal* distribution with mean 43200 seconds and standard deviation 21600 seconds.

4.6 Life Span

Since temporal correlation is observed in media workloads, an independent reference model combined with a global popularity distribution [9] is insufficient for a synthetic workload generator to generate a file access stream. SURGE [6] uses a stack distance model to generate web reference streams with reference locality. Both the independent reference model [9] and the stack distance model [6, 4] assume that each file’s popularity is stationary over the entire trace period and that each file is introduced at the start of the trace. Since we observe non-stationary popularity in streaming media workloads, such models are unsuitable for generating session arrivals in streaming media workloads.

To accurately model the non-stationarity of file popularity, we use the new file introduction process to mimic how media files are introduced at the media sites, as we described above. In addition, each file has its own life span, which characterizes its changing popularity after the file’s introduction. Thus, the global file popularity distribution, the file introduction process and life spans of individual files, all together capture the popularity change of media files over the entire trace.

We define the *relative access time* of a file as a random variable whose value is the time measured in days when the file is accessed by a client after the file is introduced. The distribution of a file’s *relative access times* describes the temporal correlation of all accesses to the file. We also call this distribution the life span distribution of the file. In our traces, we observe two types of life span distributions as illustrated in Figure 21 and Figure 22 respectively. Since most files in our traces have life spans similar to Figure 21, we call this type of life span a *regular life span*.

News-like streaming contents typically have life span distributions similar to Figure 22, where most accesses occur shortly after the file introduction and the access frequency diminishes relatively quickly. So we refer to this kind of life span as *news-like life span*.

We experimented with gamma, Pareto, exponential and lognormal distributions to fit the *relative access times* of our traces. Although gamma distributions can somehow capture both news-like and regular life spans, the combination of Pareto and lognormal distributions can better fit them. Thus, news-like life spans follow *Pareto* distributions, and regular life spans follow *lognormal* distributions.

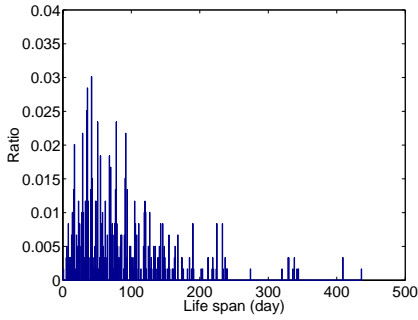


Figure 21: A *regular* lifespan.

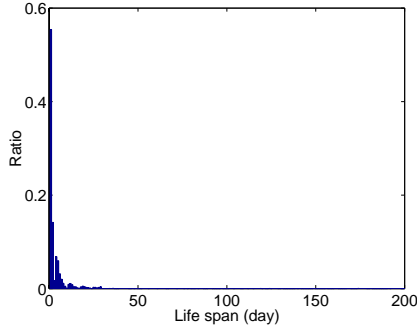


Figure 22: A *news-like* lifespan.

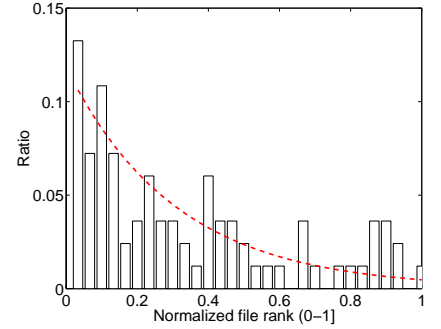


Figure 23: PDF that a file at a certain rank has a Pareto life span distribution.

To generate a sequence of regular life spans and news-like life spans, we need to model the distributions of the mean (μ) and the standard deviation (σ) for regular life spans, and the distributions of α for news-like life spans. Our analysis of the HPC and HPL traces shows that these parameters follow *normal* distributions. Table 8 shows the parameters for these normal distributions derived from the HPC log to capture the parameters of regular life spans (μ and σ) and news-like life spans (α).

Normal distribution parameters	lognormal μ	lognormal σ	Pareto α
μ	3.0935	1.1417	0.7023
σ	0.9612	0.3067	0.2092

Table 8: The parameters for the distributions (normal distributions) of the parameters in lognormal and pareto life span distributions.

There is a strong correlation between file popularity and life span shape. A file with a higher popularity rank tends to have a higher probability to have a *news-like life span*. Figure 23 shows the PDF for this probability. The distribution can be captured by an *exponential* distribution. File ranks have been transformed between 0 and 1 so that μ for the exponential distribution is independent of the number of media files generated. In the HPC trace, we observed 82 *news-like life spans* out of the 400 most popular files. Users of MediSyn can specify their own ratio according to the workload they want to generate. A workload including more files with *news-like life spans* has a more bursty access pattern.

4.7 Diurnal Access Pattern

After determining the life span and the global popularity of every file, MediSyn can generate the number of accesses for every day of a file’s life span. Distributing these accesses over a day is challenging because we wish to model both session interarrival time and diurnal access patterns.

Figure 24 shows a typical session interarrival time distribution for a file measured in a day. It is a heavy-tail distribution and can be fitted by a Pareto distribution better than an exponential distribution. However, if we generate all interarrival times within a day based on this Pareto distribution, it is difficult to simultaneously ensure diurnal pattern. Figure 25 shows the interarrival time distribution for the same file within one hour of the same day. This distribution is not a heavy-tail distribution and can be captured by an exponential distribution. Thus, if we can determine the number of accesses in each hour of a day according to a certain diurnal pattern, we can use an exponential distribution to generate the interarrival times of the accesses in this hour. Thus, we

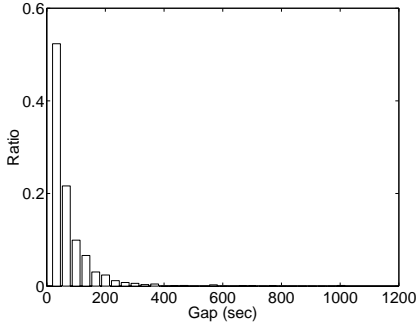


Figure 24: The PDF of session access interarrival time gaps for a file measured in a day.

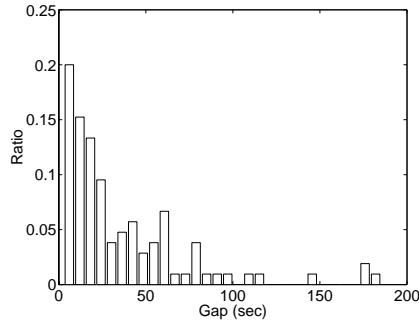


Figure 25: The PDF of session access interarrival time gaps for a file measured in an hour.

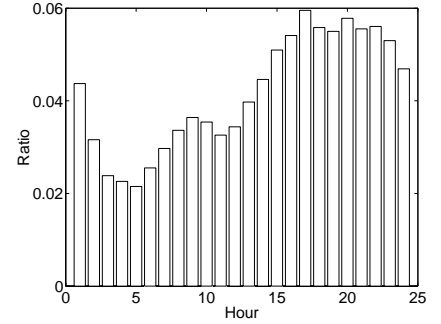


Figure 26: The session access diurnal pattern for the HPC trace. Each bin is an hour.

can both generate the diurnal pattern and satisfy the observed exponential distribution for interarrival times.

Diurnal access patterns are universally observed by other streaming workload analyses. But we do not explicitly find diurnal patterns for single files. We only observe an aggregate diurnal access pattern for all file accesses. Figure 26 shows the average ratios of accesses in each hour for all files in the HPC trace.

In MediSyn, a user can specify a global diurnal pattern like Figure 26, which contains a set of bins. Each bin specifies a time period and the ratio of accesses in this bin. Since we believe there is no temporal correlation among file accesses within a day (i.e., the temporal locality within a day is entirely determined by file popularities), we can make every file follow the diurnal pattern. Essentially, each file’s session arrival process within a given day is modeled as a *nonhomogeneous Poisson* process [22], where only the session arrivals within each bin can be modeled by a Poisson process. The session arrival rate of the file for a given bin is computed based on the diurnal pattern specified by the user and the number of accesses within a day determined by the file life span. MediSyn generates the interarrival time gaps within each bin and constructs a sequence of sessions for the file on the scale of seconds.

5 Related Work

Accurate workload characterization lays down a foundation for a successful synthesis of realistic workloads. A number of studies on multimedia workload analysis have been reported in literature [2, 3, 13, 12, 18, 21].

Acharya et al [2], presented the analysis of the six-month trace data from mMOD system (the multicast Media on Demand) which had a mix of educational and entertainment videos. They observed high temporal locality of accesses, the special client browsing pattern showing clients preference to preview the initial portion of the videos, and that rankings of video titles by popularity do not fit a Zipfian distribution.

Almeida et al [3] performed an analysis of two educational media server workloads. The authors provide a detailed study of client session arrival process: the client session arrivals in one workload can be characterized as a Poisson process, and the interarrival times in the second workload follow a heavy-tail Pareto distribution. They also observed that media delivered per session depends on the media file length.

The study by Chesire et al [13] analyzed the media proxy workload at a large university. The authors presented a detailed characterization of session duration (most of the media streams are less than 10 minutes), object popularity (78% of objects are accessed only once), sharing patterns of streaming media among the clients, and that popularity distribution follows a Zipf-like distribution (trace duration covers one week).

Two enterprise media server workloads have been extensively studied in [12]. The data was collected over significant period of time. Thus authors concentrated on the analysis of media server access trends, access locality, dynamics and evolution of the media workload over time. They reported non-Zipfian and non-stationary popularity of files observed in their data.

In our work, we attempt to summarize findings from the earlier work, and build a general, unified model for workload characteristics capturing unique properties of streaming media workloads as well as the dynamics in media workloads observed over long period of time.

Since HTTP requests and streaming media sessions are very different, streaming media workloads exhibit many new properties relative to traditional web workloads. Thus existing synthetic web workload generators [6] are not suitable for generating streaming media workloads.

The only synthetic workload generator for streaming media reported in literature is GISMO [17]. MediSyn adopts similar approach chosen in GISMO to organize the synthetic trace generation in two steps: *i)* defining the individual session characteristics, and *ii)* determining the media session arrival process. GISMO operates over a “fixed” set of media files already “introduced” at a media site, with the assumption that object popularity follows a Zipf-like distribution and remains the same over the entire duration of the experiment. Since we pursue the goal of developing a synthetic workload generator which reflects the dynamics and evolution of media workloads over time, we propose a set of new models to reflect these new temporal properties of streaming media workloads in MediSyn.

6 Conclusion and Future Work

Development of efficient resource allocation mechanisms for Internet hosting centers and CDNs, serving streaming media content, requires performing the experiments with realistic streaming media workloads which need to be scaled, parametrized, and mixed in a controllable and desirable way.

In this work, we present a synthetic streaming media workload generator, MediSyn, which is specially designed to accomplish this goal. In MediSyn, we develop a number of novel models to capture a set of characteristics critical to streaming media services, including file duration, file access prefix, non-stationary file popularity, new file introduction process, and diurnal access pattern. Among the primary features of our synthetic generator is the ability to reflect the dynamics and evolution of content at media sites and the change of access rate to the sites over time. We introduce a novel *generalized* Zipf-like distribution that captures recently-observed popularity of both web objects and streaming media not captured by existing Zipf-like distributions. Our evaluation, based on two long-term traces of streaming media services, demonstrates that MediSyn accurately captures the essential properties of media workloads, which are chosen to represent the unique (while generic) properties of streaming media workloads and their dynamics over time.

MediSyn implementation is based on a modular design allowing the particular system properties to be customized, enhanced or extended to reflect the requirements of individual scenarios. In a future work, we plan to extend MediSyn with implementation of client interactivities within media sessions.

As part of MediSyn, we plan to release a workload analysis tool reflecting the property profiles generated by MediSyn. These profiles can be conveniently used for tuning the workload generator parameters to specify the desired properties.

References

- [1] General Pareto Distribution. <http://www.math.uah.edu/stat/special/special12.html>.
- [2] Soam Acharya, Brian Smith, and Peter Parnes. Characterizing User Access to Videos on the World Wide Web. In *Proceedings of ACM/SPIE Multimedia Computing and Networking*, January 2000.
- [3] Jussara Almeida, Jeffrey Krueger, Derek Eager, and Mary Vernon. Analysis of Educational Media Server Workloads. In *Proceedings of NOSSDAV*, June 2001.
- [4] Virgílio Almeida, Azer Bestavros, Mark Crovella, and Adriana de Oliveira. Characterizing Reference Locality in the WWW. In *Proceedings of PDIS*, December 1996.
- [5] Virgilio Augusto Almeida, Marcio Cesirio, Rodrigo Fonseca, Wagner Meira Jr., and Cristina Murta. Analyzing the behavior of a proxy server in the light of regional and cultural issues. In *Proceedings of WCW*, June 1998.
- [6] Paul Barford and Mark Crovella. Generating Representative Web Workloads for Network and Server Performance Evaluation. In *Proceedings of SIGMETRICS*, June 1998.
- [7] Dave Bianchi. CNN.com: Facing A World Crisis. <http://www.tcsa.org/lisa2001/cnn.txt>.
- [8] Rebecca Braynard, Dejan Kostić, Adolfo Rodriguez, Jeffrey Chase, and Amin Vahdat. Opus: an Overlay Peer Utility Service. In *Proceedings of the 5th International Conference on Open Architectures and Network Programming (OPENARCH)*, June 2002.
- [9] Lee Breslau, Pei Cao, Li Fan, Graham Phillips, and Scott Shenker. Web Caching and Zipf-like Distributions: Evidence, and Implications. In *Proceedings of INFOCOM*, March 1999.
- [10] Jeffrey Chase, Darrell Anderson, Prachi Thakar, Amin Vahdat, and Ronald Doyle. Managing Energy and Server Resources in Hosting Centers. In *Proceedings of the 18th ACM SOSR*, October 2001.
- [11] Ludmila Cherkasova and Gianfranco Ciardo. Characterizing Temporal Locality and its Impact on Web Server Performance. In *Proceedings of ICCCN*, October 2000.
- [12] Ludmila Cherkasova and Minaxi Gupta. Characterizing Locality, Evolution, and Life Span of Accesses in Enterprise Media Server Workloads. In *Proceedings of NOSSDAV*, May 2002.
- [13] Maureen Chesire, Alec Wolman, Geoffrey Voelker, and Henry Levy. Measurement and Analysis of a Streaming-Media Workload. In *Proceedings of USITS*, March 2001.
- [14] Morris DeGroot and Mark Schervish. *Probability and Statistics*. Addison-Wesley, 3rd edition, 2002.
- [15] Raj Jain. *The art of computer systems performance analysis: technique for experimental design, measurement, simulation and modeling*. John Wiley & Sons, 1992.
- [16] Shudong Jin and Azer Bestavros. Temporal Locality in Web Requests Streams: Sources, Characteristics, and Caching Implications. Technical Report BUCS-TR-1999-009, Department of Computer Science, Boston University, August 1999.
- [17] Shudong Jin and Azer Bestavros. GISMO: A Generator of Internet Streaming Media Objects and Workloads. Technical Report BUCS-TR-2001-020, Department of Computer Science, Boston University, October 2001.
- [18] Dario Luperello, Sarit Mukherjee, and Sanjoy Paul. Streaming Media Traffic: an Empirical Study. In *Proceedings of Web Caching Workshop*, June 2002.

- [19] Real Networks. Real Producer 8 plus. <http://www.realnetworks.com>.
- [20] Hewlett Packard. Utility Data Center. <http://www.hp.com/go/hpudc>.
- [21] Jitendra Padhye and Jim Kurose. An Empirical Study of Client Interactions with a Continuous-Media Courseware Server. In *Proceedings of NOSSDAV*, June 1998.
- [22] Sheldon Ross. *Introduction to probability models*. Academic Press, 1997.
- [23] Subhabrata Sen, Jennifer Rexford, and Don Towsley. Proxy Prefix Caching for Multimedia Streams. In *Proceedings of INFOCOM*, March 1999.