

A Hierarchical and Multiscale Approach to Analyze E-Business Workloads[★]

D. A. Menascé^a V. A. F. Almeida^b R. Riedi^c F. Ribeiro^b
R. Fonseca^b W. Meira Jr.^b

^a*Dept. of Computer Science, MS 4A5, George Mason University, Fairfax, VA
22030, USA.*

^b*Dept. of Computer Science, Federal University of Minas Gerais, Belo Horizonte,
MG 31270, Brazil.*

^c*Dept. of Electrical and Comp. Engineering, Rice University, Houston TX 77251,
USA.*

Abstract

Understanding the characteristics of Electronic business (E-business) workloads is a crucial step to improve the quality of service offered to customers in E-business environments. This paper proposes a hierarchical and multiple time scale approach to characterize E-business workloads. The three levels of the hierarchy are user, application, and protocol, and are associated with customer sessions, functions requested, and HTTP requests, respectively. Within each layer, an analysis across several time scales is conducted. The approach is illustrated by presenting a detailed characterization of two actual E-business sites: an online bookstore and an electronic auction site. Our analysis of the workloads showed that the session length, measured in number of requests to execute E-business functions, is heavy-tailed, especially for sites subject to requests generated by robots. An overwhelming majority of the sessions consist of only a handful requests, which seems to suggest that most customers are human (as opposed to robots). A significant fraction of the functions requested by customers were found to be product selection functions as opposed to product ordering. An analysis of the popularity of search terms revealed that it follows a Zipf distribution. However, Zipf's law as applied to E-business is time scale dependent due to the shift in popularity of search terms. We also found that requests to execute frequent E-business functions exhibit a pattern similar to the HTTP request arrival process. Finally, we demonstrated that there is a strong correlation in the arrival process at the HTTP request level. These correlations are particularly stronger at intermediate time scales of a few minutes.

Key words: E-business, WWW, workload characterization, performance modeling, heavy-tailed distribution

1 Introduction

E-business sites are very complex, composed of several tiers of servers of different types (e.g., Web servers, application servers, and database servers), and are subject to workloads that vary in ways hard to predict. The quality of service requirements for E-business sites are strict since customers demand fast response times and high availability or else they turn to competitors. Understanding the nature and characteristics of E-business workloads is a crucial step to improve the quality of service offered to customers in electronic business environments. E-business workload characterization can lead to a better understanding of the interaction between customers and Web sites and can also help design systems with better performance and availability. This paper presents a hierarchical and multiple scale approach to the characterization of E-business workloads.

E-business workloads are composed of sessions. A *session* is a sequence of requests of different types made by a single customer during a single visit to a site. During a session, a customer requests the execution of various E-business functions such as browse, search, select, add to the shopping cart, register, and pay. A request to execute an E-business function may generate many HTTP requests to the site. For example, several images may have to be retrieved to display the page that contains the results of the execution of an E-business function.

Past studies of WWW workloads concentrated on information provider sites and found several characteristics common to them [5,7,11,20]. Some of these characteristics deal with file size distributions, file popularity distribution, self-similarity in Web traffic, reference locality, and user request patterns. A number of studies of different Web sites found file sizes to exhibit heavy-tailed distributions and object popularity to be Zipf-like. Other studies of different Web site environments demonstrated long-range dependencies in the user request process, in other words, strong correlations in the user requests. In particular, [7] identified ten workload properties, called *invariants*, across six different data sets, which included different types of information provider Web sites. Some of the most relevant invariants are: i) images and HTML files account for 90-100% of the files transferred; ii) 10% of the documents account for 90% of all requests and bytes transferred; iii) file sizes follow the Pareto distribution, and iv) file inter-reference times are independent and exponentially distributed. Shortly after, [5] discovered that the popularity of documents served by Web sites dedicated to information dissemination follows a Zipf's law. In [11], the authors pointed to the self-similar nature of Web server traffic. All these studies were performed almost five years ago. Since

* This is an expanded and revised version of [16].

then, several major changes have been observed in the WWW. The most important are: clients now have much larger bandwidth, the number of users has grown exponentially, and E-business became one of the major applications on the Web.

In [14], the authors introduce the notion of session, consisting of many individual HTTP requests. However, they do not characterize the workload of E-business sites, which is composed of typical requests such as browse, search, select, add, and pay. The analysis focuses only on the throughput gains obtained by an admission control mechanism that aims at guaranteeing the completion of any accepted session. The work in [21] proposes a workload characterization for E-business servers, where customers follow typical sequences of URLs as they move towards the completion of transactions. The authors though do not present any characterization or properties of actual E-business workloads.

There are very few published studies [6,16,19] of E-business workloads because of the difficulty in obtaining actual logs from electronic companies. Most companies consider Web logs to be very sensitive data. In [19], the authors propose a graph-based methodology for characterizing E-business workloads and apply it to an actual workload to obtain metrics related to the interaction of customers with a site. For example, the paper shows how to obtain information such as the number of sessions, average session length, and buy-to-visit ratio. Reference [17] presents several models (e.g., customer behavior model graph and customer visit model) for workload characterization of E-business sites. It also shows how workload models can be obtained from HTTP logs. Our previous work [16], extended here, discussed the issue of how to obtain invariants for E-business workloads. In [6], Arlitt et al. characterize the workload of an actual e-commerce site for the purpose of analyzing its scalability. They use performance-related criteria to cluster requests into similar groups. They then use multiclass queuing models to carry out a capacity planning study for the site. In [3], the authors study the impact of time scale on operational analysis for a large Web-based shopping system. They show that time-related service level agreements and input parameters for predictive queuing models are sensitive to time scale.

A question that naturally arises is: are the characteristics and invariants found in information provider Web sites still valid for E-business workloads? To answer this question, we define a hierarchical and multiscale approach to characterize the workload of E-business sites. The three layers of the hierarchy are: session, function, and HTTP request, as defined in Section 2. Within each layer, an analysis across several time scales is conducted. The approach is illustrated by presenting a detailed characterization of two actual E-business sites: an online bookstore and an electronic auction site. This paper extends our previous work [16] and examines statistical and distributional properties of the E-business workloads and compare these properties across the two

datasets. As much as possible, we compare the features of these workloads with the invariants that were discovered for information dissemination Web sites and provide an extended multiscale analysis of the workload. The same hierarchical approach was used by the authors to study the presence of robots in Web workloads [4].

The rest of the paper is organized as follows. Section two shows the approach used to characterize E-business workloads. The next section describes the data collection process. Section four analyzes two logs from actual E-business sites and characterizes the workload at the HTTP request level. Characterizations at the E-business function and session levels are provided in sections five and six, respectively. Finally, section seven presents concluding remarks.

2 Hierarchical Multiscale Approach

Workload characterization can be accomplished at many levels: user level, application level, protocol level, and network level. An E-business workload can be viewed in a multi-layer hierarchical way, as shown in Fig. 1. This paper focuses on the characterization of three levels, represented by the HTTP request layer (protocol level), function layer (application level), and session layer (user level). This hierarchy can be used to capture changes in user behavior and map the effects of these changes to the lower layers of the model.

Our approach is to analyze each layer individually in order to obtain a characterization of the arrival process and usage statistics. We perform multi-scale statistical analysis, study long range dependence (LRD), and burstiness. Our analysis covers properties such as: session inter-arrival times, inter-arrival times for specific E-business function requests, search term popularity distribution, session length distribution, E-business function distribution per session, and number of active sessions and initiated sessions.

More specifically, our approach can be summarized in the following steps for

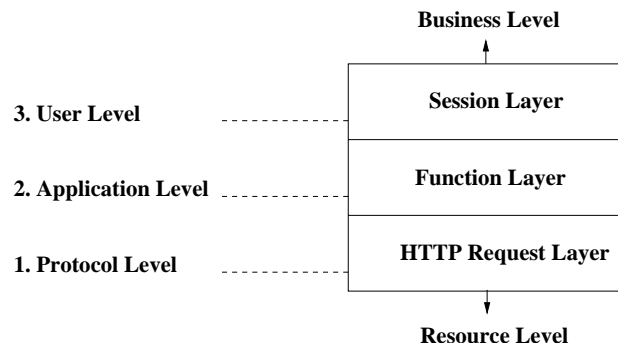


Fig. 1. A hierarchical workload model.

each level:

- *HTTP request level characterization (see Section 4).*
 - (1) Log preparation: Collect and merge by timestamp the HTTP logs of all Web servers of a site to produce a single log denoted by \mathcal{L} .
 - (2) Visual inspection: Plot the number of HTTP requests in the log \mathcal{L} at different time scales (e.g., one hour, five minutes, and five seconds) in order to conduct a visual inspection of the arrival process. Strong dependencies are characterized by long sequences of increase or decrease of traffic intensity at intermediate time scales (e.g., five minutes).
 - (3) Multiscale quantification of the dependence of the arrival process: Draw a variance time plot (VTP) for different time scales. The VTP plot is a log-log plot of the sample variance against the time scale (see Section 4 for a more formal definition of the VTP plot) and can be used to detect and quantify self-similarity and to compute the Hurst parameter.
 - (4) Multiscale next-neighbor dependence analysis: Draw neighbor-to-neighbor (N2N) plots of the number of arrivals in consecutive time slots for different time scales. Clustering of data along the diagonal indicates higher next-neighbor correlation.
 - (5) Inter-arrival time (IAT) analysis: Draw an IAT graph that plots a spike for each request in \mathcal{L} . The height of the spike is proportional to the time interval between a request and its predecessor. High peaks indicate extreme interarrival times, which results in light server load or even in server idleness.
- *Function level characterization (see Section 5)*
 - (1) Log preparation: Generate a function log \mathcal{L}_f by removing all the entries in \mathcal{L} that correspond to HTTP requests to image files or to errors and by converting the remaining entries to pairs of the type (ts, f) where ts is the timestamp of the log entry and f is the E-business function requested by the HTTP request. A lookup table that maps URLs or URL prefixes to E-business functions is generally needed to determine f .
 - (2) Function frequency determination: Determine the frequency of execution of each E-business function.
 - (3) Multiscale analysis of E-business function execution requests: Plot for each type of E-business function f the number of requests to execute f at different time scales. Patterns for finer time scales (one hour or less) should be inspected more closely. Patterns for frequently executed and non-frequently executed functions should be compared with that of the HTTP request arrival process at similar time scales. If the site is also used by robots, the analysis may need to take into account functions more likely to be requested by robots (e.g., search).
 - (4) Multiscale analysis of the popularity of search terms: One of the most popular functions of E-business sites is the search function. It is useful to understand which search terms are more popular and if there is a Zipf's law [5,7] relationship between the relative frequency $f(r)$ of occurrence of

terms in search requests and the rank r or popularity of the term. Draw, for various time scales, a plot of $\log f(r)$ vs. $\log r$. If the graph is a straight line with a negative slope close to -1, the relationship between $f(r)$ and r follows a Zipf's law. This indicates that there may be performance gains in caching query results for the most popular terms. A comparison of the popularity plots for different time scales may reveal changes in the popularity of products and services over time.

- *Session level characterization (see Section 6)*
 - (1) Log preparation: Generate a temporary log \mathcal{L}_t by removing all the entries in \mathcal{L} that correspond to HTTP requests to image files and errors. Convert the remaining entries in \mathcal{L}_t to tuples of the type (uid, ts, f) where uid is a user identification (determined by either IP address, cookie, or any other method), ts is the timestamp of the log entry, and f is the E-business function requested by the HTTP request. Generate a session log \mathcal{L}_s from \mathcal{L}_t by putting together all tuples for the same uid in increasing order of timestamp and delimited by session boundaries—generated by either implicit login/logout requests or periods of inactivity thresholds.
 - (2) Session length analysis: Draw a log-log graph of the tail of the session length distribution, defined as the number of E-business functions requested per session. A straight line indicates a heavy-tailed distribution, which may be caused by strong robot activity in some cases [4].
 - (3) Multiscale session initiation analysis: Plot a graph of the number of sessions initiated per time unit, for various time scales. This analysis is important because site resources are usually allocated on a session basis.

We applied this approach to two actual e-businesses: an online bookstore and an online auction site. Next section provides more detailed information about the data used to characterize the workload of each of these sites.

3 Data Collection for Case Studies

The online bookstore sells exclusively on the Internet. The auction site sells Internet domains. In both cases, the data consist of access logs recorded by the WWW server of each E-business.

The data comprises two weeks of accesses to each of these sites. The bookstore logs were collected from August 1st to August 15th, 1999, while the auction server logs are from March, 28th to April 11th, 2000.

During these two weeks, the bookstore handled 3,630,964 requests (242,064 daily requests on average), transferring a total of 13,711 megabytes of data (914 MB/day on average). The auction server has a smaller load, and answered 466,058 requests (31,071 requests/day) which amounts to 1,863 megabytes of

data (124 MB/day). Most of these requests are for embedded images in the response pages. In the case of the bookstore, images account for 71% of the requests, while in the auction workload they represent 85.3% of the requests.

E-business function-related requests amounted to 26.3% and 14.7% of the requests received by the bookstore and auction sites, respectively. The difference in percentage between the two sites is explained by the larger number of images used by the auction site. Thus, the bookstore executed 63,711 E-business functions per day, and each service response had 12,618 bytes, on average. We should note that service-related requests are responsible for most of the network traffic, comprising 84.6% of the data sent by the bookstore server and 92.2% of the data sent by the auction server. This is explained through the fact that most of the image files embedded in pages are usually already cached and are not transmitted back to the client. Although the auction pages contain more images than the bookstore pages, the auction site employs a smaller array of images, typically banner advertisements and page layout. The bookstore uses a larger number of different images, such as book covers, and can therefore benefit less from the advantages of caching.

4 Request-layer Characterization

In this section, we study the statistical nature of the arrival process of HTTP requests to allow for the extraction of statistically significant features towards classification, understanding, and modeling of request workload.

4.1 *Dependence and Prediction*

It is now a well accepted fact that strong correlations are present in various aspects of the World Wide Web, from request arrivals on servers to packet arrivals on the network. These correlations express considerable dependencies that lead to “burstiness” or high variability and may degrade performance and throughput if not accounted for. We carry out a statistical analysis across various time scales to detect correlations and assess their strength.

The fact that statistical analysis and modeling has to incorporate different methods according to time scale is most apparent as we attempt to accommodate various trends. On the largest time scale of days (in our study), the weekend produces somewhat less volume, while on the scale of hours the presence of a periodic sleep-wake pattern per day is visually obvious. It is not our intention to explore these patterns. On finer time scales, structure is much less obvious and it is our goal to present a simple analytical tool that distinguishes

scales of “noisy (non-predictable) oscillations” from scales with strong correlations (that foster prediction). Thereby, care is needed to avoid bias from large scale trends and non-stationarities that could manifest in both, small scale analysis and prediction.

A visual inspection of the number of requests arriving at the bookstore on different time scales, i.e., in time intervals of varying length (see Fig. 2) reveals, even to the inexperienced eye, an apparent strong dependence that shows long sequences of increase or decrease of volume (trends), particularly pointed at intermediate time scales. The purpose of our analysis is to decide whether these trends are purely due to changes in traffic volume during the day and week or whether there is predictable behavior beyond these cycles. It is important to be able to detect strong dependencies since they degrade estimation by increasing the variance of the estimation error. On the positive side, by detecting strong dependencies one can foresee not only mean behavior but also temporary phases of increase or decrease in volume and variability in workloads leading to a more accurate assessment of performance.

4.2 Overview of Findings

Before going into details we summarize our findings at a high level. We find reliable estimates of correlation and, thus, dependence in stationary periods that typically range from noon to evening on each of the fifteen days contained in the traces. The degree of dependence, measured by the LRD parameter H , amounts to $H = .73$, a value that is quite common in natural phenomena and that indicates high correlation. Furthermore, the dependence appears strongest in intermediate time scales of the order of minutes. An analysis using neighbor-to-neighbor plots (see Fig. 4) confirms this finding of the scaling analysis via the VTP. This dependence leads to superior predictability which is able to forecast not only the mean of future workloads but also its trend, i.e., whether it is increasing or decreasing (see Fig. 5)

A possible explanation for particularly strong correlations on the time scale of few hundred seconds may be human think time and human distractedness. The overall self-similarity, at least “asymptotically,” may be argued for by invoking the well-known on-off process that was crucial in explaining self-similarity in network traffic loads [13]: The number of requests per session follows a heavy-tailed distribution. Since the number of requests sent per time unit is limited, sessions are thus sending requests over on-times that are heavy-tailed. Numerical support for this claim comes from our analysis of session duration in Sec. 6, which shows that the distribution of session length follows a power law. The on-off model then relates the exponent of this heavy-tailed distribution directly with H .

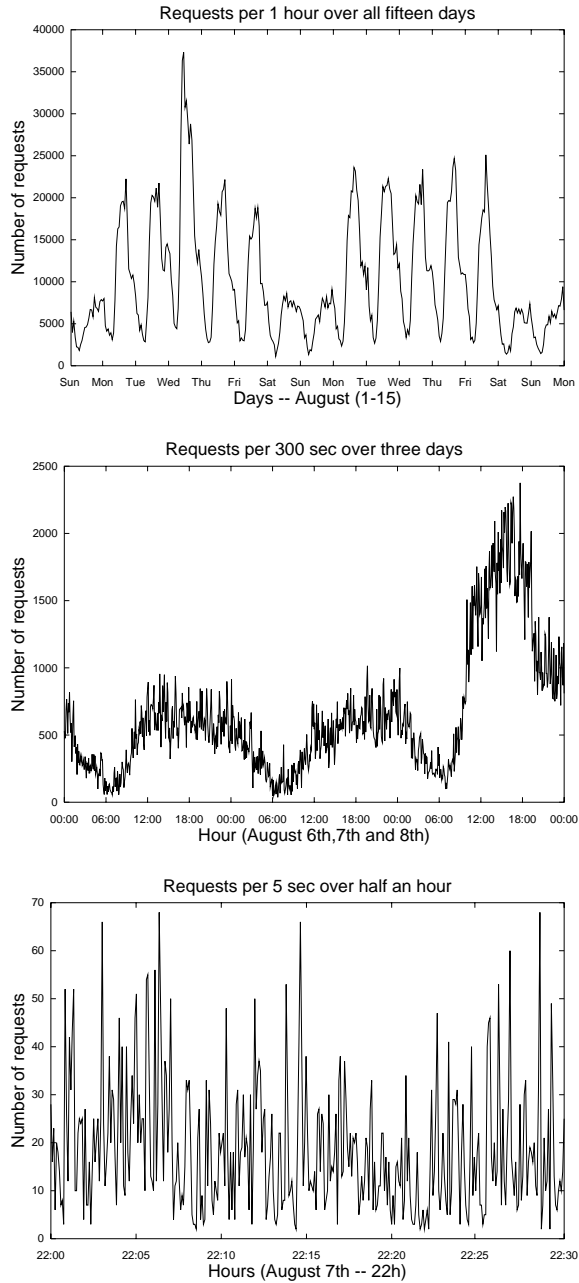


Fig. 2. Total number of requests arriving at the bookstore site at various time scales. Top: Arrivals per hour over the full fifteen days captured; Center: at a resolution of 300 sec over three days; Bottom: at full resolution (5 sec), over half an hour.

4.3 Detailed Analysis

A simple quantification of dependence over various scales is achieved by computing the sample variance by time scale: if arrivals occur independently of each other, the sample variance doubles if the length of the interval doubles. The variance exceeds twice the variance of the original interval if the arrivals

are positively correlated and will not reach twice the variance if the arrivals are negatively correlated. Indeed, $\text{var}(A + B) = \text{var}(A) + 2\text{cov}(A, B) + \text{var}(B)$. More specifically, if X_k denotes the number of arrivals in time interval $[k\delta, (k+1)\delta]$, where δ is the finest time resolution one is interested in, then $X_k^{(n)} = 2^{-n} \times (X_{k2^n} + X_{k2^n+1} + \dots + X_{k2^n+2^n-1})$ averages the arrivals in $[k2^n\delta, (k+1)2^n\delta]$ and can be computed efficiently through the recursion $X_k^{(n)} = (X_{2k}^{(n-1)} + X_{2k+1}^{(n-1)})/2$. The log-log plot of the variance against scale, i.e., $\log_2 \text{var} X^{(n)}$ versus n , is called variance time plot (VTP). This plot has the slope -1 for independent data (recall the normalization factor $1/2$ necessary to provide averages instead of total counts) and a different behavior for dependent data: The slower the VTP decays at a certain scale, the stronger the next-neighbor correlation within that scale.

The extreme case of positive correlation is a constant series X_k with a flat (horizontal) VTP. A more interesting case of dependent behavior constitutes the so-called “statistical self-similarity,” which is defined by the requirement that $\text{var} X^{(n)} = \sigma^2 2^{n(2H-2)}$. Here H denotes the Hurst parameter and lies between 0 and 1. This case is of interest due to the existence of appealing, simple, Gaussian processes with such properties, such as the fractional Gaussian noise and the auto-regressive FARIMA processes [23]. For $H = 1/2$ we find ourselves back in the case of independent data where $\text{var} X^{(n)} = \sigma^2 2^{-n}$ for all time scales n . On the other hand, if the VTP decays at a slower rate, i.e., with slope $2H - 2$ where $H > 1/2$, then we have positive correlations.

The VTP is a crude measure of the correlation structure with known bias and poor performance as an estimator of the LRD parameter H and is particularly sensitive to non-stationarities such as changing mean. However, when properly applied, the VTP is completely valid as a tool for a first look (see [1,2,23] and references therein).

The VTP plot of the number of arrivals at the online bookstore (see Fig. 3) shows a decay of particular strength corresponding to $H = 0.98$ at intermediate time scales from 80 to 5120 sec, corresponding to aggregation 4-10 in Fig. 3 (a) (there δ corresponds to 5 sec). Due to the presence of large scale trends (or non-stationarity) this number has to be considered with caution since the estimate could be highly biased. Indeed, the scaling is not optimal, and it is wise to perform local scaling tests over regions where the data shows stationarity. Indeed, over periods of several hours (see Fig. 3 (b) for an analysis of twelve hours in the afternoon of the sixth day) the VTP becomes more straight and the measured Hurst exponent falls into the region generally observed in natural phenomena (.7 to .85). The scaling we found—typically noon to evening—over the fifteen days averaged to about .73. Also, in this local analysis, the dependence seems to be strongest at intermediate time scales.

This strong dependence on intermediate scales is further confirmed by the

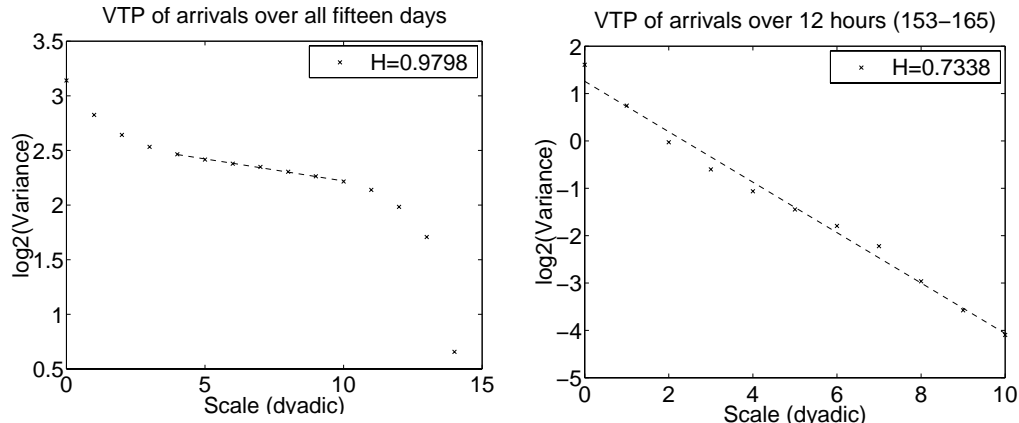


Fig. 3. Variance time plot for the online bookstore. Left (a): over the full fifteen days ($H = 0.9798$; this value is to be taken with caution since it is most likely affected by non-stationarity in the data). Right (b): over twelve hours (from 9AM to 9PM on Friday, August 6, of approximate stationarity of the mean arrival rate ($H = 0.7338$)).

“neighbor-to-neighbor” plots of the arrival process on various time scales: In Fig. 4 we display graphs of $X_{k-1}^{(n)}$ versus $X_k^{(n)}$, for three fixed values of n , where $X_k^{(n)}$ denotes the total number of requests arriving at the online bookstore in the time interval $[k2^n \delta, (k+1)2^n \delta]$, δ being 5 seconds in our data. These plots give an idea of the next-neighbor dependence on the time scales of Fig. 3. Note that the more the data is clustered along the diagonal, the higher is the predictability: large values are most likely followed by large values, small values by small values. For illustration purposes, we also show in Fig. 4 the “correlation” plot of a series of independent random variables. In this case, no structure and no clear clustering is visible.

On an intermediate time scale (Fig. 4 (b)) we find the closest clustering along the diagonal while there is a clear spread on the fine scale (Fig. 4 (a)). On the coarsest scale (Fig. 4 (c)), the data still displays dependency, though not as pronounced as on the intermediate scales. This indicates superior predictability on intermediate scales from many minutes to several hours.

The difficulty in interpreting these “neighbor-to-neighbor” plots resides in the presence of the predominant cyclic trends on the largest time scales. A critical observer could rightfully claim that the concentration along the diagonal is purely caused by these cycles, meaning that the data could be well approximated by a quasi-periodic (cyclic) mean superimposed with independent random fluctuations. Whether such an interpretation is valid cannot be decided from Fig. 4 because time information is lost: the plot does not indicate how X_n relates to X_{n-2} or any data point more than 2 steps in the past and, thus, it provides no insight with regard to stationarity.

In order to clarify this issue, a wavelet analysis could be beneficial, as wavelets allow an analysis insensitive to trends (due to vanishing moments [12]) and

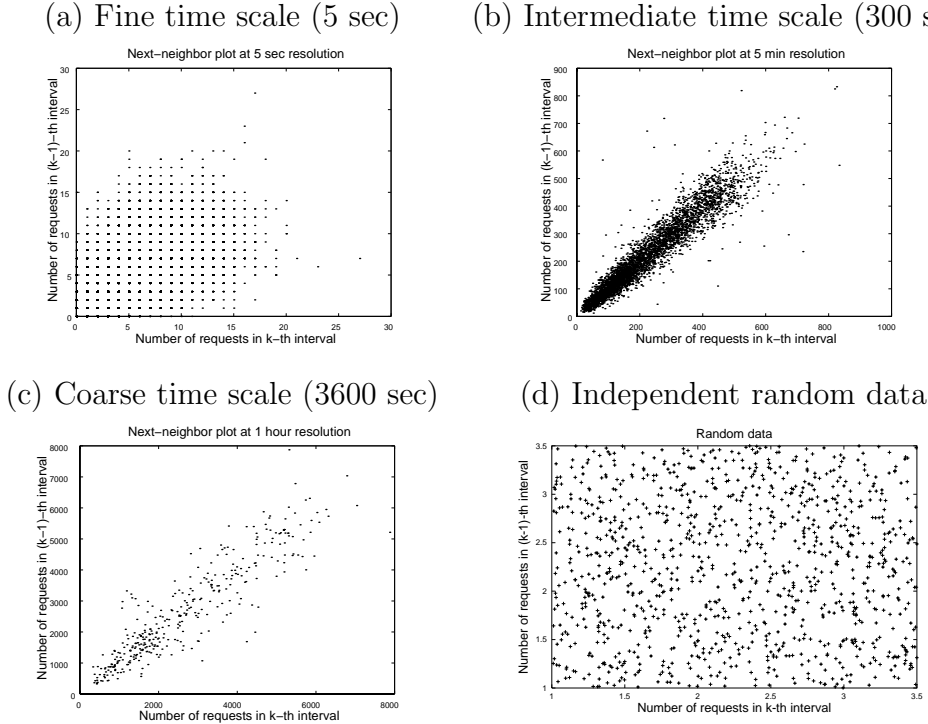


Fig. 4. Neighbor-to-neighbor plots for the number of requests arriving at the online bookstore.

provide an un-biased estimation of the Hurst parameter H [1,2]. Such an approach, however, is beyond the scope of this paper and we favor a more direct and simple approach. To test for predictability and at the same time remove bias from the changing arrival rate we study

$$Z_k^{(n)} = X_k^{(n)} - (1/8) * (X_{k-8}^{(n)} + \dots + X_{k-1}^{(n)})$$

which is, in fact, the difference between the current number of arrivals and the average of the last eight arrivals at time scale $2^n * \delta$. The choice of averaging eight is arbitrary. As a matter of fact, an auto-regressive model with better adapted coefficients for X_{k-m}, \dots, X_{k-1} (as to match the auto-correlation structure of X) would provide yet more accurate predictions. We display the next-neighbor correlation at three fixed values of n in Fig. 5, i.e. $Z_{k-1}^{(n)}$ versus $Z_k^{(n)}$, and again, we note the presence of correlations at intermediate time scales of hundreds of seconds, indicating that an increase in volume against past average volume is likely to be followed by yet another increase.

Having studied the data from the online bookstore in detail, let us now compare our two data sets. Figure 6 presents the average number of requests per day in our two-week logs. The figure clearly displays the traffic reduction during the weekends. We can confirm this behavior at the time scale of one hour by checking the graph in Fig. 7, where each point represents the hourly request average. As we can see, there are fourteen peaks, almost one per log day.

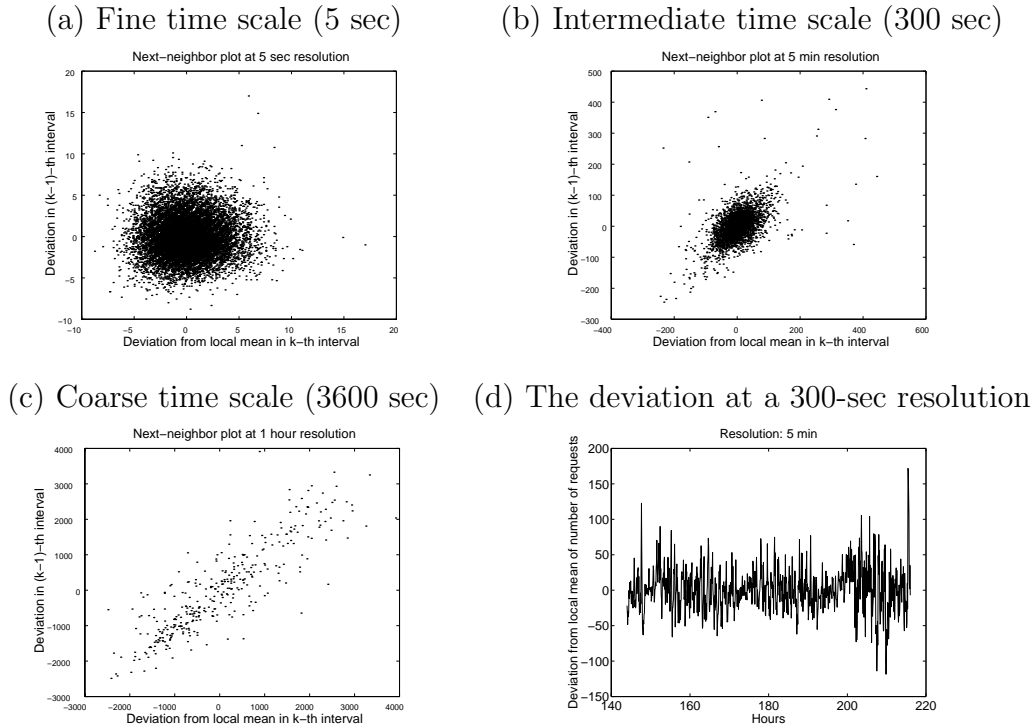


Fig. 5. Neighbor-to-neighbor plots for the deviation from the local mean of the number of requests arriving at the online bookstore.

Fig. 8 plots the inter-arrival times (IAT) graph for the bookstore and auction sites. The lighter weekend traffic may also be observed in the graph by looking at the highest peaks.

5 Function Characterization

In this section, we characterize the workload at the level of E-business functions. Our first criterion is the nature of the function. When considering an

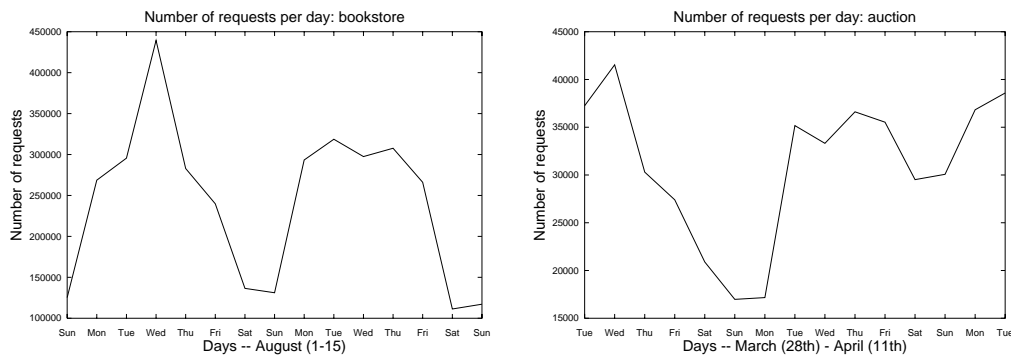


Fig. 6. Number of requests per day for the bookstore and auction sites

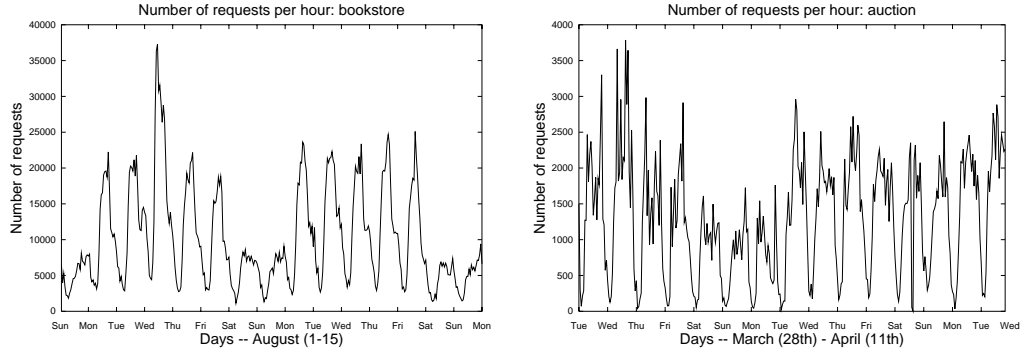


Fig. 7. Number of requests per hour for the bookstore and auction sites

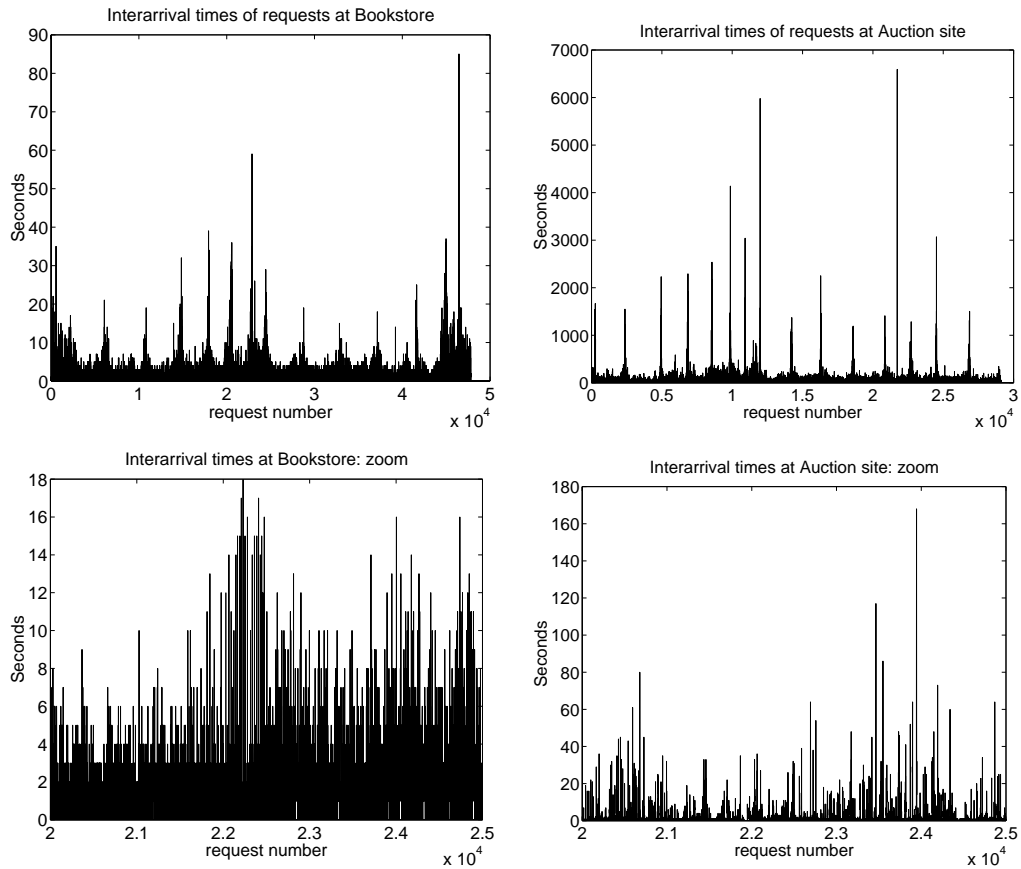


Fig. 8. Inter-arrival times (IAT) of requests at bookstore (left) and auction site (right) for the full trace (top) and a selection of 5,000 requests (bottom). The IATs at the auction site are more ‘explosive’ indicating sudden unexpected long periods of silence while the bookstore seems to see long periods of silence which are less abrupt and which could be anticipated.

online store, we may divide the functions into four groups: static, product selection, purchase, and other. Static functions comprise the home and informational pages about the store. Product selection includes all functions that allow a client to find and verify a product they are looking for: browse, search,

and view. Purchase functions indicate a desire to buy, either by selecting a product for later acquisition (e.g., add to cart) or by ordering it (e.g., pay). One interesting invariant in the logs we analyzed is that more than 70% of the functions performed are product selection functions. Table 1 presents a distribution of E-business function requests for both sites.

In the auction site, there are functions that relate to the process of posting items for sale. Similarly to the bookstore, though not as large in percentage, the majority of requests at the auction site concerns selection of products. On both sites, the functions directly related to spending money have a very low frequency.

When we split requests according to the E-business functions they invoke, i.e., search, browse, add, and pay, we find two clearly distinct classes. While the behavior on large time scales of hours and days of all functions follow the already observed human behavior, their small scale behavior is quite different. For example, Fig. 9 shows the number of requests per hour to execute searches at the bookstore and to retrieve the home page of the auction site for several days. If we compare Fig. 9 to Fig. 7, we see a similar pattern. This indicates that requests to execute frequent E-business functions exhibit a similar pattern of behavior as observed for the total number of HTTP requests.

Table 1
Distribution of E-business functions.

Bookstore		Auction	
Function	Frequency	Function	Frequency
Home	11.92%	Home	20.70%
Browse	17.72%	Browse	14.66%
Search	36.30%	Search	16.74%
View	19.99%	View	4.87%
Add	5.44%	Bid	0.08%
Pay	1.19%	Sell	7.99%
Account	2.44%	Account	5.99%
Robot	0.04%	Robot	0.06%
Info	3.66%	Info	9.44%
Other	1.31%	Other	2.31%
		Auth	9.18%
		Register	7.29%
		Admin	0.71%

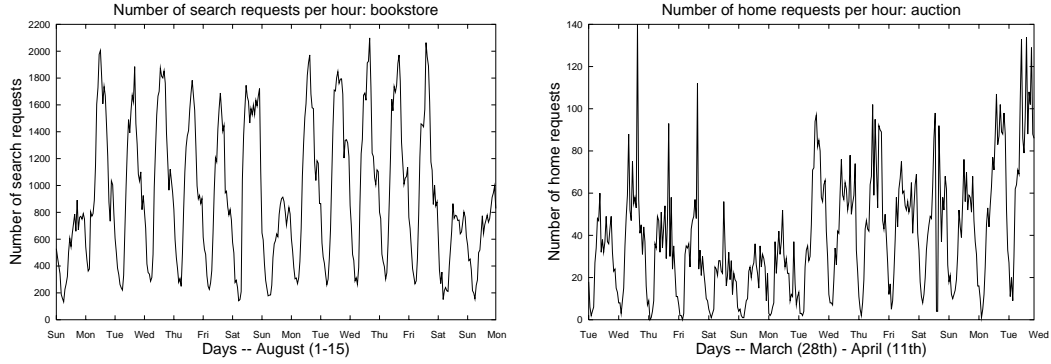


Fig. 9. Number of arriving requests per hour to execute frequent E-business functions. Bookstore: search; Auction site: home page.

The same is not true for less frequent functions such as pay and view, as indicated in Fig. 10. This figure shows clear bursts and a very different behavior from Fig. 7. Here, a more advanced statistical analysis revealing the multifractal scaling would be in place [22] and prediction is harder. In contrast, the more frequent functions such as “search” and “home” show statistics similar to the overall load of requests and are—apart from the cyclic trends—well described by Gaussian LRD processes.

This difference in small scale behavior is similar to the one we saw in the IAT process at the bookstore and the auction site (see Fig. 8). It is best understood when thinking in terms of doubly stochastic Poisson processes where Poisson arrivals are driven by a *varying intensity* which is itself random. As intensities are low, the spikyness of Poisson arrivals are apparent; as intensities grow, the Poisson distributions are well approximated by the Gaussian. In a unifying approach one would aim at measuring the “hidden” intensity, thus capturing the driving stochastics of request arrivals and allowing for a deeper understanding and more control. This is left for future investigation.

Figure 11 shows the number of search requests for the bookstore on a daily

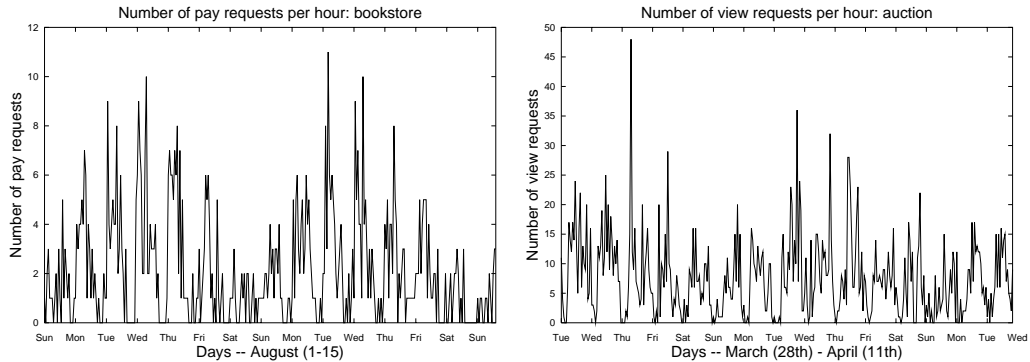


Fig. 10. Number of arriving requests to execute less frequent E-business functions, The time resolution is similar to that of Fig. 7. Bookstore: pay; Auction site: view.

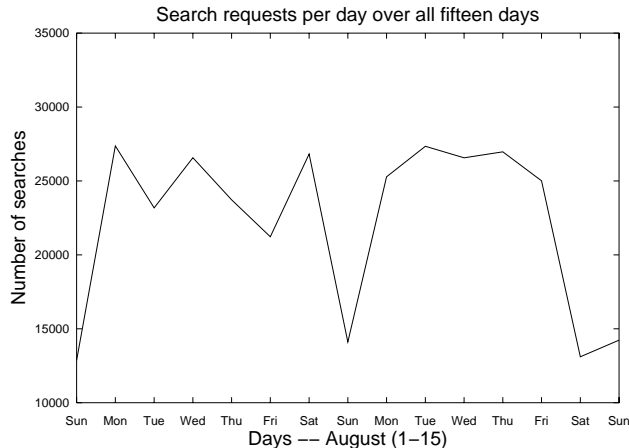


Fig. 11. Number of searches arriving per day at the bookstore over two weeks.

basis. We can see that this graph exhibits, in the first week, a behavior different from the overall number of request per day (Figure 6). We attribute the difference to the fact that the search function is used by robots, which behave differently from human users. For instance, the spike observed on Tuesday of the first week results from an unexpected number of requests for the home page. Such behavior could also be indicative of a denial of service attack and understanding such dynamics could be of advantage for security purposes.

5.1 Popularity of Search Terms

Prior studies of Web traffic have found that the popularity of static pages (i.e., documents) served by information provider Web sites follows Zipf's law [5,7]. In E-business sites, customers look for product information instead of documents or static pages. Product information is usually generated by dynamic pages, based on keywords provided by customers. A common way of finding product information in an online store is through query-based search functions, which are the central part of product seeking in E-business sites. Customers use keyword search functions to discover products and services. To improve the efficiency of search functions it is important to understand the behavior of customers when they are looking for information. So, we want to examine the frequency of specific queries and find out the underlying distribution of these queries. We conjecture that a small set of queries, which refer to popular items of the store, are repeated many times over the course of a day.

Reference [15] shows that surfing patterns on the Web display strong statistical regularities, that can be described by universal laws. Zipf's law has been extensively used to explain the patterns of access to Web servers and proxies. We investigate this issue further by studying the patterns of keywords used by customers during their interaction with an E-business site.

Zipf’s law [24] is a relationship between the frequency of occurrence of an event and its rank, when the events are ranked with respect to the frequency of occurrence. Zipf’s law [24] was originally applied to the relationship between words in a text and their frequency of use. It states that if one ranks the popularity of words used in a given text (denoted by r) by their frequency of use (denoted by $f(r)$) then $f(r) \sim 1/r$. This expression can be generalized as $f(r) = C/r^\alpha$, where C is a constant and α a positive parameter equal to one. This law describes phenomena where large events are rare, but small ones are quite common. Relationships such as Zipf’s law can be used to facilitate both cache resource planning and strategies for distributing E-business functions.

In Fig. 12, we plot the relative percent frequency of a given query term versus its popularity rank for both sites and for the entire log. The figure shows that Zipf’s law applies quite strongly to *the terms used for search functions*. This result is similar to the one found in [5], which showed that Web documents returned by Web servers also follow a Zipf’s law. The figure displays three curves: one for the bookstore, one for the auction site, and other for Zipf’s law. As it can be seen, there is a good match with Zipf’s law over an extremely wide range of popularity, except for the most popular keywords in the bookstore site. This fact is represented by the relatively flat part of the bookstore curve for small values of the term rank (i.e., popular search terms). Let us examine this fact in more detail.

At first sight, it appears that the most popular keywords for the bookstore do not follow Zipf’s Law. Let us use the multiple time scale approach to investigate why the bookstore curve has an accentuated flat region. In other words, let us look at the popular terms at different time scales. For example, in our analysis we used a two-week period log. In these two weeks, there is a kind of “shift” in popularity for the bookstore. The most popular keyword in the first week may be different from the most popular one in the second week. However, cumulatively both keywords may get the same number of requests. If this phenomenon happens for several keywords, the result can be seen as a

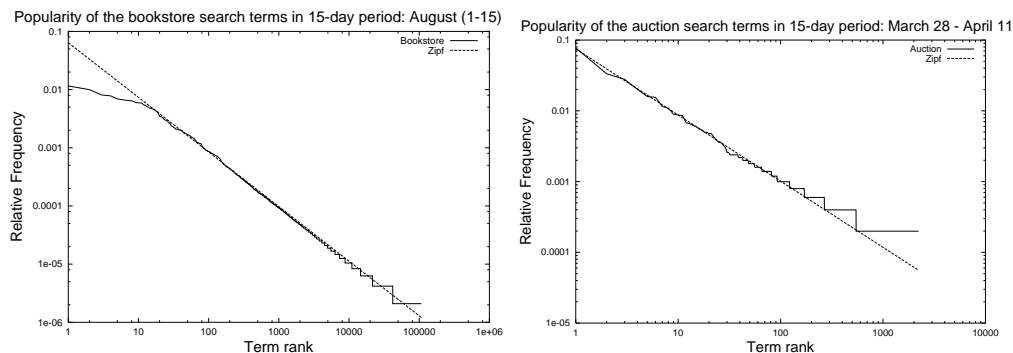


Fig. 12. Popularity of search terms for the bookstore and auction sites for a 15-day period.

flat region in the leftmost part of the frequency versus rank plot. In order to verify this conjecture, we plot the popularity graphs for different time scales (1 day, 3 days, and 7 days) for the bookstore and auction sites (see Fig. 13).

The top graphs of Fig. 13 correspond to logs of one-day period of time for the bookstore and auction sites. In this case, the flat part of the curve was clearly reduced for the bookstore. Our explanation is that one day is too short a period for a significant shift in popularity to occur. As we increase the period of analysis, we notice that the flat part of the curve increases accordingly for the bookstore as seen in the two remaining graphs for the bookstore in Fig. 13 and the one in Fig. 12. On the other hand, measuring Zipf's law over too short time intervals could bias the power law since the most popular keywords could have not enough "opportunity" to be requested in order to follow the exact power law. Figure 13 also shows that for the auction site, the popularity

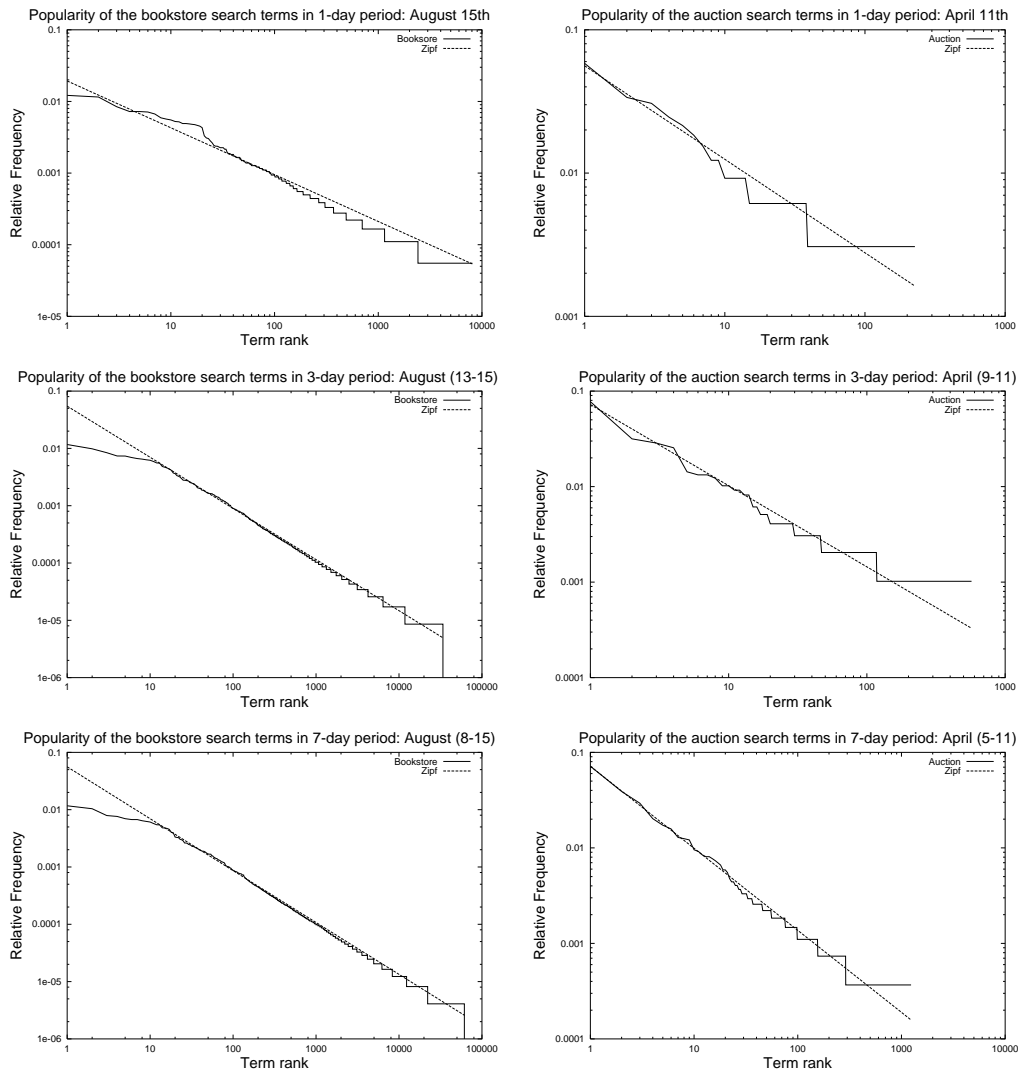


Fig. 13. Popularity of search terms for the bookstore (left) and auction sites (right) for different time scales.

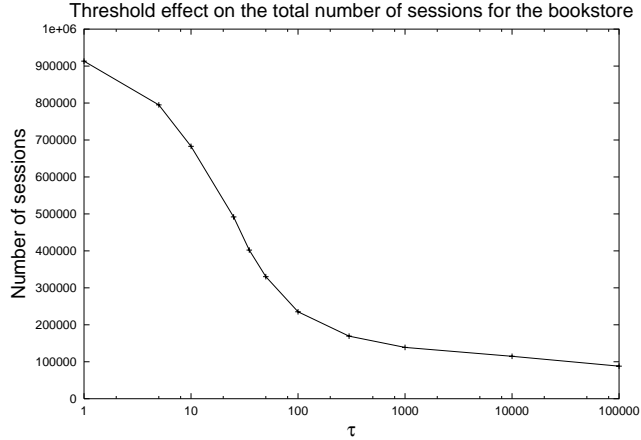


Fig. 14. Influence of the threshold τ on the total number of sessions initiated at the bookstore site.

curves do not exhibit the same temporal shift seen in the bookstore case. This can be explained by the fact that our auction site auctions domain names, which are not expected to exhibit a significant change in popularity over short periods of time.

6 Session Characterization

Session boundaries are delimited by a period of inactivity by a customer. In other words, if a customer has not issued any request for a period longer than a threshold τ , his session is considered finished. Usually, sites enforce this threshold and close inactive sessions to save resources allocated to these sessions. For the auction site, we know that the HTTP server enforced a threshold of twenty minutes. Since we do not have this information for the bookstore site, we had to estimate the threshold from the log. The value of τ has an influence on the number of sessions being handled by the site.

We discuss the effect of τ in what follows. Figure 14 shows the effect of the value of τ in the total number of sessions initiated for the bookstore site. As we can see, as the threshold increases from 1 to 100 sec, the number of sessions initiated decreases very rapidly. From 1000 sec on, the decrease is very small. This indicates that most sessions last less than 1,000 sec. A de facto industry-standard has been that 30 minutes (i.e., 1,800 sec) should be used to delimit sessions.

Figure 15 shows the distribution of session lengths, measured in number of requests to execute E-business functions, for both sites. The threshold τ used for the bookstore is 1,800 seconds while there is no threshold for delimiting the sessions at the auction site, since it implements timeouts for its sessions.

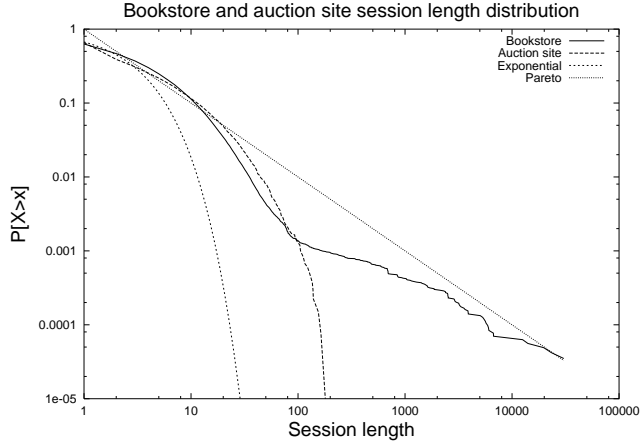


Fig. 15. Session length distribution.

The graphs of Fig. 15 show the empirical tail of the distribution of the session length X , i.e., $P\{X > x\}$ for the bookstore and auction sites, as well as the tail of the exponential and Pareto distributions. A random variable X , such as Pareto, that has a heavy-tailed distribution is characterized by $P\{X > x\} \sim x^{-a}$, $0 < a < 2$. Among other implications, a heavy-tailed distribution presents a great degree of variability, and a non-negligible probability of high sample values. The exponential distribution decays much faster than a heavy-tailed distribution. In a log-log plot, x^{-a} is a straight line with inclination $-a$. We can distinguish two regions in the plot of Fig. 15. The first one comprises session lengths of up to 100 requests, in which the curves for both sites are similar. In particular, in the region from about 5 to 100, they are fit by a straight line (not shown for clarity) with inclination ~ -2.05 . For sessions longer than 100, the behavior changes. We can see that for the auction site, the probability for longer sessions falls abruptly, whereas for the bookstore it remains close to the straight-line plot of a Pareto-like distribution with $a = 1$. This “very” heavy tail is most likely due to the accesses by robots, which tend to exhibit long sessions. The auction site was not accessed by any detectable robot, and this explains why one does not see sessions much longer than 100 requests. We can also notice that most sessions are small (about 90% of the sessions for both workloads have less than 10 requests).

6.1 Usage Analysis

The left part of Fig. 16 shows the number of sessions initiated per day at the bookstore site for various values of the threshold τ . A small value of τ corresponds to the extreme case of considering each request as a session. The picture clearly shows that there is very little difference in the number of sessions as τ is increased from 1,000 sec to 10,000 sec. This is a strong argument in favor of the 30-minute standard. A similar behavior is seen in

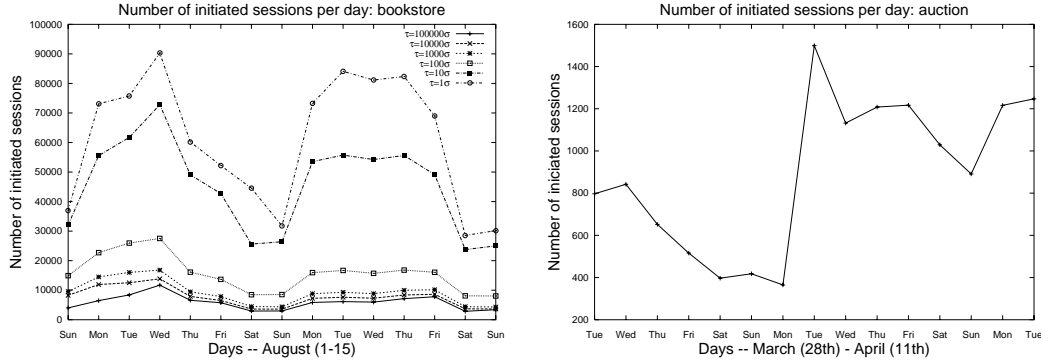


Fig. 16. Number of initiated sessions per day for the bookstore and the auction site.

the left part of Fig. 17. The right part of Figs. 16 and 17 indicate the number of sessions initiated per day and per hour for the auction site. If we compare the shape of the graph of the number of initiated sessions for the bookstore site for $\tau = 1000$ and for the number of initiated sessions for the auction site with the corresponding graphs of Fig. 6, for number of arriving requests, we see some degree of similarity.

Figure 18 displays the number of active sessions on an hourly basis for various values of the threshold τ . Again, very little variation is seen for $\tau > 1000$ sec. At a time scale of one hour, we observe a high variability in the number of active sessions per hour since the session timeout for the auction site or the threshold of 1,000 sec for the bookstore are of the same order of magnitude as the time scale.

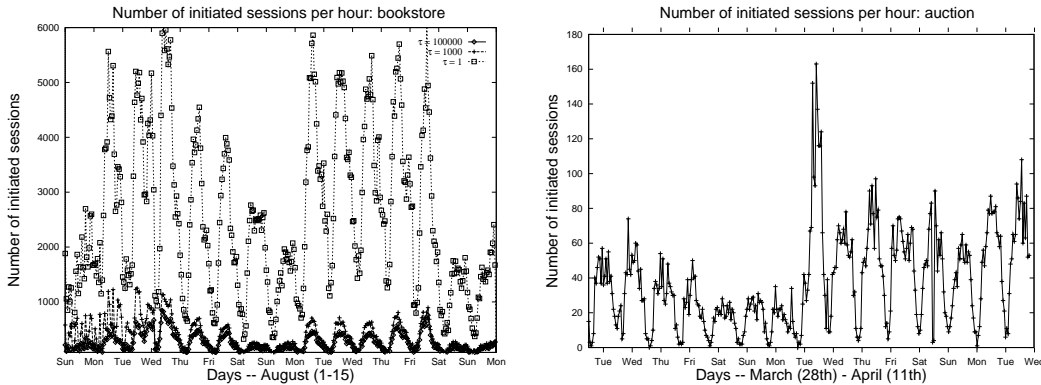


Fig. 17. Number of initiated sessions per hour for the bookstore and the auction site.

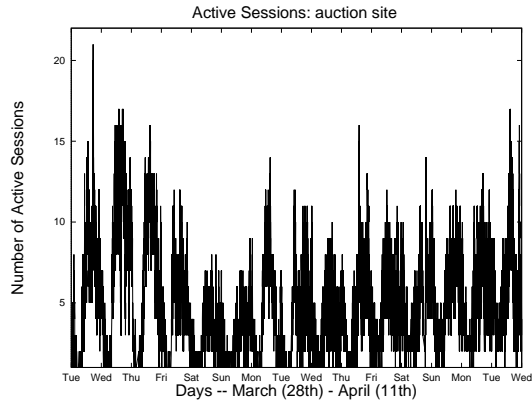
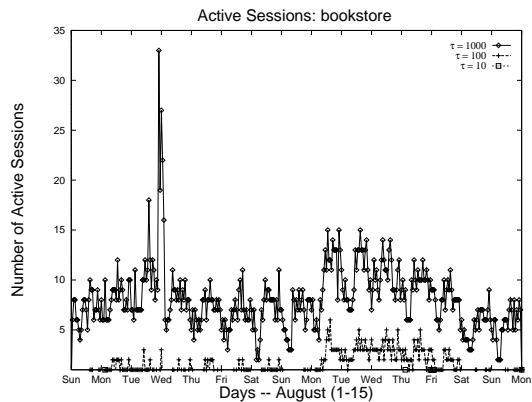
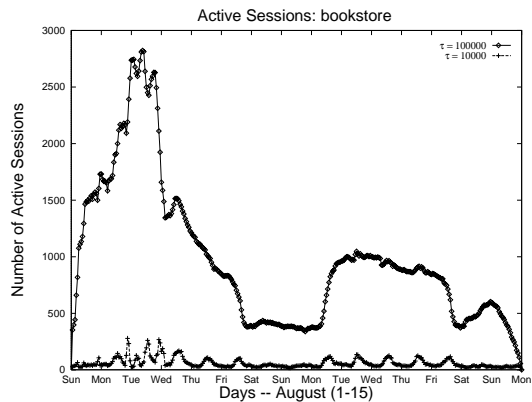


Fig. 18. Number of active sessions on an hourly basis. The top two figures illustrate the influence of τ in the bookstore.

7 Concluding Remarks

Several studies have been published regarding the workload of information provider sites. However, very few studies are available for E-business sites. This paper presented a hierarchical and multiscale approach for workload characterization of E-business sites. The characterization was done at the session, E-business function, and request levels. The approach was applied to two ac-

tual E-businesses: an online bookstore and online auction site.

The hierarchical and multiscale characterization approach allowed us to identify several new characteristics in the workload of the two sites analyzed. Some of the findings are: i) most sessions last less than 1,000 sec. ii) 88% of the sessions have less than 10 requests. iii) the session length, measured in number of requests to execute E-business functions, is heavy-tailed, especially for sites subject to requests generated by robots. iv) more than 70% of the functions performed are product selection functions as opposed to product ordering functions. v) requests to execute frequent E-business functions exhibit a similar pattern of behavior as observed for the total number of HTTP requests. vi) the popularity of search terms follows a Zipf distribution. However, Zipf's law as applied to E-business is time scale dependent, especially for sites that exhibit a shift in popularity of search terms. A similar observation in the context of media servers was discussed in [10]. There, the authors determined that file access frequencies for media workloads can be approximated by Zipf-like distributions, which vary with the time scale. In that study the authors use a much coarser time scale variation (1-month, 6-month, 1-year, and 2.5-year) than the one used here (1-day, 3-day, and 7-day). vii) there is a strong correlation in the arrival process at the request level. This correlation is given by an average Hurst parameter value of 0.73. viii) correlations in the arrival process are particularly stronger at intermediate time scales of a few minutes. We also noted that the inter-arrival time pattern at the auction site exhibits sudden unexpected periods of inactivity while the bookstore seems to see long periods of silence which are less abrupt and which could be anticipated.

It is recognized by many that one of the major challenges in carrying out experimental work in E-business is the lack of data. Most companies regard their logs as sensitive information that should not be made public. The methodology presented in this paper can certainly be applied to logs of other E-business sites and constitutes an important step for capacity planning and performance tuning. The type of workload statistics one may find when studying other E-business sites may vary as a function of the types of products and services offered and as a function of the business model implemented by the site.

One of the main advantages of our methodology is that it provides a characterization at multiple levels of abstraction, which is useful for the understanding of user behavior, site functionality, and workload intensity and arrival process at the protocol level. The multiple time scale analysis we used proved to be quite useful at uncovering aspects of the workload that one would miss by looking at a single time scale. We have not characterized the workload for capacity planning purposes as done in [6].

Acknowledgements

The authors would like to thank the anonymous reviewers for their detailed and helpful comments, which greatly improved the quality of this paper. The work of D. Menascé was partially supported by the sponsors of the E-Center for E-Business at GMU. The work of V. Almeida was partially supported by the Brazilian Research Council (CNPq) and by a grant from SIAM 76.97.1016.00. R. Riedi's support comes in part from an NSF grant no. ANI-00099148 and from Texas Instruments. He acknowledges the support of CNPq for his visit at UFMG, Belo Horizonte, Brazil, in 2000. F. Ribeiro was supported by CNPq's grant no. 14.1252/01-4 and R. Fonseca was supported by a CNPq grant no. 133149/00-5.

References

- [1] P. Abry, P. Flandrin, M. Taqqu, and D. Veitch. Wavelets for the analysis, estimation and synthesis of scaling data. In *Self-similar Network Traffic and Performance Evaluation*. Wiley, Spring 2000.
- [2] P. Abry, P. Gonçalves, and P. Flandrin. Wavelets, spectrum analysis and $1/f$ processes. In A. Antoniadis and G. Oppenheim, editors, *Lecture Notes in Statistics: Wavelets and Statistics*, volume 103, pages 15–29, 1995.
- [3] V. Almeida, M. Arlitt, and J. Rolia, “Analyzing a Web-based System's Performance at Multiple Time Scales,” *Proc. 2002 ACM SIGMETRICS Practical Aspects of Performance Analysis Workshop (PAPA)*. Also published at the ACM Sigmetrics Performance Evaluation Review, Vol. 30, No. 2, September 2002.
- [4] V. Almeida, D. Menascé, R. Riedi, F. Ribeiro, R. Fonseca, and W. Meira, Jr., “Analyzing Web Robots and their Impact on Caching,” *Proc. Sixth Workshop on Web Caching and Content Distribution*, Boston, MA, June 20-22, 2001.
- [5] V. Almeida, M. Crovella, A. Bestavros, and A. Oliveira “Characterizing Reference Locality in the WWW,” *Proc. IEEE/ACM International Conference on Parallel and Distributed System (PDIS)*, Dec. 1996.
- [6] M. Arlitt, D. Krishnamurthy, and J. Rolia, “Characterizing the Scalability of a Large Web-based Shopping System,” *ACM Transactions on Internet Technology*, Vol. 1, No. 1, August 2001, pp. 44–69.
- [7] M. Arlitt and C. Williamson, “Web Server Workload Characterization,” *Proc. 1996 SIGMETRICS Conference on Measurement of Computer Systems*, ACM, May 1996.
- [8] J. Brown and P. Duguid, *The Social Life of Information*, Harvard Business School Press, 2000.

- [9] M. Calzarossa and G. Serazzi, "Workload Characterization: A Survey," *Proc. of the IEEE*, Vol. 81, No. 8, August 1993.
- [10] L. Cherkasova and M. Gupta, "Characterizing Locality, Evolution, and Life Span of Accesses in Enterprise Media Server Workloads," *Proc. 12th Int'l. Workshop on Network and Operating System Support for Digital Audio and Video (ACM NOSSDAV 2002)*, Miami Beach, FL, May 2002.
- [11] M. Crovella and A. Bestavros, "Self-Similarity in World Wide Web Traffic: Evidence and Possible Causes," *IEEE/ACM Transactions on Networking*, 5(6), pp. 835-846, December 1997.
- [12] I. Daubechies, "Ten Lectures on Wavelets," SIAM, New York, NY, 1992.
- [13] W. Leland, M. Taqqu, W. Willinger, and D. Wilson, "On the self-similar nature of Ethernet traffic (extended version)," *IEEE/ACM Trans. Networking*, pp. 1-15, 1994.
- [14] L. Cherkasova and P. Phaal, "Session Based Admission Control: A Mechanism for Improving the Performance of an Overloaded Web Server," HPL-98-119, HP Labs Technical Reports, 1998.
- [15] B. Huberman, P. L. T. Pirolli, J. E. Pitkow, and Rajan M. Lukose, "Strong Regularities in World Wide Web Surfing," *Science*, Vol. 280, April 3, 1998.
- [16] D. A. Menascé, V. A. F. Almeida, R. Riedi, F. Pelegrinelli, R. Fonseca, and W. Meira Jr., "In Search of Invariants for E-Business Workloads," *Proc. 2000 ACM Conf. in E-commerce*, Minneapolis, MN, October 17-20, 2000.
- [17] D. A. Menascé and V. A. F. Almeida, *Scaling for E-Business: technologies, models, performance and capacity planning*, Prentice Hall, Upper Saddle River, NJ, May 2000.
- [18] D. A. Menascé, V. A. F. Almeida, R. Fonseca, and M. A. Mendes, "Business-oriented Resource Management Policies for E-Commerce Servers," *Performance Evaluation*, Vol. 42, September 2000, pp. 223-239.
- [19] D. A. Menascé, V. Almeida, R. Fonseca, and M. Mendes, "A Methodology for Workload Characterization for E-Commerce Servers," *Proc. 1999 ACM Conference in Electronic Commerce*, Denver, CO, Nov. 3-5, pp. 119-128.
- [20] Pitkow, J., Summary of WWW characterizations, *World Wide Web*, No. 2, 1999.
- [21] D. Krishnamurthy and J. Rolia, "Predicting the Performance of an E-Commerce Server: Those Mean Percentiles," in *Proc. First Workshop on Internet Server Performance*, ACM SIGMETRICS 98, June 1998.
- [22] R. Riedi, M. S. Crouse, V. Ribeiro, and R. G. Baraniuk. "A multifractal wavelet model with application to TCP network traffic," *IEEE Trans. Info. Theory, Special issue on multiscale statistical signal analysis and its applications*, Vol. 45, pp. 992-1018, April 1999.

- [23] M. Taqqu, V. Teverovsky, and W. Willinger. “Estimators for long-range dependence: An empirical study,” *Fractals*, Vol. 3, pp. 785–798, 1995.
- [24] G. Zipf, *Human Behavior and the Principle of Least Effort*, Addison-Wesley, Cambridge, MA, 1949.