

# Using Libception to Understand and Improve HTTP Streaming Video Server Throughput

Tyler Szepesi  
University of Waterloo  
Waterloo, ON, Canada  
stszepes@cs.uwaterloo.ca

Derek Eager  
University of Saskatchewan  
Saskatoon, SK, Canada  
eager@cs.usask.ca

Benjamin Cassell  
University of Waterloo  
Waterloo, ON, Canada  
becassel@cs.uwaterloo.ca

Jim Summers  
University of Waterloo  
Waterloo, ON, Canada  
jasummer@cs.uwaterloo.ca

Tim Brecht  
University of Waterloo  
Waterloo, ON, Canada  
brecht@cs.uwaterloo.ca

Bernard Wong  
University of Waterloo  
Waterloo, ON, Canada  
bernard@cs.uwaterloo.ca

## Abstract

Video streaming applications generate a large fraction of Internet traffic. Much of this content is delivered over HTTP using standard web servers. Unlike other types of web workloads, HTTP video streaming workloads are typically disk bound, and therefore an important problem is that of optimizing disk access.

In this paper we design, implement and evaluate Libception, an application-level shim library that implements techniques for improving disk I/O efficiency. Web servers can achieve the benefits of these techniques simply by linking with Libception, without the need to modify source code. In contrast to making kernel changes or attempting to optimize kernel tuning, Libception provides a portable and relatively simple setting in which techniques for optimizing I/O in HTTP video streaming servers can be implemented and evaluated.

We report experimental results evaluating the efficacy of the aggressive prefetching and disk I/O serialization techniques currently implemented in Libception, for three web servers (Apache, nginx and the userver) and two operating systems (FreeBSD, Linux). We find that on FreeBSD, video streaming throughput with all three web servers can be doubled by linking with Libception. On Linux, performance similar to that provided with Libception was eventually obtained by examining the kernel source to understand and tune kernel parameters. With the default kernel parameter settings, however, and regardless of which Linux disk scheduler is selected, we find that use of Libception can approximately double throughput. We find that both aggressive prefetching and serialization are necessary to achieve these benefits.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

ICPE'17, April 22 - 26, 2017, L'Aquila, Italy

© 2017 Copyright held by the owner/author(s). Publication rights licensed to ACM. ISBN 978-1-4503-4404-3/17/04...\$15.00

DOI: <http://dx.doi.org/10.1145/3030207.3030231>

## 1. INTRODUCTION

Video streaming over HTTP is now the largest contributor to Internet traffic. The catalogue of available content from popular video streaming services has also been growing rapidly, and although memory caching and SSD can be effective for the most popular content, HTTP streaming video server workloads are often disk-bound [13]. Techniques for improving disk throughput in such systems are therefore of considerable interest.

There has been much past work on improving disk access efficiency. In the context of HTTP video streaming, however, there are two complicating factors. First, unlike early video streaming systems, HTTP-based video streaming is pull-based: the server responds to client requests for video chunks, rather than pushing video data to the client at some server-determined rate. Second, contemporary servers may be highly concurrent, responding to video chunk requests from hundreds or thousands of clients concurrently.

A well-known approach to making disk access more efficient for applications that access files sequentially is to perform larger reads, by prefetching data before it has been requested. In prior work we observed that the *serialization* of reads may also be important in some contexts [33, 34]. In particular, we found that a modification to a web server (the userver) to both aggressively prefetch, by performing large reads, *and* to serialize its reads, yielded large performance improvements on FreeBSD. Requests are serialized by not permitting new `read` or `sendfile` system calls for data from video files for one client until the previous request (from another client) has completed. Serialization can be beneficial if the operating system breaks large reads into smaller I/O buffer-sized reads, and if sequences of these smaller reads from concurrently-issued large read requests can get interleaved inside the operating system. This prior work did not investigate whether the approach of combining aggressive prefetching with serialization could yield similarly large benefits for other web servers or on other operating systems.

Implementing techniques for improving disk access efficiency inside the application requires detailed knowledge of the application code, and must be repeated for each application of interest. This could be quite difficult for applications with large code bases like Apache and nginx. On the other hand, a kernel implementation is operating system specific, requires detailed knowledge of the relevant pieces of kernel

code, and has the additional problem of potential adverse impacts on other types of applications.

In this paper, we address the problem of improving disk access efficiency in HTTP video streaming servers by developing an application-level shim library in which applicable techniques can be implemented. We then apply this library to evaluate the performance improvements provided by aggressive prefetching and serialization, for the widely used Apache and nginx web servers as well as the userver, and two operating systems (FreeBSD, Linux).

Our contributions are as follows:

- We design and implement Libception, a portable, application-level shim library that implements techniques for improving disk I/O efficiency. We demonstrate that web servers can achieve the benefits of these techniques simply by linking with Libception, without the need for source code changes. Comparing a web server that we had modified to incorporate the techniques directly, and the unmodified server linked with Libception, we find essentially identical performance.
- We show that the aggressive prefetching and disk I/O serialization techniques currently implemented in Libception can approximately double the peak HTTP video streaming throughput of a variety of web servers (Apache, nginx, and the userver), both on FreeBSD, and on Linux when using the default kernel parameter settings, regardless of which Linux disk scheduler is chosen.
- We apply Libception to address the question of whether aggressive prefetching by itself is sufficient on Linux, or whether combining aggressive prefetching with serialization provides substantial further performance benefits. We find that combining the techniques yields substantially better performance than aggressive prefetching (or serialization) alone.
- We discover that there is great scope for improving HTTP video streaming performance on Linux, when not using Libception, by tuning kernel parameters. In particular, by tuning parameters to improve both prefetching and serialization, we find that peak throughput can be more than doubled, yielding peak throughput slightly higher than that obtained with Libception. In contrast, when using Libception, kernel parameter tuning yields only marginal additional improvements.

## 2. BACKGROUND AND RELATED WORK

HTTP streaming video workloads are typically disk-bound due to the large size of video files and the tendency of the popularity distributions of these files to have a very long tail (meaning the large majority of content is accessed infrequently) [13]. There are three general techniques for improving the bottleneck of disk performance: file caching, disk scheduling and prefetching. Our research is primarily interested in the effects of request scheduling and prefetching. File caching plays a significant role in improving server throughput for our workload, but competes with prefetching for system memory resources [5]. For the purposes of our experiments, we simply use the kernel file caching algorithm as-is.

In the following sections we describe prior work for disk scheduling and prefetching, and discuss how the results apply to an HTTP streaming video workload. We also specifically discuss handling concurrent I/O streams, because the bulk of the research into scheduling and prefetching assumes a single-threaded workload (which is not consistent with our workload).

### 2.1 Block I/O Scheduling

Block I/O schedulers play a major role in most modern operating systems, and typically act as a layer between user I/O requests and requests to block device drivers. As of kernel version 2.6.33, Linux provides three different default options for scheduling block I/O devices: NOOP, deadline, and completely fair queuing (CFQ). All three schedulers attempt to take advantage of different aspects of temporal and spatial locality between requests in order to yield higher throughput for disk-bound workloads.

The NOOP scheduler is the least complicated of the Linux schedulers. It provides a simple first-in-first-out (FIFO) queue for requests, and also performs basic request merging [2]. The deadline scheduler maintains sector-sorted read and write queues, as well as queues which are organized by expiration times. The deadline scheduler gives priority to expired requests in the secondary queues, and otherwise batches requests from the sector-sorted queues. It is tailored towards workloads that require latency guarantees on I/O requests [2]. CFQ divides access to the disk into time slices which are allocated between groups of per-process request queues and sized by process priority. CFQ idles shortly on empty queues whose time-slices have not expired, even if other queues contain outstanding requests [2].

Linux previously offered an anticipatory scheduler (AS), which was introduced as a means to eliminate “deceptive idleness”. Deceptive idleness occurs when processes leave a small data processing gap between I/O requests. During this time, naive schedulers may switch to servicing other processes, introducing seeks that can degrade system performance [11]. The abilities of AS are mostly a subset of CFQ, and as such AS was removed in version 2.6.33 of the Linux kernel [3].

The schedulers built into the Linux kernel are necessarily designed to handle a wide range of workloads. However, there are characteristics of streaming video workloads that can be exploited by specialized scheduling algorithms. There are many studies, aimed at broadcast streaming scenarios in which servers push data to clients, where scheduling is used as a means to maximize throughput [29, 9, 7]. These studies contain valuable insights, but are not directly applicable to HTTP streaming video servers, where clients individually pull requests from the server. As we will demonstrate experimentally, the choice of block I/O scheduler is not very important for HTTP streaming video workloads. Regardless of which scheduler is chosen, significant throughput benefits are gained from the use of other techniques to improve I/O.

### 2.2 Prefetching

Prefetching is a well-studied technique, and refers to reading data from the disk into memory before it is actually requested by the user. This allows subsequent read requests to return immediately instead of blocking on disk operations. Prior research has shown the effectiveness of using prefetching as a means of off-setting latency and CPU stalls [22, 6,

36]. Papatthasiou and Scott demonstrated that aggressive prefetching (prefetching far beyond what is requested by the user) can be used to offset request latency [20]. However, latency of disk access is not a significant concern for HTTP streaming video workloads. Clients use buffering to cope with potentially high network latencies, so lower latencies incurred by disk I/O are unlikely to affect the quality of service experienced by end users.

Instead, HTTP streaming video would benefit from using prefetching as a vehicle for increasing disk throughput. Examples of systems that have studied prefetching within this context are DiskSeen, a system that modifies the Linux kernel to introduce history-aware prefetching into the operating system [12], and libprefetch, which uses both kernel and application modifications to provide application-directed prefetching [35]. Unlike previous work, our research considers disk request serialization in combination with highly aggressive prefetching, and does not require code changes at the user level nor in the kernel.

Prefetching is commonly performed at the hardware level, in addition to the software level. We refer to prefetching done by the disk itself as “lookahead”. Ruemmler and Wilkes demonstrated that lookahead could improve read times by up to 42% on Unix-based systems [27], and subsequently showed that effective use of on-disk caches for lookahead can yield optimal results for sequential workloads by eliminating unnecessary rotational delays [26].

An important issue in prefetching concerns how much data to request with each read operation. A larger prefetch amortizes the overhead cost of a disk access over more bytes of data but can have adverse consequences, such as the eviction of useful data from the cache. Panagiotakis, et al. [19] demonstrate that using a large fixed prefetch size to service 100 sequential streams improves throughput by up to 4 times compared to not prefetching. Li, et al. [14] provide a 2-competitive algorithm that uses hard drive performance specifications to determine a prefetch size. In prior work, we found that the best prefetch size depends on both available system resources and specific workload characteristics [32]. We provide an automated algorithm for dynamically determining a good prefetch size [32], and we demonstrate that it is possible to further increase the efficiency of servicing a streaming video workload by exploiting knowledge of specific workload characteristics to implement a prefetch algorithm [31]. For this paper, we implement a simple prefetch algorithm in Libception that uses fixed-size prefetches, with a default size of 2 MB (we experiment with different fixed prefetch sizes in Sections 5.4 and 5.5). In the future, for use in a production server or for different workloads, we could implement an automated workload-specific prefetch algorithm in Libception.

### 2.3 Concurrent I/O Streams

It is important to consider the trade-off between reading from disk efficiently using large prefetches and servicing concurrent disk requests fairly. Panagiotakis et. al demonstrated rapid degradation of I/O throughput as additional I/O streams are introduced across a variety of Linux schedulers [19].

The Argon system [37] focused on meeting service-level agreements by using disk-head time-slicing to ensure minimum levels of throughput to competing applications. Unlike Argon’s target workload, HTTP streaming video workloads

only require that clients avoid re-buffering (it is not uncommon for clients to buffer about 5 – 30 seconds of video [1, 4]). Because HTTP streaming clients are insensitive to disk latency, there is scope to reduce fairness between individual clients in exchange for higher disk throughput.

We will demonstrate that, by using request serialization in combination with aggressive prefetching, we can achieve significantly higher HTTP streaming video server throughput than if we were to allow the system to process the I/O requests in parallel.

## 3. DESIGN AND IMPLEMENTATION

We developed a library shim, Libception, that provides applications with the necessary capabilities for both I/O serialization and aggressive prefetching without the need for source code or kernel modifications.

Our library is portable, operates in user space, and has been tested on FreeBSD, Linux and Mac OS X. It is comprised of two components: The first component is Deception, a dynamically linked shared object which inserts itself between the application and libc calls. Deception intercepts I/O requests from the application and forwards them to the system’s second component, Reception. Reception is a server process that runs separately from applications using Deception. Reception receives, services, and responds to requests generated by applications using Deception (including prefetching data when necessary). Figure 1 shows a high-level overview of the components of Libception.

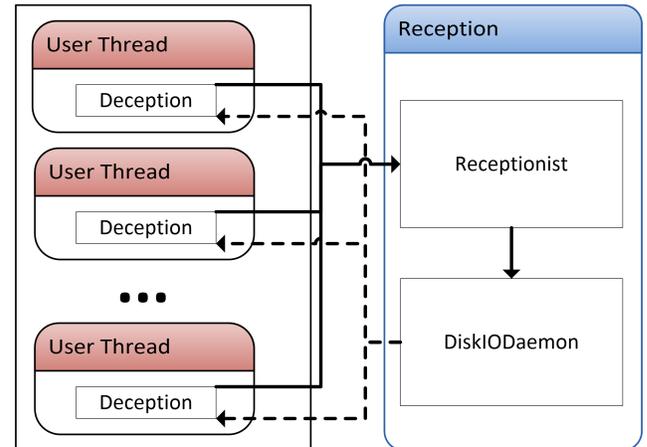


Figure 1: Libception design

### 3.1 Deception

Deception is the primary interface for communication between user applications and the Libception library. It is implemented as a shared object that is dynamically linked at launch-time by the application that wishes to make use of it (for example by using the LD\_PRELOAD environment variable on most Unix-based operating systems). Once loaded, Deception begins silently intercepting calls to a variety of libc I/O-related functions including `open`, `read` and `sendfile`. This technique allows applications to take advantage of the Libception library’s benefits without requiring any source code modifications. Furthermore, because Deception is implemented in user space, it requires no changes to the underlying operating system.

Most calls to Deception perform several validity checks, and then determine whether or not the application’s I/O request is already resident in the system’s file cache. This is done by using the system calls `mmap` and `mincore` to check all memory blocks of the request (excluding blocks that would be prefetched if the request went to disk). If a request is found to be contained entirely in memory, Deception passes the call to the original libc function (which returns without needing to go to disk). Otherwise, disk I/O is required, and a Libception I/O request is constructed and sent to Reception using Unix sockets (which duplicates any necessary file descriptors).

Should a disk read be necessary, Deception waits on a response from Reception before it proceeds. To ensure transparency, Deception always finishes by executing the application’s original libc call and returning to it the appropriate return codes (even if a parameter input fails sanity checks). This allows Deception to run invisibly, without affecting the guarantees of the API for the associated libc call. This in turn means that the application can expect the same control and user-functionality for libc I/O calls that it would be afforded if Libception was not being employed.

Some additional non-I/O libc functions are also intercepted by Deception (such as `getpid` and `fork`). These functions are intercepted for functionality purposes, do not communicate with the Reception layer, and are also invisible to the user (as they terminate by making the original libc calls as well). All initialization for Deception is handled transparently at launch by GCC constructor functions and all cleanup is likewise handled unobtrusively at termination by GCC deconstructor functions which are run when the Libception shared object is loaded and unloaded, respectively.

### 3.2 Reception

In our experiments, Reception is launched as a separate user space process by the user, prior to running applications with the Deception shim. It could also be launched as a daemon process in a production environment. Reception is primarily responsible for accepting, serializing and servicing incoming I/O requests from one or more Deception shims. Reception furthermore modifies requests as necessary (for example, by enlarging read sizes to introduce readahead), executes the prefetch itself, and finally responds to Deception. These tasks are divided between a single Receptionist thread, and one or more DiskIODaemon threads.

The Receptionist thread acts as a server, accepting requests over Unix sockets. I/O requests received by the Receptionist thread are sorted by device, and are placed in one of multiple queues to be serviced by the appropriate DiskIODaemon thread. The Receptionist server thread also may delegate requests to a separate maintenance thread that performs utility tasks, such as statistics collection and aggregation.

DiskIODaemon threads are responsible for performing I/O for different devices. By default, Reception runs in a single DiskIODaemon mode, with all requests being serviced sequentially in the system by one thread. Alternatively, the user may set Reception to create individual DiskIODaemon threads for individual devices in the system. In this case, the Receptionist de-multiplexes incoming I/O requests to the DiskIODaemon threads based on the underlying device for the file descriptor in the request.

Regardless of which mode is selected, DiskIODaemon threads each use their own lock-protected request queue, which is filled by the Receptionist, and drained by the DiskIODaemon thread. Only one request is removed from this queue at a time by the DiskIODaemon thread, ensuring that I/O to whichever device it is servicing is serialized. As requests are pulled from this queue, they are expanded to include any necessary prefetch information, and are then sent to disk. Once finished servicing the request, the DiskIODaemon thread sends a message to the Deception shim that made the request, allowing it to unblock and proceed.

Libception additionally contains options that allow it to perform prefetch-free serialization, serialization-free prefetching, and simple request tracking without either serialization or prefetching (useful, for instance, in latency profiling or other statistics gathering). Serialization-only mode indicates to Libception not to extend reads provided by the user and is identical to prefetching with a prefetch size of zero. Prefetch-only mode transfers the burden of extending and performing requests from Reception to Deception. In this mode, instead of communicating with Reception, Deception shims simply immediately extend their requests and perform application-side `pread` calls before completing the original I/O request.

## 4. EXPERIMENTAL METHODOLOGY

We used the methodology described in [33] to generate experimental workloads and benchmarks. Our workload represents a large number of HTTP streaming video clients requesting videos with characteristics similar to requests for YouTube videos in 2011. We use a small number of client machines to generate traffic simulating thousands of concurrent sessions. Each session represents an end user viewing video. The video selected for each session is chosen using a Zipf distribution with an  $\alpha$  value of 0.8. This video is watched for some fraction of its duration. It is an important characteristic of video workloads that users do not typically watch to the end of a video, so this property is reflected in our workloads: The percentage of a video that any given user requests is in line with the watching patterns of a typical YouTube-like workload.

Another important characteristic of our workload is that the network is not the primary bottleneck. Netflix, for example, uses Open Connect Appliance (OCA) servers that have operational throughputs between 9 Gbps and 36 Gbps [16], but are provisioned with up to 40 Gbps of network capacity [18]. Instead, the workload is heavily disk-bound. This can be seen in the findings of our recent work characterizing the properties of a Netflix video workload [31]. Video providers maintain a large catalogue of content in order to appeal to a broad audience. Netflix, for instance, has a catalogue which is approximately 2 Petabytes in size [31]. Large video libraries have long tails [13], meaning they contain a large amount of infrequently viewed content. Cost effective storage is achieved by storing videos on large, inexpensive HDDs and the infrequent access means that the data being requested must be serviced from disk.

Videos are stored on the server hard drives by storing each video in a separate file. We stored data for a video in a single file rather than multiple chunks because we found this approach is more efficient in prior experiments [34]. The disk is populated with 20,000 video files which have an average duration of 265 seconds, with a similar distribution as

YouTube videos [8]. The files represent videos with a fixed bit rate of 420 kbps. This bit rate was chosen based on information available at the time of creation of the benchmark [10]. With these duration and bit rate characteristics, our average file size is 13 MB.

Each client session consists of a sequence of requests for 10 second intervals of video data, which are 0.5 MB in size. The first three requests in a session are made as quickly as the server can deliver the results, then subsequent requests are made on a fixed 10 second interval. This represents the filling of a playout buffer at the beginning of a session, followed by requests to refill the buffer as it is consumed at the bit rate of the video. This is a simplified model of a pull-based video client; it does not attempt to represent user actions like pausing the video, skipping to different points in the playback, or changing the quality level of the video. These actions were rare in the YouTube workload we modelled, but our methodology is flexible enough that we could represent these actions for workloads where they are significant.

The average duration of a video session is 160 seconds, using a distribution derived from real-world measurements. During experiments, video sessions are started at a chosen rate, using a Poisson distribution for session initiation. Experiments consist of 14,400 sessions, with a maximum of about 650 concurrent sessions. A total of 118 GB of video data is requested from 6581 different videos. Each video is viewed 2.2 times on average and 67.8% of videos are requested a single time during the experiment.

The clients monitor the service time for each request, and if it takes longer than 10 seconds to completely receive the data from a request, the client terminates the session and stops making further requests. For our experiments, we are interested in determining the highest aggregate client request rate that can be serviced, so that we can compare different web servers and configurations. To determine this rate, the *maximum failure-free rate*, we conduct a number of benchmark runs with a range of different aggregate rates of requests. From this, we determine the highest rate that results in fewer than 0.3% session failures. This value was chosen to permit clients to perform a small but very limited amount of re-buffering.

The clients are connected to the server over a local-area network with high bandwidth and low delay. To better represent the conditions available to real-world users, we use dummynet [23], which allows us to simulate different network types. We throttle 50% of client sessions in the workload to 3.5 Mbps, and the other 50% of client sessions to 10.0 Mbps, in order to represent a mix of end-user cable and DSL access speeds. Furthermore, we add 50 ms of delay to the network in each direction in order to model more realistic wide-area network conditions.

The equipment and environment we use to conduct our experiments were selected to ensure that network and processor resources are not a limiting factor in the experiments. We use two server machines, one for FreeBSD experiments and the other for Linux experiments. Both are HP DL380 G5 systems which contain two four-core Intel E5400 2.8 GHz processors and 8 GB of RAM. The Linux system uses Ubuntu 12.04 with a Linux 3.2.0 kernel, and a Western Digital Red (WDC WD10EFRX) 1.0 TB 5,400 RPM 3.5 inch SATA3 disk to store video files (chosen for its combination of relatively high throughput and low power consumption).

The FreeBSD system uses FreeBSD 8.0 and stores videos on an HP 146 GB 10,000 RPM 2.5 inch SAS disk. We use FreeBSD 8.0 so we can try to match the performance obtained previously with a modified web server [33]. Video files accessed during experiments are stored on a separate disk from the operating system.

Four client machines are used to generate the load of thousands of viewers on the servers. Each of these systems contains either dual 2.4 or dual 2.8 GHz Xeon processors and 2 GB or 3 GB of memory. Client machines run Ubuntu 10.04 on top of version 2.6.32-30 of the Linux kernel. They also use a version of httpperf [15] that was modified locally to support new features in a workload generation module named wssesslog. These modifications also allow the clients to track additional statistics. All clients are connected to the server via multiple 1 Gbps network links through multiple 24-port switches, helping to further ensure that the network is not a bottleneck.

## 5. PERFORMANCE EVALUATION

We now evaluate the maximum failure free throughput (often referred to henceforth as throughput) of the Apache, nginx, and userver web servers while utilizing Libception on FreeBSD and Linux servers. In all cases, we have tuned the web server to the best of our ability so that it provides the greatest maximum failure free throughput. Unless otherwise specified, the prefetch size used by Libception is 2 MB (we examine other sizes in Sections 5.4 and 5.5).

### 5.1 Evaluation on FreeBSD

In previous work [33, 34] we demonstrated how modifications to the userver web server to perform asynchronous serialized aggressive prefetching (ASAP) within the application significantly increased server throughput when servicing streaming video workloads. Unfortunately, these benefits rely on modifying the web server to use the `SF_NODISKIO` option [25] to the `sendfile` system call which is only available on FreeBSD. This flag causes `sendfile` calls that would block on disk I/O to instead return `EBUSY`.

The basic architecture of the userver using ASAP is to have a separate thread which performs large disk reads (thus implementing asynchronous, serialized, aggressive prefetching). This was relatively straightforward in the userver because it integrates well with its event-driven architecture and we were very familiar with the relatively small code base of the userver.

In this section we are interested in providing similar benefits to the more widely used Apache and nginx web servers without directly modifying either application. Note that nginx running on FreeBSD is of interest because the servers in the Netflix Open Connect Content Delivery Network use nginx on FreeBSD [17]. This is particularly relevant because Netflix currently accounts for a large fraction of peak Internet traffic in the United States [28].

We avoid making code modifications to the web servers because they each have a much larger code base than the userver and because Apache uses a significantly different software architecture (thread-per connection) [21]. Additionally, we want to determine if Libception can improve web server performance without using the non-portable `SF_NODISKIO` option to the `sendfile` system call.

Figure 2 shows the maximum failure free throughput obtained when using each of the Apache (labelled “A”), ng-

inx (labelled “N”) and userver (labelled “U”) web servers. Throughput is shown without Libception (labelled “Vanilla”), when using Libception (labelled “Libception”), and for the modified version of the userver that uses the `SF_NODISKIO` option (labelled “ASAP”). Additionally, this graph shows both the disk throughput and the web server throughput as observed by all of the client machines.

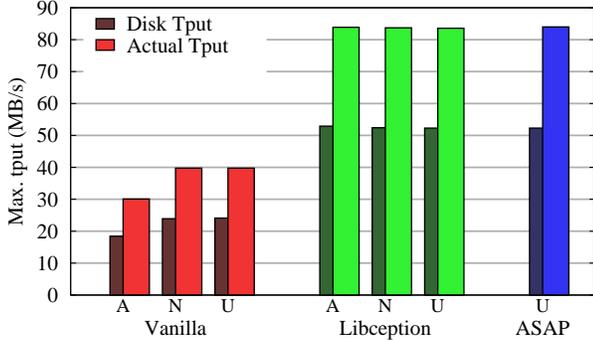


Figure 2: FreeBSD throughput without and with Libception and with ASAP

These results show that Libception is able to more than double disk throughput and total server throughput for all three web servers. As well, when using Libception, the maximum failure free throughput obtained by each server is equal to that obtained by the modified ASAP userver. This is despite Libception’s use of `mincore` prior to each call to `sendfile` (rather than relying on the non-portable `SF_NODISKIO` option) to determine whether or not data needs to be prefetched before calling `sendfile`. Although previous work reports that `mincore` call overhead can be significant [24], these video server workloads are disk-bound and can therefore easily tolerate the increase in CPU overhead.

It is worth pointing out that, in these experiments, not all of the data that is requested needs to be read from disk. For example, two clients requesting the same video content in quick succession will only require the data to be read from disk one time. Therefore, the difference between the total server throughput and the disk throughput is due to file system cache hits.

To the best of our knowledge, the version of FreeBSD used for these experiments does not provide any options to control the block I/O scheduler. There were also relatively few options available to influence kernel prefetching decisions and we were not able to significantly improve throughput using kernel parameters.

## 5.2 Evaluation on Linux

As noted previously, one of the key goals of Libception is to provide improved throughput for HTTP video web servers using techniques that are portable across different Unix-based operating systems. As a result, we now examine the performance of Apache, nginx and the userver on Linux. Recall that the disks used on the FreeBSD and Linux systems are different so we can not compare performance across the different operating systems.

Figure 3 shows the disk throughput and maximum failure free throughput obtained using each of the different web servers on Linux running with and without Libception. As was the case for FreeBSD, Libception again provides signifi-

cant improvements in disk and server throughput. On Linux server throughput is increased by a factor of about 2.5 times when using Libception.

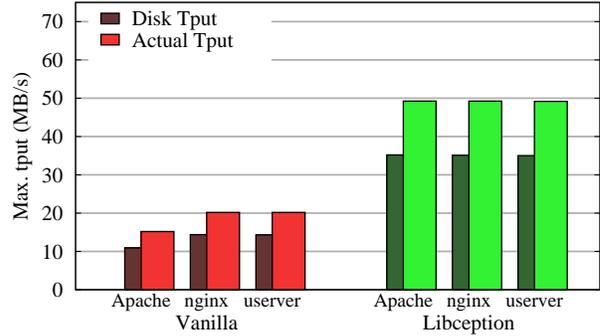


Figure 3: Linux throughput without and with Libception

To our surprise, despite years of research on prefetching techniques in operating systems, these applications and workloads do not appear to perform very well on either FreeBSD or Linux.

## 5.3 Evaluating Linux Block I/O Schedulers

The default Linux configuration on our system uses the CFQ block I/O scheduler [2]. We expected that the anticipatory nature of the Linux CFQ scheduler might be well-suited to this workload. Because web servers simultaneously process requests from thousands of clients, we expected that blocks from different requests might provide reordering opportunities that could be exploited by CFQ to improve disk and server throughput. For completeness we now examine the performance of all web servers with each of the three block I/O schedulers available in Linux. Figure 4 shows the throughput obtained without Libception while using the CFQ (labelled “C”), deadline (labelled “D”) and NOOP (labelled “N”) schedulers. Figure 5 shows the results obtained using the same schedulers but this time while using Libception.

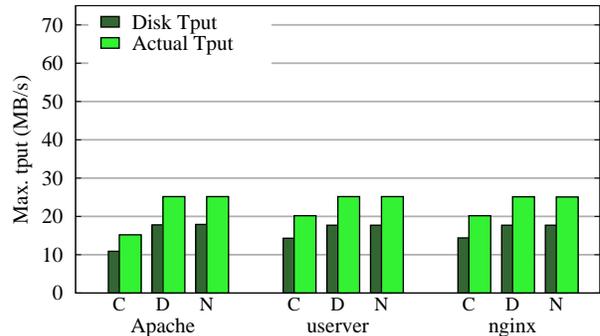


Figure 4: Servers without Libception using different block I/O schedulers

Interestingly, Figure 4 shows that without Libception the server throughput is slightly higher with the deadline and NOOP schedulers than with the CFQ scheduler. On the other hand, when using Libception (see Figure 5) all web

servers obtain the same maximum failure free throughput of nearly 50 MB/second regardless of the block I/O scheduler used. We believe that this is because Libception is serializing all of the reads that go to disk and as a result the schedulers only ever see a single outstanding request, which leaves no room for scheduling policies to make a difference. Note that we have spent some time attempting to adjust the parameters designed to control the behaviour of the deadline and CFQ block I/O schedulers. We did not see throughput improvements when compared with the default values.

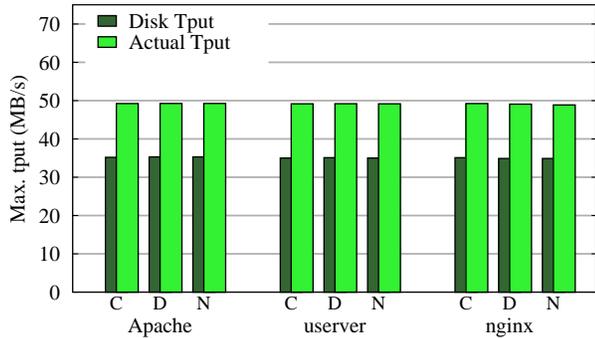


Figure 5: Servers with Libception using different block I/O schedulers

In summary, the different schedulers do not significantly improve web server throughput on this workload. More importantly, Libception provides significant increases in throughput that are not possible using the block I/O scheduling algorithms.

## 5.4 Evaluation Insights

In this section we first conduct a sequence of experiments designed to understand the relative importance of the prefetching and serialization components of Libception. For the remainder of the paper we focus solely on the nginx server and the CFQ block I/O scheduler. We chose nginx because it is used by Netflix for serving HTTP video streaming workloads. We chose CFQ because it is the default block I/O scheduler on our Linux server and because the schedulers did not significantly affect performance when using Libception (see Figure 5).

Figure 6 shows the maximum failure free throughput obtained using nginx without Libception (labelled “Vanilla”), with Libception using only serialization (labelled “Serialized”), with Libception using only prefetching (labelled “Prefetching”), and with Libception using both serialization and prefetching (labelled “Libception”). As can be seen in this figure, serialization alone actually reduces throughput. We believe that this is primarily due to the relatively small size of reads that the application performs. These small reads, in conjunction with serialization, result in small requests being issued one at a time, which causes very poor performance. On the other hand, using prefetching without serialization does significantly increase both disk and server throughput when compared with the “Vanilla” server. Finally, by combining both aggressive prefetching and serialization, a further increase of approximately 25% beyond that of prefetching is alone is achieved. These experiments demonstrate that, while aggressive prefetches are essential,

the full potential of Libception is not realized unless the requests to the disk are serialized.

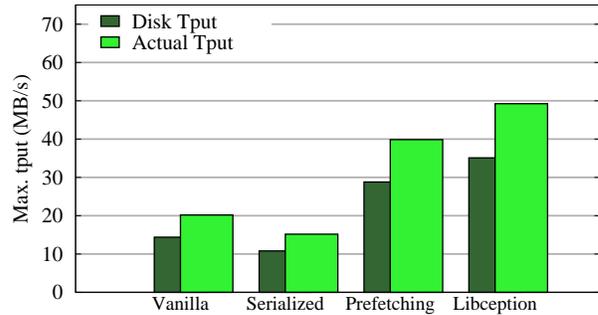


Figure 6: nginx Vanilla, and nginx with Libception using serialization only, prefetching only, and both

In previous work [33, 34] and in all experiments in this paper up to this point, the prefetch size was set to 2 MB. We now examine a range of prefetch sizes and study the impact on the throughput of nginx while using Libception with prefetching but without serialization (see Figure 7) and with both prefetching and serialization (see Figure 8)

Figure 7 shows that, without serialization, throughput does not improve until a prefetch size of 2 MB or greater is used. It also demonstrates that prefetch sizes of 3 and 4 MB provide slightly better throughput than a prefetch size of 2 MB.

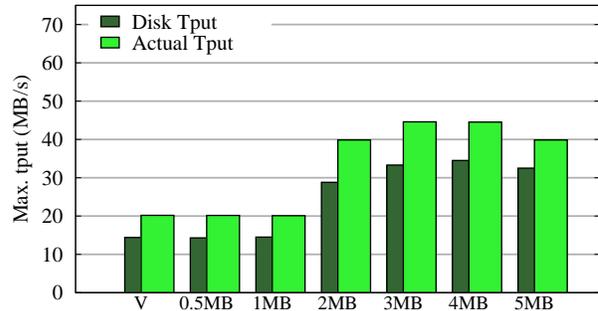


Figure 7: nginx with Libception using various prefetch sizes (no serialization)

Figure 8 shows that when using Libception, with both serialization and prefetching, small prefetch reads actually slightly reduce server throughput. However a prefetch size of 1 MB does significantly improve server throughput, which is not the case when prefetching is used without serialization. When using both serialization and prefetching, throughput peaks with prefetch sizes of 2 – 4 MB, and also shows that serialization provides additional benefits (about 10%) when compared with prefetching alone.

## 5.5 Using Insights to Improve Linux

We now use the insights obtained from the previous section to modify Linux kernel parameters in an attempt to improve the performance of nginx when running on Linux without the use of Libception. The question being examined

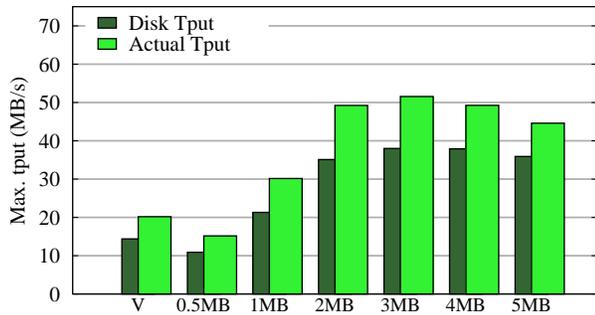


Figure 8: nginx with Libception using serialization and various prefetch sizes

is: can we find and tune appropriate Linux kernel parameters in order to obtain throughput that equals that obtained using Libception?

It was not too difficult to find a Linux kernel parameter that we were able to adjust to increase the amount of data being read from disk by the kernel. When obtaining data from disk the Linux kernel may (depending on several factors the details of which aren't relevant to this discussion) extend the read beyond what the user has requested. Roughly speaking, a readahead size is tracked per file and under proper conditions grows for each successive read. However, the maximum readahead size for all files on a device is limited by the kernel parameter `read_ahead_kb`, which can be set differently for each device. The default setting for this parameter for the version of Linux used in our experiments is 128 KB. This is quite small compared to the aggressive prefetches we use in Libception (e.g., 2 MB in many cases) in order to obtain significant increases in throughput.

Figure 9 demonstrates how the throughput of nginx changes as the value of the readahead parameter is increased. As an example of how we set the readahead size to 1024 KB on disk drive used to store video files (`/dev/sdb1`), we use the command `blockdev --setra 1024 /dev/sdb1`. As can be seen in this graph, `read_ahead_kb` values of 0.5 MB and 1 MB provide significant improvements over the default value of 128 KB (labelled "V" for Vanilla). Larger values for `read_ahead_kb` do not perform as well as 1 MB. While throughput obtained with 1 MB is about 45 MB/second it is not as high as that obtained using Libception (50 MB/second). Understanding why this is the case involved significantly more work.

Although the size of the prefetch for files were actually reaching the limit imposed by `read_ahead_kb`, using the Linux `blktrace` facility we were able to determine that requests to the disk were being limited to 512 KB. We believe that, as a result, some requests for reads to different files were being interleaved (this is similar to the Libception case where prefetching is used but serialization is not).

By examining the Linux kernel source we were eventually able to determine that another kernel parameter was placing further limits on the size of reads. This value `/sys/block/sdb/queue/max_sectors_kb` (for the disk `/dev/sdb`) uses the default size of 512 KB. We expect that the default values for these two limits (`read_ahead_kb` and

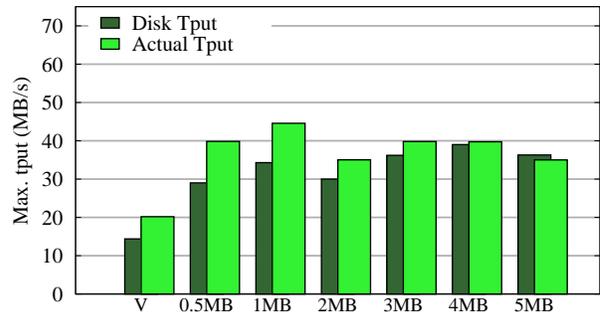


Figure 9: nginx without Libception using various `read_ahead_kb` sizes

`max_sectors_kb`) are chosen in order to ensure fairness across different disk requests and to keep latencies low.

To ensure that `max_sectors_kb` does not limit the size of disk reads we set its value to 16 MB. We then adjust `read_ahead_kb` and examine the throughput obtained by nginx. Figure 10 shows the results of these experiments and demonstrates the importance of setting both kernel parameters to proper values in order to obtain good throughput on this workload. In this case a `read_ahead_kb` size of 1 or 2 MB, now results in throughput of about 58 MB/second, which is slightly better than that obtained using Libception.

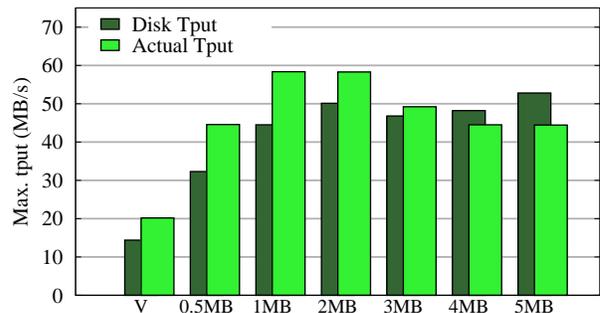


Figure 10: nginx without Libception using various `read_ahead_kb` sizes, with `max_sectors_kb = 16 MB`

Figure 11 now shows results obtained using Libception in addition to using modified Linux kernel parameters. In this case `max_sectors_kb` is set to 16 MB to ensure that it does not limit read sizes and the Libception prefetch size and `read_ahead_kb` are adjusted together using the values shown along the x-axis of the graphs. The column labelled "V" is showing the vanilla case where nginx runs without Libception and the default kernel parameter values are used. Although these results show a slight improvement in maximum failure free throughput when compared with just using Libception, they are still lower than adjusting the kernel parameters alone and not using Libception.

We observe that for some configurations of the experiments conducted in this section, the disk throughput exceeds the total server throughput. In these cases data that is being prefetched is actually being evicted before it is requested by and sent to the client. This means that, for the amount of memory in the current system, prefetching has become too

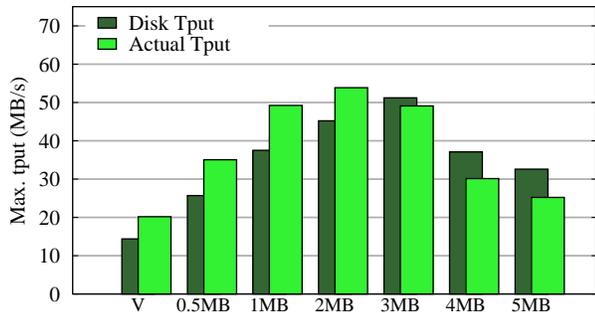


Figure 11: nginx with Libception using various prefetch and readAhead\_kb sizes, with max\_sectors\_kb = 16 MB

aggressive and is causing evictions. We expect that in these cases, increasing the amount of memory in the system will likely increase total server throughput because there will be less memory pressure created by aggressive prefetching. It seems that, in HTTP video streaming workloads, it may be more important to have memory acting as a buffer for large aggressive prefetching (in order to obtain high disk throughput) than as a file system cache.

## 5.6 Evaluating Latencies

When servicing HTTP video server workloads, server throughput is strongly influenced by disk throughput. This is because video services like YouTube and Netflix have large numbers of videos that are viewed infrequently (i.e., the popularity distribution of videos has a long tail). In order to improve disk throughput Libception prefetches relatively large amounts of data and serializes access to each disk. This naturally increases disk throughput while potentially increasing latencies for some requests. The potential for increased latencies should be relatively harmless when servicing video server workloads because clients are designed to be able to tolerate fairly significant latencies by using a play out buffer. The play out buffer is filled before play begins and is used to seamlessly continue playing the video even during periods where the client may experience latencies due to the network or HTTP server. For example, a client with a 10 second play out buffer can tolerate nearly 10 seconds of latency for some requests. As long as the data being requested arrives (and can be decoded before) within 10 seconds of when it is requested, the user will not experience any problems.

Before delving into further empirical analysis, it is important to note that all of our results place an implicit bound on latency. Client timeouts occur when a request has not received a corresponding response after 10 seconds. This latency includes the time to transmit the request, service the request at the server, and transmit the response. Therefore, if a client completes a session without errors, all of the latencies must have been acceptable.

In order to better understand the latencies incurred by using Libception and the aggressive tunings of Linux kernel parameters that support high server throughput, we now examine the latencies experienced by server processes. We use an existing system call tracing facility that exists in the userver to record the time required for every call to `sendfile` in memory and print those times to disk after the server

has finished servicing all requests. While client-side latency measurement would give an indication of end-user quality of experience, latency measurement at the server simplifies the collection process. Furthermore, as noted in Section 4, in our workload the network does not serve as a bottleneck. Therefore, server-side latencies should be reflective of client-side latencies, minus network transfer time.

Figure 12 shows the cumulative distribution function of `sendfile` call times without Libception and with default kernel parameters (labelled “Vanilla”), with Libception, and with the aggressive kernel parameter tunings. The “Vanilla” and “Libception” lines are created using data from one run with the userver using the configuration that obtains the highest error free rate (i.e., the two userver configurations used in Figure 3). The “Aggressive Kernel Params” line was created by setting `readAhead_kb` to 2 MB and `max_sectors_kb` to 16 MB and using the same request rate as used for the configuration of nginx shown in Figure 10.

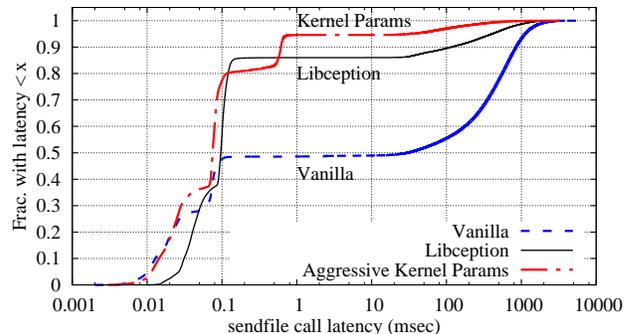


Figure 12: userver sendfile latency CDF

Figure 12 shows that Libception has the smallest density of requests that are serviced in the lowest range of latencies (i.e., below 0.1 ms). This is due to the extra overhead incurred from the `mincore` system call, which is used to determine if the data being requested is in memory or if it should be prefetched from disk. In cases where the data is already in memory, the lowest latencies are achieved by calling `sendfile` directly. For requests where the desired data is not in memory at the time of the `sendfile` call, the “Vanilla” configuration results in a blocking call to `sendfile` because the data being requested needs to be read from disk. When a disk read occurs using either Libception or aggressively tuned kernel parameters, the application-initiated read is serialized and increased in size (to 2 MB in this case) to implement prefetching. Although this incurs some overhead for that individual request, subsequent reads will often be served directly from the file system cache with comparatively low latencies. The net result is a significant reduction in latencies for a large number of `sendfile` calls. While about 49% of `sendfile` calls take less than 1 ms under the “Vanilla” configuration, about 85% of `sendfile` calls take less than 1 ms for the Libception configuration, and about 95% of `sendfile` calls take less than 1 ms for the aggressive kernel parameters configuration.

The key observation from this experiment is that, although one might expect latencies to increase because of large serialized prefetches, for this workload they actually decrease for a large number of data requests.

## 5.7 Evaluation Summary

Table 1 summarizes the main results from our experiments. Entries marked with an asterisk are not possible to obtain (e.g., the ASAP userver requires an option to `sendfile` that is only available on FreeBSD so it can't be run on Linux). Entries marked with “-” are not included due to space and/or because they are unlikely to provide new insights. On FreeBSD, the use of Libception more than doubled the peak server throughput for Apache, nginx and userver. We were not able to improve the poor performance achieved by these web servers without Libception by tuning FreeBSD kernel parameters. Comparing the results for the userver modified to directly incorporate aggressive prefetching and serialization (ASAP) shows that the overhead of implementing these techniques in a shim library, rather than directly, is negligible for our HTTP video streaming workload (on FreeBSD the throughput obtained with Libception for all servers matches that obtained with ASAP).

	Default	Libception	Kernel Tuning	Libception + Kernel Tuning
<b>FreeBSD</b>				
nginx	39.75	83.71	*	*
Apache	30.10	83.88	*	*
userver	39.74	83.56	*	*
ASAP	83.99	-	*	*
<b>Linux</b>				
nginx	25.11	49.06	58.29	53.85
Apache	25.18	49.26	-	-
userver	25.16	49.18	-	-
ASAP	*	*	*	*

**Table 1: Summary of Results: Throughput in MB/sec.**

These FreeBSD results raise the question of whether they might reflect some deficiency in the block I/O scheduler in that system. This prompted us to take a careful look at performance on Linux, which supports three different block I/O schedulers. We found that use of Libception approximately doubled the peak server throughput on Linux, for all three web servers, regardless of the choice of block I/O scheduler. However, using our insights from Libception as a guide and in some cases examining the Linux kernel source code, we were able to discover kernel parameters that we could tune to obtain slightly better performance than that with Libception. In contrast to Libception which is highly portable, we could not run our modified version of the userver on Linux since it makes use of a system call option that is not available on that system. Finally, a potential concern might be that these throughput improvements have significant cost with respect to latency, but as shown in Section 5.6 this does not appear to be the case.

## 6. DISCUSSION

Perhaps surprisingly, our results show that, without Libception, none of the web servers we investigated yield good performance for HTTP video streaming workloads, on either FreeBSD, or Linux with default kernel parameter settings. This finding suggests that web server implementations are still optimized for more traditional web workloads,

despite the rapid growth of HTTP video streaming. Also, web server developers might assume that, after many years of research and development, prefetching (or readahead) in operating systems, to efficiently use a magnetic disk is a “solved problem”, but evidently this is not the case. Although with Linux, it was eventually possible to achieve good performance after kernel parameter tuning, it is noteworthy that such manual tuning was required. For services like Netflix and YouTube, due to the large amounts of video available, the wide variety of bit rates at which each video is encoded, and because of the large number of videos that are viewed infrequently, it is not economically viable to store it all on SSDs. As a result, obtaining good performance when servicing video from disks is still important in these settings. We have found Libception to be a relatively simple and portable platform for implementing and evaluating techniques for improving disk I/O efficiency. In particular, using Libception we were able to readily evaluate the benefits of aggressive prefetching and I/O serialization for HTTP video streaming workloads, using multiple web servers and operating systems.

As demonstrated in Section 5.4 and Section 5.5, selecting a prefetch size which is too small results in low disk throughput. Likewise, selecting a prefetch size which is too large leads to a drop in system throughput. As a result obtaining peak server throughput requires choosing the most appropriate prefetch size. In previous work we have demonstrated how the best prefetch size and the benefits obtained from prefetching are sensitive to workload and system properties [32]. We show how the best prefetch size can be affected by the amount of available system memory, the distribution of the popularity of videos requested, hard drive characteristics, and the bit rates of files being served.

In order to avoid having to exhaustively and repeatedly determine the best prefetch size when workload or system characteristics change, we have designed an algorithm for dynamically and automatically adjusting the prefetch size with the goal of obtaining peak server throughput [32]. We demonstrate that a gradient descent algorithm, which minimizes a score based on a combination of disk transfer times and file system cache miss ratio, is effective at selecting prefetching sizes that result in high server throughput. This adaptive algorithm results in throughputs that rival those obtained by exhaustive manual tuning across a variety of different workload and system characteristics. We believe that such a strategy could be added to Libception, allowing it to continue providing high throughput while also eliminating the need to manually set a prefetch size.

In more recent work [30, 31] we have analyzed log files of HTTP requests to characterize the workloads of two different types of production Netflix servers. We use simulation to show that workload-specific adjustments can be used to increase server throughput. However, some of these improvements require knowledge of request streams and other information about the workload that may be difficult to infer in Libception. An interesting question is how much information can be provided to Libception, without requiring changes to web server code and whether or not it such changes can compete with the performance that can be obtained by directly modifying the web server. We intend to explore such questions in future work.

## 7. CONCLUSIONS

HTTP video streaming has become an important class of web server workloads. Such workloads have quite different characteristics than the web server workloads that have been the focus of most prior work on web server performance. In particular, the requests to HTTP video streaming servers are for large chunks of data commonly stored on disk, so these servers are frequently disk-bound.

In this work, we have designed, implemented and evaluated Libception, an application-level shim library incorporating techniques for improving disk access efficiency. HTTP video streaming servers can achieve the benefits of these techniques simply by linking with Libception, without the need for any source code modification. Experiments with three web servers (Apache, nginx and the userver) and two operating systems (FreeBSD and Linux) showed that with the aggressive prefetching and disk I/O serialization techniques currently implemented in Libception, peak server throughput can be doubled. Only on Linux with kernel parameters tuned for this workload was it possible to achieve performance competitive with Libception. We believe that Libception could be fruitfully applied to investigate other techniques for improving HTTP video streaming performance, and possibly for improving the performance of other disk-intensive applications.

## Acknowledgments

Brecht, Eager, and Wong thank the Natural Sciences and Engineering Research Council (NSERC) of Canada for partial support for this project through Discovery Grants. Brecht has also received an NSERC Discovery Accelerator Supplement in support of this work. Cassell, Summers and Szepesi were partially supported by NSERC graduate scholarships and Cassell and Summers were partially supported by a University of Waterloo Cheriton Scholarship. The authors would also like to thank the anonymous reviewers for their comments.

## 8. REFERENCES

- [1] ADHIKARI, V. K., GUO, Y., HAO, F., VARVELLO, M., HILT, V., STEINER, M., AND ZHANG, Z.-L. Unreeling Netflix: Understanding and improving multi-CDN movie delivery. In *Proc. IEEE International Conference on Computer Communications (INFOCOM)* (2012), pp. 1620–1628.
- [2] AXBOE, J. Linux block IO – present and future. In *Proc. Ottawa Linux Symposium (OLS)* (2004), pp. 51–61.
- [3] AXBOE, J. Linux kernel Git commit. <http://git.kernel.org/cgiit/linux/kernel/git/torvalds/linux.git/commit/?id=492af6350a5ccf087e4964104a276ed358811458>, 2009.
- [4] BEGEN, A., AKGUL, T., AND BAUGHER, M. Watching video over the web: Part 1: Streaming protocols. *IEEE Internet Computing* 15, 2 (2011), 54–63.
- [5] BUTT, A. R., GNIADY, C., AND HU, Y. C. The performance impact of kernel prefetching on buffer cache replacement algorithms. In *Proc. ACM SIGMETRICS International Conference on Measurement and Modeling of Computer Systems* (2005), pp. 157–168.
- [6] CAO, P., FELTEN, E. W., KARLIN, A. R., AND LI, K. Implementation and performance of integrated application-controlled file caching, prefetching, and disk scheduling. *ACM Transactions on Computer Systems (TOCS)* 14, 4 (1996), 311–343.
- [7] DHAGE, S. N., PATIL, S. K., AND MESHAM, B. B. Survey on: Interactive video-on-demand (VoD) systems. In *Proc. IEEE International Conference on Circuits, Systems Communication and Information Technology Applications (CSCITA)* (2014), pp. 435–440.
- [8] FINAMORE, A., MELLIA, M., MUNAFÒ, M. M., TORRES, R., AND RAO, S. G. YouTube everywhere: Impact of device and infrastructure synergies on user experience. In *Proc. ACM SIGCOMM Internet Measurement Conference (IMC)* (2011), pp. 345–360.
- [9] GHOSE, D., AND KIM, H. J. Scheduling video streams in video-on-demand systems: A survey. *Springer Multimedia Tools and Applications* 11, 2 (2000), 167–195.
- [10] GILL, P., ARLITT, M., LI, Z., AND MAHANTI, A. Youtube traffic characterization: a view from the edge. In *Proc. ACM SIGCOMM Internet Measurement Conference (IMC)* (2007), pp. 15–28.
- [11] IYER, S., AND DRUSCHEL, P. Anticipatory scheduling: A disk scheduling framework to overcome deceptive idleness in synchronous I/O. In *ACM Symposium on Operating Systems Principles (SOSP)* (2001), pp. 117–130.
- [12] JIANG, S., DING, X., XU, Y., AND DAVIS, K. A prefetching scheme exploiting both data layout and access history on disk. *ACM Transactions on Storage (TOS)* 9, 3 (2013), 10:1–10:23.
- [13] KASBEKAR, M. On efficient delivery of web content. GreenMetrics Keynote Talk, 2010.
- [14] LI, C., SHEN, K., AND PAPATHANASIOU, A. E. Competitive prefetching for concurrent sequential I/O. In *Proc. ACM European Conference on Computer Systems (EuroSys)* (2007), pp. 189–202.
- [15] MOSBERGER, D., AND JIN, T. `httperf` – a tool for measuring web server performance. *ACM SIGMETRICS Performance Evaluation Review (PER)* 26, 3 (1998), 31–37.
- [16] NETFLIX. Appliance hardware. [https://openconnect.netflix.com/en\\_gb/hardware](https://openconnect.netflix.com/en_gb/hardware), 2017.
- [17] NETFLIX. Appliance software. [https://openconnect.netflix.com/en\\_gb/software](https://openconnect.netflix.com/en_gb/software), 2017.
- [18] NETFLIX. Requirements for deploying embedded appliances. [https://openconnect.netflix.com/en\\_gb/requirements-for-deploying](https://openconnect.netflix.com/en_gb/requirements-for-deploying), 2017.
- [19] PANAGIOTAKIS, G., FLOURIS, M. D., AND BILAS, A. Reducing disk I/O performance sensitivity for large numbers of sequential streams. In *Proc. IEEE International Conference on Distributed Computing Systems (ICDCS)* (2009), pp. 22–31.
- [20] PAPATHANASIOU, A. E., AND SCOTT, M. L. Aggressive prefetching: An idea whose time has come. In *Proc. USENIX Workshop on Hot Topics in Operating Systems (HotOS)* (2005).
- [21] PARIAG, D., BRECHT, T., HARJI, A., BUHR, P., SHUKLA, A., AND CHERITON, D. R. Comparing the

- performance of web server architectures. In *Proc. ACM European Conference on Computer Systems (EuroSys)* (2007), pp. 231–243.
- [22] PATTERSON, R. H., GIBSON, G. A., GINTING, E., STODOLSKY, D., AND ZELENKA, J. Informed prefetching and caching. In *Proc. ACM Symposium on Operating Systems Principles (SOSP)* (1995), pp. 79–95.
- [23] RIZZO, L. Dummynet: a simple approach to the evaluation of network protocols. *ACM SIGCOMM Computer Communication Review (CCR)* 27, 1 (1997), 31–41.
- [24] RUAN, Y., AND PAI, V. S. Making the “Box” transparent: System call performance as a first-class result. In *Proc. USENIX Annual Technical Conference (ATC)* (2004).
- [25] RUAN, Y., AND PAI, V. S. Understanding and addressing blocking-induced network server latency. In *Proc. USENIX Annual Technical Conference (ATC)* (2006).
- [26] RUEMMLER, C., , AND WILKES, J. An introduction to disk drive modeling. *IEEE Computer* 27, 3 (1994), 17–28.
- [27] RUEMMLER, C., AND WILKES, J. UNIX disk access patterns. In *Proc. USENIX Winter Conference* (1993), pp. 405–420.
- [28] SANDVINE. Global Internet phenomena report, 2012.
- [29] STEINMETZ, R. Multimedia file systems survey: Approaches for continuous media disk scheduling. *Elsevier Computer Communications* 18, 3 (1995), 133–144.
- [30] SUMMERS, J. *Understanding and Efficiently Servicing HTTP Streaming Video Workloads*. PhD thesis, University of Waterloo, 2016. <https://uwspace.uwaterloo.ca/handle/10012/10956>.
- [31] SUMMERS, J., BRECHT, T., EAGER, D., AND GUTARIN, A. Characterizing the workload of a Netflix streaming video server. In *Proc. IEEE International Symposium on Workload Characterization (IISWC)* (2016).
- [32] SUMMERS, J., BRECHT, T., EAGER, D., SZEPESI, T., CASSELL, B., AND WONG, B. Automated control of aggressive prefetching for HTTP streaming video servers. In *Proc. ACM International Conference on Systems and Storage (SYSTOR)* (2014), pp. 5:1–5:11.
- [33] SUMMERS, J., BRECHT, T., EAGER, D., AND WONG, B. Methodologies for generating HTTP streaming video workloads to evaluate web server performance. In *Proc. ACM International Systems and Storage Conference (SYSTOR)* (2012), pp. 2:1–2:12.
- [34] SUMMERS, J., BRECHT, T., EAGER, D., AND WONG, B. To chunk or not to chunk: Implications for HTTP streaming video server performance. In *Proc. ACM International Workshop on Network and Operating System Support for Digital Audio and Video (NOSSDAV)* (2012), pp. 15–20.
- [35] VANDEBOGART, S., FROST, C., AND KOHLER, E. Reducing seek overhead with application-directed prefetching. In *Proc. USENIX Annual Technical Conference (ATC)* (2009).
- [36] VARKI, E., HUBBE, A., AND MERCHANT, A. Improve prefetch performance by splitting the cache replacement queue. In *Proc. IEEE International Conference on Advanced Infocomm Technology (ICAIT)* (2012), pp. 98–108.
- [37] WACHS, M., XU, L., KANEVSKY, A., AND GANGER, G. R. Exertion-based billing for cloud storage access. In *Proc. USENIX Workshop on Hot Topics in Cloud Computing (HotCloud)* (2011).