

Developing a High-Performance Web Server in Concurrent Haskell

Simon Marlow

Microsoft Research Ltd., Cambridge
simonmar@microsoft.com

Abstract

Server applications, and in particular network-based server applications, place a unique combination of demands on a programming language: lightweight concurrency, high I/O throughput, and fault tolerance are all important.

This paper describes a prototype web server written in Concurrent Haskell (with extensions), and presents two useful results: firstly, a conforming server could be written with minimal effort, leading to an implementation in less than 1500 lines of code, and secondly the naive implementation produced reasonable performance. Furthermore, making minor modifications to a few time-critical components improved performance to a level acceptable for anything but the most heavily loaded web servers.

1 Introduction

The Internet has spawned its own application domain: multithreaded server applications, capable of interacting with hundreds or thousands of clients simultaneously, are becoming increasingly important. Examples include FTP (File Transfer Protocol), E-Mail transport, DNS (name servers), Usenet News, chat servers, distributed file-sharing, and the most popular of all: web servers.

The basic function of a web server is to service requests for files on the local file system. A client contacts the web server over the network, sends a request containing the name of the required file, and the server responds with the contents of the file, if it exists. In reality web servers may provide a great deal more functionality, such as allowing the web site manager to design pages with dynamic content with scripts that are run by the web server, but in this paper we are only concerned with the basic file-serving facilities.

Nevertheless, even providing this basic service leads to a number of problems on the implementation front, if we are to produce a server that is capable of realistically running a real web site:

- The server has to be able to communicate with several clients simultaneously, because we can't afford for the server to be unavailable for long periods while it waits for the current client to complete its transaction. Moreover, the server should scale well to a large number of simultaneous clients. In implementation terms, this usually means that we have to employ some form of concurrency;

the choice of concurrency model tends to have the largest impact on a server's performance, so we discuss the relative merits of several concurrency models in Section 2.

- For small requests, the server must be able to process the request quickly to avoid tying up resources for too long. This means the latency for a new connection should be as low as possible. This aspect of a server's performance can be measured by firing small requests at it at an increasing rate; at some point, the server's performance will start to drop off as the rate increases beyond the latency and a backlog of connections in progress builds up. Another advantage of low latency is that the server is more resistant to denial-of-service type attacks. If it can throw out bogus requests as quickly as possible, then the overall server performance will be impacted less when the server is flooded with requests from a malevolent client.
- Fault tolerance is as important as performance: the server should be able to recover gracefully from errors, and never crash (or if it does, arrange that it restarts itself so that there is a minimal interruption in service). It should also be possible to re-configure the server, perhaps to add or remove content, without taking it down.

So why write a web server in Haskell? Firstly, because we can! It's no bad thing for Haskell if we can compete effectively in important application domains such as web serving. Secondly, several recent extensions to Haskell are particularly useful for this problem domain:

- Concurrent Haskell (Peyton Jones *et al.*, 1996) provides a lightweight concurrency model that helps to provide the low latency and low overhead for multiple simultaneous clients that is essential for good server performance.
- Recent extensions to support exceptions (Peyton Jones *et al.*, 1999) provide useful facilities for coping with run-time errors.
- Asynchronous Exceptions (Marlow *et al.*, 2001) allow concise implementations of important features such as timeouts, as we shall see in Section 4. Asynchronous exceptions are also useful for allowing the server to respond to external stimulus, such as when the administrator wishes to make changes to the server's configuration (see Section 6).
- We used a wide range of libraries in constructing the web server, including a networking library, a parsing combinator library, an HTML generation library and a POSIX system interface library. None of these libraries are specified as part of the Haskell standard as yet, but all are distributed with the GHC compiler.

This paper describes our web server implementation, focussing on the parts which required extensions to Haskell, in particular the concurrency and exception support. The server implementation performs well, as we shall show in Section 7. It is also reliable - we tested it as an alternate server for the `haskell.org` site, where it ran for a week, collected about 2000 hits during that time, and kept within a memory footprint of 3M. This test uncovered one bug in the HTTP protocol implementation, which has since been fixed. We hope to replace the main `haskell.org` web

server (currently Apache) with the Haskell implementation at some point in the future.

2 Concurrency Model

The choice of concurrency model is crucial in the design of a web server. In this section we briefly examine the common models in use and list their advantages and disadvantages, and compare them with Concurrent Haskell's approach.

2.1 *Separate Processes*

This is the model used by Apache, where each new connection is handled by a new operating system process on the machine. A single top-level process monitors incoming connections and spawns a new worker-process to talk to each client.

In terms of our requirements, the separate process model shapes up thus:

- **Scalability and Latency.** By virtue of the fact that the operating system is handling the concurrency, the separate process model will automatically take advantage of multiple processors if the operating system supports them. However, processes are rather heavyweight beasts in terms of the startup and shutdown cost, context-switch overhead, and memory overhead for each process, so latency is likely to be high. Systems which use this model, like Apache, tend to keep a cache of available processes and use a single process to serve multiple serial requests, in order to alleviate the high startup costs.
- **Fault tolerance.** The operating system's memory protection provides a certain amount of fault tolerance, in that a single crashing process can't affect any other process.
- **Programming model.** The separate process model is fairly simple to implement, but communication between processes tends to be inconvenient because of the explicit nature of the interprocess communication methods provided by most operating systems. Interprocess communication is required for certain aspects of a web server's operation. For example: there is normally a single log file which records transaction events (see Section 5), and hence multiple threads must cooperate for access to the log file.

2.2 *Operating-System Threads*

Depending on the particular operating system, OS threads may map onto processes, in which case the overhead will be similar as for separate processes (modern operating systems will share page tables between forked processes in any case), or they may be implemented as light-weight kernel threads, in which case the overhead will be lower. A light-weight kernel threads implementation has slightly different characteristics compared to processes:

- **Scalability and latency.** Performance is expected to be better than processes, because threads generally have less overhead than full processes. Again,

the operating system may be able to take advantage of multiple processors without any work on the part of the programmer.

- **Fault tolerance.** There is no memory protection, and no explicit support for fault tolerance in this model.
- **Programming model.** Inter-thread communication is somewhat easier, but dealing with thread-local state is harder (data is thread-global by default).

2.3 Monolithic Process with I/O multiplexing

Another approach is to implement the desired concurrency directly, foregoing any time-sharing facilities provided by the operating system itself. The single requirement for this approach is to be able to multiplex several I/O channels.

The existing methods for multiplexing I/O in a single process include:

- Use POSIX's `select()` (or equivalently `poll()`) functions. These functions tests multiple file descriptors simultaneously, returning information on which of the descriptors are available for reading or writing. The idea is that the application then performs any available reads and writes (using non-blocking I/O), and then returns to call `select()` on the list of open file descriptors again.

This approach suffers from the problem that `select()` is $O(n)$, where n is the number of file descriptors being tested, because the application must build up a list of length n to pass to the `select` function, and the OS must traverse this list to build up the results.

Another problem with non-blocking I/O is that it doesn't normally apply to disk I/O: so a web server using non-blocking I/O and `select()` could become effectively single-threaded while reading from disk.

- Asynchronous I/O, POSIX real-time signals, and kernel event queues. These are all methods of alleviating the aforementioned problems with `select()`. They are relatively new, non-standard features which are not supported by all operating systems. However, any implementation which uses `select()` can be converted to use one of these alternatives with relatively little effort.

The characteristics of these methods are:

- **Scalability and latency.** Very low latency, because of the lack of threading overhead. Scalability can suffer if the `select` method is used.
- **Fault tolerance.** There is no explicit support for fault tolerance in this model.
- **Programming model.** A web server is an inherently multi-threaded application, so programming without a concurrency abstraction is bound to be painful.

There exist several web servers which use these methods, and they are currently the fastest servers around.

2.4 User-space threads

User-space threads are essentially an implementation of a thread abstraction inside a single process (or possibly on top of a small number of operating system threads; see later). The programmer gets to write his/her application using the concurrency primitives provided by the language, and the user-space threads implementation will provide the low-level time-sharing and I/O multiplexing support. Several implementations of user-space POSIX threads exist for Unix.

Concurrent Haskell (or at least the implementation in GHC) is also an instance of this model; the Haskell runtime system runs in a single operating system process and multiplexes many Haskell threads. To support multiple Haskell threads performing I/O simultaneously, the runtime system may choose between the I/O multiplexing options described in the previous section. However this is implemented, the choice is invisible to the programmer.

The characteristics of this model combine the performance of the monolithic process model with the programming model of your chosen threading library. The concurrency is lightweight, and the programmer doesn't need to be concerned with the details of I/O multiplexing. There is one disadvantage, though: a user-space threads package won't normally be able to take advantage of multiple processors on the host machine. The GHC development team are currently working on an implementation of Concurrent Haskell that doesn't suffer from this deficiency, by using a small number of operating-system threads to share the load.

3 Implementing a web server in Haskell

A web server can be thought of as two components: an implementation of the HTTP protocol, and a top-level loop that continually accepts new connections and initiates new transactions. In this section we describe the structure of the Haskell implementation of the web server including code snippets from the implementation, beginning with the top-level loop and then going on to describe the implementation of the HTTP protocol.

3.1 The main loop

We will write the code to use GHC's networking library, `Socket`, which provides an interface similar to the traditional socket API present on most operating systems. The functions we need to know about are:

```
listenOn :: PortId -> IO Socket
accept   :: Socket -> IO (Handle, HostAddress)
```

The general pattern followed by a server application is to open a listening socket using `listenOn`, which returns an abstract value of type `Socket`. We can wait for a new connection request on this socket using the `accept` function; `accept` waits until a connection request is received, then returns a `Handle` for communication with the remote host (a normal Haskell read/write handle), and the address of the remote host.

Here is the main loop of the web server:

```
acceptConnections :: Config -> Socket -> IO ()
acceptConnections conf sock = do
    (handle, host_addr) <- accept sock
    forkIO (catch
        (talk conf handle host_addr 'finally' hClose handle)
        (\e -> logError e)
    )
    acceptConnections conf sock
```

`acceptConnections` takes a server configuration of type `Config` and a listening socket, and calls `accept` to wait for a new connection request on the socket. When a connection request is received, a new worker thread is forked using `forkIO`, and the main loop goes back to waiting for connections.

The worker thread calls `talk` (the definition of `talk` is given in the next section), which is the main function for communicating in HTTP with a client. The interesting part here is what happens if an exception is raised during `talk`. The `finally` combinator allows strict sequencing to be specified, independent of exceptions:

```
finally :: IO a -> IO b -> IO a
```

This combinator behaves much like `finally` in Java. It performs its first argument, then performs its second argument (even if the first argument raised an exception), then returns the value of the first argument (or re-raises the exception).

In the main loop above, we're using `finally` to ensure that the socket to the client is properly closed down if we encounter an error of any kind, including a bug in our code. Although the Haskell runtime system will automatically close files which are determined to be unused, it is beneficial to close them down as early as possible in order to free up the resources associated with the file.

The call to `talk` is also enclosed in a `catch` combinator:

```
catch :: IO a -> (Exception -> IO a) -> IO a
```

This combinator performs its first argument, and if an exception is raised, passes it to the second argument (the *exception handler*), otherwise it returns the result. In contrast to `finally`, `catch` specifies an action to be performed only when an exception is raised, whereas `finally` specifies an action which is always to be executed.

The code for `acceptConnections` uses `catch` to catch any errors and log them to the error log file, which we describe in Section 5.

3.2 HTTP protocol implementation

The HTTP protocol is essentially transaction-based. The client first opens a connection to the server and sends a request message. The server interprets the request and, if the request refers to a valid document on the server, replies to the client sending it the contents of the document. In early versions of the HTTP protocol,

the connection between client and server would be closed at this point, requiring the client to open a new connection for each document request. The latest revision of the HTTP protocol allows the connection to remain open for further transfers (known as a keep-alive connection).

So, serving a request is a simple pipeline:

1. read the request from the socket,
2. parse the request,
3. generate the response,
4. send the response back to the client,
5. if the connection is to be kept alive, return to step 1.

Reading the request from the socket is performed by `getRequest`:

```
getRequest :: Handle -> IO [String]
```

which takes the file handle representing the socket on which communication with the client is taking place, and returns a list of strings, each one being a single line of the request. A typical request looks something like this:

```
GET /index.html HTTP/1.1
Host: www.haskell.org
Date: Wed May 31 11:08:40 GMT 2000
```

The first line gives the command (GET in this case), the name of the object requested, and the version of the HTTP protocol being used by the client. Subsequent lines, termed *headers*, give additional information, and are mostly optional. The server is required to ignore any headers it doesn't understand.

The next stage is to parse the request into a `Request`:

```
data Request = Request {
    reqCmd      :: RequestCmd,
    reqURI      :: ReqURI,
    reqHTTPVer  :: HTTPVersion,
    reqHeaders  :: [RequestHeader]
}
```

The `Request` record contains elements for the command name (GET in the above example), the requested URI¹, the HTTP protocol version being used by the client, and a list of optional headers. The server is required to interpret requests differently depending on the protocol version being used by the client, although it can respond using its native protocol version; the protocol is designed to be upwards-compatible.

Requests are parsed by `parseRequest`:

```
parseRequest :: Config -> [String] -> Either Response Request
```

¹ Universal Resource Indicator, a more general form of URL

Note that parsing the request may return a response: this indicates failure, and the response will in most cases be a “Bad Request” response, but may be something more specific.

Next, we generate the response:

```
data Response = Response {
    respCode      :: Int,
    respHeaders  :: [ResponseHeader],
    respCoding    :: [TransferCoding],
    respBody      :: ResponseBody,
    respSendBody :: Bool
}
```

```
data ResponseBody
  = NoBody
  | FileBody Integer{-size-} FilePath
  | HTMLBody HTML
```

```
genResponse :: Config -> Request -> IO Response
```

`genResponse` performs a number of checks on the validity of the request, and generates an appropriate response. A valid GET request will result in a response with a `FileBody`. An invalid request will result in an error response of some description, with some automatically generated HTML describing the error, in an `HTMLBody`. If we later extend the server to support dynamically generated pages for example, then the `ResponseBody` datatype could be extended with further constructors for different types of content.

In the common case where the response body consists of an entire file verbatim, the `respBody` component of the `Response` structure doesn't contain the entire file body as a string, rather it contains just the path to the file. This is so that we can use more efficient methods for sending the file to the client than simply converting the contents to and from a `String`.

The final step is to send the response to the client:

```
sendResponse :: Config -> Handle -> Response -> IO ()
```

Pulling all this together, the top-level talk function looks like this:

```
talk :: Config -> Handle -> HostAddress -> IO ()
talk conf handle host_addr
  = do strs <- getRequest handle
      case parseRequest strs of
        Left resp -> sendResponse conf handle resp
        Right req -> do
          resp <- genResponse conf req
          sendResponse conf handle resp
          logAccess req resp host_addr
```



```
if (isKeepAlive req)
  then talk conf handle haddr
  else return ()
```

In reality, there is some extra code to deal with catching and logging of errors (Section 5.1) and timeouts (Section 4) in there too.

The call to `logAccess` causes an entry to be written to the log file describing the transaction, see Section 5. The real implementation also takes timestamps before and after the transaction, and passes the time difference to `logAccess` (timing facilities are available from the standard Haskell Time library).

4 Timeouts

A web server needs some form of timeout mechanism, so that clients which hang or take an inordinately long time to respond can be disconnected, and the resources associated with the connection freed back to the system.

We need a generic time-out combinator, with the following type:

```
timeout :: Int    -- timeout in seconds
        -> IO a   -- action to run
        -> IO a   -- action to run on timeout
        -> IO a
```

The application `timeout t a b` should behave as follows: `a` is run until it either completes, or `t` seconds passes. If it completes in time `t`, `timeout` returns the result immediately, otherwise `a` is terminated with an exception and `b` is executed. If `a` (or `b`, in the case of a timeout) raises an exception, then the exception will be propagated by `timeout`. The `timeout` function has no other side effects, so timeouts can be nested arbitrarily.

Thanks to asynchronous exceptions (Marlow *et al.*, 2001), we can implement a `timeout` combinator with the above properties. Note that because the action `a` can be terminated with an asynchronous exception at any time, it should be *exception safe*, that is it must be sure not to leave any mutable data structures in an inconsistent state or leak any resources. In fact, all our code should be written to be exception-safe, because exceptions like stack overflow and heap overflow are delivered asynchronously.

There are two primitives which are important when writing exception-safe code:

```
block   :: IO a -> IO a
unblock :: IO a -> IO a
```

where `block a` executes `a` with asynchronous exceptions *blocked*, that is any thread wishing to raise an asynchronous exception in the current thread must wait until exceptions are unblocked again. Similarly, `unblock a` unblocks asynchronous exceptions during the execution of `a`. Applications of `block` and `unblock` can be arbitrarily nested. Here's an example of acquiring a lock, where the lock is represented by an `MVar`, `m`, such that the lock will always be released safely if we receive an exception:

```

block (do a <- takeMVar m
        unblock (...)
        'catch' (\e -> do putMVar m a; throw e)
        putMVar m a
    )

```

Use of combinators such as `finally` (described in the previous section), and `bracket`:

```
bracket :: IO a -> (a -> IO b) -> (a -> IO c) -> IO c
```

are also helpful in writing exception-safe code. For example, a simpler way of writing the above locking sequence is

```
bracket (takeMVar m) (putMVar m) (...)
```

The full story on asynchronous exceptions in Haskell, including the implementations of the `timeout`, `finally` and `bracket` combinators above, can be found in (Marlow *et al.*, 2001).

5 Logging

A web server normally produces log files listing all the requests made and certain information about the response sent by the server. Each entry in the log normally records

- the time the request was received,
- the requestor's address,
- the URL requested,
- the response code (either success or some kind of failure),
- the number of bytes transferred,
- the time taken for the request to complete,
- the client software's type & version,
- the referring URL.

The format of log entries is configurable, and may include other fields from the request or response. There exist standard log entry formats produced by the popular servers, and software available which processes the log files to produce reports. For this reason, we decided that the Haskell web server should be able to produce compatible logs.

A worker thread causes a log entry to be recorded by calling the function `logAccess`:

```
logAccess :: Request -> Response -> HostAddress -> TimeDiff -> IO ()
```

passing the request, the response, the address of the client and the time difference between the request being received and completion of the response.

The actual generation of the log entries and writing of log entries to the file is done by a separate thread². Worker threads communicate with the logging subsystem

² or threads, if resources are abundant!

via a global unbounded channel, by calling `logAccess`. The logging thread removes items from the channel and manufactures log entries which are then written to a log file.

Placing the logging in a separate thread is a good idea for several reasons:

- It helps reduce the total load on the system, because the worker threads can finish, and hence be garbage collected, before the log entry has been written.
- It means that a thread serving multiple requests can proceed immediately with the next request without waiting for the log entry of the first request to be written.
- The logging thread can batch multiple requests and write them out in one go, which is likely to be more efficient than writing them one at a time.

The logging thread is designed to be fault tolerant: if it receives an exception of any kind, it attempts to restart itself by re-opening the log file for writing, and continuing with the next log request. This behaviour is (ab)used by the main loop, which needs to restart the logging thread whenever it receives a request to re-read its configuration file.

5.1 Error logging

Logging of server errors is handled in a similar way to logging of requests. A separate thread writes log entries to an error log file, taking log requests from a global channel. Exception handlers are scattered around the main request/response handling code, which catch exceptions and log them with an informative message indicating where the error occurred, before passing the exception on to be handled at the top level.

The error logging thread also restarts itself if it receives an exception (and logs this event to the log file).

5.2 Global variables

In the above description of the logging threads, we mentioned that communication between a worker thread and a logging thread was via a “global channel”. How does one define a global channel in Haskell? This is one instance of a global mutable variable, a concept familiar to those who program in imperative languages, but until recently not available to Haskell programmers.

Global variables let you avoid the loss of modularity (and performance) that results from passing around extra parameters that are rarely needed. A cleaner alternative to using global variables would be to use implicit parameters (Lewis *et al.*, 2000), although implicit parameters still carry the performance cost of passing the additional parameters around.

A global `MVar` can be defined in Haskell as:

```
global_mvar :: MVar String
global_mvar = unsafePerformIO newEmptyMVar
```

Although we have declared `global_mvar` using `unsafePerformIO`, the declaration is normally perfectly safe (however, see later for caveats). We still access the `MVar` using the standard `putMVar` and `takeMVar` operations.

One can think of the `newEmptyMVar` as being executed at program initialisation time, but in fact it doesn't matter when the action is executed, as long as it happens before `global_mvar` is first accessed. In fact, the action will most probably happen lazily, being deferred until the first time `global_mvar` is demanded.

Global mutable objects are particularly useful when a single instance of a mutable variable is required, since they avoid the need to pass the object around explicitly. However, there are a few points to bear in mind:

- This use of `unsafePerformIO` is strictly speaking unsafe, because the program now behaves differently if occurrences of `global_mvar` are replaced with their values, namely `(unsafePerformIO newEmptyMVar)`. In order to stop this happening, we have to circumvent any optimisations in the compiler which may replace `global_var` with its value. In GHC, this amounts to adding the pragma

```
{-# NoInline global_mvar #-}
```

somewhere in the source code.

This stops the variable from being *unshared*, but the compiler may also decide to increase sharing by commoning up several global variable declarations with the same definitions (which is perfectly legal, of course). In GHC, we also have to disable this optimisation when compiling code with global variables.

- Care must be taken to give a type signature for the global variable and not to declare global mutable variables with polymorphically-typed contents. Type safety is in danger if this rule is broken, because the contents of a polymorphically-typed variable could be extracted and used at any type (this problem is described in more detail in (Launchbury *et al.*, 1999)).
- In a concurrent program, it is important to use `MVars` instead of just `IORefs` when multiple threads may have access to the variable (an `IORef` is a plain mutable variable, whereas an `MVar` adds synchronisation).

If we observe these rules, however, global mutable variables are a useful concept. We use global mutable variables in the web server in the following places:

- To store the channels by which the worker threads can communicate with the logging threads.
- To store the `ThreadId`s of the logging threads, so that the main thread can send them an exception to restart them.
- To store the command line options. The program is only started once, so this is a write-once mutable variable. Write-once mutable variables can semi-safely be read from pure, non-IO, code: if the initial value given to the mutable variable is \perp , then the program will fail immediately if it tries to access the variable before its value has been written.

6 Run-time configuration

Our web server is configured by editing a text file, in a similar way to other popular web servers. The syntax of the configuration file is similar to that of Apache's. When the server starts up, it parses the configuration file, and if there are no errors found, immediately starts serving requests.

In the interests of high availability, a web server should be run-time configurable. For example, when new content is placed on the server, the administrator somehow needs to inform the server that the new content is available and where to find it. It is occasionally necessary to change certain options, or tweak security settings, on a running server.

To take the server down and restart it with the new configuration would be unsatisfactory, because the site would be off-line during the restart. So a running server should be able to re-read its configuration file without interrupting operations. But what about transactions that are already in progress? Should they see the new configuration immediately?

In our server, we take the approach that the new configuration should only take effect for new connections, and existing connections should be allowed to continue using the old settings. This approach avoids a number of problems with changing configuration settings while a request is in progress: for example, if the security settings are changed such that a file being transmitted is no longer available to the client that requested it, should the transfer be terminated? We don't attempt to tackle this problem; our implementation requires you to restart the server if security settings change and/or existing connections need to be terminated.

Implementing run-time configuration updates in the Haskell web server turned out to be straightforward: as we've already seen, the key functions in the inner pipeline all take an argument of type `Config`, which contains the current configuration. The configuration is passed into the worker thread when it is created, so when the configuration changes all we need to do is ensure that any new threads receive the new configuration.

The approach we took is to send the main thread an asynchronous exception when it should re-read the configuration file. This gives us the option of having several ways to force a configuration change:

- A signal on Unix-like operating systems. This is the traditional way to kick a process into re-reading its configuration file, and consists of sending the process a signal from the command line. This method is implemented in our web server as follows: the incoming signal causes a new thread to start, which immediately sends an exception to the main thread. The main thread catches the signal and re-reads the configuration file.
- Implementing a proprietary HTTP command, which the administrator can use to re-configure the server on-line and remotely. Secure authorisation would certainly be needed if this method were to be used.
- Any other type of inter-process communication provided by the host operating system.

7 Performance Results

In this section we present our preliminary performance results for the Haskell web server.

7.1 Performance tweaks

We made several tweaks to the initial implementation of the server to remove some of the larger performance bottlenecks.

- We replaced the naive file transfer code which used `getContents` and `hPutStr`, with a version which does I/O directly to and from an array of bytes. GHC's `IOExts` library provides simple primitives for doing this.
- By default, GHC's scheduler context switches about 5000 times a second. We reduced this to something more reasonable, 50 times/sec, which made a substantial difference to the results. The reason is that GHC's scheduler currently does a `select` on every context switch to determine which I/O bound threads can be woken up. As discussed in Section 2, `select` is $O(n)$, so reducing the number of times we do it is a win when the system becomes heavily loaded.
- Tweaking the garbage collection settings had some effect, in particular increasing the allocation area size. GHC by default increases its heap usage in line with the program's demands, but giving the program more memory from the outset is usually a win, and was in this case.
- Reading a request turned out to be expensive, due to an inefficient, non-tail-recursive, implementation of `hGetLine` in GHC's I/O library. Rewriting this function improved performance by 10% or so.
- GHC's I/O library uses a system of finalizers to ensure that the buffers and other resources associated with a file descriptor are freed when the program releases the file handle. Finalizers are normally run in a thread by themselves, but this turned out to be expensive for the web server, since most connections give rise to two finalizers: one for the handle to the socket itself, and one for the file being transferred. We changed the finalisation mechanism to batch finalizers in a single thread after each garbage collection, which led to a small overall performance improvement.

Note that only one of these tweaks, namely the optimisation of the file transfer code, was made to the web server code itself, the rest were tweaks to GHC's runtime system and libraries. Indeed, the web server has been a useful source of insight into performance bottlenecks in GHC's concurrency and I/O support.

7.2 Connection latency

These measurements were made using `httperf` (Mosberger & Jin, n.d.), a tool which can be used for generating requests at a specific rate. It is used primarily for determining the rate of connection requests a web server can sustain before performance starts to drop off.

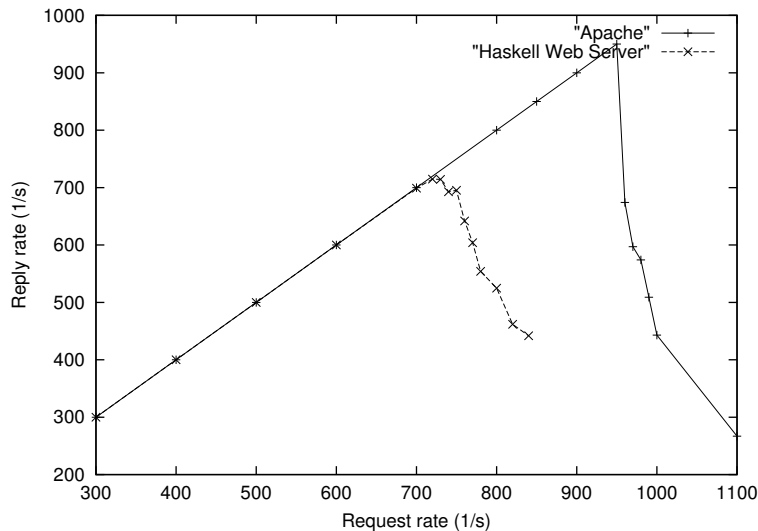


Figure 1. Connection latency results

In this test, the server machine was a single-processor PII/450 running Linux 2.2. The client was a separate machine on a local 100Mbit ethernet connection. The total number of requests sent was 4000 in each test. All the requests were for the same 1k file. The timeout on the client was set at 1 second.

A graph of reply rate against requests issued per second is given in Figure 1. This shows clearly how the server keeps up with the client until the request rate rises above the rate that the server can handle without accumulating a backlog (about 710 requests/second), at which point performance begins to decrease sharply. Why does performance decrease so dramatically? Two possible factors are:

- As connections in progress accumulate on the server, the $O(n)$ behaviour of `select()` as used by GHC's scheduler comes into play.
- As the number of threads in the system increases, thus the cost of garbage collection also increases. Garbage collection is necessarily $O(n)$ in the number of live threads, since it must traverse the active thread queues to determine which threads are live.

The server doesn't currently limit the number of connections in progress, and in fact at a rate of 850 connections/sec the number of concurrent connections observed on the server peaked at over 700 during the test. Setting a limit on the number of concurrent connections would help to flatten the graph after the drop-off point.

On the same hardware, Apache (the most commonly used web server software) tops out at 950 requests/second, and the drop-off is just as sharp. One way to flatten out the drop-off is to limit the number of active connections that the server

can process, and to stop listening for new connections when the limit is reached. Apache uses this technique but sets the limit high by default.

To put these figures into perspective, the most heavily loaded web servers on the net (eg. <http://www.yahoo.com/>) take an average of about 5000 hits/second, with peaks of probably 10000 hits/second. These sites use collections of identically configured servers with a load-balancing arrangement to spread the requests between the available machines.

However, for most sites on the net the performance turned in by our Haskell Web Server is more than adequate, and there's still plenty of opportunities for improvement: we haven't really made any attempt to optimise the code of the server itself, beyond fixing the slow file transfer.

8 Conclusions

The primary result presented here is that we constructed a web server in Haskell which conforms to the HTTP/1.1 standard (and more) in less than 1500 lines of Haskell (not including library code), and the resulting server performs admirably in real-world conditions. Furthermore, it is fault-tolerant and runs in a constant, and small, amount of memory over a sustained period.

In order to achieve this, we had to make use of a number of extensions to Haskell, the main ones being concurrency and exceptions. We also made use of a large amount of library code, all of which is part of GHC's library collection. The libraries we used include a networking library, a parsing combinator library, an HTML generation library and a POSIX interface library.

References

- Launchbury, J, Lewis, J, & Cook, B. (1999). On embedding a microarchitectural design language within haskell. *Pages 60–69 of: Acm sigplan international conference on functional programming (icfp'99)*. Paris: ACM Press.
- Lewis, J., Shields, M., Meijer, E., & Launchbury, J. 2000 (January 19–21). Implicit parameters: Dynamic scoping with static types. *27th annual acm sigplan-sigact symposium on principles of programming languages (popl'00)*.
- Marlow, Simon, Jones, Simon Peyton, Moran, Andrew, & Reppy, John. 2001 (June 20–22). Asynchronous exceptions in haskell. *Acm sigplan conference on programming language design and implementation*.
- Mosberger, David, & Jin, Tai. *httperf—a tool for measuring web server performance*. http://www.hpl.hp.com/personal/David_Mosberger/httperf.html. Hewlett-Packard Research Labs.
- Peyton Jones, S. L., Gordon, A.D., & Finne, S. (1996). Concurrent Haskell. *Pages 295–308 of: Proc. popl'96*. St. Petersburg, Florida: ACM Press.
- Peyton Jones, S. L., Reid, A., Hoare, T., Marlow, S., & Henderson, F. (1999). A semantics for imprecise exceptions. *Pages 25–36 of: Proc. pldi'99*. ACM SIGPLAN Notices, vol. 34(5). ACM Press.