

Flux: A Language for Programming High-Performance Servers

Brendan Burns Kevin Grimaldi Alexander Kostadinov Emery D. Berger Mark D. Corner

*Department of Computer Science
University of Massachusetts Amherst
Amherst, MA 01003*

{bburns, kgrimald, akostadi, emery, mcorner}@cs.umass.edu

Abstract

Programming high-performance server applications is challenging. It is both complicated and error-prone to write the concurrent code required to deliver high performance and scalability. Server performance bottlenecks are difficult to identify and correct. Finally, it is difficult to predict server performance prior to deployment.

This paper presents Flux, a language that dramatically simplifies the construction of scalable high-performance server applications. Flux lets programmers compose off-the-shelf, sequential C or C++ functions into concurrent servers. Flux programs are type-checked and guaranteed to be deadlock-free. We have built a number of servers in Flux, including a web server with PHP support, an image-rendering server, a BitTorrent peer, and a game server. These Flux servers match or exceed the performance of their counterparts written entirely in C. By tracking hot paths through a running server, Flux simplifies the identification of performance bottlenecks. The Flux compiler also automatically generates discrete event simulators that accurately predict actual server performance under load and with different hardware resources. Flux is both easy to use and unusually compact, allowing entire servers to be specified in tens of lines of code.

1 Introduction

Server applications need to provide high performance while handling large numbers of simultaneous requests. Programming servers remains a daunting task. Concurrency is required for high performance but introduces errors like race conditions and deadlock that are difficult to debug. The mingling of server logic with low-level systems programming complicates development and makes it difficult to understand and debug server applications. Consequently, the resulting implementations are often either lacking in performance, buggy or both. At the same time, the interleaving of multiple threads of server logic makes it difficult to identify performance bottlenecks or predict server performance prior to deployment.

This paper introduces *Flux*, a domain-specific language that addresses these problems¹. A Flux program

describes two things: (1) the flow of data from client requests through nodes, typically off-the-shelf C or C++ functions, and (2) mutual exclusion requirements for these nodes, expressed as high-level *concurrency constraints*. Flux requires no other typical programming language constructs like variables or loops – a Flux program executes inside an implicit infinite loop. The Flux compiler combines the C/C++ components into a high performance server using just the flow connectivity and concurrency constraints.

Flux captures a programming pattern common to server applications: concurrent executions, each based on a client request from the network and a subsequent response. This focus enables numerous advantages over conventional server programming:

- **Ease of use.** Flux is a declarative, implicitly-parallel language that eliminates the error-prone management of concurrency via threads or locks. A typical Flux server consists of just tens of lines of code.
- **Reuse.** By design, Flux directly supports the incorporation of unmodified existing code. There is no “Flux API” that a component must adhere to; as long as components follow the standard UNIX conventions, they can be incorporated unchanged. As an example, we were able to add PHP support to our web server just by implementing a required PHP interface layer.
- **Runtime independence.** Because Flux is not tied to any particular runtime model, it is possible to deploy Flux programs on a wide variety of runtime systems. Section 3 describes three runtimes we have implemented: thread-based, thread pool, and event-driven.
- **Correctness.** Flux programs are type-checked to ensure their compositions make sense. The concurrency constraints eliminate deadlock by enforcing a canonical ordering for lock acquisitions.
- **Performance prediction.** The Flux compiler optionally outputs a discrete event simulator. As we show in Section 5, this simulator accurately predicts actual server performance.

- **Bottleneck analysis.** Flux servers include light-weight instrumentation that identifies the most-frequently executed or most expensive paths in a running Flux application.

Our experience with Flux has been positive. We have implemented a wide range of server applications in Flux: a web server with PHP support, a BitTorrent peer, an image server, and a multi-player online game server. The longest of these consists of fewer than 100 lines of code, with the majority of the code devoted to type signatures. In every case, the performance of these Flux servers matches or exceeds that of their hand-written counterparts.

The remainder of this paper is organized as follows. Section 2 presents the semantics and syntax of the Flux language. Section 3 describes the Flux compiler and runtime systems. Section 4 presents our experimental methodology and compares the performance of Flux servers to their hand-written counterparts. Section 5 demonstrates the use of path profiling and discrete-event simulation. Section 6 reports our experience using Flux to build several servers. Section 7 presents related work, and Section 8 concludes with a discussion of planned future work.

2 Language Description

To introduce Flux, we develop a sample application that exercises most of Flux's features. This sample application is an image server that receives HTTP requests for images that are stored in the PPM format and compresses them into JPEGs. Recently-requested images are stored in a cache managed with a least-frequently used (LFU) replacement policy.

This section describes the entire Flux language, including the **concrete nodes** that implement the server logic, **abstract nodes** that represent a flow through multiple nodes, **semantic types** that implement conditional data flow, **error handlers** that deal with exceptional conditions, and the **concurrency constraints** that control simultaneous access to shared state.

2.1 Concrete Nodes

The first step in designing a Flux program is describing the *concrete nodes* that correspond to C and C++ implementations.

Flux requires type signatures for each node. In the current syntax, the name of the node is followed by the input arguments in parentheses, followed by an arrow and the output arguments.

Below are the signatures for three of the concrete nodes in the image server: `ReadRequest` parses client input, `Compress` compresses images, and `Write` outputs the compressed image to the client.

```
ReadRequest (int socket)
  -> (int socket, char* data);

Compress (int socket, char* raw, int size)
  -> (int socket, char* jpeg, int size);

Write (int socket, char* data, int size)
  -> (int socket)
```

While most concrete nodes both receive input data and produce output, *source* nodes only produce output to initiate a data flow. The statement below indicates that `Listen` is a source node, which Flux executes inside an infinite loop. Whenever `Listen` receives a connection, it transfers control to the `Image` node.

```
source Listen => Image;
```

2.2 Abstract Nodes

In Flux, concrete nodes can be composed to form *abstract nodes*. These abstract nodes represent a flow of data from concrete nodes to concrete nodes or other abstract nodes. Arrows connect nodes, and Flux checks to ensure that these connections make sense. The output type of the node on the left side of the arrow must match the input type of the node on the right side. For example, the abstract node `Image` in the image server corresponds to a flow from client input that checks the cache for the requested image, handles the result, writes the output, and completes.

```
Image =
  ReadRequest -> CheckCache -> Handler
  -> Write -> Complete;
```

2.3 Semantic Types

A client request for an image may result in either a cache hit or a cache miss. These need to be handled differently. Instead of exposing control flow directly, Flux lets programmers use the *semantic type* of a node's output to direct the flow of data to the appropriate subsequent node. A semantic type is a Boolean function supplied by the Flux programmer that is applied to the node's output.

Using semantic types, a Flux programmer can express multiple possible paths for data through the server. Semantic type dispatch is processed in order of the tests in the Flux program. The `typedef` statement binds the type `hit` to the Boolean function `TestInCache`. The node `Handler` below checks to see if its first argument is of type `hit`; in other words, it applies the function `TestInCache` to the third argument. The underscores are wildcards that match any type. `Handler` does nothing for a hit, but if there is a miss in the cache, the image

server fetches the PPM file, compresses it, and stores it in the cache.

```
typedef hit TestInCache;

Handler:[_, _, hit] = ;
Handler:[_, _, _] =
    ReadInFromDisk -> Compress
    -> StoreInCache;
```

2.4 Error Handling

Any server must handle error conditions. Flux expects nodes to follow the standard UNIX convention of returning error codes. Whenever a node returns a non-zero value, Flux checks to see if an error handler has been declared for the node. If no error handler exists, the current data flow is simply terminated.

In the image server, if the function to read an image from disk discovers that the image does not exist, it signals an error. We handle this error by directing the flow to a node `FourOhFour` that outputs a 404 page to the client:

```
handle error ReadInFromDisk -> FourOhFour;
```

2.5 Concurrency Constraints

All flows through the image server access a single shared image cache. Access to this shared resource must be controlled to ensure that two different data flows do not interfere with each other's operation.

The Flux programmer specifies such *concurrency constraints* in Flux rather than inside the component implementation. The programmer specifies concurrency constraints by using arbitrary symbolic names. These concurrency constraints can be thought of as locks, although this is not necessarily how they are implemented. A node only runs when it has “acquired” all of the constraints. This acquisition follows a two-phase locking protocol: the node acquires (“locks”) all of the constraints in order, executes the node, and then releases them in reverse order.

The image server maintains cache consistency by ensuring that only one node can modify the cache at a time.

```
constraint CheckCache:{cache};
constraint StoreInCache:{cache};
constraint Complete:{cache};
```

Readers/Writers

Concurrency constraints can be specified as either *readers* or *writers*. Using these constraints allows multiple readers to execute a node at the same time, supporting greater efficiency when most nodes read shared data rather than update it. Reader constraints have a question

mark appended to them (“?”). Although constraints are considered writers by default, a programmer can append an exclamation point (“!”) for added documentation.

Scoped Constraints

While flows generally represent independent clients, in some server applications, multiple flows may constitute a single *session*. For example, a file transfer to one client may take the form of multiple simultaneous flows. In this case, the state of the session (such as the status of transferred chunks) only needs to be protected from concurrent access in that session.

In addition to *program-wide constraints* that apply across the entire server (the default), Flux supports *per-session constraints* that apply only to particular sessions. Using session-scoped concurrency constraints increases efficiency by eliminating contention across sessions. Sessions are implemented as hash functions on the output of each source node. The Flux programmer implements this session id function that takes the source node's output as its parameter and returns a unique session identifier, and then adds (`session`) to a constraint name to indicate that it applies only per-session.

Discussion

Specifying concurrency constraints in Flux rather than placing locking operations inside implementation code has a number of advantages, beyond the fact that it allows the use of libraries whose source code is unavailable.

Safety: Declaring concurrency constraints in Flux allows Flux to prevent deadlock. The Flux compiler enforces a canonical order for lock acquisition (corresponding to concurrency constraints). Since data flows are acyclic, deadlock is thus impossible. Concurrency constraints in Flux are also reentrant, preventing deadlock that could occur due to multiple locking. Flux programs that only use Flux level concurrency constraints are thus guaranteed to not deadlock.

Efficiency: Exposing concurrency constraints also enables the Flux compiler to generate more efficient code. In particular, it provides implementations of the concurrency constraints tailored to particular runtimes. For example, a multithreaded runtime require locks, while a single-threaded event-driven runtime does not. The Flux compiler thus generates locks or other mutual exclusion operations only when needed.

Granularity selection: Finally, concurrency constraints let programmers easily find the appropriate granularity of locking. Grain selection is often difficult: too coarse a grain results in contention, while too fine a grain can impose excessive locking overhead. As we describe in Section 5.1, Flux can generate a discrete event simulator for the Flux program. This simulator can let a

developer identify the appropriate granularity of locking before actual server deployment.

3 Compiler and Runtime Systems

A Flux program is transformed into a working server by a multi-stage process. The compiler first reads in the Flux source and constructs a representation of the program graph. It then processes the internal representation to type-check the program. Once the code has been verified, the runtime code generator processes the graph and outputs C code that implements the server's data flow for a specific runtime. Finally, this code is linked with the implementation of the server logic into an operational server. We first describe the compilation process in detail. We then describe the three runtime systems that Flux currently supports.

3.1 The Flux Compiler

The Flux compiler is a three-pass compiler implemented in Java, and uses the JLex lexer [5] in conjunction with the CUP LALR parser generator [3].

The first pass parses the Flux program text and builds a graph-based internal representation. During this pass, the compiler links nodes referenced in the program's data flows. All of the conditional flows are merged, with an edge corresponding to each conditional flow.

The second pass decorates edges with types, connects error handlers to their respective nodes, and verifies that the program is correct. First, each node mentioned in a data flow is looked up in the table of function definitions. From these definitions, each node is labelled with its input and output types. Each semantic type used by a conditional node is looked up in the type definitions. Finally, the error handlers and concurrency constraints are attached to each node. If any of the referenced nodes or semantic types are undefined, the compiler signals an error and exits. Otherwise, the program graph is completely instantiated. The final step of program graph construction checks that the output types of each node match the inputs of the nodes that they are connected to. If all type tests pass, then the compiler has a valid program graph.

The third pass generates the intermediate code that implements the data flow of the server. Flux supports generating code for arbitrary runtime systems. The compiler defines an object-oriented interface for code generation. New runtimes can easily be plugged into the Flux compiler by implementing this code generator interface.

The current Flux compiler supports several different runtimes, described below. In addition to the runtime-specific intermediate code, the Flux compiler generates a Makefile and stubs for all of the functions that provide the server logic. These stubs ensure that the programmer uses the appropriate signatures for these meth-

ods. When appropriate, the code generator outputs locks corresponding to the concurrency constraints, but in a canonical order (alphabetically by name) that eliminates deadlock.

3.2 Runtime Systems

The current Flux compiler supports three different runtime systems: one thread per connection, a thread-pool system, and an event-driven runtime.

In the thread-based runtimes, each request handled by the server is dispatched to a thread function that handles all possible paths through the server's data flows. In the one-to-one thread server, a thread is created for every different data flow. In the thread-pool runtime, a fixed number of threads are allocated to service data flows. If all threads are occupied when a new data flow is created, the data flow is queued and handled in first-in first-out order.

The event-driven runtime operates differently. In this runtime, every input to a functional node is seen as an event. Each event is placed into a queue and handled in turn by a single thread. Additionally, each source node (a node with no input) is repeatedly placed on the queue to originate each new data flow. The transformation of input to output by a node generates a new event corresponding to the output data being propagated to the subsequent node.

The implementation of the event-based runtime is complicated by the fact that node implementations may perform blocking function calls. If blocking function calls like `read` and `write` were allowed to run unmodified, the operation of the entire server would block until the function returned.

Instead, the event-based runtime intercepts all calls to blocking functions using a handler that is pre-loaded via the `LD.PRELOAD` environment variable. This handler captures the state of the node at the blocking call and moves onto the next event in the queue. The formerly-blocking call is then executed asynchronously. When the event-based runtime receives a signal that the call has completed, the event is reactivated and re-queued for completion. Because the mainstream Linux kernel does not currently support callback-driven asynchronous I/O, the current Flux event-based runtime uses a separate thread to simulate callbacks for asynchronous I/O using the `select` function. A programmer is thus free to use synchronous I/O primitives without interfering with the operation of the event-based runtime.

Each of these runtimes was implemented in the C using POSIX threads and locks. Flux can also generate code for different programming languages. We have also implemented a prototype that targets Java, using both SEDA [16] and a custom runtime implementation, though we do not evaluate the Java systems here.

In addition to these runtimes, we have implemented a code generator that transforms a Flux program graph into code for the discrete event simulator CSIM [13]. This simulator can predict the performance of the server under varying conditions, even prior to the actual implementation of the core server logic. This process is described in greater detail in Section 5.1.

4 Experimental Evaluation

To demonstrate its effectiveness for building high-performance server applications, we have implemented a number of servers in Flux. We summarize these in Table 1. We chose these servers specifically to span the space of possible server applications. Most server applications can be broadly classified into one of the following categories: request-response client/server, “heartbeat” client/server and peer-to-peer. What differentiates these categories is their pattern of interactions. We implemented a server in Flux for each of these categories and compared their performance under varying load with existing hand-tuned server applications written in conventional programming languages.

4.1 Methodology

We evaluate all server applications by measuring their throughput and latency in response to realistic workloads.

All testing was performed with a server and client machine, both running Linux version 2.4.20. The server machine was a Pentium 4 (2.4Ghz, 1GB RAM), connected via gigabit Ethernet on a dedicated switched network to the client machine, a Xeon-based machine (2.4Ghz, 1GB RAM). All server and client applications were compiled using GCC version 3.2.2. During testing, both machines were running in multi-user mode with only standard services running. All results are for a run of two minutes, ignoring the first twenty seconds to allow the cache to warm up.

4.2 Request-Response: Web Server

Request-response based client/server applications are among the most common examples of network servers. This style of server includes most major Internet protocols including FTP, SMTP, POP, IMAP and HTTP. As an example of this application class, we implemented a web server in Flux. The Flux web server implements the HTTP/1.1 protocol and can serve both static and dynamic PHP web pages.

We implemented a benchmark to load test the Flux webserver that is similar to SPECweb99 [14]. The benchmark simulates a number of clients requesting files from the server. Each simulated client sends five requests over a single HTTP/1.1 TCP connection using keep-alives. When one file is retrieved, the next file is

immediately requested. After the five files are retrieved, the client disconnects and reconnects over a new TCP connection. The files requested by each simulated client follow the static portion of the SPECweb benchmark and each file is selected using the Zipf distribution.

We compare the performance of the Flux webserver against the latest versions of the *knot* webserver distributed with Capriccio [15] and the *Haboob* webserver distributed with the SEDA runtime system [16]. Figure 1 presents the throughput and latency for a range of simultaneous clients. These graphs represent the average of five different runs for each number of clients.

The results show that the Flux web server provides comparable performance to the fastest webserver (*knot*), regardless of whether the event-based or thread-based runtime is used. All three of these servers (*knot*, *flux-threadpool* and *flux-event-based*) significantly outperform *Haboob*, the event-based server distributed with SEDA. As expected, the naïve one-thread, one-client server generated by Flux has significantly worse performance due to the overhead of creating and destroying threads.

The results for the event-based server highlight one drawback of running on a system without true asynchronous I/O. With small numbers of clients, the event-based server suffers from increased latency that initially decreases and then follows the behavior of the other servers. This hiccup is an artifact of the interaction between the webserver’s implementation and the event-driven runtime, which must simulate asynchronous I/O. The first node in the webserver uses the `select` function with a timeout to wait for network activity. In the absence of other network activity, this node will block for a relatively long period of time. Because the event-based runtime only reactivates nodes that make blocking I/O calls after the completion of the currently-operating node, in the absence of other network activity, the call to `select` imposes a minimum latency on all blocking I/O. As the number of clients increases, there is sufficient network activity that `select` never reaches its timeout and frozen nodes are reactivated at the appropriate time. In the absence of true asynchronous I/O, the only solution to this problem would be to decrease the timeout call to `select`, which would increase the CPU usage of an otherwise idle server.

4.3 Peer-to-Peer: BitTorrent

Peer-to-peer applications act as both a server and a client. Unlike a request-response server, they both receive and initiate requests.

We implemented a BitTorrent server in Flux as a representative peer-to-peer application. BitTorrent uses a scatter-gather protocol for file sharing. BitTorrent peers exchange pieces of a shared file until all participants

Server	Style	Description	Lines of Flux code
Web server	request-response	a basic HTTP/1.1 server with PHP	36
Image server	request-response	image compression server	23
BitTorrent	peer-to-peer	a file-sharing server	84
Game server	heartbeat client-server	multiplayer game of “Tag”	54

Table 1: Servers implemented in Flux, described in Section 4.

have a complete copy. Network load is balanced by randomly requesting different pieces of the file from different peers.

To facilitate benchmarking, we changed the behavior of both of the BitTorrent peers we test here (the Flux version and CTorrent). First, all client peers are *unchoked* by default. Choking is an internal BitTorrent state that blocks certain clients from downloading data. This protocol restriction prevents real-world servers from being overwhelmed by too many client requests. We also allow an unlimited number of unchoked client peers to operate simultaneously, while the real BitTorrent server only unchokes clients who upload content.

We are unaware of any existing BitTorrent benchmarks, so we developed our own. Our BitTorrent benchmark mimics the traffic encountered by a busy BitTorrent peer. It simulates a series of clients continuously sending requests for randomly distributed pieces of a 54MB test file to a BitTorrent peer with a complete copy of the file. When a peer finishes downloading a piece of the file, it immediately requests another random piece of the file from those still missing. Once a client has obtained the entire file, it disconnects. This benchmark does not simulate the “scatter-gather” nature of the BitTorrent protocol – all requests go to a single peer. Using single peers has the effect of maximizing load, since obtaining data from a different source would lessen the load on the peer being tested.

Figure 2 compares the latency, throughput (completions per second) and network throughput to CTorrent, an implementation of the BitTorrent protocol written in C. Because BitTorrent is network-bound, there is little difference between the various server implementations. Nonetheless, prior to saturating the network, all of the Flux servers perform slightly worse than the CTorrent server. We are investigating the cause of this small performance gap.

4.4 Heartbeat Client-Server: Game Server

Unlike request-response client/server applications and most peer-to-peer applications, certain server applications are subject to real-time deadlines. An example of such a server is an online multi-player game. In these applications, the server maintains the shared state of the

game and distributes this state to all of the players at “heartbeat” intervals. There are two important conditions that must be met by this communication: the state possessed by all clients must be the same at each instant in time, and the inter-arrival time between states can not be too great. If either of these conditions is violated, the game will be unplayable or susceptible to cheating. These requirements place an important real-time constraint on the server’s performance.

We have implemented an online multi-player game of Tag in Flux. The Flux game server enforces the rules of tag. Players can not move beyond the boundaries of the game world. When a player is tagged by the player who is “it”, that player becomes the new “it” and is teleported to a new random location on the board. All communication between clients and server occurs over UDP at 10Hz, a rate comparable to other real-world online games. While simple, this game has all of the important characteristics of servers for first person shooter or real-time strategy games.

Benchmarking the gameserver is significantly different than load-testing either the webserver or BitTorrent peer. Throughput is not a consideration since only small pieces of data are transmitted. The primary concern is the latency of the server as the number of clients increases. The server must receive every player’s move, compute the new game state, and broadcast it within a fixed window of time.

To load-test the game server, we measured the effect of increasing the number of players. The performance of the gameserver is largely based upon the length of time it takes the server to update the game state given the moves received from all of the players, and this computation time is identical across the servers. The latency of the gameserver is largely a product of the rate of game turns, which stays constant at 10Hz. We found no appreciable differences between the traditional implementation of the gameserver and the various Flux versions. These results show that Flux is capable of producing a server with sufficient performance for multi-player online gaming.

5 Performance

In addition to its programming language support for writing server applications, Flux provides support for

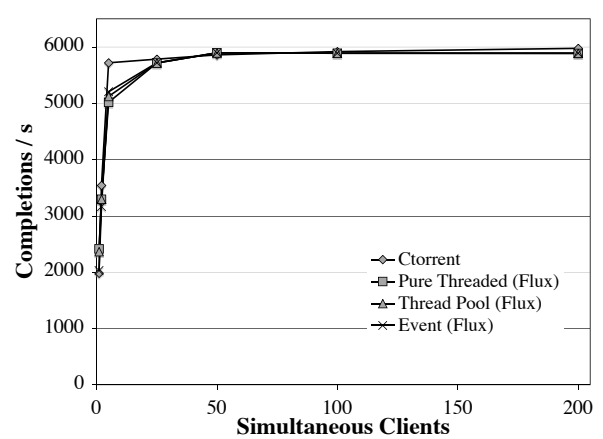
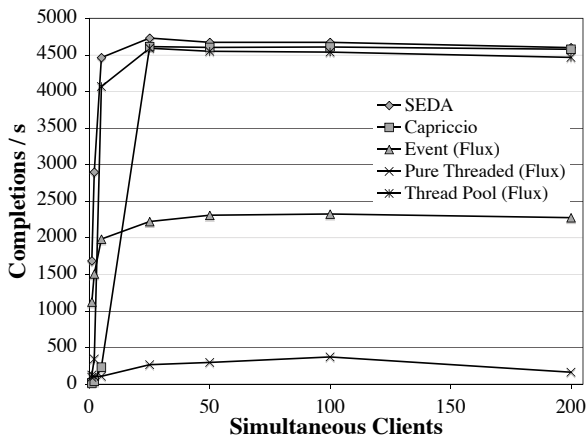
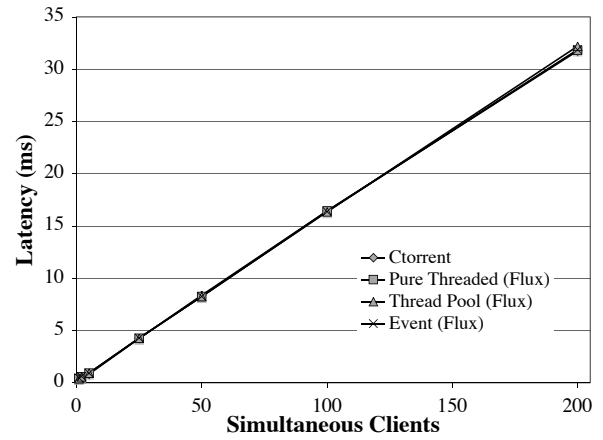
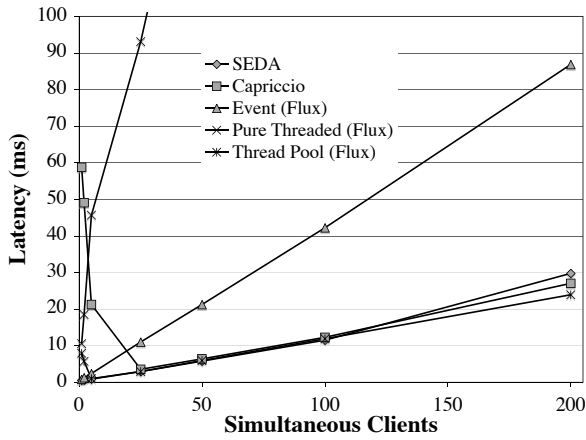
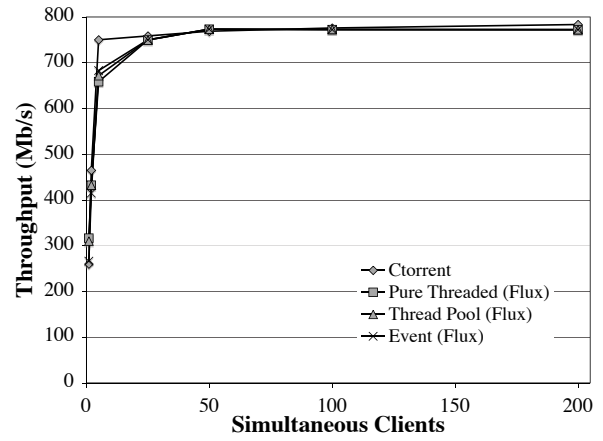
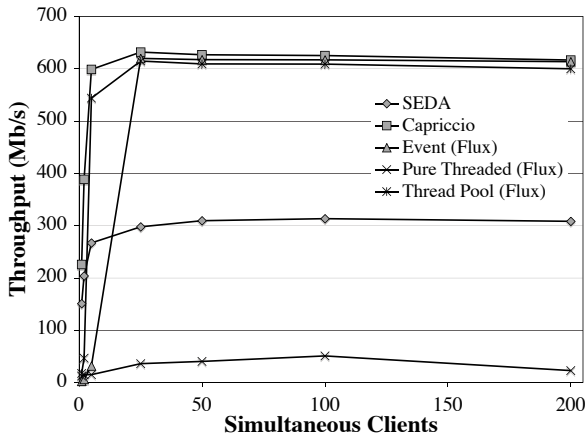


Figure 1: Comparison of Flux web servers with other high-performance implementations (see Section 4.2).

Figure 2: Comparison of Flux BitTorrent servers with CTorrent (see Section 4.3).

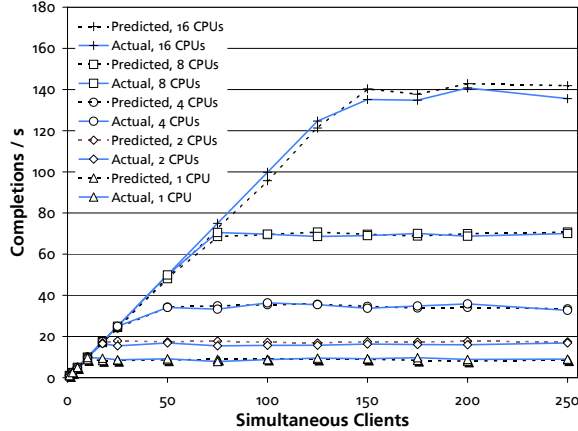


Figure 3: Predicted performance of the image server (derived from a single-processor run) versus observed performance for varying numbers of processors and load.

predicting and measuring the performance of server applications. The Flux system can generate **discrete-event simulators** that predict server performance for synthetic workloads and on different hardware and perform. It can also perform **path profiling** to identify server performance bottlenecks on a deployed system.

5.1 Performance Prediction

Predicting the performance of a server prior to deployment is important but often difficult. For example, performance bottlenecks due to contention may not appear during testing because the load placed on the system is insufficient. In addition, system testing on a small-scale system may not reveal problems that arise when the system is deployed on an enterprise-scale multiprocessor.

In addition to generating executable server code, the Flux code generator can transform a Flux program directly into a discrete-event simulator that models the performance of the server. We use CSIM as the implementation language for the simulator [13].

In the simulator, each node acquires a shared CPU resource for the period of time observed in the real world. Increasing the number of nodes that can simultaneously acquire the CPU resource simulates the addition of processors to the system. Each concurrency constraint becomes a shared resource. Every node using a particular concurrency constraint acquires that resource for the duration of the node’s execution. The simulator conservatively treats session-level constraints as globals, and reader constraints as writers.

It is important to note that this simulation does not model disk or network resources. While this is a realistic assumption for CPU-bound servers (such as dynamic web-servers), other servers may require more complete modeling.

The simulator can either use observed parameters from a running system on a uniprocessor (per-node execution times, source node inter-arrival times, and observed branching probabilities), or the Flux programmer can supply estimates for these parameters. The latter approach allows server performance to be estimated prior to any actual implementation of the server logic.

To demonstrate that the generated simulations accurately predict actual performance, we tested the image server described in Section 2. To simulate load on the machine, we made requests at increasingly small interarrival times. The images requested were selected using a uniform random distribution. The image server is CPU-bound, with each image taking a half second to compress on average.

We first measured the performance of this server on a 16-processor SunFire 6800, but with only a single CPU enabled. We then used the observed node runtime and branching probabilities to parameterize the generated CSIM simulator. We compare the predicted and actual performance of the server by making more processors available to the system. As Figure 3 shows, the predicted results (dotted lines) and actual results (solid lines) match closely, demonstrating the effectiveness of the simulator at predicting performance.

5.2 Path Profiling

The Flux compiler optionally instruments generated servers to simplify the identification of performance bottlenecks. This profiling information takes the form of “hot paths”, the most frequent or most time-consuming paths in the server. Flux identifies these hot paths using the Ball-Larus path profiling algorithm [4]. Because Flux graphs are acyclic, the Ball-Larus algorithm identifies each unique path through the server’s data-flow graph.

The overhead of path profiling is low enough that hot path information can be maintained by a production server. Profiling adds just one arithmetic operation and two high-resolution timer calls to each node. A performance analyst can obtain path profiles from a running Flux server by connecting to a dedicated socket.

To demonstrate the use of path profiling, we compiled a version of the BitTorrent peer with profiling enabled. For the experiments, we used a patched version of Linux that supports per-thread time gathering. The BitTorrent peer was load-tested with the same tester as in the performance experiments. For profiling, we used loads of 25, 50, and 100 clients. All profiling information was automatically generated from a running Flux server.

In BitTorrent, the most time-consuming path identified by Flux was, unsurprisingly, the file transfer path (Listen → GetClients → SelectSockets → CheckSockets → Message → ReadMessage →

HandleMessage \rightarrow Request \rightarrow MessageDone, 0.295 ms). However, the second most expensive path was the path that finds no outstanding chunk requests (Listen \rightarrow GetClients \rightarrow SelectSockets \rightarrow CheckSockets \rightarrow ERROR, 0.016ms). While this path is relatively cheap compared to the file transfer path, it also turns out to be the most frequently executed path (780,510 times, compared to 313,994 for the file transfer path). Since this path accounts for 13% of BitTorrent’s execution time, it is a reasonable candidate for optimization efforts.

Hot paths not only aid understanding of server performance characteristics but also identify places where optimization would be most effective. Because profiling information can be obtained from an operating server and is linked directly to paths in the program graph, a performance analyst can easily understand the performance characteristics of deployed servers.

6 Developer Experience

In this section, we examine the experience of programmers implementing Flux applications. In particular, we focus on the implementation of the Flux BitTorrent peer.

The Flux BitTorrent peer was implemented by two undergraduate students in less than one week. The students began with no knowledge of the technical details of the BitTorrent protocol or the Flux language. The design of the Flux program for the BitTorrent peer was entirely their original work. The implementation of the functional nodes in BitTorrent is loosely derived from the CTorrent source code. The program graph for the BitTorrent server is shown in Figure 4 at the end of this document.

The students had a generally positive reaction to programming in Flux. Primarily, they felt that organizing the application into a Flux program graph prior to implementation helped modularize their application design and debug server data flow prior to programming. They also found that the exposure of concurrency constraints at the Flux language level allowed for easy identification of the appropriate locations for mutual exclusion. Flux’s immunity to deadlock and the simplicity of the concurrency constraints increased their confidence in the correctness of the resulting server.

Though this is only anecdotal evidence, this experience suggests that programmers can quickly gain enough expertise in Flux to build reasonably complex server applications.

7 Related Work

In this section, we discuss the most closely related work to Flux.

Several previous domain-specific languages allow the integration of off-the-shelf code into data flow graphs,

though for different domains. The Click modular router is a domain-specific language for building network routers out of existing C components [11]. The Flux OSKit (no relation) is a domain-specific language for building operating systems, with rich support for integrating code implementing COM interfaces [8]. In addition to its linguistic and tool support for programming server applications, Flux ensures deadlock-freedom by enforcing a canonical lock ordering; this is not possible in Click and OSKit because they permit cyclic program graphs.

Flux is an example of a *coordination language* [9] that combines existing code into a larger program in a data flow setting. There have been numerous data flow languages proposed in the literature; Johnston et al.’s recent survey includes over one hundred references [10]. Most dataflow languages focus on extracting parallelism from individual programs, while Flux describes parallelism across multiple clients or event streams. Most of these languages also operate at the level of fundamental operations rather than functional granularity, although some medium-grained dataflow languages exist (e.g., CODE 2 [7]).

In recent years, there have been a number of papers proposing a wide variety of runtime systems, including SEDA [16], Hood [6, 1], Capriccio [15], libasync/mp [17], Fibers [2], and cohort scheduling [12]. Users of these runtimes are forced to implement a server using a particular API. Once implemented, the server logic is generally inextricably linked to the runtime. By contrast, Flux programs are independent of any particular choice of runtime system, so advanced runtime systems can be integrated straightforwardly into Flux’s code generation pass.

8 Future Work

We plan to build on this work in several directions. First, we are actively porting Flux to other architectures, especially multicore systems. We are also planning to extend Flux to operate on clusters. Because concurrency constraints identify nodes that share state, we plan to use these constraints to guide the placement of nodes across a cluster to minimize communication.

To gain more experience with Flux, we are adding further functionality to the web server. In particular, we plan to build an Apache compatibility layer so we can easily incorporate Apache modules. We also plan to enhance the simulator framework to support per-session constraints and to distinguish between reader and writer constraints.

We plan to release the entire Flux system by publication time, via the Flux-based BitTorrent and web servers described in this paper.

9 Acknowledgments

The authors would like to thank Gene Novark for helping to design of the discrete event simulation generator, and Vitaliy Lvin for assisting in experimental setup and data gathering.

Notes

1. The name Flux comes from the Latin *fluxus*, past participle of *fluere* = “flow”.

References

- [1] U. A. Acar, G. E. Blelloch, and R. D. Blumofe. The data locality of work stealing. In *SPAA '00: Proceedings of the twelfth annual ACM symposium on Parallel algorithms and architectures*, pages 1–12, New York, NY, USA, 2000. ACM Press.
- [2] A. Adya, J. Howell, M. Theimer, W. J. Bolosky, and J. R. Douceur. Cooperative task management without manual stack management. In *Proceedings of the General Track: 2002 USENIX Annual Technical Conference*, pages 289–302, Berkeley, CA, USA, 2002. USENIX Association.
- [3] A. W. Appel, F. Flannery, and S. E. Hudson. CUP parser generator for Java. <http://www.cs.princeton.edu/~appel/modern/java/CUP/>.
- [4] T. Ball and J. R. Larus. Optimally profiling and tracing programs. *ACM Transactions on Programming Languages and Systems*, 16(4):1319–1360, July 1994.
- [5] E. Berk and C. S. Ananian. Jlex: A lexical analyzer generator for Java. <http://www.cs.princeton.edu/~appel/modern/java/JLex/>.
- [6] R. D. Blumofe and D. Papadopoulos. The performance of work stealing in multiprogrammed environments (extended abstract). In *SIGMETRICS '98/PERFORMANCE '98: Proceedings of the 1998 ACM SIGMETRICS joint international conference on Measurement and modeling of computer systems*, pages 266–267, New York, NY, USA, 1998. ACM Press.
- [7] J. C. Browne, E. D. Berger, and A. Dube. Compositional development of performance models in POEMS. *The International Journal of High Performance Computing Applications*, 14(4):283–291, Winter 2000.
- [8] B. Ford, G. Back, G. Benson, J. Lepreau, A. Lin, and O. Shivers. The Flux OSKit: A substrate for OS and language research. In *16th ACM Symposium on Operating System Principles*, October 1997.
- [9] D. Gelernter and N. Carriero. Coordination languages and their significance. *Commun. ACM*, 35(2):96, 1992.
- [10] W. M. Johnston, J. R. P. Hanna, and R. J. Milar. Advances in dataflow programming languages. *ACM Comput. Surv.*, 36(1):1–34, 2004.
- [11] E. Kohler, R. Morris, B. Chen, J. Jannotti, and M. F. Kaashoek. The Click modular router. *ACM Transactions on Computer Systems*, 18(3):263–297, August 2000.
- [12] J. R. Larus and M. Parkes. Using cohort-scheduling to enhance server performance. In *Proceedings of the General Track: 2002 USENIX Annual Technical Conference*, pages 103–114, Berkeley, CA, USA, 2002. USENIX Association.
- [13] M. Software. The CSIM simulator. <http://www.mesquite.com>.
- [14] Standard Performance Evaluation Corporation. SPECweb99. <http://www.spec.org/osg/web99/>.
- [15] R. von Behren, J. Condit, F. Zhou, G. C. Necula, and E. Brewer. Capriccio: scalable threads for internet services. In *SOSP '03: Proceedings of the nineteenth ACM symposium on Operating systems principles*, pages 268–281, New York, NY, USA, 2003. ACM Press.
- [16] M. Welsh, D. Culler, and E. Brewer. Seda: an architecture for well-conditioned, scalable internet services. In *SOSP '01: Proceedings of the eighteenth ACM symposium on Operating systems principles*, pages 230–243, New York, NY, USA, 2001. ACM Press.
- [17] N. Zeldovich, A. Yip, F. Dabek, R. Morris, D. Mazières, and F. Kaashoek. Multiprocessor support for event-driven programs. In *Proceedings of the 2003 USENIX Annual Technical Conference (USENIX '03)*, San Antonio, Texas, June 2003.

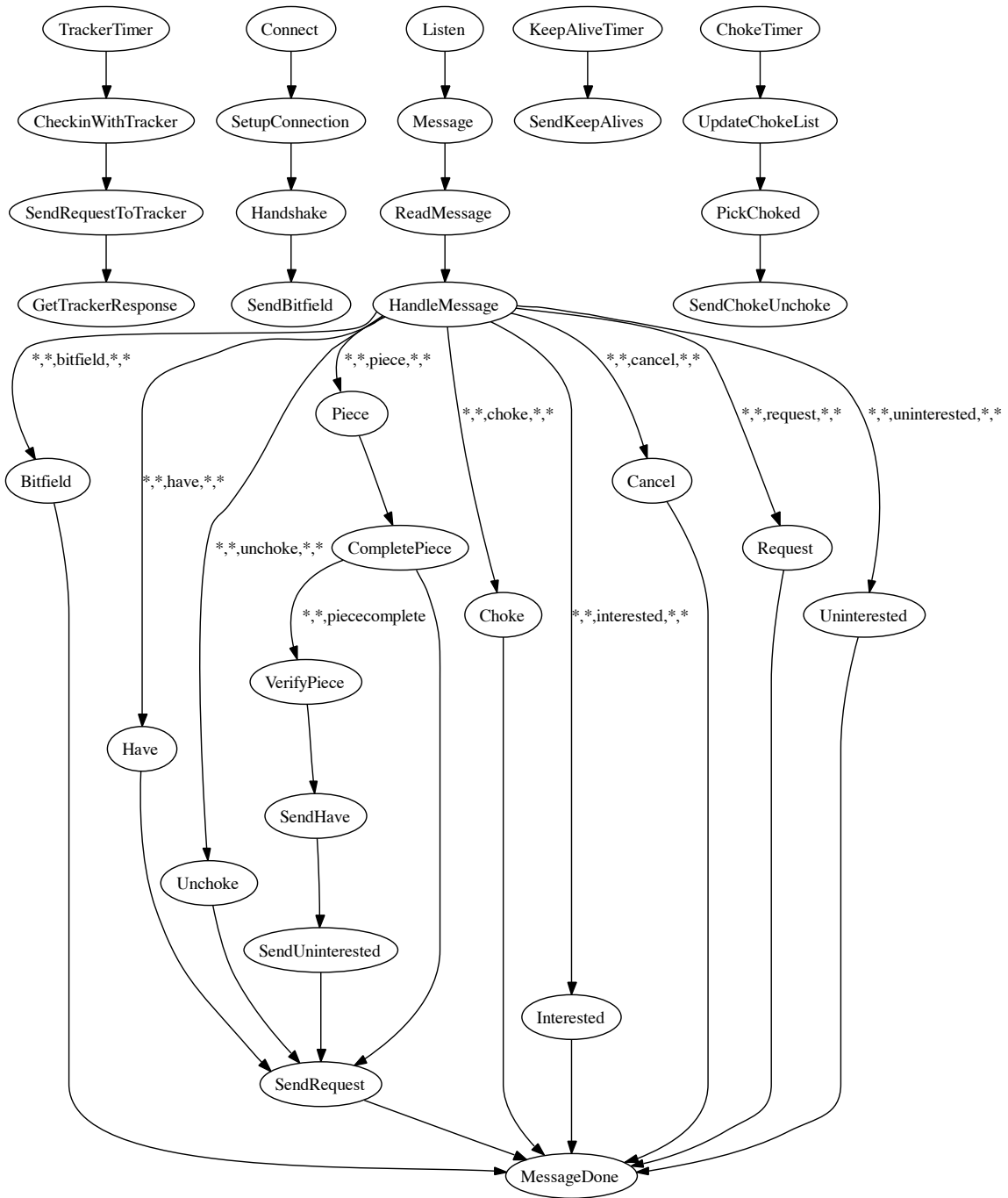


Figure 4: The Flux program graph for the example BitTorrent server.