EXPLOITING NETWORK PROCESSORS FOR LOW LATENCY,

HIGH THROUGHPUT, RATE-BASED SENSOR UPDATE DELIVERY

By

Kim Christian Swenson

A thesis submitted in partial fulfillment of
the requirements for the degree of

Master of Science in Computer Science

WASHINGTON STATE UNIVERSITY
School of Electrical Engineering and Computer Science

DECEMBER 2009

To the Faculty of Washington State University:

   The members of the Committee appointed to examine the thesis of KIM CHRISTIAN SWENSON find it satisfactory and recommend that it be accepted.

<br>

            ————————————————————————

            David E. Bakken, Ph.D., Chair

<br>

            ————————————————————————

            Carl Hauser, Ph.D.

<br>

            ————————————————————————

            Min Sik Kim, Ph.D.

ACKNOWLEDGMENT

EXPLOITING NETWORK PROCESSORS FOR LOW LATENCY,

HIGH THROUGHPUT, RATE-BASED SENSOR UPDATE DELIVERY

Abstract

by Kim Christian Swenson, M.S.
Washington State University
December 2009

Chair: David E. Bakken

GridStat is a novel publish-subscribe middleware framework which delivers rate-based streams of sensor updates for the power grid. Its novelty includes providing different subscribers to the same sensor update stream with different guarantees in terms of end-to-end delay, rate, and number of disjoint delivery paths. GridStat implementations on a general-purpose, dual-core CPU have very good per-hop forwarding latency on the order of 0.1 milliseconds. However, they can only support approximately 20 thousand forward/sec. In this thesis, we consider the problem of how to exploit high-speed network processors to increase the forwarding throughput for a GridStat status router (publish-subscribe forwarding engine). This thesis presents an analysis of the problem, and a design for GridStats routing algorithm on top of the Intel IXP 2400 series network processor. Experimental analysis shows how it is possible to support up to approximately 1 million forwards/second with a forwarding latency of approximately 0.01 milliseconds.

# Table of Contents

# List of Figures

# List of Tables

# Chapter 1

# Introduction

Third generation network processors are high performance network solutions designed for efficient packet processing and deep packet inspection. A third generation network processor generally include a physical network interface with multiple links, several RISC processors to process packets along a fast data path, and an embedded general processor to perform higher level protocol management [4]. Network processors can take on different networking roles, including real-time encryption, protocol translation, or multicast routing. Network processors are a relatively new technology aimed to combine the high performance of expensive but less flexible ASIC solutions with the low cost and flexibility of software based packet processing. A reduced instruction set tailored for packet processing and scalability through parallelism and pipelining are the main design ideas that enable the high performance of a network processor. Although the most demanding backbone networks still require custom hardware, network processors deliver high enough performance to meet many networking demands and

by far outperforms any regular computer based packet processing system.

GridStat is a status dissemination middleware framework designed to manage the flow of power grid status data in real-time [15, 10, 22, 6, 2, 7, 5]. GridStat use publish and subscribe concepts to hide the details of an advanced management, rate-filtering and routing middleware system. GridStat is designed ground-up to support the real-time delivery of very high bandwidth data streams, achieved among other things through early rate-filtering and efficient multicasting. GridStat supports quality of service (QoS) solutions like reliable timely delivery and secure communication.

In this thesis we explore some of the performance potential and scalability challenges of implementing a GridStat status router using network processors. Of particular interest is how the current GridStat rate-filtering and routing algorithm [15] scales in a network processor as throughput and routing-table complexity increases. The design and implementation of a prototype of a GridStat status router using the Intel Internet Exchange Architecture (IXA) is explained in detail in chapter 3. Chapter 4 present our findings from simulation and analysis, and chapter 5 suggests some subscription constraints that would enable GridStat status routers to provide latency and throughput guarantees.

# Chapter 2

# Background

## 2.1 GridStat

GridStat is a flexible publish-subscribe middleware designed for real-time delivery of power grid status information [19, 14, 10, 22, 6, 7]. Users of the middleware need not concern themselves with the details of the data delivery plane, rather they need only relate to either writing updates to variables (publisher role) or receiving updates of selected variables (subscriber role). Figure 2.1 show publishers writing updates to variables X, Y, and Z, and subscribers receiving updates of selected subsets of these variables. Each node in the directed graph of the data delivery plane is called a status router, it's the performance of such a GridStat status router implemented on network processors that is the focus of this thesis.

GridStat has more than one mode of operation, allowing the middleware to gracefully degrade service in the event of congested or damaged communication channels. Subscribers of

Figure 2.1: GridStat middleware overview

GridStat still receive useful and correct status information when the network runs at a lower level of operation, though likely receiving updates less often. A GridStat status router is a leaf node in the GridStat management hierarchy where all status routers make an interconnected graph to compose the possible data paths of power grid status data. These status routers all run the same algorithm for rate filtering and multicast routing, and differ only by where in the management topology they reside, therefore using different routing table information.

### 2.1.1 Routing algorithm pseudo code

1: **loop**

2:     $P \leftarrow GetNextPacket()$

3:     $V \leftarrow StatusVariableID(P)$

4:     $Pub \leftarrow PublishInterval(V)$

5:     **for** $L$ **in** $RouterLinks()$ **where** $L \neq IncomingLink(P)$ **do**

6:         $S \leftarrow SubscriberList(V, L)$

7:         **for** $s$ **in** $S$ **do**

8:             $R \leftarrow (Timestamp(P) + Pub/2) \; modulo \; SubscriberInterval(s)$

9:             **if** $R < Pub$ **then**

10:                 Forward packet $P$ on link $L$

11:                 continue for loop at line 5

12:             **end if**

13:         **end for**

14:     **end for**

15: **end loop**

Throughout this document lines 5 to 14 are referred to as the algorithm's outer loop and lines 7 to 13 as the algorithm's inner loop. The subscriber list at line 6 is ordered ascending on subscriber interval values.

## 2.2  Hardware platform

### 2.2.1  Network Processors in general

In principle, a network processors is a network capable device that specializes in processing packets. Network processors are generally categorized in three generations [4]. The first generation of network processors are packet processing systems built with conventional computer systems. The second generation of network processors added specialized processors more directly connected to the physical network interface itself allowing a fast data path and much higher response times and throughput capacity. Third generation network processors extend the second generation by also including an embedded general processor to manage higher level of protocols and to interact with the other processors along the fast data path through shared memory.

### 2.2.2  The Intel IXP2400 architecture

Details on the IXP2400 network processor can be found in [1].

**Intel XScale Core**

The Intel XScale core is a 32-bit general purpose processor. The core is compatible with the ARM Version 5 (V5) Architecture with the exception that the XScale core does not provide hardware support of floating point instructions. Features of the XScale include Memory Management (MMU), Instruction Cache, Data Cache, and Interrupt Controller. The core

7

Figure 2.2: Schematic overview of the IXP24xx network processor

has 4 GB of addressable memory, where the 2 GB low addresses map onto DRAM and the 2 GB high addresses map onto other targets like SRAM, PCI, and CSR's. Typical operating system choices are VxWorks® and Linux®.

**Microengines**

Microengines allow fast programmable packet processing. The IXP2400 has 8 microengines. Each microengine is its own processor with a range of registers available. Key to the speed and flexibility of the Microengines is their fast access to most special purpose subsystems of the IXP2400 chip through buses.

**Control Store** Microengines in the IXP2400 each has a static RAM control store. It can hold 4096 instructions each 40-bits wide. Before any execution can be done, code must be loaded into this control store. The code store is typically not changed after execution starts, even though it's possible, however patching of symbols is a common operation after code is loaded but before execution starts and allows more dynamic software architectures.

**Contexts** Each Microengine has 8 hardware contexts available. Each context has its own set of registers, a program counter, and context specific local registers. By default registers are allocated in a per context exclusive manner giving each context their fair share. Naturally this limits the amount of registers available per context but it allows very fast context switches which is key to achieve good performance. A context can only be in 1 of 4 states:

- Inactive. Applications may not require all 8 contexts, when that is the case, select

contexts are marked as inactive and will not participate in the context switch rotation. The corresponding available registers may be shared among the remaining contexts.

- Executing. At most one context may be in the executing state at any time. When in the executing state, that context's program counter is used to fetch instructions from the control store. A context will stay in the executing state until it executes an instruction that tells it to sleep.

- Ready. This state indicates that the context is ready to execute. The only reason it cannot execute is because another context is already in the execute state. The next context to execute is selected round robin among the Ready state contexts once the executing context goes to sleep.

- Sleep. Waiting on external event(s). Once all events the context is waiting on has occurred the context is put in the Ready state.

**Registers**  The Microengines of the IXP2400 has several types of registers:

- Local Control and Status Registers (CSRs). These registers hold special purpose status information and is not part of the direct data path registers. There is a delay of 3-4 cycles to read/write Local CSRs. These registers reside in the SHaC unit.

- 256 General purpose registers organized in 2 banks A and B. Most instructions require operands from separate banks, the allocation of registers to banks in software is done by the compiler and not by the programmer.

- 512 Transfer registers. There are 4 types of transfer registers: `S_Transfer_In`, `S_Transfer_Out`, `D_Transfer_In`, `D_Transfer_Out`. These registers are used when accessing data external to the microengine like SRAM, DRAM, and the SHaC.

- 128 Next Neighbor (NN) Registers. Microengines have dedicated data paths (NN registers) between them. These data paths are predefined in which Microengine can send/read data to/from another Microengine.

- 640 32-bit words of Local Memory (LM). Microengine LM is accessed through context copies of the `LM_ADDR` registers. Each context have 2 such registers: `LM_ADDR_0` and `LM_ADDR_1`.

**SRAM Controller**

There are 2 independent SRAM controllers, each supporting pipelined QDR synchronous static RAM (SRAM). SRAM is ideal for representing packet queues and meta-data and the addresses here are typically mapped directly to packet data in DRAM through some simple arithmetic function. The SRAM is logically 32-bit wide, so special care must be taken by the programmer when access outside this alignment occur.

**DRAM Controller**

The DRAM controller support a single 64-bit channel of DRAM. 2 GB of address space is allocated to DRAM. The DRAM controller may receive read and write requests from the Intel XScale core, Microengines, and the PCI bus master. These connect to the DRAM

11

controller via the Command Bus and the Push and Pull buses. DRAM on the IXP2400 also support Error Correcting Code (ECC). The DRAM is logically 64-bit wide, so special care must be taken by the programmer when access outside this alignment occur.

## Media and Switch Fabric Interface (MSF)

The MSF connect the IXP2400 to a physical layer device (PHY) and/or to a switch fabric. The MSF acts as the bus for SRAM push/pull, DRAM push/pull, and the physical network hardware layer. The MSF allows transparent and uniform access to different king of physical network layers without the need for application code changes except some configuration options.

## PCI Controller

The IXP2400 PCI Controller may be used to communicate between the Intel XScale core and its host (any computer with a compatible PCI interface).

## SHaC unit

The SHaC unit consist of Scratchpad, Hash Unit, and CAP. These three units all share the Command, S_Push, S_Pull bus interfaces and only one command to either unit can be processed per cycle.

**Scratchpad Memory**    Scratchpad Memory has 16KB general purpose storage. The memory is organized as 4K 32-bit words accessible by the Intel XScale core and Microengines.

Scratchpad is not byte-writable, all read and write operations on Scratchpad must work with 4 byte aligned words of 4 bytes each. Max 16 32-bit words (total of 64 bytes) can be accessed in one read/write operation. Scratchpad support 16 hardware assisted rings. These rings are fixed size FIFO structures using head and tail pointers and the hardware operations on them are atomic. The size of the rings must be decided at compile or load time and must be one of 128, 256, 512, or 1024 32-bit words. Scratchpad allows much faster memory access than SRAM or DRAM, and is therefore ideal for small data structures accessed frequently.

**Hash Unit**   The Hash Unit take 48, 64, and 128-bit data and produce a 48, 64, and 128-bit hash index respectively. The Hash Unit can be accessed by Microengines and the Intel XScale core. The hash unit is ideal when doing searches with large keys like L2 addresses.

**Control and Status Register Access Proxy (CAP)**   The CAP has many control and status registers (CSRs) configuring network processor behavior. In addition to configuring the hardware, the CSRs provide mechanisms for inter-thread communication across microengines and ways to communicate between the XScale core and microengines. Note that Scratchpad rings and SRAM memory can also be used to communicate between the XScale core and the microengines.

**Performance monitor**

**Intel XScale core peripherals (XPI)**

The following peripherals fall under this category: Interrupt Controller, Timers, GPIO, UART, and SlowPort.

**Interrupt Controller**   The Interrupt Controller provide ability to enable or mask interrupts from sources like timers, microengine software, PCI device, hardware error conditions and more. Reading interrupt status and enabling/masking interrupts are done as read and write operations on memory mapped registers.

**Timers**   4 programmable 32-bit timers are provided. They can be clocked by the internal clock, a divided version of the clock, or by a signal on an external GPIO. Periodic timers can generate interrupts if programmed so.

**GPIO**   The IXP2400 has 8 General Purpose IO (GPIO) pins. Usage of GPIO pins range from output LED indicators or input switches. GPIO pins can provide interrupts to the Intel XScale core or clocking the timers.

**UART**   There is a RS-232 compatible Universal Asynchronous Receiver/Transmitter (UART) usable by a maintenance console or a debugger.

**SlowPort**   The SlowPort allows access to the Flash ROM and asynchronous external devices at 8, 16, or 32-bit access rate.

Figure 2.3: General packet flow

## 2.3 Related work

### 2.3.1 GridStat

An initial prototype of the GridStat framework, including a status router written in Java, was developed in 2000 [5, 13]. Since then, the GridStat framework has seen many enhancements and extensions, particularly related to security, quality of service (QoS) and ease of use. In 2005, it was shown that GridStat running on a regular computer with Java based status routers could satisfy the industry end-to-end latency requirement for delivering phasor measurement unit data through the GridStat middleware at relatively low throughput rates [14]. In 2008, a C implementation of a GridStat status router proved capable of more than twice the throughput than current Java implementations and with lower and far more

predictable delays [19]. While computer based status routers seem to be capable of thousands and even tens of thousands of forwards per second, network processors (of comparable vintage) are expected to be capable of many million forwards per second and at much lower latencies and better predictability than regular computer based implementations.

## 2.3.2   Network processors

Third generation network processors have been widely available to industry and academia for almost ten years now, and a lot of interesting research has been done with them.

The IXP2400 network processor has been used for studies on various topics. Bursty traffic and a multi threaded caching algorithm [3]. Packet reordering algorithms that are more efficient than those that come bundled with the SDK [8]. Automated application design and mapping of code onto different processing elements [20]. Network processor runtime operating system to dynamically reconfigure and switch code running in the data path of the network processor based on traffic patterns [25].

The IXP2800 and IXP2850 are other very popular network processors for studying specific applications or general methods of network processor architecture. Common applications to study are packet and flow classification algorithms [16, 27, 21] which must eventually be part of a GridStat status router if it is to handle unexpected input properly. Such classification algorithms were also studied along with fast TCP session creation and packet reordering. A fast IPv6 forwarding and routing algorithm have been demonstrated [11], while at the same time analyzing some general IXP network processor architecture pro-

gramming challenges and how to best maintain packet ordering and hide internal IO latency. There are also other studies on hiding the internal latency of memory and other subsystems through multi-threading techniques [9].

Parallelization and multi-threading is key to performance and scalability of network processors in general [26]. There exists methods and applications to compare different network processor platforms, which was tested to work on the IXP1200 [18].

There is also much research done on automatic and dynamic mapping of network processor code onto microengine as runtime requirements change. Methods to do this has been developed through both pure analysis [24] and statistical analysis [17]. There even exists methods for comparing different dynamic runtime systems with each other [12].

### 2.3.3 Multicast IP routing

Multicast IP is a facility that send an IP packet to a group of hosts identified by a group address [23]. Multicast routers are a subset of IP routers that manage and provide the multicast IP facility. In a router with $N$ physical network ports of equal capacity, the outgoing throughput is potentially up to $N - 1$ times bigger than the incoming throughput. This is because that the most expensive case of multicast is a broadcast, where each incoming packet is sent to every outgoing port expect through the port it arrived. In the case of sustained data streams where all packets are broadcasted, the router will not be able to accept an incoming throughput of more than the capacity of a single port. A simplified explanation proving this is as follows: Observe that when a router receives incoming data streams at rates higher

than a single port, sending every such incoming packet to every other port would exceed any other port's capacity. The real multicast world is generally more complex than such an ideal simplification of the worst case. It may be expected that a real-world multicast traffic is bursty and would not necessarily trigger a full broadcast in any router. Such environments can be very complex and the throughput capacity and packet processing time are dynamic and can only be described through advanced network traffic and routing buffering analysis and would most likely also require simulation or real-world observation to be well described.

# Chapter 3

# Prototype design and implementation

## 3.1 Implementation Challenges

Implementing the GridStat Status Router on an Intel IXP network processor had some challenges:

- Complexity of such a highly parallel programming model as the one presented through the IXA SDK.

- In fixed rate streams, predictable packet delay is important.

- Exposure to network processor hardware complexity, e.g. many layers of memory but no MMU, some special instructions share resources across microengines.

- No hardware support for division or modulo operations used by the routing algorithm.

### 3.1.1   Router tailored hardware

The Intel IXP network processors are designed with routing applications in mind.

The IXP network processors have a dedicated hardware/software layer called fast path or data-plane. The hardware at this layer is very fast and designed with fast packet processing and forwarding in mind. The efficient way packet data is transferred from physical links to DRAM and sent out again is one of the best examples of dedicated hardware for packet processing. Hardware supported atomic ring buffer operations on SRAM, fast atomic synchronization primitives in scratch memory, and next-neighbor registers are all examples of hardware support for a fast and reliable software pipeline at the data-plane.

Perhaps the most obvious benefit of using a network processor chip like the IXP family is the presence of many processors running in parallel each with multiple threads making it possible to hide the waiting on internal hardware IO. Some of the benefits and drawbacks of the parallel nature of the IXP network processor are discussed in section 3.1.2.

Typical to any routing application is the presence of a semi-static routing table data structure. Navigating the routing table data structure has to be fast at the data-plane where table lookup is performed. Hashing has proven a fast method to navigate large data structures and the IXP2xxx chips provide hardware accelerated hashing using either 48-bit, 64-bit, and 128-bit keys and hashes. Hashing is the method of choice in navigating the routing tables of all the sample routing applications including IP and multicast IP. Hashing remain the obvious choice for navigating the routing table in a GridStat status router, and the one-to-one relationship between variable-ID (32-bit flow membership identifier) and the corresponding

complete routing table entry makes mapping a hashing algorithm onto navigating the routing table in GridStat straightforward.

The IXP2xxx chip provide a small but fast mechanism to support a microengine-local cache available at all microengines. By combining the 640 4-byte words of microengine local memory and the 16 entry Content Addressable Memory (CAM) it's easy to store 16 entries of a large data structure (like a routing table) in local memory and using the fast CAM to navigate it. Although the cache often prove an effective performance enhancement to an application like IP routing and forwarding, it's not given that it will in the GridStat application. The benefit of a cache is much greater in an application with a somewhat unpredictable and bursty traffic environment. General Internet IP traffic has this exact nature, due to it being used by such a diverse number of applications often with unpredictable users dictating traffic and bursts. The GridStat application however, is neither bursty nor unpredictable in its steady state. The publish/subscribe architecture is designed with reliable traffic streams that are not bursty, rather they are fixed rate streams. Only when alerts in GridStat are issued will we see bursty behavior, alerts are not discussed in this thesis. Given a reasonably large amount of publish streams where many of them operate with the same or similar publish rates, we can not expect a cache to give much benefit. On the contrary, In GridStat, we might expect the routing table cache to be a lookup overhead for most publish streams. It should be noted that in the case of a few (less than some fraction of the cache size) very fast (high rate) publish streams, these fast streams might always get a cache-hit while most other streams get a cache-miss, hence a cache will favor the fast streams. This caching problem in

21

GridStat could have other and better solutions, one of which would be to dedicate certain microblocks to some high rate streams giving it the benefit of always having the relevant routing table entries in local memory.

## 3.1.2 Parallelism and throughput

Parallelism in a network processor when applied to an application like the GridStat status router, is always about processing more than one packet at a time and at the same section of the packet processing pipeline. Given a simple packet processing pipeline consisting of `Receive`, `Route`, `Queue Manager`, and `Transmit` microblocks. The fact that these microblocks run in parallel give no speedup benefit when we look at the processing delay of a single packet. Rather, the packet processing pipeline allows high throughput as we can have many packets in the pipeline each being processed by a different part at the same time. It's important to note that having one processor doing all steps on every packet would yield the same single packet delay as that of a parallel pipeline, however the throughput when using multiple processors each with a separate task is much higher. Parallelism can also be extended within part of the pipeline by having one of the microblocks (steps) use more than one microengine (processor). An example would be having the `Route` microblock use four microengines. Each microengine would run the same code but always process different packets allowing more packets to pass through this step of the packet processing pipeline over the same time interval. Any given packet would only visit one of the four instances of the `Route` microblock.

Ideally, the ability to process packets in parallel will give a speedup corresponding with the degree of parallelism. e.g. Running the `Route` microblock on three microengines simultaneously allow three times the amount of packets to be routed over the same time interval than running `Route` on one microengine only. If we assume that the task of dividing and combining work takes no time and that synchronization between the microengines running the same microblock is not needed, we get a throughput scalability where adding one microengine give one speedup. In GridStat, the routing tables might change at runtime and packet order must be maintained, hence processing many packets in parallel require some synchronization to access the routing table and to maintain packet order.

There are two levels of parallelism within a microblock running on multiple microengines. The most obvious one is the fact that when you have more than one microengine, these can execute instructions at the same physical time providing true parallelism. The other not so obvious one comes from the fact that each microengine have eight hardware supported threads (contexts) each with its own program counter and fair share of the registers. As the threads share the same physical processor they can't execute instructions at the same physical time, however some part of packet processing don't require executing instructions. Waiting on IO is something that most processing code must do, and while waiting on IO, a thread will yield, allowing other threads to execute while it's waiting. This is how "hiding" IO delay is achieved, and although any "waiting on IO" will add to the delay of the corresponding packet being processed, the delay does not "steal" throughput capability of the application as long as the total amount of "waiting on IO" delay in a microengine per

packet is no more than eight times (as there are eight threads) that of the delay added by executing instructions for that packet in the microengine.

### 3.1.3   Predictable packet delay

Data passed through the GridStat framework is of a "time-critical" nature as some users / systems expect to act quickly based on what the data describes. Ideally all packets would have zero delay allowing systems to respond immediately. Clearly we can't obtain zero delay in such a framework so the question becomes "how well can we do." The fixed-rate nature of traffic streams intended for the GridStat framework would also suggest that we want predictable delays as the streams will only stay as "fixed-rate" as the variance in delay from packet to packet is. As the users of GridStat will need to obtain "global-snapshots" using data from multiple streams, the snapshots can only be viewed as soon as all the packets belonging within that certain time-window defined as a single snapshot have reached the subscriber.

### 3.1.4   Scalability of throughput

Scalability within the status router itself depend on the routing algorithm used. The GridStat status router routing and rate-filtering algorithm (see section 2.1.1) is simply referred to as "algorithm" or "routing algorithm". Each published variable in the GridStat architecture is called a "flow" and each flow has an associated unique "flow identifier." Every flow has a routing table entry which contains a two-dimensional list of links and values indicating

subscriber interval periods. For a given flow, each link has a set of 0 or more subscriber values. This set we call the "collapsed subscriber interval list," see 3.3.4 for a more detailed explanation of this set.

The original routing algorithm scale based on the number of links on the status router and on the collapsed subscriber interval list. Changing the number of links on a status router falls into the category of managing the status router hierarchy as a whole. From a single status router's point of view, the number of physical links to use is dictated by the management layer. The original routing algorithm assume that subscribers may register to receive published streams at any rate (using millisecond resolution.) The nature of these constraints will in the worst-case cause packet delay to scale at $O(mn)$ where $m$ is the number of links and $n$ is the number of subscribers. This is because the algorithm will always iterate over all outgoing links and for each link iterate through every item in the associated collapsed subscriber interval list. The constant time required in each iteration is high as well, given that each operation requires a modular (division) operation which is not even supported by IXP2xxx hardware.

## 3.2 Mapping the GridStat algorithm on network processor hardware

Given the Java version of the GridStat algorithm, mapping this algorithm onto IXP hardware has proven quite the challenge. Many things about network processor hardware forces

specific design on aspects of the software. The most obvious are that all the processing units (microengines) of the network processor runs in parallel. Secondly is the exposure of 4 layers of RAM to the programmer, DRAM, SRAM, Scratch, and Microengine local memory. And the last major challenge is the use of a a programming language called "Microcode" which is a kind of macro'ed assembly with some simple block constructs.

## Choosing the programming language

There exist a compiler for a higher level programming language called "Micro C" which is very close to standard "C" but with some special added constructs to support the IXP hardware. In the prototype used in this thesis Microcode was chosen since most of the current SDK versions were already written using Microcode. Most of the development tools also seem to support Microcode better than Micro C which also made the decision easier. Examining sample projects written in both languages reveal that the amount of work on the programmers part is somewhat lighter in Micro C than it is in Microcode, but not as significant a difference as usually found between assembly and C-code in general programming. Micro C is easier to read and also hide as much hardware exposure as possible giving it a more abstract feel than Microcode. Ultimately the compatibility with libraries and development tools along with help material available made the choice fall on Microcode. However, in hindsight it's easy to feel that the implementation would have been easier to do in Micro C, and perhaps that should be the choice of a production implementation in order to favor code clarity over complete control of every microengine instruction.

**Working with multiple processors**

The IXP2400 platform has 8 microengines per network processor. Each microengine is its own processor running in perfect clock synchronization with all the other microengines. As the packet processing done in GridStat is fairly expensive the code per microblock should utilize more than one microengine to give a significant performance advantage over using only one microengine. This performance advantage only comes to its right when traffic loads are high, at low traffic rates there is no benefit of using more than one microengine per microblock. The implementation used for simulations in this thesis (see chapter 4) utilize two microengines for the routing algorithm, although this is probably not optimal, it suffices and ease the performance analysis a little since a slower moving system is easier to measure.

**Designing the GridStat status router**

The IXP hardware platform provides the programmer with full freedom as to how software should run. Intel however, in designing the IXP features and structure as they have, did so optimizing its features to support some special ways of designing software. Given all the example software provided in the IXA SDK, it's easily demonstrated that this specific way of designing software works well for most common network applications. We have chosen for the GridStat status router implementation to use the software design principles demonstrated in the IXP/IXA documentation that follows the network processor's software development kit.

The introduction of a "microblock" appear as a chief characteristic of applications

on the IXP. The IXA support the concept of microblocks in many ways to streamline implementing them. A microblock is a software component in the packet processing pipeline performing a specific operation on packets. Conceptually the core of each microblock is an infinite loop called a "dispatch loop" like:

```
// dispatch loop start:
(1) Source packet P
(2) Process packet P
(3) Sink packet P
// dispatch loop end.
```

A packet processing pipeline consists of many microblocks chained together. Consider a pipeline with microblock $A$ and microblock $B$ each with a dispatch loop as the one above. In order to chain A and B together we can create a shared data structure $AB\_RING$ with FIFO operation logic. Now $A$ will sink packets onto the $AB\_RING$ and $B$ will source packets off the $AB\_RING$. The data structure used to link microblock together in this way is usually a ring buffer that reside in scratch memory called a "scratch ring". These rings are an excellent choice as they reside in the fastest shared RAM available big enough to support such structures for large ring buffers. Scratch rings are easily accessed by all microengines using fast atomic hardware instructions to get and put data on them solving all synchronization issues with shared ring structures. As these rings are fixed size they might become full or empty. It's the responsibility of the microblock sinking a packet on a ring to check whether the ring has available space before performing a put operation. If the ring is full, the packet is either dropped or the microblock stall waiting until there is space. In order to support reusing microblocks across applications the specific assignment of scratch rings in between

28

Figure 3.1: Prototype packet flow

microblocks can be set at load time which is after microblock compilation but before the code starts executing. This enables an application designer to reuse pre-compiled microblocks in their own packet pipeline by patching scratch ring assignments into the code after it's loaded into the microengine hardware but before code execution starts.

The GridStat status router application use two microblocks taken from the IXA and one microblock specific to the application. The `RX` and `TX` microblocks taken from the IXA provide endpoints in the packet processing pipeline as they are responsible for receiving and sending packets on the physical medium. The GridStat features are all implemented within the `Route` microblock which provide the core route and filtering logic used in the GridStat routing algorithm. This is a simple design that can be improved on, but suffice to demonstrate some performance properties of using network processors to implement status routers. See figure 3.1 for a prototype architecture overview.

## 3.3   Microblocks

### 3.3.1   General packet flow

Packets arrive through physical links and all enter the GridStat software pipeline through the `Ethernet Receive` driver (`RX`) running on a single microengine. `RX` put packets on an outgoing scratch ring (`PACKET_RX`) in the order they are received. The `Route` microblocks take turns in an ordered fashion getting packets off the `PACKET_RX`. Once gotten, the packet is immediately prepared for an eventual multicast. `Route` then runs the GridStat filtering algorithm 2.1.1 on each outgoing port to determine which ports the packet is to be forwarded on. When a forwarding decision for each port has been made `Route` then iterates through all outgoing ports and forward packets on those meant to have the packet. Forwarding packets from the perspective of `Route` means putting the packet on a scratch ring corresponding to the link it's meant for. Packets prepared for a link but not meant to be forwarded should be dropped by `Route` but are allowed to pass through anyhow in order that performance may be more easily measured after simulation. The `Ethernet Transmit` (`TX`) microblock handles as many scratch rings are there are hardware links on the network processor hardware. `TX` will get packets from a ring then do what it must to send this packet out on the physical medium corresponding to the scratch ring the packet was found on.

### 3.3.2 Packet Receive (RX) microblock

This software component is a driver running on a single microengine. Although configurable to run on two microengines if the hardware line capacity demands it, something we will not discuss here as we currently have no such demand. RX is responsible for queuing all packets that arrive at any physical link on a single scratch ring (`PACKET_RX`). RX maintains logical packet ordering by making sure that all packets arriving on the same physical link are queued on `PACKET_RX` in the order they arrived on that link. However, order between packets arriving on separate physical links need not be maintained, and it's up to RX to choose which packet to queue on `PACKET_RX` first.

**Input**

RX is responsible for getting packets from the physical medium and into the software packet pipeline. This microblock is part of the IXA written by Intel and it handles all the configuration setup needed to initialize the hardware and logic for pulling packets off the physical media.

**Output**

All incoming packets are put on the `PACKET_RX` scratch ring.

**Concurrency**

RX can be configured to run on 1 or 2 microengines. We use the configuration of 1 microengine in our project.

**Packet ordering**

Packet order is maintained on a per link basis. So any number of packets arriving on the same physical link will remain their relative packet order when the packets are output on PACKET_RX. Packets arriving on different links have no ordering guarantees between them.

**Performance**

RX can keep up with line rates on all IXP hardware it's bundled with given proper configuration. Hence its performance is at least as good as we need.

### 3.3.3   Packet Transmit (TX) microblock

TX handles $N$ scratch rings and $N$ outgoing physical links. Each scratch ring has a fixed outgoing link that it corresponds to. TX is responsible for sending packets from all scratch rings out on its corresponding link in the order packets were queued on the scratch ring. TX maintains packet order within one scratch ring to physical link stream, and uses a simple round robin scheduling for the inter queue packet selection. TX can be configured with other scheduling algorithms like weighted round robin and negative deficits, but this is not important for the purpose of this project.

**Input**

TX has a number of scratch rings where packets are pulled as input. Each of these scratch rings correspond to an outgoing physical link in a 1 to 1 fashion. Packets are pulled from

the scratch rings round robin.

**Output**

TX send packets out on the physical links of the network processor. The logic of TX is built to control the hardware needed to send packets that reside in DRAM out on the physical network medium.

**Concurrency**

TX can be configured to run on 1 or 2 microengines. We use the configuration of 1 micro-engine in our project. It utilizes network processor hardware optimized for transferring data directly from DRAM into the transfer buffers used by the hardware to generate packet data on the physical medium.

**Packet ordering**

Packet order is maintained per scratch ring. So any number of packets pulled from the same scratch ring will remain their relative packet order when the packets are output on the corresponding physical link. Packets pulled from different scratch rings have no ordering guarantees, however a simple round robin policy is used in this project's configuration.

**Performance**

TX can keep up with line rates on all IXP hardware it's bundled with given proper configuration. Hence its performance is at least as good as we need.

### 3.3.4 `Route` microblock

`Route` implements the GridStat routing and rate-filtering algorithm from sectino 2.1.1, and is the software component which determines whether and where a packet is to be forwarded. This microblock get packets from the `PACKET RX` scratch ring and put packets on one or more of the `TX`'s scratch rings.

**Input**

Packets for route processing are pulled from the `PACKET RX` scratch ring.

**Output**

Depending on whether and where a packet is to be forwarded, `Route` put each packet on the scratch ring corresponding to which outgoing link it belongs. If a packet is not to be forwarded anywhere, it should simply be dropped, but is allowed to pass through the router so that its performance can be measured from simulation packet logs.

**Concurrency**

`Route` can run on 1, 2, 3, or 4 microengines.

**Packet ordering**

Packet ordering is inherit when using the ordered threading model since the order in which threads do scratch ring put and get operations (interface between microblocks) is always the same. E.g. Thread 1 gets a packet, then thread 2 gets a packet, and etc. until thread 8 gets

34

a packet. Then Thread 1 puts a packet, then thread 2 puts a packet, and etc. until thread 8 puts a packet. And only when that entire sequence is complete it start again with thread 1 getting another packet.

**Routing algorithm**

The `Route` microblock implements the algorithm described in section 2.1.1 After getting a packet $P$ from the `PACKET_RX` scratch ring, `Route` makes a meta-data and reference copy of $P$ for each link except the one $P$ came from and one other where $P$ itself will be used. It's not a full packet copy, since only the meta data is copied. In a production implementation this should be changed to use proper reference counting for multicast support. In the configuration used here, there are only 3 links, hence only 1 copy is needed. e.g. If $P$ arrived at link 0, $P$ would be the packet to send out on link 1, and the one copy of $P$ would be the packet to send out on link 2. Even though packets are copied and prepared for sending this early on in the algorithm, a forwarding decision must still be made on a per link basis. After starting a packet pseudo copy, the packet data needed to decide whether and where it's to be forwarded is downloaded from the packet buffer in DRAM. For GridStat this is currently the variable-ID $P.varid$ and time-stamp $P.timestamp$ , a total of 12 bytes which can be extracted in a single DRAM read operation. After initializing some registers used to store forwarding decisions, preparing packet copy to link assignment, and getting all needed packet data from DRAM, execution enters a loop which iterates through all links (skipping the link where $P$ arrived). For each link $L$, we call the `dl_gridstat_filter` macro, which will decide

whether $P$ is to be forwarded on $L$. The per link decision is marked on a bit-array, where a set bit signifies a decision to forward the packet on the corresponding link. When the loop has completed, the result is the updated decision bit-array. `Route` then enters a final loop iterating through all links sending or dropping packets according to the decision bit-array. After iterating through all links and forwarding/dropping packets, the algorithm starts over by getting another packet from the `PACKET_RX)` scratch ring.

The `dl_gridstat_filter` macro uses a link number, and the following fields of $P$ as input: variable-ID, time-stamp, publication rate. As output a single register called `forward` that is given value false if the packet is to be dropped, and true if the packet is to be forwarded. First in the macro we ensure that the routing table entry of $P.variableid$ is in local memory, then we perform a filtering operation and use the truth value to determine whether $P$ is to be forwarded. The operation is as follows: $(P.timestamp + P.pubint/2) \ modulo \ divisor) < P.pubint$, and corresponds to lines 8 and 9 of the routing algorithm in section 2.1.1. The *modulo* operator is implemented with the `remainder64x32` macro as the IXP hardware does not provide any general division or modular operator. The macro does not implement general integer division or modulo in software either, rather it uses reciprocal multiplication to compute the 32-bit remainder of a 64-bit dividend which is exactly what the algorithm requires. Reciprocal multiplication exploits some properties of manipulating numbers that are represented in binary, i.e. the fact that $x/2^n = x >> n$ and $x * 2^n = x << n$. In order to understand how reciprocal multiplication works, consider the following. By definition:

$$Remainder = Dividend \ modulo \ Divisor$$

and

$$Dividend = Quotient * Divisor + Remainder$$

If we only need to find the $Quotient$ (like in division), then the value of $Remainder$ does not matter and we may assume that $Remainder = 0$ allowing the previous equation to transform into:

$$Quotient = \frac{Dividend * \frac{2^n}{Divisor}}{2^n}$$

Let $R_n$ be the $n$-bit reciprocal $\frac{2^n}{Divisor}$ of $Divisor$. Appropriate values for $n$ and $R_n$ are computed based on $Divisor$ before stored in the routing table, see the `computeReciprocal` method in A.2 for details on how to compute them. Given $n$ and $R_n$ from routing table lookup at runtime, the $Quotient$ may be computed through only multiplication and shift operators as:

$$Quotient = (Dividend * R_n) >> n$$

and once the $Quotient$ is known, then $Remainder$ may be computed through only multiplication and subtraction operators as:

$$Remainder = Dividend - (Quotient * Divisor)$$

The value of the dividend is given by the relevant packet's time-stamp, the divisors are known in the algorithm's inner loop as the subscriber interval for which $n$ and $R_n$ are pre-computed and to be found in the routing table.

**Collapsing subscriber rate list**   See figures 3.4 and 3.5 on how to collapse subscriber lists for the routing table.

Figure 3.2: Routing table data structure

The routing table data structure consists of a single linked list of routing table entries. Each routing table entry has an associated variable-ID, publication rate, and a linked list of subscriber intervals for each port (link) of the hardware. A routing table entry can also be found by a hash lookup using the variable-ID as key. The IXP2xxx platform have hardware capable of doing very fast hash generation.

Figure 3.3: Routing table entry details

Dynamic data structures like linked lists are cumbersome and slow to work with from the data-plane. Microengines prefer data structures that are small and sequential in memory as to minimize memory read operations. Routing table entries therefore use a fixed size array for each port which contains a collapsed list of subscriber intervals.



Figure 3.4: Original subscriber list

As an example consider the above list of subscriber intervals. Due to the modular arithmetic in the `dl_gridstat_filter` macro, we need only check a subset of the subscriber intervals. This subset is found by eliminating any subscriber interval value $SI_i$ when there exist another subscriber interval of lower value $SI_j$ such that $SI_i \ MOD \ SI_j = 0$ holds.

| 20 | 30 | 50 | 110 |

Figure 3.5: Collapsed subscriber list

Above we see the list of subscriber intervals after elimination. Note how 40, 80, 90, and 120 were eliminated due to $40 \ MOD \ 20 = 0$, $80 \ MOD \ 20 = 0$, $90 \ MOD \ 30 = 0$, $120 \ MOD \ 20 = 0$.

**Cache** `Route` gets routing table entry information from a cache that resides in microengine local memory. If the cache does not hold the desired routing table entry, it's downloaded to the local memory cache from SRAM. Only one routing table entry is ever used by a single packet, and the same entry is often used very frequently across packets. The IXP2400 has a per-microengine Content Addressable Memory (CAM) providing a fast and easy to use 16 entry cache.

We can think of the CAM as a 16 row table, where each row has the following fields: A 32-bit tag value that it written by the programmer. A 4-bit state field also written by the programmer. An entry of the time-ordered Least Recently Used (LRU) list. The LRU list keeps track of which row was most and least recently used at all times. Using the CAM is simple and fast. In order to use the CAM as a cache in the `Route` microblock we do as follows. The `cam_clear` instruction is called once per microblock running `Route` upon initialization. Whenever a packet $P$ is processed by one of the 8 threads in a microengine, before we enter the loop iterating through all links calling `dl_gridstat_filter`, a `cam_lookup` on *P.variableid* is performed to check whether the current microengine has $P$'s routing table entry in local memory (cache). If we get a cache hit (the routing table entry is present in

| | | | |
|---|---|---|---|
| Entry  0 | LRU list | 4-bit state | 32-bit tag |
| Entry  1 | LRU list | 4-bit state | 32-bit tag |
| Entry  2 | LRU list | 4-bit state | 32-bit tag |
| Entry  3 | LRU list | 4-bit state | 32-bit tag |
| Entry  4 | LRU list | 4-bit state | 32-bit tag |
| Entry  5 | LRU list | 4-bit state | 32-bit tag |
| Entry  6 | LRU list | 4-bit state | 32-bit tag |
| Entry  7 | LRU list | 4-bit state | 32-bit tag |
| Entry  8 | LRU list | 4-bit state | 32-bit tag |
| Entry  9 | LRU list | 4-bit state | 32-bit tag |
| Entry 10 | LRU list | 4-bit state | 32-bit tag |
| Entry 11 | LRU list | 4-bit state | 32-bit tag |
| Entry 12 | LRU list | 4-bit state | 32-bit tag |
| Entry 13 | LRU list | 4-bit state | 32-bit tag |
| Entry 14 | LRU list | 4-bit state | 32-bit tag |
| Entry 15 | LRU list | 4-bit state | 32-bit tag |

Figure 3.6: CAM structure

Each microengine has it's own local Content Addressable Memory. The IXP hardware provide very fast CAM compare and lookup instruction with automatic LRU list updates making the CAM the perfect cache support mechanism.

local memory), then nothing is explicitly done to update the CAM but the hardware will mark the entry that matched as Most Recently Used (MRU). If we get a cache miss, then the routing table entry must be downloaded from SRAM to local memory before it can be used. Before the download process starts the CAM tag field of the LRU item is written with the variable-ID of the relevant routing table entry. However, the state field is marked with a value indicating it being invalid, so that if another thread $T_x$ try a lookup for the same variable-ID before the routing table entry is downloaded from SRAM, $T_x$ will get a cache hit but with associated state value indicating that it's being downloaded. The `dl_gridstat_filter` macro being called inside the loop iterating through all links, performs a second `cam_lookup`. This will always produce a cache hit, but the data might be either valid or invalid (being downloaded). If the data is valid the macro continues the filtering process using that data, If the data is invalid then $T_x$ will yield (force context switch) and try the lookup again until the result of the lookup is a hit with valid data. A cache miss is structurally impossible for $T_x$ since there are fewer other packets processed on this microengine (at most 7) before $P$ gets processed than there are entries available in the CAM, hence the entry which has the matching variable-ID but invalid data will never become the LRU item (the LRU item is the item to always overwrite when cache is full) while P is being processed.

**Performance**

Several issues make good performance hard to achieve in the `Route` microblock. Accessing DRAM for packet data extraction, accessing the SRAM for routing table entry lookups.

Performance of the routing algorithm itself is also very important. There are four kinds of RAM access done within `Route`. Scratch, DRAM, SRAM, and local memory. Scratch access is used getting packets from the `PACKET_RX` scratch ring, and putting packets on the transmit rings. Local memory access is used by the routing table entry cache and is discussed in a paragraph below. DRAM and SRAM access are both discussed in their own separate paragraphs below.

**Routing algorithm**   Performance on the routing algorithm itself depends on several variables. Most importantly it depends on how many entries in the subinterval list that need be checked. There are two nested loops, the outer loop iterates through all outgoing links, and the inner loop iterates through all subintervals in the routing table entry lookup table. The outer loop need worst case consider all outgoing links (except 1) as they all have different lookup tables. The inner loop has several complications to overcome. It uses modular arithmetic for every subinterval in the subinterval list to decide whether the packet can be forwarded. A separate macro `remainder64x32` implements the modular operation. This macro does at least 60 instructions making it very expensive.

**DRAM access**   DRAM access, this is only to extract information about the packet needed by the routing logic. It's a single read-only access done once per packet pulled from the `PACKET_RX` scratch ring. DRAM lookup is never needed on copies of packets, as everything (except destination address) about the copy is the same as the original.

**SRAM access**   SRAM access is used to look up routing table entries.

**Cache** The cache in `Route` is used to store routing table entries that have been recently used by the microengine running `Route`. On a cache hit we see that routing table entry lookup is very fast, reading from local memory takes at most 3 cycles per long word whereas a register read takes 1 cycle. On a cache miss we download the routing table entry from SRAM, which takes takes 3 SRAM read operations of 8 long words each which implies waiting on IO.

## 3.4 IXP2xxx synchronization and locking mechanisms

This prototype implementation does not use synchronization mechanisms other than the atomic scratch ring put and get and sram queue and dequeue instructions. Those instructions automatically handle synchronization, and so the microcode of the `Route` microblock has no critical regions from a programmer's point of view. The hardware does however support many different synchronization mechanisms, of which some probably must be used in order to implement features like routing table updates at runtime. In particular, the sram and scratchpad memory have supporting atomic operations to do read-modify-write instructions that are guaranteed to appear atomic to all other atomic operations (but not to non-atomic operations). We already mentioned the special ring and queue instructions which are used in this implementation, but there are also general atomic operations like test-and-set, test-and-clear, test-and-incr, test-and-decr, test-and-add, swap, and a few others. These instructions are more than sufficient to implement any kind of synchronization between microengines, but care must be taken so that synchronization contention does not prevent microengines

from achieving expected degree of parallelism.

## 3.5   Constraints on the IXP2xxx hardware

### 3.5.1   Soft constraints

**Division and Modular operators**

The Intel IXP 2400 and 2800 series of hardware has no hardware support providing division or modular (remainder) operation. So these operations must either be implemented in software, or avoided. The GridStat router algorithm requires modular computation (remainder operator) where all the divisors are known in advance. The straightforward solution of implementing a remainder function in software suffers from very poor performance. As the amount of remainder operations needed per packet can be very many, then even at traffic rates at far less than half of line capacity the software would be unable to keep up with traffic rates resulting in massive packet drops once the packet pipeline buffers (scratch rings) fill up. A small but significant improvement over using general software division implementations is to use precomputed divisor reciprocals to speed up the remainder function. See section 3.3.4 for details.

**Precomputed reciprocals**   The solution using precomputed reciprocals requires extra work at the routing table setup and manipulation stages. In the implementation example we only do this work at routing table setup as we omit changing the routing table at runtime

for simplicity.

## 3.5.2  Hard constraints

**Microengines**

The IXP2400 has 8 microengines. A minimum of 2 of them are always used to control the physical receive and transmit hardware, and are usually occupied by the `Ethernet Receive` and `Ethernet Transmit` microblocks provided by the IXA. This leaves 6 microengines to use in the application. As demonstrated by applications like IP Forwarding, Diffserv, and others, this suffice for most applications needing typical IXP2400 line speed throughput. The GridStat routing algorithm (see section 2.1.1) have one major scalability challenge, it allows many unique subscriber intervals (no common factors except the value 1) per status variable. This challenge cannot be completely overcome by adding more microengines to the hardware, as the time complexity of the algorithm grows much faster than any reasonable hardware upgrade in amount and speed of microengines could mitigate. With a frequency of 600 MHz the microengines should not spend more than about 100 cycles per packet in order to keep up with line rates of 2-3Gbps. Applications that need more instructions than this must run their microblock on several microengines in parallel to mitigate this. However, as the amount of microengines increase, the contention of critical regions get worse and requires careful design and implementation to fully utilize all CPUs which isn't always possible.

**Memory**

The 256MB of DRAM currently available on typical IXP2400 is by far sufficient in order to implement a complete GridStat status router. Part of the DRAM is used as system memory of the local network processor OS, the rest is used to hold packet data. The SRAM of 64MB should suffice for routing tables in the GridStat project with proper data structure design. If we assume routing table entries of 64 byte, this adds up to a routing table capacity of one million flows. The scratch memory area of 8MB is used only for scratch rings in this project and the size is plenty sufficient. Microengine local memory of 640 long words of 4 bytes each suffice for a medium size (16 entries) cache. The microengine CAM only has 16 entries on the IXP2400 hardware, which is small but could be utilized for larger groups like routing table pages and not only status variable ids.

## 3.6 Features not included in the status router prototype

In order that the network processor prototype described in this chapter may become a full-featured GridStat status router, the following functionality must be added:

- Add ability for the GridStat management-plane to interface with the status router at runtime. This is generally about updating the SRAM routing table as publishers are added/removed and subscribers add/remove or change their subscriptions or as relevant network topology changes. This would also require proper routing table synchronization

between management-plane and the microengines who share the routing table data structure.

- Mode switching must be implemented, with appropriate interface for controlling it either via management packages arriving on the data-path or other management messages arriving through host computer PCI bus or Ethernet debug interface.

- Support for packets with more than one status variable contained in it and/or more than one sample of the same status variable in the same packet. Currently, the prototype expects only packets with a single sample of a single status variable.

- The mapping of variable-IDs to routing table entry SRAM-addresses should be done through hash lookups rather than be hard-coded in order to support numerous values. This would also require a change in the routing table initialization script.

- The routing table must be used for publish interval value lookups, rather than the hard-coded value of 10 millisecond used. This is a very simple task.

# Chapter 4

# Analysis and Simulations

This chapter will demonstrate how the GridStat status routing algorithm (see section 2.1.1) may perform on an IXP2400 chip. Analysis form the basis of expected and ideal performance while simulation is used to estimate real performance. The SDK used to implement the status router includes a simulator to run microengine code, this is a true simulator that interprets the same compiled code that would be used with real hardware. The simulator is deterministic and cycle accurate, which means that running microengine code in the simulator is predictable and will produce results as if it was run in real hardware, but in slow motion. The simulator is most easily used with fixed rate incoming packet streams, hence testing many different throughputs requires many different simulations. Each experiment in this chapter consists of several independent simulations where hardware configuration, code, and scripts remain the same while the throughput of the incoming packet stream varies from one simulation to another. The throughput capacity and expected packet processing times

are estimated for each experiment and then compared across experiments in a summary. The analysis and simulation results also reveal some scalability properties that are summarized in the end.

## 4.1   Experiment configuration details

Two configurable aspects of simulation determine the outcome of the experiments:

- The stream of packets arriving at the various ports, or "input" for short. The definition of input also includes a simulation configuration option of the packet arrival rate for each port measured in Mbps. Note that the actual timestamps in the status variables described by packets are synthetic and that the difference in timestamps from one packet to the next do not reflect the time passed between their respective packet arrival. The input data is on the GridStat packet format where the variables relevant to routing are the status variable id and the publish time-stamp of the packet. These variables must be inspected for each packet at runtime in order to determine whether and where a packet is to be forwarded. In a single simulation run, the input stream is a fixed rate, so in order to see the performance of different throughputs the simulation is rerun within the same experiment with different input rates. Variable (non-fixed) input rates are not analyzed or simulated in these experiments. The packet data stream is a file generated by a Java program (see code listing A.1) with a set of configurable options. The bytes of the packet data stream used for all experiments are created by invoking the following parameters to the PacketStreamGenerator:

- "experiments-p0 1 500 244813791234"

- "experiments-p1 2 500 244813667203"

- "experiments-p2 3 500 244813641339"

- The routing table. The routing table holds for each status variable and outgoing port a collapsed list (see 3.3.4 for details) of subscriber intervals (i.e. the computed reciprocals of the intervals). The routing table setup script is generated by a Java program (see code listing A.2) with a set of configurable options. The routing table is kept exactly the same across all simulation runs within a single experiment, and is never the same across experiments.

Other aspects that could significantly change the result of the experiments are kept the same across all experiments. These aspects include:

- The routing algorithm used. One single algorithm implementation is used.

- Microengine code other than the routing algorithm.

- Hardware simulation configuration with the exception of throttling input rates done by changing incoming line speed capacity. This is set to be an IXP2400 B0 chip with a microengine clock frequency of 400MHz. The incoming and outgoing network devices is set to be a x32MPHY4 with 3 ports each and the external media bus clock runs at 104MHz.

- Scripts for simulation setup and SRAM/DRAM memory layout, with the exception

of routing table setup in SRAM. This includes things like packet buffers and packet meta-data memory area.

A simulation configuration is a specific combination of hardware configuration, initialization scripts (includes routing table setup), microengine code, and input data streams. These parameters control the runtime behavior and output of the simulated network processor and suffice to allow the simulator to deterministically process the input and produce output in the form of outgoing packet streams. The simulator is deterministic, which means that every time the same simulation configuration is simulated it will always produce the exact same result. A simulation run, means the simulator running a specific simulation configuration starting with hardware initialization and running some number of microengine clock cycles while input is injected through one of the simulators input devices. Exactly how many microengine clock cycles to run is either controlled manually or by predefining some number of packets that should be injected or transmitted by the simulator.

Simulations are generally run so that at least 1000 incoming packets have completed processing. However, in some simulation runs where the incoming throughput is either very low or very high compared with the capacity of the router, the results are obvious after much fewer packets have been processed and such simulations may be stopped before 1000 packets have been processed.

## 4.2 Performance measurements

Performance of networks in general and also routers are chiefly measured through three measurements with sometimes very dynamic properties. They are throughput, latency, and jitter. Properties of throughput receive particularly much attention here because it relates so strongly to scalability of the GridStat algorithm. Throughput, latency and jitter challenges can often be satisfied with hardware and/or software improvements because most routing applications have well understood algorithms with performance relatively independent of packet content. The GridStat routing challenge is special because the routing algorithm itself has runtime dependencies to a management plane (routing table) and packet stream data contents that could significantly affect its performance and scalability properties.

Statistics of packet processing time (delay) at different throughputs are key when measuring performance of these experiments. In general the throughput capacity of a network processor router is dependent on runtime features like complexity of incoming packets, routing table management and current internal thread utilization. These experiments are designed so that the current internal thread utilization is the only significant runtime aspect that is directly related to the throughput capacity. Statistics like mean, minimum, maximum, and standard deviation of the packet delays are relevant only as long as the input rate does not exceed the throughput capacity of the router. Each experiment attempt to estimate the router's throughput capacity by the same input stream at different throughputs and observing the dynamics of packet processing times. Each experiment estimates a routing table runtime complexity for its input when processing exactly 1000 packets.

In general the routing table throughput capacity is said to be exceeded if the incoming packet stream enters faster than the router is capable of processing, which must eventually lead to unwanted packet drops. A method of estimating the routing table throughput capacity is needed. Let $Delay(s, t)$ be a function so that its output is the expected processing time of a packet that enters the router at time $t$ and where $s$ represents a simulation run. Delay samples are observed by analyzing the output logs of the simulation run of $s$, and form the data set used to estimate $Delay(s, t)$ through a least squares linear regression model. $Delay(s, t)$ is on the linear form $Delay(s, t) = intercept + (slope * t)$. We say that the incoming throughput rate of a simulation exceeds the router's capacity if the *slope* part of the corresponding $Delay(s, t)$ is greater than 0. Approximating $Delay(s, t)$ is done by a Java program (See code listing A.5) that inspects the packet logs after running $s$.

Packet receive and transmit "start of packet" (SOP) and "end of packet" (EOP) is measured in units of media bus cycles since beginning of simulation, simply referred to as cycle X. Packet processing time is also measured in media bus cycles, referred to as cycles. Packet processing time is the difference between packet transmit SOP and packet receive EOP, in other words the time it spends in the actual packet processing software pipeline along the fast data path. Note that a certain number of media bus cycles do not reflect the number of instructions that the microengines carry out, since media bus cycles are much slower (about 4 times) than microengine cycles and some instructions take more than one microengine cycle to carry out due to internal IO latency.

For a single simulation the following statistics are compiled by analyzing its packet

logs and shown through tables and graphs of the experiments:

- Number of packets entering and leaving the system.

- Inter packet arrival time (in cycles). A mean measurement with a standard deviation generally so low that it can be thought of as an exact measurement. In practice this means we have a fixed rate input stream. However, the simulator will start sending packets right away even before the network processor has been initialized, this will cause the observed inter packet arrival time to be smaller at the beginning of a simulation then quickly even out to a fixed rate.

- Standard deviation of inter packet arrival time (in cycles). This give a handle on the jitter added to the packet stream passing through the router.

- Mean, minimum, maximum and standard deviation of the observations used to estimate $Delay(s, t)$ (in cycles). This is only relevant as long as the simulation does not exceed the router's throughput capacity.

## 4.3 Experiments overview and expectations

A total of 5 experiments is conducted where the only significant change in simulation configuration between any two experiments is the routing table. When we consider a packet being processed by the routing algorithm, the combined length of the collapsed subscriber lists of the packet's corresponding status variable and target ports define the algorithm's

processing time for that packet. In fact these subscriber list lengths (i.e. routing table contents) is the only free variable in the routing algorithm and hence the variable that defines the algorithm's asymptotic time complexity. More than one such subscriber lists are used (due to multicasting) by the algorithm when processing a packet, and part of the payload data (status variable id and time-stamp) of the packets decide how many of the elements in each relevant subscriber list need be considered. Hence it's impossible to determine the average number of iterations of the inner loop of the routing algorithm until the routing table and input packet stream is known. All experiments use the same pre-generated packet stream as input, and only the first 1000 packets passing through is considered in performance evaluation.

Let $e$ be an experiment. Let $RoutingTable$ be a function of $e$ such that $RoutingTable(e)$ defines the routing table of $e$. Let $RuntimeComplexity$ be a function of routing tables such that $RuntimeComplexity(RoutingTable(e))$ defines the runtime complexity of the routing table of $e$. $RuntimeComplexity(RoutingTable(e))$ is the per packet average number of routing table lookups required to route the first 1000 packets from the combined input streams of $e$. Table 4.1 maps $e$ to $RuntimeComplexity(RoutingTable(e))$ which corresponds to the number of inner loop iterations of the algorithm:

Let $K_1$ be the minimum cost of routing a packet (i.e. when there are no subscribers for the corresponding status variable). Let $K_2$ be the cost added per iteration of the algorithm's inner loop.

Let $ThroughputCapacity$ be a function of $RuntimeComplexity$ such that

Table 4.1: Mapping $e$ to $RuntimeComplexity(RoutingTable(e))$

| $e$ | $runtimecomplexity$ |
|---|---|
| 1 | 0 |
| 2 | 2 |
| 3 | 7.08 |
| 4 | 14.16 |
| 5 | 21.24 |

$ThroughputCapacity(RuntimeComplexity(RoutingTable(e)))$ defines the throughput capacity of experiment $e$ measured in megabits per second (Mbps). Intuitively the throughput formula to find the packets per second capacity is:

$$pps_{\frac{packets}{second}} = \frac{computationcapacity_{\frac{instructions}{second}}}{packetcost_{\frac{instructions}{packet}}}$$

The relationship between packets per second and megabits per second is:

$$throughput_{\frac{megabits}{second}} = \frac{pps_{\frac{packets}{second}} * bpp_{\frac{bits}{packet}}}{10^6}$$

Furthermore

$$packetcost_{\frac{instructions}{packet}} = K_{1instructions} + K_{2instructions} * RuntimeComplexity(RoutingTable(e))$$

and

$$computationcapacity_{\frac{instructions}{second}} = me_{microengines} * clock_{\frac{\frac{tick}{second}}{microengine}} * ipt_{\frac{instructions}{tick}}$$

where $me$ is the number of microengines running the routing algorithm in parallel.

Consider the three following assumptions:

Assumption (1): All internal IO of microengines is successfully hidden through microengine concurrency.

Assumption (2): Each instruction take exactly 1 clock cycle to complete, implying $ipt = 1$.

Assumption (3): Running the algorithm in $me$ parallel microengines (where $me > 1$) is $me$ times faster than running the algorithm with only 1 microengine.

By assuming that (1), (2), and (3) holds, the expected throughput capacity is upper bound by the following function:

$$IdealThroughputCapacity(x) = \frac{me * clock * ipt * bpp}{10^6 * (K_1 + K_2 * x)}$$

where $x$ is the number of iterations of the inner loop of the algorithm.

In these experiments the following applies:

$K_1$ and $K_2$ can be estimated by counting the number of instructions in the microengine code that implements the routing algorithm. A minimal (thereby ideal) instruction path of a single packet through our implementation yield $K_1 \geq 320$ and $K_2 \geq 88$.

The microengine clock is set to 400Mhz for all simulations.

The ideal throughput capacity (in Mbps) function for these experiments is therefore:

$$IdealThroughputCapacity(x) = \frac{me * clock * ipt * bpp}{10^6 * (K_1 + K_2 * x)} = \frac{2 * 400 * 10^6 * 1 * 528}{10^6 * (320 + 88 * x)} = \frac{422400}{320 + 88 * x}$$

Figure 4.1 shows the $IdealThroughputCapacity(x)$ function drawn as a graph for $x$ such that $0 \leq x \leq 30$.

Figure 4.1: $IdealThroughputCapacity(x)$

A realistic estimation of the throughput capacity would not assume that any of (1), (2), or (3) holds. Though network processors are designed to get close to these ideals, it does take careful implementation of code. Such careful implementation is outside the scope of this thesis, so it must be expected that the implementation used here suffer a significant penalty compared with a more ideal implementation.

For our prototype, let $F$ be the ideal throughput capacity function and let $F'$ be the real throughput capacity function. $F'$ is expected to have the same general shape as $F$ and $F'(x) < F(x)$ is expected to hold for all positive values of x. Exactly how great the difference between $F$ and $F'$ is best explored by estimating $F'$ through simulation.

## 4.4   Experiment 1

### 4.4.1   Setup

Routing table with no subscribers hence a runtime complexity of 0 no matter what the input is. RTSetupCodeGenerator run with no parameters.

### 4.4.2   Analysis

This experiment estimates the real throughput capacity $TC(x)$ when $x = 0$. A more realistic $K_1$ can then be estimated by the equation:

$$K_1 = \frac{me * clock * ipt * bpp}{10^6 * TC(0)} = \frac{2 * 400 * 10^6 * 1 * 528}{10^6 * TC(0)} = \frac{422400}{TC(0)}$$

Observe that the above equation assumes that assumptions (2) and (3) still holds. Though this is not completely accurate, it suffices to quantify the difference between the ideal throughput capacity and an estimation of the real throughput capacity.

A baseline throughput capacity of 692 Mbps (see figure 4.3) might seem like a very low and bad capacity of hardware supposedly capable of almost 4 times this throughput. The reason for this low capacity is that the routing algorithm implementation is not at all optimized, nor designed to utilize the full potential of the network processor. Having such a "slow" implementation is in fact a benefit when attempting to uncover scalability properties of the algorithm. When the algorithm implementation performs way below the capabilities of every other part of the routers software and hardware even for simple inputs and routing tables, we can be sure that it will always be the bottleneck of throughput capacity. This makes every change affecting the routing algorithm's packet processing times in turn affect the observed total packet processing times proportionally. The algorithm as a slow bottleneck also makes the router very sensitive to variables like routing table runtime complexity which in turn becomes very visible in packet processing times.

## 4.5   Experiment 2

### 4.5.1   Setup

Routing table with a runtime complexity of exactly 2 no matter what the input is.

Figure 4.2: Simulations for experiment 1

Figure 4.3: Slopes of experiment 1

Throughput capacity is approximately 692 Mbps.

Table 4.2: Incoming and outgoing inter-packet-delay per simulation of experiment 1

| $s$ | Input (Mbps) | Input Average Inter-packet delay (media-bus cycles) | StdDev | Output Average Inter-packet delay (media-bus cycles) | StdDev |
|---|---|---|---|---|---|
| 1 | 600 | 274.713 | 0.452456 | 265.303 | 122.719 |
| 2 | 690 | 238.887 | 0.317176 | 236.583 | 112.661 |
| 3 | 699 | 235.807 | 0.395132 | 236.519 | 116.301 |
| 4 | 705 | 233.802 | 0.398545 | 235.637 | 69.2877 |
| 5 | 720 | 228.930 | 0.255289 | 235.676 | 122.151 |
| 6 | 750 | 219.773 | 0.418908 | 235.581 | 111.718 |

## 4.5.2 Analysis

The router has a throughput capacity of about 508 Mbps in this experiment, see figure 4.6. The inner loop is run exactly 2 times for every packet, and already there is a noticeable drop in throughput capacity. This experiment represents the "cheapest" routing we expect to get, even if subscribers were strictly constrained in their subscriber intervals.

# 4.6   Experiment 3

## 4.6.1   Setup

Routing table with a runtime complexity of 7.08 for the first 1000 packets of the data stream.

Figure 4.4: Simulations for experiment 2 part 1

Figure 4.5: Simulations for experiment 2 part 2

Figure 4.6: Slopes of experiment 2

Throughput capacity is approximately 508 Mbps.

Table 4.3: Incoming and outgoing inter-packet-delay per simulation of experiment 2

| $s$ | Input (Mbps) | Input Average Inter-packet delay (media-bus cycles) | StdDev | Output Average Inter-packet delay (media-bus cycles) | StdDev |
|---|---|---|---|---|---|
| 7 | 60 | 2747.71 | 0.455071 | 2747.01 | 113.559 |
| 8 | 180 | 915.777 | 0.416712 | 915.599 | 145.635 |
| 9 | 300 | 549.452 | 0.497989 | 541.927 | 119.405 |
| 10 | 330 | 499.509 | 0.500199 | 487.382 | 148.353 |
| 11 | 450 | 366.302 | 0.459476 | 357.191 | 122.251 |
| 12 | 480 | 343.393 | 0.488761 | 332.606 | 165.961 |
| 13 | 495 | 332.997 | 0.0576708 | 325.114 | 75.7650 |
| 14 | 501 | 329.013 | 0.114761 | 318.797 | 160.637 |
| 15 | 507 | 325.107 | 0.308861 | 316.513 | 162.167 |
| 16 | 510 | 323.197 | 0.397699 | 340.587 | 152.215 |
| 17 | 540 | 305.243 | 0.429333 | 338.316 | 150.626 |
| 18 | 600 | 274.713 | 0.452456 | 324.726 | 160.500 |

Table 4.4: Incoming and outgoing inter-packet-delay per simulation of experiment 3

| $s$ | Input (Mbps) | Input Average Inter-packet delay (media-bus cycles) | StdDev | Output Average Inter-packet delay (media-bus cycles) | StdDev |
|---|---|---|---|---|---|
| 19 | 180 | 1026.78 | 0.415519 | 1029.92 | 242.030 |
| 20 | 255 | 646.423 | 0.494362 | 648.530 | 440.270 |
| 21 | 270 | 610.488 | 0.500129 | 602.042 | 429.324 |
| 22 | 276 | 597.231 | 0.421778 | 601.459 | 367.205 |
| 23 | 285 | 578.390 | 0.488021 | 571.173 | 256.472 |
| 24 | 291 | 566.424 | 0.494533 | 561.737 | 260.637 |
| 25 | 294 | 560.652 | 0.476530 | 556.111 | 325.264 |
| 26 | 297 | 554.997 | 0.0576708 | 550.822 | 150.288 |
| 27 | 300 | 549.452 | 0.497989 | 630.566 | 254.958 |
| 28 | 303 | 544.017 | 0.128090 | 623.421 | 325.256 |

## 4.6.2   Analysis

The router seem to be at its limit around 297 Mbps, see figure 4.9. This is a severe decrease in throughput capacity compared with a simple routing table complexity of 2, indicating that the routing algorithm implementation have an expensive "inner loop" relative to the constant time cost.

Figure 4.7: Simulations for experiment 3 part 1

Figure 4.8: Simulations for experiment 3 part 2

Figure 4.9: Slopes of experiment 3

Throughput capacity is approximately 297 Mbps.

## 4.7 Experiment 4

### 4.7.1 Setup

In order to simulate a routing table runtime complexity above 5, the code was slightly changed in this experiment so that the inner loop of the algorithm is run twice rather than once per port per packet. This adds 5 microengine clock cycles to the constant time and doubles the cost of using the routing table. The average routing table runtime complexity is in this experiment 14.16.

### 4.7.2 Analysis

Jitter is getting noticeably higher in this experiment, see table 4.5, this is also expected due to the variance in inner loop iterations of the algorithm being higher.

## 4.8 Experiment 5

### 4.8.1 Setup

In order to simulate a routing table runtime complexity above 5, the code was slightly changed in this experiment so that the inner loop of the algorithm is run three times rather than once per port per packet. This adds 10 microengine clock cycles to the constant time and triples the cost of using the routing table. The routing table runtime complexity is 21.24 for this experiment.
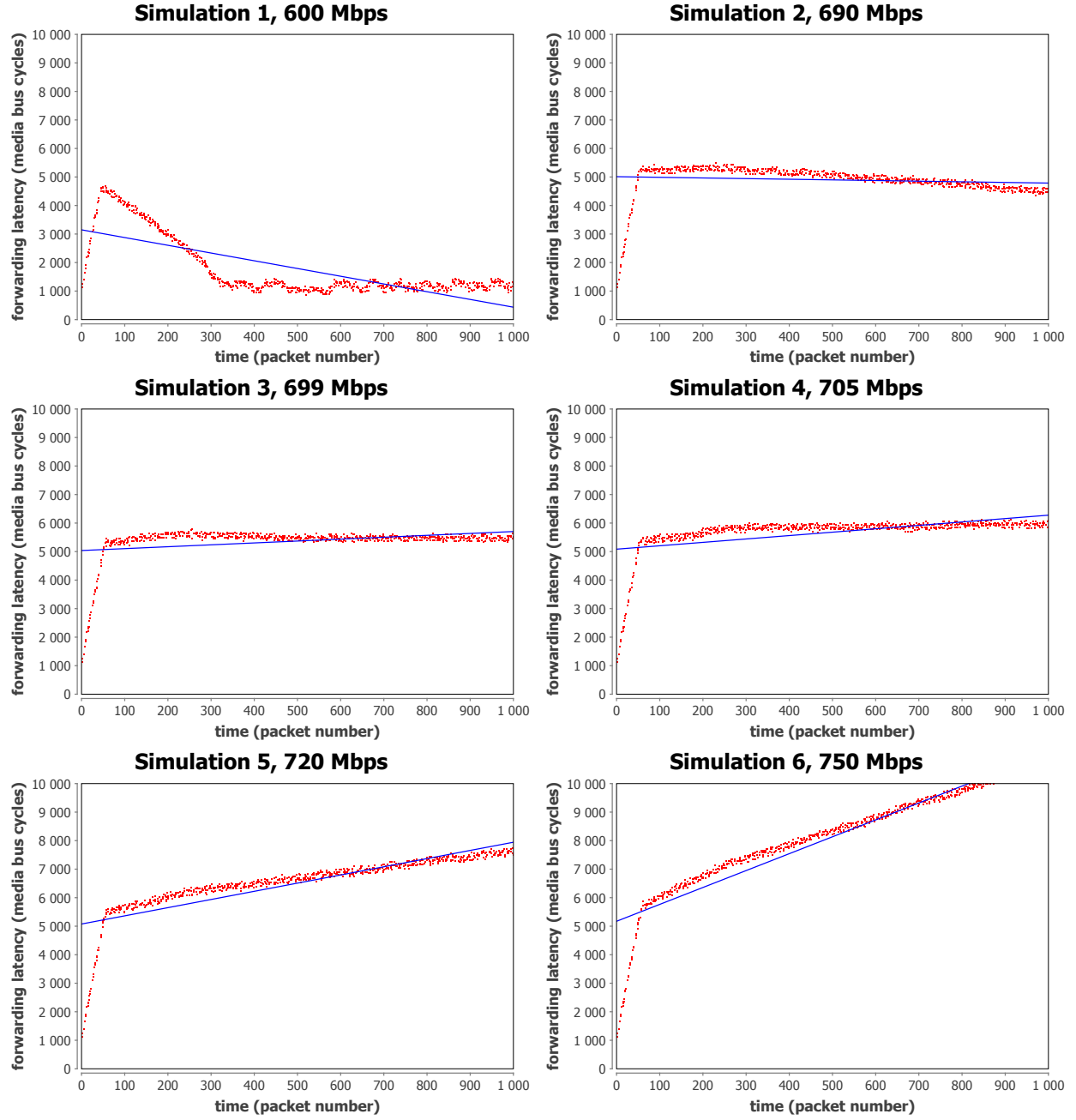
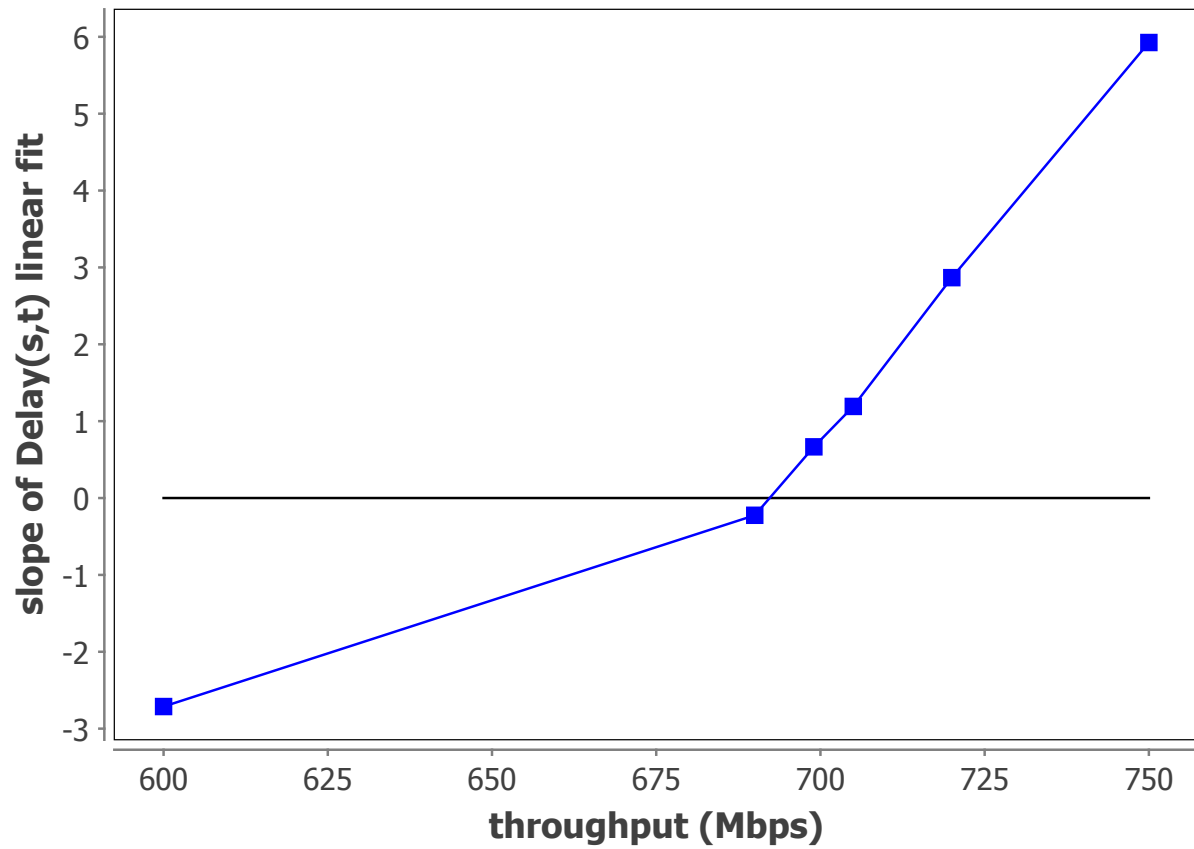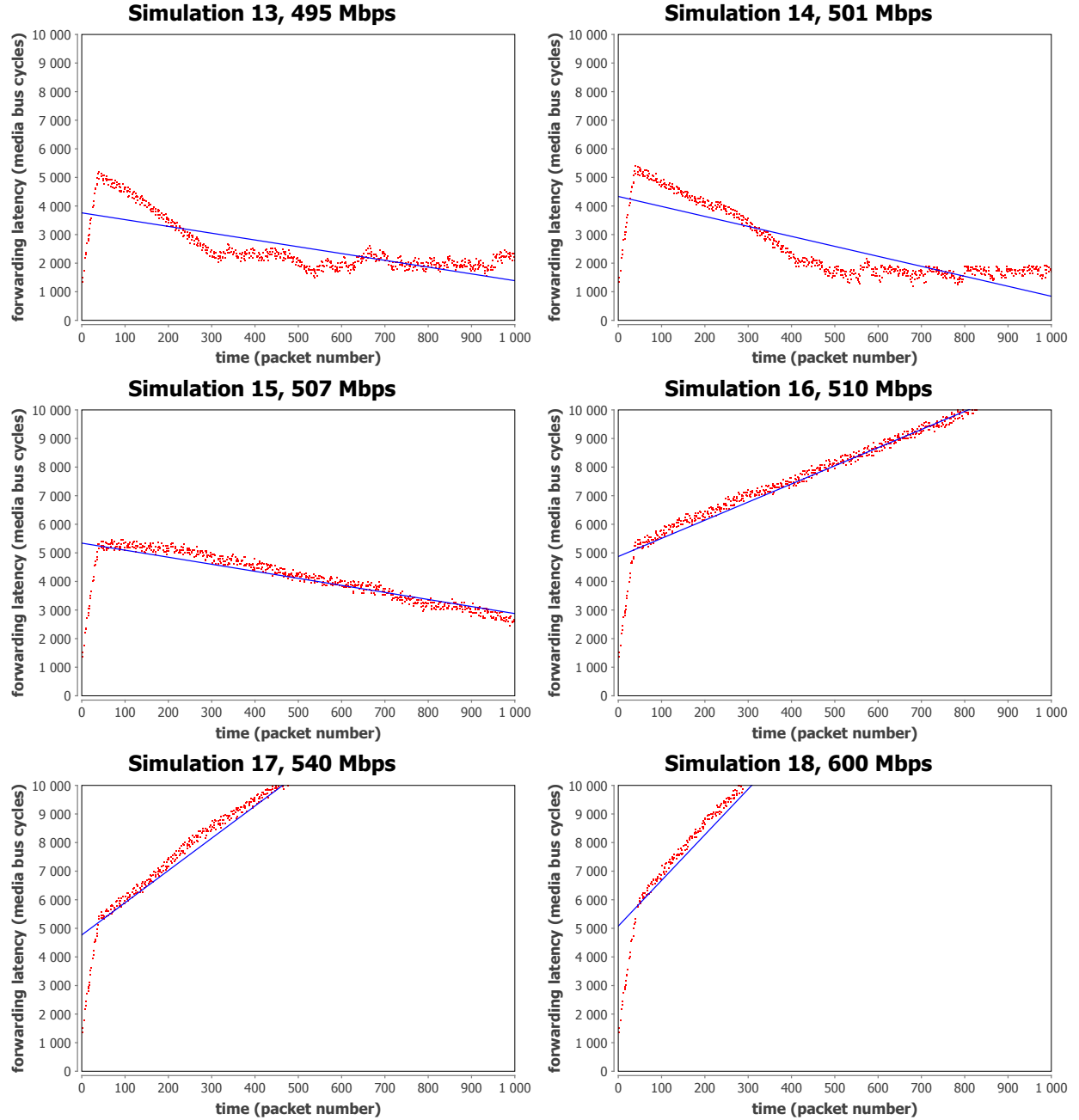Figure 4.10: Simulations for experiment 4

Figure 4.11: Slopes of experiment 4

Throughput capacity is approximately 186 Mbps.

Table 4.5: Incoming and outgoing inter-packet-delay per simulation of experiment 4

| $s$ | Input (Mbps) | Input Average Inter-packet delay (media-bus cycles) | StdDev | Output Average Inter-packet delay (media-bus cycles) | StdDev |
|---|---|---|---|---|---|
| 29 | 150 | 1098.90 | 0.297179 | 1098.85 | 888.507 |
| 30 | 180 | 915.777 | 0.416712 | 927.063 | 664.934 |
| 31 | 186 | 886.263 | 0.440687 | 881.416 | 561.297 |
| 32 | 189 | 872.209 | 0.406741 | 1023.90 | 486.746 |
| 33 | 192 | 858.482 | 0.499962 | 996.034 | 432.342 |
| 34 | 225 | 732.602 | 0.489711 | 802.936 | 352.152 |

### 4.8.2 Analysis

Fewer simulations were done in this experiment, but the results are still clear. The throughput capacity for experiment 5 is about 168 Mbps.

## 4.9 Summary

### 4.9.1 Throughput

The algorithm implementation is not optimized and suffer accordingly. With a simple routing table it can keep up with an input rate of almost 510 Mbps which translates into 1020 Mbps outgoing rate when every packet is forwarded to every other port. As the routing table

Figure 4.12: Simulations for experiment 5

The Y-axis ranges from 5000 to 15000 in these graphs to better show the results, this differ from all the other experiments that show these graphs with Y-axis from 0 to 10000.
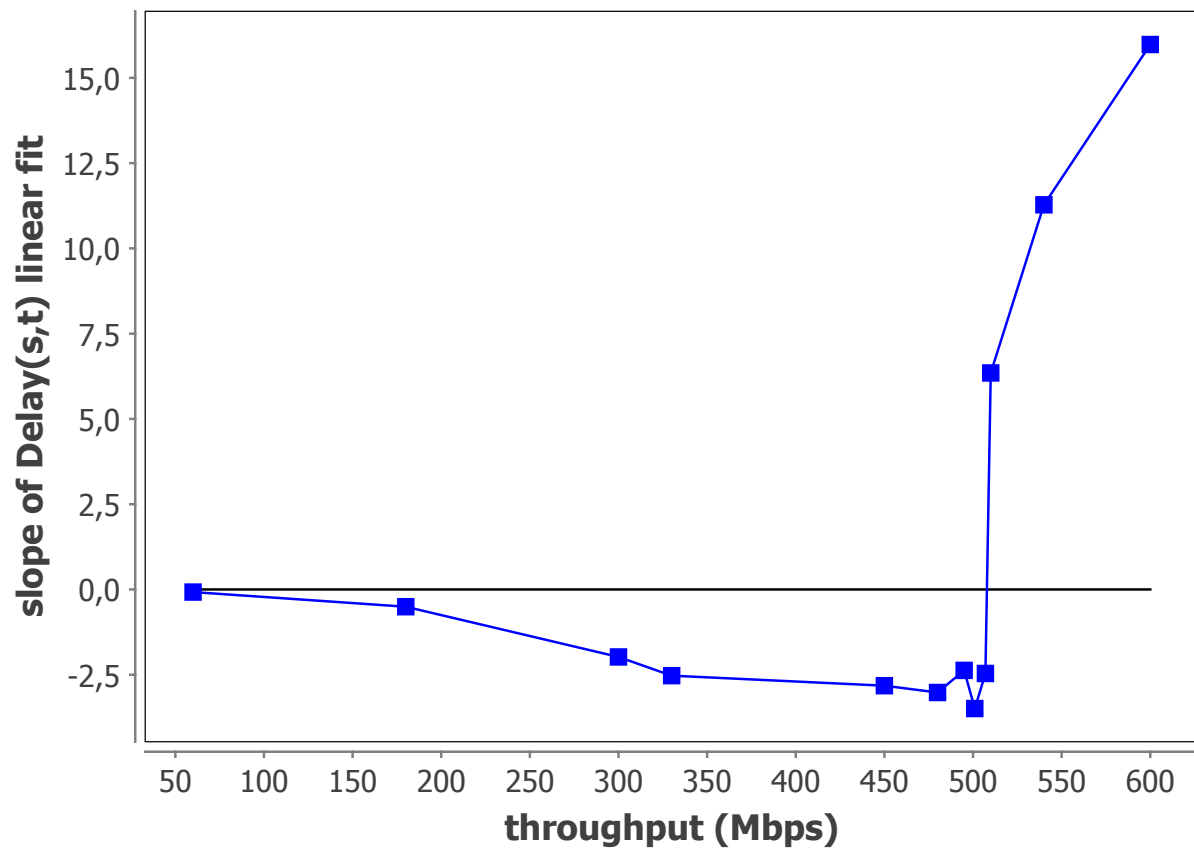
Figure 4.13: Slopes of experiment 5

Throughput capacity is approximately 168 Mbps.

Table 4.6: Incoming and outgoing inter-packet-delay per simulation of experiment 5

| $s$ | Input (Mbps) | Input Average Inter-packet delay (media-bus cycles) | StdDev | Output Average Inter-packet delay (media-bus cycles) | StdDev |
|---|---|---|---|---|---|
| 35 | 165 | 999.090 | 0.286341 | 993.453 | 370.152 |
| 36 | 168 | 981.122 | 0.327724 | 977.416 | 398.056 |
| 37 | 171 | 963.927 | 0.260828 | 1296.24 | 658.474 |
| 38 | 180 | 915.777 | 0.416712 | 1012.33 | 577.386 |

runtime complexity increases, so does the number of filtering operation checks which include a software division, and the throughput capabilities of the router decreases accordingly. With a upper-bound routing table runtime complexity it's possible given enough hardware and a well designed and optimized implementation of the routing algorithm to reach whatever throughput desirable. Even though one can only get so far with code optimization, throughput capacity may be increased through more parallelism usually at the cost of a little higher packet processing times. In figure 4.14 the observed throughput capacity and its fitted function is well below the ideal line. This is implies that our three assumptions of an ideal model were far from holding in these experiments. However, the observed graph has the same general shape, and it's a matter of code design and optimization to bring both the ideal and the observed lines up to production quality which is expected to keep up with line rates of the hardware the router runs.

Figure 4.14: Throughput capacity as a function of inner loop iterations

### 4.9.2 Latency

Theoretically the latency added by the routing microblock is upper bound by 8 times the number of ME cycles along the routing microblock code's longest path assuming that we are able to hide all IO by computing 8 packets in each ME concurrently. It has been shown through simulation in the experiments that the expected latencies depend on how full of packets the microengine pipeline is, though latency is not proportionally higher the fuller the pipeline since a full pipeline will hide a lot more internal IO and a pipeline with a single packet won't hide any of the internal IO related to processing it. It's also important to realize that the longest path in the routing microblock can be very long given a complex routing table. Ideally, the longest path should be static and predictable and not dependent on such runtime variables. It is in other words very hard to predict packet latencies with the current algorithm. These experiments show latencies between 600 and 5000 media bus cycles for stable routers. This translates to approximately 6 and 50 microseconds as the media bus clock runs at 104MHz. In a more optimal implementation, this number should be improved significantly even on such old network processors as IXP2400. It would not be surprising if the improvements yielded merely half the delay on average and worst case within throughput capacity.

### 4.9.3 Jitter

The tables showing inter packet delay of incoming and outgoing packets reveal that the router does indeed impose a lot of jitter on the fixed rate incoming stream. However, such jitter

is easily mitigated through a leaky-bucket algorithm applied after packet routing decisions are made. This would enable more stable outgoing inter packet delays at the cost of a small added total packet processing delay.

## 4.9.4   Algorithm throughput scalability

The algorithm implementation does suffer from not being production optimized but still exhibit the expected scalability properties. It's worth noting that increasing the number of outgoing ports of the system will raise the requirements of the hardware accordingly since potentially as many more packets are going out of the system as there are ports. It's clear from figure 4.14 that the algorithm does not scale well at all with regards to the runtime dependencies. I.e. the number of iterations of the inner loop of the algorithm required on average per incoming packet.

## 4.9.5   Conclusions

In order to support a large amount of status variables (flows) and their subscribers we can't use an algorithm that scales on the order of $O(mn)$, but need an $O(m)$ algorithm, where $m$ is the number of links and $n$ the number of subscribers. The point here being that $m$ is fixed and the hardware can be tuned/bought to fit the need of $m$'s size. $n$ however is a dynamic and potentially very large number that if allowed to dictate the performance of the routing algorithm will severely decrease the algorithms scaling potential. However, it's possible to apply a similar logic to $n$ as was done to $m$. By recognizing that different subscribers may

share at most $N_{max}$ number of non-compatible subscriber rates and buying hardware to fit $N_{max}$ and $m$ together. Ideally $n$ should be upper bound by a relatively small number so that the router may exhibit predictable packet processing times and throughput capacities regards of runtime conditions like routing table contents.

An attempt was made to improve the performance by running the algorithm on additional microengines. However, for unknown reasons it did not yield the expected speedups. Going from running the `Route` microblock on 1 microengine to 2 microengines gave speedup close to expectations while adding a third or fourth microengine yielded almost no speedup at all. Why remains an open question, and is possibly due to implementation weaknesses in the prototype.

Overall, the performance of the network processor is much better than that of comparable computer based implementations. Figure 4.15 shows how experiment 2 compares with the computer based C-implementation of a GridStat status router as demonstrated in [19]. Note that both the throughput and the delay axis are using 10-base logarithmic scales, which implies that both throughput and latency performance of the network processor implementation are 1-2 orders of magnitude better than that of a efficient C implementation which was designed to achieve high performance.

Figure 4.15: Comparing status router performance

# Chapter 5

# Conclusions

In this thesis we have explained in detail and demonstrated a prototype of a GridStat status router implemented on the IXP2400 network processors. This prototype revealed some challenges of network processor programming in general, like the exposure of complex and parallelized hardware and the need for well optimized code in order to keep the performance of code on par with the expected networking line-speeds. In particular, the GridStat rate-filtering algorithm does not inherently provide packet processing time guarantees, and the potential adverse impact of this has been demonstrated in chapter 4. A solution to provide routing latency guarantees by constraining subscriber options is suggested in the following section which meets the one part of GridStat routing that can not be solved with regards to worst-case merely by upgrading hardware. It is also clear after observing a non-optimal performing prototype, that network processors represent packet processing power that far exceeds that of a regular computer based system of comparable vintage.

## 5.1 Performance and constraints

GridStat leaf status routers can be implemented on third-generation network processor architectures, and perform well given some constraints. Consider a status router $R$ with $M$ physical links each capable of $K$ Mbps and with a routing table $T$.

Some scalability challenges of GridStat routing that are also challenges of multicast routing in general can be met by upgrading the hardware of $R$.

- If $T$ dictate that all incoming packets should be forwarded to every other link ($M - 1$ total), the maximum incoming throughput the router could tolerate without dropping packets is $K$ Mbps regardless of what $M$ is. Therefore, $K$ (and probably also the computational capacity) must be increased in order to increase the incoming throughput capacity of $R$ for all possible cases of $T$, this requires a hardware upgrade of $R$.

- $M$ corresponds to the number of iterations in the outer loop of the original routing algorithm. This implies that a larger $M$ requires $R$ to have more computational capacity in order to maintain the same throughput capacity. $M$ is not a free variable, but dictated by the topology of the corresponding GridStat pub-sub routing architecture which in turn is set by design. The GridStat pub-sub topology can therefore be designed in such a way that $M$ is sufficiently small enough for $R$'s computational capacity, and $R$ can be hardware upgraded to match the requirements of a larger but sufficiently small $M$.

The routing algorithm has one chief scalability challenge. The innermost loop requires

an integer modulo operation (which is at least as expensive as a division) where the dividend and divisor are both free variables within the positive integer domain. For each outgoing link in $M$ this operation must in the worst case be executed as many times as there are elements in the corresponding collapsed subscriber list. The number of elements in any collapsed subscriber list is, given sufficiently many subscribers, potentially as many as there are primes in $2^{32}$. Even though any collapsed subscriber list size is finite, it's still too large in a bad case that a status router implementation could have its throughput capacity and/or packet processing time suffer unacceptably compared with a trivial subscriber list size. More formally, let $v_1$ be a status variable and let $n$ be the size of a collapsed subscriber list of $v_1$. The computational complexity of reaching a forwarding decision for a packet describing only a single sample of $v_1$ is $O(Mn)$, which can be simplified to $O(n)$ since $M$ is expected to be a sufficiently small fixed value set design time. As seen by simulating the original algorithm on a IXP2400 chip (figure 4.14) increasing $n$ only from 1 to 5 almost halves the status router's throughput capacity due to the high cost of software division. Network processors usually have a computational capacity only a little higher than sufficient to process packets at the hardware's configured network line speed. This implies that increasing $n$ only by a little could decrease the throughput capacity of a network processor well below the line speed of a single physical link. Even with an ideal implementation, where division (and multiplication) was relatively cheap, increasing $n$ would still at some point decrease throughput capacity below line speeds.

One solution to the scalability challenge of the routing algorithm is to introduce

subscriber rate restrictions at the management level of the GridStat middleware. Let $S_1$ be a set of positive integers so that any integer $i$ in $S_1$ can be evenly divided by all integers in $S_1$ smaller than $i$. Let $v_1$ be a status variable. By restricting all $v_1$ subscriptions to choose subscriber rates from $S_1$, any collapsed subscriber list of $v_1$ would have a size of 0 or 1. If a single set like $S_1$ is insufficient to meet subscriber rate demands, more such sets could be defined. Let $S_2$ be a set of positive integers so that any integer $i$ in $S_2$ can be evenly divided by all integers in $S_2$ smaller than $i$. Allowing subscribers of $v_1$ to choose subscriber rates from $S_1 \cup S_2$ would make any collapsed subscriber list of $v_1$ have a size of 0, 1, or 2. And so it's possible to define $x$ such sets like $S_1$ and $S_2$. By allowing subscribers of $v_1$ to choose subscribers rates from the union of all $x$ sets, any collapsed subscriber list of $v_1$ would have a size no larger than $x$. As long as $x$ is kept small enough, the computational complexity of routing a sample of $v_1$ is guaranteed $O(1)$ with a constant time that depends on implementation and the size of $x$ which should be chosen at GridStat management level.

A more dynamic approach to solve the scalability challenge of the original routing algorithm could let the GridStat management plane use the knowledge of all the currently registered publishers and subscribers to present new subscribers with restrictions in subscriber rates such that it would not incur routing complexity above a certain threshold in any status router.

In addition to restricting subscribers' possibilities to subscribe at any rates they want, it might be possible to predict the remainders (see line 8 in section 2.1.1) in the routing algorithm. This because the router know the publish intervals and subscriber rates, and

should be able to predict the next remainder of any new incoming packet of a stream where the previous packet has already been processed. This could lead to substantial reduction of the inner loop constant time cost and allow far more iterations of it for a given throughput capacity.

All in all even older network processors like the IXP2400 may be expected to do a lot better than current regular computer based systems. In this thesis, with a far from optimal implementation of the routing algorithm, the packet processing times through the status router is between 5 and 50 microseconds and expected to be maybe as little as half that with an optimized implementation. This is better than a relatively fast C implementation [19], and has the benefit of higher throughput because network processor based software can be upscaled well in throughput capacity through increased parallelization.

These are some aspects of network processor hardware and software architecture that might affect the throughput and/or latency performance of a GridStat status router:

- The design and implementation optimization of algorithm to process packets. This is clear after demonstrating a prototype that only utilizes less than 1/3 of the hardware's potential. This requires expert developer knowledge of both the network processor hardware and the accompanying software development kit. As with general programming, the developer of code for a network processor have many options and ways to design and implement the same algorithm, many of which result in performance far below satisfactory when considering expectations of the hardware.

- Average and max code path length taken through the routing algorithm implemen-

tation. As this length increases, latency will increase and throughput capacity will decrease. The important aspect here is to be as small as possible, and if more instructions are required, then more microengines must be used to run the code in parallel in order to maintain throughput capacity at acceptable rates.

- Number of microengines $N$ running the routing algorithm in parallel. as $N$ increases, we may generally expect throughput capacity to be proportionally higher, at the cost of a likely negligible latency hit.

- As parallelism increases, the shared buses for access to memory and other subsystems could become congested and have a strong adverse effect on both latency and throughput capacity.

- The instruction set of microengines can potentially make a huge difference. In the case of GridStat routing, hardware support for modulo, or even remainder using reciprocal division could potentially bring the instruction count down by a significant amount.

- Clock frequency of microengines. In general, performance is proportional with clock frequencies. However, unless all subsystems (e.g. memory, and SHaC) also have proportionally higher clock frequencies, it's possible that faster microengines congest these subsystems more and won't get the expected speedup.

## 5.2 Future work

There are many possibilities to continue research of utilizing network processors in the GridStat middleware data delivery plane, some ideas are:

- Implementing a status router on a "modern" network processor, there are already network processor hardware available that can do more than ten times the throughput of the hardware simulated in this thesis, with at least 40 microengines potentially allowing massively parallelized designs.

- Implement a full-feature status router, including integration with the management layer of GridStat and ensure performance can be upheld while doing lots of management control at runtime.

- Investigate possibilities for line-speed encryption

- Given powerful enough network processors, one could even think of running security applications like intrusion detection in parallel with the normal packet processing pipeline and implementing the appropriate management actions to be taken in GridStat when an intrusion is detected.

# References

[1] July 2005. Intel IXP2400 Network Processor Hardware Reference Manual.

[2] Stian F. Abelsen, Harald Gjermundrød, David E. Bakken, and Carl H. Hauser. Adaptive data stream mechanism for control and monitoring applications. In *Proceedings of 1st International Conference on Adaptive and Self-adaptive Systems and Applications (ADAPTIVE'09)*, pages 86–91, Athens, Greece, November 2009.

[3] Yeim-Kuan Chang and Fang-Chen Kuo. Packet processing with blocking for bursty traffic on multi-thread network processor. In *High Performance Switching and Routing, 2009. HPSR 2009. International Conference on*, pages 1–6, June 2009.

[4] Douglas E. Comer. *Network Systems Design using Network Processors*. Pearson Prentice Hall, 2004.

[5] C. Hauser I. Dionysiou H. Gjermundrd L Xu D Bakken, A. Bose and S. Bhowmik. Towards more extensible and resilient real-time information dissemination for the electric power grid. In *Proceedings of Power Systems and Communications Systems for the Future*, September 2002.

[6] Harald Gjermundrød, David E. Bakken, and Carl H. Hauser. Integrating an event pattern mechanism in a status dissemination middleware. In *Proceedings of 1st International Conferences on Pervasive Patterns and Applications (PATTERN'09)*, pages 259–264, Athens, Greece, November 2009.

[7] Harald Gjermundrød, David E. Bakken, Carl H. Hauser, and Anjan Bose. GridStat: A flexible QoS-managed data dissemination framework for the power grid. *IEEE Transactions on Power Delivery*, 24(1):136–143, January 2009.

[8] S. Govind, R. Govindarajan, and J. Kuri. Packet reordering in network processors. In *Parallel and Distributed Processing Symposium, 2007. IPDPS 2007. IEEE International*, pages 1–10, March 2007.

[9] Xiaofeng Guo, Jinquan Dai, Long Li, Zhiyuan Lv, and Prashant R. Chandra. Latency hiding through multithreading on a network processor. In *PPoPP '07: Proceedings of the 12th ACM SIGPLAN symposium on Principles and practice of parallel programming*, pages 130–131, New York, NY, USA, 2007. ACM.

[10] Carl H. Hauser, David E. Bakken, Ioanna Dionysiou, K. Harald Gjermundrød, Venkata Irava, Joel Helkey, and Anjan Bose. Security, trust and qos in next-generation control and communication for large power systems. *International Journal of Critical Infrastructures (Interscience)*, 4(1-2):3–16, 2008.

[11] Xianghui Hu, Xinan Tang, and Bei Hua. High-performance ipv6 forwarding algorithm for multi-core and multithreaded network processor. In *PPoPP '06: Proceedings of the eleventh ACM SIGPLAN symposium on Principles and practice of parallel programming*, pages 168–177, New York, NY, USA, 2006. ACM.

[12] Xin Huang and Tilman Wolf. A methodology for evaluating runtime support in network processors. In *ANCS '06: Proceedings of the 2006 ACM/IEEE symposium on Architecture for networking and communications systems*, pages 113–122, New York, NY, USA, 2006. ACM.

[13] H. Gjermundrød I. Dionysiou and D. Bakken. Fault tolerance issues in publish-subscribe status dissemination middleware for the electric power grid. In *Supplement of the International Conference on Dependable Systems and Networks (DSN-2002)*, pages B–62–63, June 2002.

[14] Ryan Johnson. Obtaining high performance phasor measurements in a geographically distributed status dissemination network. Master's thesis, Washington State University, August 2005.

[15] Carl Hauser Dave Bakken K. Harald Gjermundrød, Ioanna Dionysiou and Anjan Bose. Flexible and robust status dissemination middleware for the electronic power grid. September 2003.

[16] Duo Liu, Zheng Chen, Bei Hua, Nenghai Yu, and Xinan Tang. High-performance packet classification algorithm for multithreaded ixp network processor. *ACM Trans. Embed. Comput. Syst.*, 7(2):1–25, 2008.

[17] Arindam Mallik, Yu Zhang, and Gokhan Memik. Automated task distribution in multi-core network processors using statistical analysis. In *ANCS '07: Proceedings of the 3rd ACM/IEEE Symposium on Architecture for networking and communications systems*, pages 67–76, New York, NY, USA, 2007. ACM.

[18] Gokhan Memik and William H. Mangione-Smith. Evaluating network processors using netbench. *ACM Trans. Embed. Comput. Syst.*, 5(2):453–471, 2006.

[19] Sunil Muthuswamy. System implementation of a real-time, content based application router for a managed publish-subscribe system. Master's thesis, Washington State University, August 2008.

[20] Chris Ostler, Karam S. Chatha, Vijay Ramamurthi, and Krishnan Srinivasan. Ilp and heuristic techniques for system-level design on network processor architectures. *ACM Trans. Des. Autom. Electron. Syst.*, 12(4):48, 2007.

[21] Yaxuan Qi, Bo Xu, Fei He, Baohua Yang, Jianming Yu, and Jun Li. Towards high-performance flow-level packet processing on multi-core network processors. In *ANCS '07: Proceedings of the 3rd ACM/IEEE Symposium on Architecture for networking and communications systems*, pages 17–26, New York, NY, USA, 2007. ACM.

[22] Erlend S. Viddal, David E. Bakken, Harald Gjermundrød, and Carl H. Hauser. Wide-area actuator rpc over gridstat with timeliness, redundancy, and safety. In *International Conference on Complex, Intelligent and Software Intensive Systems (CISIS'2010*, page Accepted for Publication, Krakow, Poland, February 15-18 2010.

[23] Jean Walrand and Pravin Varaiya. *High-Performance Communication Networks*. Morgan Kaufmann Publishers, 2000.

[24] Ning Weng and Tilman Wolf. Analytic modeling of network processors for parallel workload mapping. *ACM Trans. Embed. Comput. Syst.*, 8(3):1–29, 2009.

[25] T. Wolf and Ning Weng. Runtime support for multicore packet processing systems. *Network, IEEE*, 21(4):29–37, July-August 2007.

[26] Chunqing Wu, Xiangquan Shi, Xuejun Yang, and Jinshu Su. The impact of parallel and multithread mechanism on network processor performance. In *Grid and Cooperative Computing, 2006. GCC 2006. Fifth International Conference*, pages 236–240, Oct. 2006.

[27] Bo Xu, Yaxuan Qi, Fei He, Zongwei Zhou, Yibo Xue, and Jun Li. Fast path session creation on network processors. In *Distributed Computing Systems, 2008. ICDCS '08. The 28th International Conference on*, pages 573–580, June 2008.

# Appendix A

# Code Listing

## A.1  PacketStreamGenerator.java

Listing A.1: PacketStreamGenerator.java

```java
package thesis;

import java.io.ByteArrayOutputStream;
import java.io.File;
import java.io.FileNotFoundException;
import java.io.FileOutputStream;
import java.io.IOException;
import java.io.PrintStream;
import java.util.zip.CRC32;

import org.apache.commons.codec.DecoderException;
import org.apache.commons.codec.binary.Hex;

public class PacketStreamGenerator {

  /**
   * This program produce data streams that can be used by the ↩
        Development Workbench for IXP network processors by
   * Intel. A single stream for 1 status variable is created. ↩
        The publish interval is hardcoded to 10.
   *
   * Usage: java thesis.PacketStreamGenerator <name of ↩
        datastream> <status variable id> <first timestamp in ↩
        stream>
   *
```

```
22      */
23    public static void main(String[] args) {
24      String dataStreamName = args[0];
25      int variableId = Integer.parseInt(args[1]);
26      int nPackets = Integer.parseInt(args[2]);
27      long firstTimestamp = Long.parseLong(args[3]);
28      int pubIntMillis = 10; // hard-coded for now
29
30      CRC32 crc32 = new CRC32();
31
32      File f = new File(dataStreamName + ".strm");
33      FileOutputStream fos = null;
34      try {
35        fos = new FileOutputStream(f);
36        PrintStream out = new PrintStream(fos);
37
38        // generate stream header
39        out.println("# IXP1200 Developer Workbench Data Stream ←↩
              File");
40        out.println("# Format Version 6.00");
41        out.println("#********  Do not edit this file ***********"←↩
              );
42        out.println("");
43        out.println("# Begin Data Stream " + dataStreamName);
44        out.println("STREAM_TYPE = Ethernet IP");
45
46        long timestamp = firstTimestamp;
47
48        for (int i = 0; i < nPackets; ++i) {
49
50          // Define packet data
51          String etherHeaderFixedData = "3←↩
                a3a7b7b7bcc4646464646460800";
52          String ipHeaderFixedData = "4500003000000000004063500←↩
                c0a80032c0a80046";
53          ByteArrayOutputStream baos = new ByteArrayOutputStream()←↩
                ;
54          PrintStream ps = new PrintStream(baos);
55          ps.format("%016x%08x%024x%08x", timestamp, variableId, (←↩
                int) 0, i + 1); // print data payload in hex
56          String payloadFixedData = new String(baos.toByteArray())←↩
                ;
57
58          // compute packet crc32
59          byte[] binaryPacketData;
```

96

```
60            binaryPacketData = Hex.decodeHex((etherHeaderFixedData +↩
                  ipHeaderFixedData + payloadFixedData)
61                .toCharArray());
62            crc32.reset();
63            crc32.update(binaryPacketData);
64            baos.reset();
65            ps.format("%08x", crc32.getValue());
66            String etherTrailerFixedDataWrongByteOrder = new String(↩
                  baos.toByteArray());
67            String etherTrailerFixedData = ↩
                  etherTrailerFixedDataWrongByteOrder.substring(6, 8)
68                + etherTrailerFixedDataWrongByteOrder.substring(4, ↩
                    6)
69                + etherTrailerFixedDataWrongByteOrder.substring(2, ↩
                    4)
70                + etherTrailerFixedDataWrongByteOrder.substring(0, ↩
                    2);
71
72            // generate frame in file stream
73            out.println("# Begin Ethernet Frame");
74            out.println("# Begin Ethernet Header");
75            out.println("FIXED_DATA = " + etherHeaderFixedData);
76            out.println("# End Ethernet Header");
77            out.println("# Begin Ip Packet");
78            out.println("# Begin Ip Header");
79            out.println("COMPUTE_CHECKSUM = TRUE");
80            out.println("COMPUTE_PACKET_LENGTH = TRUE");
81            out.println("FIXED_DATA = " + ipHeaderFixedData);
82            out.println("# End Ip Header");
83            out.println("# Begin Data Payload");
84            out.println("USE_FILL_PATTERN = FALSE");
85            out.println("FILL_PATTERN_TYPE = 15");
86            out.println("FILL_START_VALUE = " + (i * 14));
87            out.println("DATA_PAYLOAD_SIZE = 28");
88            out.println("DATA_PAYLOAD = " + payloadFixedData);
89            out.println("# End Data Payload");
90            out.println("# End Ip Packet");
91            out.println("# Begin Ethernet Trailer");
92            out.println("FIXED_DATA = " + etherTrailerFixedData);
93            out.println("# End Ethernet Trailer");
94            out.println("# End Ethernet Frame");
95
96            timestamp += pubIntMillis;
97        }
98
```

```
 99          out.println("# End Data Stream");
100          out.println();
101
102        } catch (FileNotFoundException e) {
103          e.printStackTrace();
104        } catch (DecoderException e) {
105          e.printStackTrace();
106        } finally {
107          if (fos != null) {
108            try {
109              fos.close();
110            } catch (IOException e) {
111              e.printStackTrace();
112            }
113          }
114        }
115      }
116
117 }
```

## A.2    RTSetupCodeGenerator.java

Listing A.2: RTSetupCodeGenerator.java

```
 1 package thesis;
 2
 3 import java.util.ArrayList;
 4 import java.util.Collections;
 5 import java.util.Comparator;
 6 import java.util.List;
 7
 8 public class RTSetupCodeGenerator {
 9
10   static class ReciprocalResult {
11     int divisor;
12     int nbits;
13     long reciprocal;
14   }
15
16   /*
17    *
18    * Input: nbits number of significant bits used in reciprocal.↩
            divisor divisor.
```

```
19      *
20      * Output: returns reciprocal with nbits significant bits.
21      */
22     static void computeReciprocal(int divisor, ReciprocalResult rr↩
           ) {
23       int z;
24       long r;
25       double f_tmp;
26       rr.divisor = divisor;
27
28       // find z so that 0.5 <= (1/d * 2^z) < 1
29       // then 1/d will have z leading zeros between the decimal ↩
             point and its
30       // most significant set bit
31       z = 0;
32       f_tmp = 1.0 / divisor;
33       while (0.5 > f_tmp || f_tmp >= 1) {
34         z++;
35         // f_tmp = 2^z / d
36         f_tmp = ((double) (2 << (z - 1))) / divisor;
37       }
38
39       // compute (int) ( 2^(z+n+1) / d )
40       r = (2l << (z + 32)) / divisor;
41
42       rr.nbits = 32 + z;
43       rr.reciprocal = (long) (r >> 1);
44     }
45
46     static void generate_init_rtent_function() {
47       System.out.println("");
48       System.out.println("int init_rtent( int addr, int varid, int↩
             pubint ) {");
49       System.out.println("  // VID_ENTRY* next");
50       System.out.println("  set_sram(addr, 0x0);");
51       System.out.println("");
52       System.out.println("  // long variableid");
53       System.out.println("  addr = addr + 4;");
54       System.out.println("  set_sram(addr, varid);");
55       System.out.println("");
56       System.out.println("  // long pub_interval, hardcoded as 10 ↩
             for all status variables");
57       System.out.println("  addr = addr + 4;");
58       System.out.println("  set_sram(addr, 10);");
59       System.out.println("");
```

```
60      System.out.println("  // long reserved");
61      System.out.println("  addr = addr + 4;");
62      System.out.println("  set_sram(addr, 0x0);");
63      System.out.println("");
64      System.out.println("  // FILTER_ENTRY subint[NO_PORTS][←
            MAX_SUB_INTS]");
65      System.out.println("  addr = addr + 4;");
66      System.out.println("");
67      System.out.println("  // SUBSCRIBER_ENTRY *sub_head[NO_PORTS←
            ]");
68      System.out.println("  addr = addr + (8*3*6);");
69      System.out.println("  set_sram(addr, 0x0);");
70      System.out.println("  addr = addr + 4;");
71      System.out.println("  set_sram(addr, 0x0);");
72      System.out.println("  addr = addr + 4;");
73      System.out.println("  set_sram(addr, 0x0);");
74      System.out.println("");
75      System.out.println("  return 1;");
76      System.out.println("}");
77      System.out.println("");
78    };
79
80    static void generate_rt_init_header() {
81      System.out.println("int rt_init() {");
82      System.out.println("  int offset;");
83      System.out.println("  int addr;");
84      System.out.println("");
85      System.out.println("  // clear entire routing table");
86      System.out.println("  addr = GRIDSTAT_RTABLE_BASE;");
87      System.out.println("  offset = addr + GRIDSTAT_RTABLE_SIZE;"←
            );
88      System.out.println("  init_sram(0x0, addr, offset);");
89      System.out.println("");
90    }
91
92    static void generate_rt_init_body(RoutingTable rt) {
93      ReciprocalResult rr = new ReciprocalResult();
94      for (StatusVariableSubscriptions sv : rt.←
            statusVariableSubscriptions) {
95
96        System.out.println("  // VARIABLEID " + sv.variableid);
97        System.out.println("");
98        System.out.println("  offset = " + sv.variableid + ";");
99        System.out.println("  addr = GRIDSTAT_RTABLE_BASE + offset←
              * GRIDSTAT_RTENT_SIZE;");
```

```java
100        System.out.println("  init_rtent(addr, " + sv.variableid +↵
              ", 10);");
101        System.out.println("  addr = addr + 16;");
102        System.out.println("");
103
104        if (sv.variableid == 0) {
105          continue;
106        }
107
108        for (int p = 0; p < NO_PORTS; ++p) {
109          System.out.println("  // port " + p);
110          System.out.println("");
111          for (int i = 0; i < MAX_DIVISORS + 1; ++i) {
112            if (sv.divisors[p][i] != 0) {
113              computeReciprocal(sv.divisors[p][i], rr);
114              long rtd = (rr.divisor & 0x07ffffff) | (rr.nbits << ↵
                  27);
115              System.out.println("  // subscriber " + (i + 1) + " ↵
                  with divisor " + sv.divisors[p][i]);
116              System.out.println("  set_sram(addr, 0x" + Long.↵
                  toHexString(rtd) + ");");
117              System.out.println("  addr = addr + 4;");
118              System.out.println("  set_sram(addr, 0x" + Long.↵
                  toHexString(rr.reciprocal) + ");");
119              System.out.println("  addr = addr + 4;");
120              System.out.println("");
121            } else {
122              System.out.println("  addr = addr + 8;");
123            }
124          }
125          System.out.println("");
126          System.out.println("");
127        }
128      }
129    }
130
131    static void generate_rt_init_footer() {
132      System.out.println("");
133      System.out.println("  return 1;");
134      System.out.println("}");
135      System.out.println("");
136    }
137
138    public static class Subscription {
139      int variableid;
```

```
140        int port;
141        int divisor;
142     };
143
144     public static final int NO_PORTS = 3;
145     public static final int MAX_DIVISORS = 5;
146
147     public static class StatusVariableSubscriptions {
148        int variableid;
149        int[] n = new int[NO_PORTS]; // number of active ↩
                subscriptions
150        int[][] divisors = new int[NO_PORTS][MAX_DIVISORS + 1];
151     }
152
153     public static class RoutingTable {
154        List<StatusVariableSubscriptions> ↩
                statusVariableSubscriptions = new ArrayList<↩
                StatusVariableSubscriptions>();
155
156        public StatusVariableSubscriptions find(int variableId) {
157           for (StatusVariableSubscriptions svs : ↩
                 statusVariableSubscriptions) {
158              if (svs.variableid == variableId) {
159                 return svs;
160              }
161           }
162           return null;
163        }
164     }
165
166     public static RoutingTable createRoutingTable(List<Integer> ↩
             args) {
167        // prepare routing table based on user input
168        List<Subscription> subscriptions = new ArrayList<↩
                Subscription>();
169        for (int i = 0; i < args.size(); i += 3) {
170           Subscription s = new Subscription();
171           s.variableid = args.get(i);
172           s.port = args.get(i + 1);
173           s.divisor = args.get(i + 2);
174           subscriptions.add(s);
175        }
176        Collections.sort(subscriptions, new Comparator<Subscription↩
             >() {
177           @Override
```

```
178        public int compare(Subscription arg0, Subscription arg1) {
179          return arg0.port - arg1.port;
180        }
181      });
182      Collections.sort(subscriptions, new Comparator<Subscription↩
           >() {
183        @Override
184        public int compare(Subscription arg0, Subscription arg1) {
185          return arg0.variableid - arg1.variableid;
186        }
187      });
188      RoutingTable rt = new RoutingTable();
189      StatusVariableSubscriptions sv = new ↩
           StatusVariableSubscriptions();
190      rt.statusVariableSubscriptions.add(sv); // variableid 0
191      for (Subscription s : subscriptions) {
192        if (sv.variableid != s.variableid) {
193          sv = new StatusVariableSubscriptions();
194          sv.variableid = s.variableid;
195          rt.statusVariableSubscriptions.add(sv);
196        }
197        if (sv.n[s.port] >= MAX_DIVISORS) {
198          throw new IndexOutOfBoundsException("too many ↩
               subscriptions for variableid " + s.variableid
199            + " and port " + s.port);
200        }
201        sv.divisors[s.port][sv.n[s.port]++] = s.divisor;
202      }
203      return rt;
204    }
205
206    /**
207     * This program will generate a script to be used by the ↩
           Development Workbench to set up the routing table in SRAM.
208     * The script is printed to standard out.
209     *
210     * Usage: java thesis.RTSetupCodeGenerator [<status variable ↩
           id> <port> <subscriber rate >]...
211     *
212     */
213    public static void main(String[] args) {
214
215      // parse arguments
216      List<Integer> list = new ArrayList<Integer>();
217      for (int i = 0; i < args.length; i += 3) {
```

```
218        list.add(new Integer(args[i]));
219        list.add(new Integer(args[i + 1]));
220        list.add(new Integer(args[i + 2]));
221      }
222
223      // setup internal routing table
224      RoutingTable rt = createRoutingTable(list);
225
226      // generate code
227      generate_init_rtent_function();
228      generate_rt_init_header();
229      generate_rt_init_body(rt);
230      generate_rt_init_footer();
231    }
232 }
```

## A.3    Packet.java

Listing A.3: Packet.java

```java
 1 package thesis;
 2
 3 /**
 4  * An instance of this class represents one incoming packet and ←
        hold performance data of that packet as computed by
 5  * looking at the original packet logs from the simulator. Each ←
        incoming packet has potentially N−1 target ports, and
 6  * the delay between each incoming packet's EOP and its ←
        corresponding outgoing packets' SOP form the basis of packet
 7  * processing durations that almost all other performance ←
        measurements are based on.
 8  *
 9  */
10 public class Packet {
11
12   // number of physical ports connected to network processor
13   static final int N = 3;
14
15   static class TransmittedPacket {
16     boolean forwarded = false;
17     long tx_sop;
18     long tx_eop;
19     Double avg5Delay = null;
```

```java
20    }
21
22    static final String[] CSV_HEADERS = { "Packet", "AVG TX Delay"↩
        , "AVG TX last 5 AVG Delay", "RX Port",
23        "RX PacketNum", "Receive SOP", "Receive EOP", "Timestamp",↩
            "Variable ID", "TX0 Forwarded", "TX0 SOP",
24        "TX0 EOP", "TX0 Delay", "TX0 last 5 AVG Delay", "TX1 ↩
            Forwarded", "TX1 SOP", "TX1 EOP", "TX1 Delay",
25        "TX1 last 5 AVG Delay", "TX2 Forwarded", "TX2 SOP", "TX2 ↩
            EOP", "TX2 Delay", "TX2 last 5 AVG Delay" };
26
27    static final String csvHeader;
28
29    static {
30      StringBuilder sb = new StringBuilder();
31      sb.append(CSV_HEADERS[0]);
32      for (int i = 1; i < CSV_HEADERS.length; ++i) {
33        sb.append("\t").append(CSV_HEADERS[i]);
34      }
35      csvHeader = sb.toString();
36    }
37
38    String packetData;
39
40    int packet;
41    int variableId;
42    long timestamp;
43
44    int rx_packetNum;
45    int rx_port;
46    long rx_sop;
47    long rx_eop;
48
49    Double avgTxDelay = null;
50    Double avgLast5AvgTxDelay = null;
51
52    TransmittedPacket[] tp = new TransmittedPacket[N];
53
54    @Override
55    public int hashCode() {
56      final int prime = 31;
57      int result = 1;
58      result = prime * result + (int) (timestamp ^ (timestamp >>> ↩
          32));
59      result = prime * result + variableId;
```

```java
60        return result;
61    }
62
63    @Override
64    public boolean equals(Object obj) {
65      if (this == obj)
66        return true;
67      if (obj == null)
68        return false;
69      if (getClass() != obj.getClass())
70        return false;
71      Packet other = (Packet) obj;
72      if (timestamp != other.timestamp)
73        return false;
74      if (variableId != other.variableId)
75        return false;
76      return true;
77    }
78
79   String toCsv() {
80      StringBuilder sb = new StringBuilder();
81      sb.append(packet).append("\t").append(avgTxDelay == null ? "←
             " : avgTxDelay).append("\t").append(
82           avgTxDelay == null ? "" : avgLast5AvgTxDelay == null ? "←
                 " : avgLast5AvgTxDelay).append("\t").append(
83           rx_port).append("\t").append(rx_packetNum).append("\t").←
                 append(rx_sop).append("\t").append(rx_eop)
84           .append("\t").append(timestamp).append("\t").append(←
                 variableId);
85      for (int i = 0; i < N; ++i) {
86        if (tp[i] != null) {
87          TransmittedPacket t = tp[i];
88          sb.append("\t").append(t.forwarded ? "1" : "").append("\←
                 t").append(t.tx_sop).append("\t").append(
89             t.tx_eop).append("\t").append(t.tx_sop - rx_eop).←
                 append("\t").append(
90             (t.avg5Delay == null) ? "" : t.avg5Delay);
91        } else {
92          sb.append("\t\t\t\t\t");
93        }
94      }
95      return sb.toString();
96    }
97 }
```

## A.4 GraphProducer.java

Listing A.4: GraphProducer.java

```java
package thesis;

import java.awt.Color;
import java.awt.Rectangle;
import java.awt.geom.Ellipse2D;
import java.io.File;
import java.io.FileNotFoundException;
import java.io.FileOutputStream;
import java.io.IOException;
import java.io.OutputStreamWriter;
import java.io.UnsupportedEncodingException;
import java.io.Writer;
import java.util.ArrayList;
import java.util.List;

import org.apache.batik.dom.GenericDOMImplementation;
import org.apache.batik.svggen.SVGGraphics2D;
import org.apache.batik.svggen.SVGGraphics2DIOException;
import org.apache.commons.math.stat.regression.SimpleRegression;
import org.jfree.chart.ChartFactory;
import org.jfree.chart.JFreeChart;
import org.jfree.chart.plot.DatasetRenderingOrder;
import org.jfree.chart.plot.PlotOrientation;
import org.jfree.chart.plot.XYPlot;
import org.jfree.chart.renderer.xy.XYDotRenderer;
import org.jfree.chart.renderer.xy.XYLineAndShapeRenderer;
import org.jfree.data.xy.XYDataItem;
import org.jfree.data.xy.XYSeries;
import org.jfree.data.xy.XYSeriesCollection;
import org.jfree.ui.RectangleInsets;
import org.w3c.dom.DOMImplementation;
import org.w3c.dom.Document;

/**
 * This class is capable of producing JFreeChart graphs and ←
      rendering these as SVG files. Three different types of
 * graphs are implemented:
 *
 * 1. One graph per simulation representing the packet ←
      processing duration (delay) as a function of time (packet
```

```
39   * number). Also in the same plot an estimated linear fit of ←↩
        that delay function. One observation per incoming packet in
40   * every simulation.
41   *
42   * 2. One graph per experiment representing the slope parts of ←↩
        the estimated linear fits from (1). One observation per
43   * simulation in experiment.
44   *
45   * 3. A single graph depicting ideal, observations of, and ←↩
        estimated fit of the throughput capacity function as seen
46   * across all experiments and simulations.
47   *
48   */
49   public class GraphProducer {
50
51     static JFreeChart createDelayWithLinearFit(List<Packet> ←↩
          packets, SimpleRegression sr, String title, int yAxisStart)←↩
          {
52       // create new chart
53       JFreeChart chart = ChartFactory.createScatterPlot(title, "←↩
            time (packet number)", "forwarding latency (media bus ←↩
            cycles)", null,
54         PlotOrientation.VERTICAL, false, true, true);
55
56       // general settings
57       chart.setPadding(new RectangleInsets(0, 0, 0, 5));
58       XYPlot xyPlot = chart.getXYPlot();
59       xyPlot.setBackgroundPaint(null);
60       xyPlot.setDatasetRenderingOrder(DatasetRenderingOrder.←↩
            FORWARD);
61
62       // X-axis
63       xyPlot.getDomainAxis().setAutoRange(false);
64       xyPlot.getDomainAxis().setRange(0, 1000);
65
66       // Y-axis
67       xyPlot.getRangeAxis().setAutoRange(false);
68       xyPlot.getRangeAxis().setRange(yAxisStart, yAxisStart + ←↩
            10000);
69
70       // delay dataset
71       XYSeries delaySeries = new XYSeries("observations");
72       for (Packet p : packets) {
73         delaySeries.add(p.packet, p.avgTxDelay);
74       }
```

```
75      xyPlot.setDataset(0, new XYSeriesCollection(delaySeries));
76
77      // delay rendering
78      XYDotRenderer dotRenderer = new XYDotRenderer();
79      dotRenderer.setDotWidth(1);
80      dotRenderer.setDotHeight(1);
81      dotRenderer.setSeriesPaint(0, Color.RED);
82      xyPlot.setRenderer(0, dotRenderer);
83
84      // regression dataset
85      XYSeries linearSeries = new XYSeries("linear fit");
86      linearSeries.add(0, sr.predict(0));
87      linearSeries.add(1000, sr.predict(1000));
88      xyPlot.setDataset(1, new XYSeriesCollection(linearSeries));
89
90      // regression rendering
91      XYLineAndShapeRenderer lineRenderer = new ↩
            XYLineAndShapeRenderer();
92      lineRenderer.setBaseShapesVisible(false);
93      lineRenderer.setSeriesPaint(0, Color.BLUE);
94      xyPlot.setRenderer(1, lineRenderer);
95
96      return chart;
97    }
98
99    static JFreeChart createSlopeOfDelay(XYSeries slopeSeries) {
100     // create new chart
101     JFreeChart chart = ChartFactory.createXYLineChart(null, "↩
            throughput (Mbps)", "slope of Delay(s,t) linear fit", ↩
            null,
102         PlotOrientation.VERTICAL, false, true, true);
103
104     /*
105      * general settings
106      */
107     XYPlot xyPlot = chart.getXYPlot();
108     xyPlot.setBackgroundPaint(null);
109
110     /*
111      * slope dataset. One sample per simulation.
112      */
113     xyPlot.setDataset(0, new XYSeriesCollection(slopeSeries));
114     // render in bluw with straight lines between markers.
115     XYLineAndShapeRenderer splineRenderer = new ↩
            XYLineAndShapeRenderer();
```

```
116        splineRenderer.setSeriesPaint(0, Color.BLUE);
117        xyPlot.setRenderer(0, splineRenderer);
118
119        /*
120         * y=0 series. slope dataset will intersect with this line ←
               approximately where the throughput capacity is.
121         */
122        XYSeries linearSeries = new XYSeries("y=0");
123        linearSeries.add(slopeSeries.getMinX(), 0);
124        linearSeries.add(slopeSeries.getMaxX(), 0);
125        xyPlot.setDataset(1, new XYSeriesCollection(linearSeries));
126        // render as black solid line with no markers
127        XYLineAndShapeRenderer lineRenderer = new ←
               XYLineAndShapeRenderer();
128        lineRenderer.setSeriesShapesVisible(0, false);
129        lineRenderer.setSeriesPaint(0, Color.BLACK);
130        xyPlot.setRenderer(1, lineRenderer);
131
132        return chart;
133    }
134
135    /*
136     * The following three functions are based on the equation:
137     *
138     * throughput capacity = (me * clock * bpp) / (k1 + (k2 * x))
139     */
140
141    static double throughputCapacity(int nMe, int clockMHz, int ←
           bpp, double k1, double k2, double x) {
142      return (nMe * clockMHz * bpp) / (k1 + k2 * x);
143    }
144
145    static double estimateK1(int nMe, int clockMHz, int bpp, ←
             double throughput, double k2, double x) {
146      return ((nMe * clockMHz * bpp) / throughput) − (k2 * x);
147    }
148
149    static double estimateK2(int nMe, int clockMHz, int bpp, ←
             double throughput, double k1, double x) {
150      return (((nMe * clockMHz * bpp) / throughput) − k1) / x;
151    }
152
153    @SuppressWarnings("unchecked")
154    static JFreeChart createThroughputCapacityFunction(boolean ←
           drawObserved, int nMe, int clockMHz, int bitsPerPacket,
```

```
155              int idealK1, int idealK2, XYSeries observations) {
156
157          // create new chart
158          JFreeChart chart = ChartFactory.createXYLineChart(null, "↩
                 iterations of algorithm's inner loop per packet",
159              "throughput capacity (Mbps)", null, PlotOrientation.↩
                     VERTICAL, true, true, true);
160
161          /*
162           * general settings
163           */
164          XYPlot xyPlot = chart.getXYPlot();
165          xyPlot.setBackgroundPaint(null);
166          xyPlot.setDatasetRenderingOrder(DatasetRenderingOrder.↩
                 FORWARD);
167          // X-axis
168          // compute xAxisRange to the closest even 10 greater than ↩
                 all x observations.
169          int xAxisRange = 10 * (((int) observations.getMaxX() + 10) /↩
                 10);
170          xyPlot.getDomainAxis().setAutoRange(false);
171          xyPlot.getDomainAxis().setRange(0, xAxisRange);
172          // Y-axis
173          // compute yAxisRange to the closest even 100 greater than y↩
                 when x=0.
174          int yAxisRange = 100 * (((int) throughputCapacity(nMe, ↩
                 clockMHz, bitsPerPacket, idealK1, idealK2, 0) + 100) / ↩
                 100);
175          xyPlot.getRangeAxis().setAutoRange(false);
176          xyPlot.getRangeAxis().setRange(0, yAxisRange);
177
178          /*
179           * ideal dataset
180           */
181          XYSeries idealSeries = new XYSeries("ideal");
182          for (double x = 0; x <= xAxisRange; x = x + 0.1) {
183            idealSeries.add(x, throughputCapacity(nMe, clockMHz, ↩
                   bitsPerPacket, idealK1, idealK2, x));
184          }
185          xyPlot.setDataset(0, new XYSeriesCollection(idealSeries));
186          // ideal rendering, as dots
187          XYDotRenderer dotRenderer = new XYDotRenderer();
188          dotRenderer.setDotWidth(1);
189          dotRenderer.setDotHeight(1);
190          dotRenderer.setSeriesPaint(0, Color.BLACK);
```

```
191        xyPlot.setRenderer(0, dotRenderer);
192
193        /*
194         * observed series
195         */
196        if (drawObserved) {
197          xyPlot.setDataset(1, new XYSeriesCollection(observations))↩
                ;
198          // render observations as large dots with no connecting ↩
                line
199          XYLineAndShapeRenderer simulatedRenderer = new ↩
                XYLineAndShapeRenderer();
200          simulatedRenderer.setSeriesLinesVisible(0, false);
201          simulatedRenderer.setSeriesPaint(0, Color.RED);
202          simulatedRenderer.setSeriesShape(0, new Ellipse2D.Double↩
                (0, 0, 5, 5));
203          xyPlot.setRenderer(1, simulatedRenderer);
204        }
205
206        /*
207         * estimate K1 and K2 based on observations
208         */
209        double k1 = 0;
210        double k2 = 0;
211        double k2Sum = 0;
212        int k2Count = 0;
213        for (XYDataItem item : (List<XYDataItem>) observations.↩
            getItems()) {
214          if (item.getXValue() == 0) {
215            // k1 is based on x=0 observation
216            k1 = estimateK1(nMe, clockMHz, bitsPerPacket, item.↩
                getYValue(), 1 /* K2 have no impact when x is 0 */, ↩
                0);
217            break;
218          }
219        }
220        for (XYDataItem item : (List<XYDataItem>) observations.↩
            getItems()) {
221          if (item.getXValue() > 0) {
222            k2Sum += estimateK2(nMe, clockMHz, bitsPerPacket, item.↩
                getYValue(), k1, item.getXValue());
223            k2Count++;
224          }
225        }
226        if (k2Count > 0) {
```

```
227          // k2 is based on the average estimate of all experiments ←
                with x>0
228          k2 = k2Sum / k2Count;
229        }
230
231      /*
232       * fit a line to observations
233       */
234      if (drawObserved) {
235        XYSeries estimatedFitSeries = new XYSeries("estimated fit"←
              );
236        for (double x = 0; x <= xAxisRange; x = x + 0.1) {
237          estimatedFitSeries.add(x, throughputCapacity(nMe, ←
                clockMHz, bitsPerPacket, k1, k2, x));
238        }
239        xyPlot.setDataset(2, new XYSeriesCollection(←
              estimatedFitSeries));
240        // fitted line rendered as dots
241        XYDotRenderer estimatedFitRenderer = new XYDotRenderer();
242        estimatedFitRenderer.setDotWidth(1);
243        estimatedFitRenderer.setDotHeight(1);
244        estimatedFitRenderer.setSeriesPaint(0, Color.BLUE);
245        xyPlot.setRenderer(2, estimatedFitRenderer);
246      }
247
248      return chart;
249    }
250
251    static final int DELAY_CHART_WIDTH = 500;
252    static final int DELAY_CHART_HEIGHT = 350;
253
254    // At most 6 graphs per page
255    static Rectangle[] BOUNDS = new Rectangle[6];
256
257    static {
258      BOUNDS[0] = new Rectangle(0, 0, DELAY_CHART_WIDTH, ←
            DELAY_CHART_HEIGHT);
259      BOUNDS[1] = new Rectangle(DELAY_CHART_WIDTH, 0, ←
            DELAY_CHART_WIDTH, DELAY_CHART_HEIGHT);
260      BOUNDS[2] = new Rectangle(0, DELAY_CHART_HEIGHT, ←
            DELAY_CHART_WIDTH, DELAY_CHART_HEIGHT);
261      BOUNDS[3] = new Rectangle(DELAY_CHART_WIDTH, ←
            DELAY_CHART_HEIGHT, DELAY_CHART_WIDTH, DELAY_CHART_HEIGHT←
            );
```

```
262      BOUNDS[4] = new Rectangle(0, 2 * DELAY_CHART_HEIGHT, ←
            DELAY_CHART_WIDTH, DELAY_CHART_HEIGHT);
263      BOUNDS[5] = new Rectangle(DELAY_CHART_WIDTH, 2 * ←
            DELAY_CHART_HEIGHT, DELAY_CHART_WIDTH, DELAY_CHART_HEIGHT←
            );
264    }
265
266    static void exportChartAsSVG(List<JFreeChart> charts, File ←
          svgFile) {
267      DOMImplementation domImpl = GenericDOMImplementation.←
            getDOMImplementation();
268      Document document = domImpl.createDocument(null, "svg", null←
            );
269      SVGGraphics2D svgGenerator = new SVGGraphics2D(document);
270      for (int i = 0; i < charts.size(); ++i) {
271        charts.get(i).draw(svgGenerator, BOUNDS[i]);
272      }
273      FileOutputStream fos = null;
274      try {
275        fos = new FileOutputStream(svgFile);
276        Writer out = new OutputStreamWriter(fos, "UTF-8");
277        svgGenerator.stream(out, true /* css */);
278        out.flush();
279      } catch (FileNotFoundException e) {
280        throw new RuntimeException(e);
281      } catch (UnsupportedEncodingException e) {
282        throw new RuntimeException(e);
283      } catch (SVGGraphics2DIOException e) {
284        throw new RuntimeException(e);
285      } catch (IOException e) {
286        throw new RuntimeException(e);
287      } finally {
288        if (fos != null) {
289          try {
290            fos.close();
291          } catch (IOException e) {
292            e.printStackTrace();
293          }
294        }
295      }
296    }
297
298    /**
299     *
```

```java
300      * This method produces the throughput capacity function based↩
             on known hardware configuration, estimated ideal K1
301      * and K2 (based on code analysis), and the observations ↩
            derived from experiment simulations (one observation per
302      * experiment).
303      *
304      * @param args
305      *           ignored.
306      */
307    public static void main(String args[]) {
308
309      /*
310       * Assumptions
311       */
312      int nMe = 2; // number of microengines running routing ↩
             algorithm in parallel
313      int clockMHz = 400; // microengine clock
314      int bitsPerPacket = 8 * 66; // each packet is 66 bytes
315      int idealK1 = 320; // ideal k1 estimate based on instruction↩
             count of route microcode
316      int idealK2 = 88; // ideal k2 estimate based on instruction ↩
            count of route microcode
317
318      /*
319       * Observations, one pair (X,Y) per experiment. X is the ↩
             runtime complexity of routing table (iterations of
320       * algorithm's inner loop per packet). Y is the throughput ↩
            capacity.
321       */
322      XYSeries observations = new XYSeries("observed");
323      observations.add(0, 692); // experiment 1
324      observations.add(2, 508); // experiment 2
325      observations.add(7.08, 297); // experiment 3
326      observations.add(14.16, 186); // experiment 4
327      observations.add(21.24, 168); // experiment 5
328
329      /*
330       * compute the throughput capacity chart
331       */
332      JFreeChart idealOnlyChart = createThroughputCapacityFunction↩
           (false, nMe, clockMHz, bitsPerPacket, idealK1,
333          idealK2, observations);
334      JFreeChart chart = createThroughputCapacityFunction(true, ↩
           nMe, clockMHz, bitsPerPacket, idealK1, idealK2,
335          observations);
```

```
336
337     /*
338      * export  chart  to  SVG
339      */
340     List<JFreeChart> charts = new ArrayList<JFreeChart>();
341     charts.add(chart);
342     exportChartAsSVG(charts, new File("../experiments/↩
            throughput_capacity.svg"));
343     charts.clear();
344     charts.add(idealOnlyChart);
345     exportChartAsSVG(charts, new File("../experiments/↩
            ideal_throughput_capacity.svg"));
346   }
347 }
```

## A.5   IXPSimulationLogAnalyser.java

Listing A.5: IXPSimulationLogAnalyser.java

```java
 1 package thesis;
 2
 3 import java.io.BufferedReader;
 4 import java.io.File;
 5 import java.io.FileFilter;
 6 import java.io.FileInputStream;
 7 import java.io.FileNotFoundException;
 8 import java.io.FileOutputStream;
 9 import java.io.IOException;
10 import java.io.InputStreamReader;
11 import java.io.PrintStream;
12 import java.nio.charset.Charset;
13 import java.util.ArrayList;
14 import java.util.Collections;
15 import java.util.Comparator;
16 import java.util.Iterator;
17 import java.util.LinkedHashMap;
18 import java.util.LinkedList;
19 import java.util.List;
20 import java.util.Map;
21 import java.util.regex.Matcher;
22 import java.util.regex.Pattern;
23
```

```java
24  import org.apache.commons.math.stat.descriptive.↩
        SummaryStatistics;
25  import org.apache.commons.math.stat.regression.SimpleRegression;
26  import org.jfree.chart.JFreeChart;
27  import org.jfree.data.xy.XYSeries;
28
29  import thesis.Packet.TransmittedPacket;
30  import thesis.RTSetupCodeGenerator.RoutingTable;
31  import thesis.RTSetupCodeGenerator.StatusVariableSubscriptions;
32
33  public class IXPSimulationLogAnalyser {
34
35    Map<String, SummaryStatistics> statistics;
36    Map<String, SimpleRegression> regressions;
37    List<Packet> packets;
38    RoutingTable routingTable;
39    double routingTableRuntimeComplexity;
40
41    static final Pattern packetLogLinePattern = Pattern.compile("[↩
          ]*([0-9]+)[ ]([0-9]+)[ ]([0-9]+)"
42        + "[ ]([0-9a-f]{8})[ ]([0-9a-f]{8})[ ]([0-9a-f]{8})[ ↩
              ]([0-9a-f]{8})"
43        + "[ ]([0-9a-f]{8})[ ]([0-9a-f]{8})[ ]([0-9a-f]{8})[ ↩
              ]([0-9a-f]{8})"
44        + "[ ]([0-9a-f]{8})[ ]([0-9a-f]{8})[ ]([0-9a-f]{8})[ ↩
              ]([0-9a-f]{8})"
45        + "[ ]([0-9a-f]{8})[ ]([0-9a-f]{8})[ ]([0-9a-f]{8})[ ↩
              ]([0-9a-f]{8})" + "[ ]([0-9a-f]{4})");
46
47    IXPSimulationLogAnalyser(File experimentFolder, RoutingTable ↩
          rt) {
48      routingTable = rt;
49      parsePacketLogs(experimentFolder);
50
51      // compute statistics
52      computeForwardingDecisionsAndRoutingTableRuntimeComplexity()↩
            ;
53      computeAvgPacketDelay();
54      computeLast5AvgDelay();
55      computeSummaryStatistics();
56      computeLinearRegressions();
57    }
58
59    void parsePacketLogs(File experimentFolder) {
60      // file handles
```

```java
61      File fileD0P0Rx = new File(experimentFolder, "↩
           Device0_Port0_Rx.log");
62      File fileD0P1Rx = new File(experimentFolder, "↩
           Device0_Port1_Rx.log");
63      File fileD0P2Rx = new File(experimentFolder, "↩
           Device0_Port2_Rx.log");
64      File fileD1P0Tx = new File(experimentFolder, "↩
           Device1_Port0_Tx.log");
65      File fileD1P1Tx = new File(experimentFolder, "↩
           Device1_Port1_Tx.log");
66      File fileD1P2Tx = new File(experimentFolder, "↩
           Device1_Port2_Tx.log");
67
68      // temporary packet lists.
69      List<Packet> port0Rx = new ArrayList<Packet>();
70      List<Packet> port1Rx = new ArrayList<Packet>();
71      List<Packet> port2Rx = new ArrayList<Packet>();
72
73      // parse files
74      if (fileD0P0Rx.isFile()) {
75        parseReceiveFile(fileD0P0Rx, 0, port0Rx);
76      }
77      if (fileD0P1Rx.isFile()) {
78        parseReceiveFile(fileD0P1Rx, 1, port1Rx);
79      }
80      if (fileD0P2Rx.isFile()) {
81        parseReceiveFile(fileD0P2Rx, 2, port2Rx);
82      }
83      if (fileD1P0Tx.isFile()) {
84        parseTransmitFile(fileD1P0Tx, 0, port2Rx, port1Rx);
85      }
86      if (fileD1P1Tx.isFile()) {
87        parseTransmitFile(fileD1P1Tx, 1, port0Rx, port2Rx);
88      }
89      if (fileD1P2Tx.isFile()) {
90        parseTransmitFile(fileD1P2Tx, 2, port1Rx, port0Rx);
91      }
92
93      createPacketList(port0Rx, port1Rx, port2Rx);
94    }
95
96    void createPacketList(List<Packet> port0Rx, List<Packet> ↩
         port1Rx, List<Packet> port2Rx) {
97      // reorganize packets into a single list sorted by packet ↩
           arrival and
```

```java
 98        // port
 99        packets = new ArrayList<Packet>();
100        for (Packet p : port0Rx) {
101          packets.add(p);
102        }
103        for (Packet p : port1Rx) {
104          packets.add(p);
105        }
106        for (Packet p : port2Rx) {
107          packets.add(p);
108        }
109        // sort by least significant sorting criteria first
110        Collections.sort(packets, new Comparator<Packet>() {
111          @Override
112          public int compare(Packet arg0, Packet arg1) {
113            return (int) (arg0.rx_port - arg1.rx_port);
114          }
115        });
116        Collections.sort(packets, new Comparator<Packet>() {
117          @Override
118          public int compare(Packet arg0, Packet arg1) {
119            return (int) (arg0.rx_eop - arg1.rx_eop);
120          }
121        });
122        // sort by most significant sorting criteria last
123        Collections.sort(packets, new Comparator<Packet>() {
124          @Override
125          public int compare(Packet arg0, Packet arg1) {
126            return (int) (arg0.rx_sop - arg1.rx_sop);
127          }
128        });
129
130        // assign packet numbers (correlating with arrival time)
131        int packet = 1;
132        for (Packet p : packets) {
133          p.packet = packet++;
134        }
135
136        // trim list of packets to contain at most 1000 packets
137        while (packets.size() > 1000) {
138          packets.remove(packets.size() - 1);
139        }
140      }
141
142      void parseReceiveFile(File f, int port, List<Packet> list) {
```

```java
143        FileInputStream fis = null;
144      try {
145        fis = new FileInputStream(f);
146        BufferedReader br = new BufferedReader(new ↩
             InputStreamReader(fis, Charset.forName("US-ASCII")));
147        while (br.ready()) {
148          String line = br.readLine();
149          Matcher m = packetLogLinePattern.matcher(line);
150          if (m.matches()) {
151            if (port == 2 && m.group(1).equals("212")) {
152              toString();
153            }
154            Packet p = new Packet();
155            p.packetData = m.group(4) + m.group(5) + m.group(6) + ↩
               m.group(7) + m.group(8) + m.group(9)
156                + m.group(10) + m.group(11) + m.group(12) + m.↩
                   group(13) + m.group(14) + m.group(15)
157                + m.group(16) + m.group(17) + m.group(18) + m.↩
                   group(19).substring(0, 4);
158            p.rx_port = port;
159            p.rx_packetNum = Integer.parseInt(m.group(1));
160            p.rx_sop = Integer.parseInt(m.group(2));
161            p.rx_eop = Integer.parseInt(m.group(3));
162            p.timestamp = Long.parseLong(m.group(12).substring(4) ↩
               + m.group(13) + m.group(14).substring(0, 4),
163                16);
164            p.variableId = Integer.parseInt(m.group(14).substring↩
               (4) + m.group(15).substring(0, 4), 16);
165            list.add(p);
166          }
167        }
168      } catch (FileNotFoundException e) {
169        e.printStackTrace();
170      } catch (IOException e) {
171        e.printStackTrace();
172      } finally {
173        if (fis != null) {
174          try {
175            fis.close();
176          } catch (IOException e) {
177            e.printStackTrace();
178          }
179        }
180      }
181    }
```

```
182
183    // Uninitialized memory byte patterns
184    static final String UIP_4 = "5a5a5a5a"; // 4-byte pattern
185    static final String UIP_2 = "5a5a"; // 2-byte pattern
186    static final String UIP_PACKET = UIP_4 + UIP_4 + UIP_4 + UIP_4↩
           + UIP_4 + UIP_4 + UIP_4 + UIP_4 + UIP_4 + UIP_4
187        + UIP_4 + UIP_4 + UIP_4 + UIP_4 + UIP_4 + UIP_2;
188
189    void parseTransmitFile(File f, int port, List<Packet> source1,↩
           List<Packet> source2) {
190      FileInputStream fis = null;
191      try {
192        Iterator<Packet> s1It = source1.iterator();
193        Iterator<Packet> s2It = source2.iterator();
194        fis = new FileInputStream(f);
195        BufferedReader br = new BufferedReader(new ↩
             InputStreamReader(fis, Charset.forName("US-ASCII")));
196        int s1 = 0;
197        int s2 = 0;
198        while (br.ready()) {
199          String line = br.readLine();
200          Matcher m = packetLogLinePattern.matcher(line);
201          if (m.matches()) {
202            if (port == 1 && m.group(1).equals("424")) {
203              toString();
204            }
205            Packet p;
206            String packetData = m.group(4) + m.group(5) + m.group↩
               (6) + m.group(7) + m.group(8) + m.group(9)
207                + m.group(10) + m.group(11) + m.group(12) + m.↩
                   group(13) + m.group(14) + m.group(15)
208                + m.group(16) + m.group(17) + m.group(18) + m.↩
                   group(19).substring(0, 4);
209          if (UIP_PACKET.equals(packetData)) {
210              /*
211               * The current implementation of multicasting ↩
                     translates copied packet handles to ↩
                     uninitialized
212               * packet buffers in DRAM. These areas can be ↩
                     recognized by repeating bytes of value 0x5a.
213               * Because we only use 3 ports we don't need the ↩
                     actual packet data to know where it came from,
214               * because any incoming packet can only be forwarded↩
                     to at most two other ports, and one of
```

121

```
215                     * those  ports  will  translate  the  packet  handle  to  ←
                           the  correct  DRAM  area  and  get  actual  packet
216                     * data  rather  than  0x5a  patterns.
217                     */
218                 p = s2It.next();
219                 s1++;
220                 int rxPort = (port + 1) % 3;
221                 if (rxPort != p.rx_port) {
222                     throw new RuntimeException("rxPort != p.rx_port, ←
                           rxPort=" + rxPort + ", p.rx_port="
223                         + p.rx_port + ", tx_packetNum=" + m.group(1));
224                 }
225                 if (s1 != p.rx_packetNum) {
226                     throw new RuntimeException("s1 != p.rx_packetNum, ←
                           s1=" + s1 + ", p.rx_packetNum="
227                         + p.rx_packetNum);
228                 }
229             } else {
230                 p = s1It.next();
231                 s2++;
232                 int rxPort = (port + 2) % 3;
233                 if (rxPort != p.rx_port) {
234                     throw new RuntimeException("rxPort != p.rx_port, ←
                           rxPort=" + rxPort + ", p.rx_port="
235                         + p.rx_port + ", tx_packetNum=" + m.group(1));
236                 }
237                 if (s2 != p.rx_packetNum) {
238                     throw new RuntimeException("s2 != p.rx_packetNum, ←
                           s2=" + s2 + ", p.rx_packetNum="
239                         + p.rx_packetNum);
240                 }
241                 if (!packetData.equals(p.packetData)) {
242                     throw new RuntimeException("tx_port=" + port + ", ←
                           rx_port=" + p.rx_port + ", tx_packetNum="
243                         + m.group(1) + ", rx_packetNum=" + p.←
                           rx_packetNum
244                         + ", !packetData.equals(p.packetData), ←
                           packetData=\"" + packetData
245                         + "\", p.packetData=\"" + p.packetData + "\"")←
                           ;
246                 }
247             }
248             if (p.tp[port] != null) {
249                 throw new RuntimeException("Packet " + p.←
                       rx_packetNum + " from port " + p.rx_port
```

```
250                    + " already has a TransmittedPacket for port " +↩
                         port
251                    + ". Current transmit packetNum is " + Integer.↩
                         parseInt(m.group(1)));
252              }
253            p.tp[port] = new TransmittedPacket();
254            p.tp[port].tx_sop = Integer.parseInt(m.group(2));
255            p.tp[port].tx_eop = Integer.parseInt(m.group(3));
256          } else {
257            toString();
258          }
259        }
260    } catch (FileNotFoundException e) {
261        e.printStackTrace();
262    } catch (IOException e) {
263        e.printStackTrace();
264    } finally {
265        if (fis != null) {
266          try {
267            fis.close();
268          } catch (IOException e) {
269            e.printStackTrace();
270          }
271        }
272      }
273    }
274
275    void ↩
          computeForwardingDecisionsAndRoutingTableRuntimeComplexity↩
          () {
276      int pubint = 10;
277      int innerLoopIterations = 0;
278      for (Packet p : packets) {
279        StatusVariableSubscriptions svs = routingTable.find(p.↩
            variableId);
280        /*
281         * this is basically an implementation of the GridStat ↩
              routing filtering algorithm with additional
282         * statistics to estimate the routing table runtime ↩
              complexity for this simulation run
283         */
284        for (int port = 0; port < 3; ++port) {
285          if (p.tp[port] != null) {
286            for (int i = 0; i < svs.n[port]; ++i) {
287              innerLoopIterations++;
```

```
288              if ((p.timestamp + (pubint / 2)) % svs.divisors[port↩
                    ][i] < pubint) {
289                p.tp[port].forwarded = true; // packet should be
290                // forwarded to port
291                break;
292              }
293            }
294          }
295        }
296      }
297      routingTableRuntimeComplexity = innerLoopIterations / (↩
            double) packets.size();
298    }
299
300    void computeAvgPacketDelay() {
301      for (Packet p : packets) {
302        int nTransmitted = 0;
303        double sumTxDelay = 0;
304        for (int i = 0; i < 3; ++i) {
305          if (p.tp[i] != null && p.tp[i].forwarded) {
306            sumTxDelay += p.tp[i].tx_sop − p.rx_eop;
307            ++nTransmitted;
308          }
309        }
310        if (nTransmitted > 0) {
311          p.avgTxDelay = sumTxDelay / nTransmitted;
312        }
313      }
314    }
315
316    void computeLast5AvgDelay() {
317      LinkedList<Packet> last5Packets = new LinkedList<Packet>();
318      for (Packet p : packets) {
319        if (p.avgTxDelay != null) {
320          last5Packets.add(p);
321          if (last5Packets.size() > 5) {
322            last5Packets.removeFirst();
323          }
324        }
325        int nTransmitted = 0;
326        double sumTxDelay = 0;
327        if (p.avgTxDelay != null && last5Packets.size() >= 5) {
328          for (Packet a : last5Packets) {
329            sumTxDelay += a.avgTxDelay;
330            ++nTransmitted;
```

```
331              }
332            p.avgLast5AvgTxDelay = sumTxDelay / nTransmitted;
333          }
334       }
335     }
336
337     void computeSummaryStatistics() {
338       statistics = new LinkedHashMap<String, SummaryStatistics>();
339       computeAvgDelaySS();
340       computeAvgLast5DelaySS();
341       computeIncomingInterPacketDelaySS();
342       computeAddedInterPacketDelayJitterSS();
343     }
344
345     void computeAvgDelaySS() {
346       SummaryStatistics avgDelaySS = new SummaryStatistics();
347       for (Packet p : packets) {
348         if (p.avgTxDelay != null) {
349           avgDelaySS.addValue(p.avgTxDelay);
350         }
351       }
352       statistics.put("AverageDelay", avgDelaySS);
353     }
354
355     void computeAvgLast5DelaySS() {
356       SummaryStatistics avglast5DelaySS = new SummaryStatistics();
357       for (Packet p : packets) {
358         if (p.avgLast5AvgTxDelay != null) {
359           avglast5DelaySS.addValue(p.avgLast5AvgTxDelay);
360         }
361       }
362       statistics.put("AverageLast5AverageDelay", avglast5DelaySS);
363     }
364
365     void computeIncomingInterPacketDelaySS() {
366       SummaryStatistics rxInterPacketDelaySS = new ←↩
             SummaryStatistics();
367       Packet prevIncomingPacket[] = new Packet[3];
368       for (Packet p : packets) {
369         if (prevIncomingPacket[p.rx_port] != null) {
370           /*
371            * only look at the packets after the first 100 because ←↩
                 the first few packets possibly arrive faster
372            * than the configured throughput dictate to compensate ←↩
                 for the time it takes to initialize the hardware
```

```
373              * and memory layout.
374              */
375            if (p.packet > 100) {
376              rxInterPacketDelaySS.addValue(p.rx_eop − ←
                     prevIncomingPacket[p.rx_port].rx_eop);
377            }
378          }
379          prevIncomingPacket[p.rx_port] = p;
380        }
381      statistics.put("IncomingInterPacketDelay", ←
             rxInterPacketDelaySS);
382    }
383
384    void computeAddedInterPacketDelayJitterSS() {
385      SummaryStatistics txInterPacketDelaySS = new ←
             SummaryStatistics();
386      Packet prevIncomingPacket[] = new Packet[3];
387      for (Packet p : packets) {
388        if (prevIncomingPacket[p.rx_port] != null) {
389          /*
390              * only look at the packets after the first 100 because ←
                     the first few packets possibly arrive faster
391              * than the configured throughput dictate to compensate ←
                     for the time it takes to initialize the hardware
392              * and memory layout.
393              */
394            if (p.packet > 100) {
395              int d1 = (p.rx_port + 1) % 3; // first destination (←
                     port)
396              int d2 = (p.rx_port + 2) % 3; // second destination (←
                     port)
397              if (p.tp[d1] != null && prevIncomingPacket[p.rx_port].←
                     tp[d1] != null) {
398                txInterPacketDelaySS.addValue(p.tp[d1].tx_eop − ←
                       prevIncomingPacket[p.rx_port].tp[d1].tx_eop);
399              }
400              if (p.tp[d2] != null && prevIncomingPacket[p.rx_port].←
                     tp[d2] != null) {
401                txInterPacketDelaySS.addValue(p.tp[d2].tx_eop − ←
                       prevIncomingPacket[p.rx_port].tp[d2].tx_eop);
402              }
403            }
404          }
405          prevIncomingPacket[p.rx_port] = p;
406        }
```

```
407      statistics.put("OutgoingInterPacketDelay", ←
            txInterPacketDelaySS);
408    }
409
410    void computeLinearRegressions() {
411      regressions = new LinkedHashMap<String, SimpleRegression>();
412      computeLinearRegressionOfPacketDelayAsFunctionOfTime();
413      computeLinearRegressionOfLastHalfPacketDelayAsFunctionOfTime←
            ();
414    }
415
416    void computeLinearRegressionOfPacketDelayAsFunctionOfTime() {
417      SimpleRegression sr = new SimpleRegression();
418      for (Packet p : packets) {
419        for (int port = 0; port < 3; ++port) {
420          if (p.tp[port] != null) {
421            double y = p.tp[port].tx_sop - p.rx_eop;
422            double x = p.packet;
423            sr.addData(x, y);
424          }
425        }
426      }
427      regressions.put("PacketDelayAsFunctionOfTime", sr);
428    }
429
430    void ←
          computeLinearRegressionOfLastHalfPacketDelayAsFunctionOfTime←
          () {
431      SimpleRegression sr = new SimpleRegression();
432      for (Packet p : packets) {
433        if (p.packet >= packets.size() / 4) {
434          for (int port = 0; port < 3; ++port) {
435            if (p.tp[port] != null) {
436              double y = p.tp[port].tx_sop - p.rx_eop;
437              double x = p.packet;
438              sr.addData(x, y);
439            }
440          }
441        }
442      }
443      regressions.put("TrimmedPacketDelayAsFunctionOfTime", sr);
444    }
445
446    void printCsv(File outputFile) {
447      FileOutputStream fos = null;
```

```
448        try {
449          fos = new FileOutputStream(outputFile);
450          PrintStream out = new PrintStream(fos);
451          out.println(Packet.csvHeader);
452          for (Packet p : packets) {
453            out.println(p.toCsv());
454          }
455          out.flush();
456        } catch (FileNotFoundException e) {
457          e.printStackTrace();
458        } finally {
459          if (fos != null) {
460            try {
461              fos.close();
462            } catch (IOException e) {
463              e.printStackTrace();
464            }
465          }
466        }
467      }
468
469    static final Pattern experimentFolderPattern = Pattern.compile↩
           ("experiment([1-9][0-9]*)");
470    static final Pattern simulationRunFolderPattern = Pattern.↩
           compile("3x([1-9][0-9]*)");
471
472    static class SingleSimulationRun {
473      int simulationId;
474      int incomingThroughput;
475      File f;
476      IXPSimulationLogAnalyser logAnalyzer;
477    }
478
479    static void printCrossSummaryTexJitterTable(File f, List<↩
           SingleSimulationRun> simulations) {
480      FileOutputStream fos = null;
481      try {
482        fos = new FileOutputStream(f);
483        PrintStream out = new PrintStream(fos);
484        out.println("\\begin{table}");
485        out.println("\\caption{Incoming and outgoing inter-packet-↩
               delay per simulation of experiment \\theex}");
486        out.println("\\centering");
487        out.println("\\begin{tabular}{|c|c|c c|c c|}");
488        out.println("\\hline");
```

128

```
489        out.println("$s$ & Input & Input Average & StdDev & Output↩
               Average & StdDev \\\\");
490        out.println(" & (Mbps) & Inter-packet delay & & Inter-↩
               packet delay & \\\\");
491        out.println(" & & (media-bus cycles) & & (media-bus cycles↩
               ) & \\\\");
492        out.println("\\hline");
493        for (SingleSimulationRun ssr : simulations) {
494          Matcher m = simulationRunFolderPattern.matcher(ssr.f.↩
               getName());
495          m.matches();
496          int fractionOfThroughput = Integer.parseInt(m.group(1));
497          SummaryStatistics incomingInterPacketDelaySS = ssr.↩
               logAnalyzer.statistics
498            .get("IncomingInterPacketDelay");
499          SummaryStatistics outgoingInterPacketDelaySS = ssr.↩
               logAnalyzer.statistics
500            .get("OutgoingInterPacketDelay");
501          out.format("\\thesimc \\addtocounter{simc}{1} & %d & %g ↩
               & %g & %g & %g\\\\", 3 * fractionOfThroughput,
502              incomingInterPacketDelaySS.getMean(), ↩
                 incomingInterPacketDelaySS.getStandardDeviation()↩
                 ,
503              outgoingInterPacketDelaySS.getMean(), ↩
                 outgoingInterPacketDelaySS.getStandardDeviation()↩
                 );
504          out.println();
505        }
506        out.println("\\hline");
507        out.println("\\end{tabular}");
508        out.println("\\label{table:se\\theex}");
509        out.println("\\end{table}");
510        out.println("\\endinput");
511        out.flush();
512      } catch (FileNotFoundException e) {
513        e.printStackTrace();
514      } finally {
515        if (fos != null) {
516          try {
517            fos.close();
518          } catch (IOException e) {
519            e.printStackTrace();
520          }
521        }
522      }
```

```
523    }
524
525    /**
526     * Usage: java thesis.IXPSimulationLogAnalyser <experiment ↩
              folder> [<status variable id> <port> <subscriber
527     * rate >]...
528     *
529     */
530    public static void main(String[] args) {
531
532      // parse program arguments
533      File experimentFolder = new File(args[0]);
534      int nextSimulationId = Integer.parseInt(args[1]);
535      List<Integer> list = new ArrayList<Integer>();
536      for (int i = 2; i < args.length; i += 3) {
537        list.add(new Integer(args[i]));
538        list.add(new Integer(args[i + 1]));
539        list.add(new Integer(args[i + 2]));
540      }
541
542      // setup internal routing table which is the same across ↩
              simulations
543      // within an experiment
544      RoutingTable rt = RTSetupCodeGenerator.createRoutingTable(↩
          list);
545
546      // find experiment folder sub directories which corresponds ↩
              with
547      // simulation runs
548      File[] subdirs = experimentFolder.listFiles(new FileFilter()↩
              {
549        @Override
550        public boolean accept(File f) {
551          Matcher m = simulationRunFolderPattern.matcher(f.getName↩
                ());
552          return m.matches() && f.isDirectory();
553        }
554      });
555
556      // parse and analyze input files
557      boolean problem = false;
558      List<SingleSimulationRun> simulations = new LinkedList<↩
              SingleSimulationRun>();
559      for (File f : subdirs) {
560        try {
```

130

```
561             SingleSimulationRun ssr = new SingleSimulationRun();
562             ssr.f = f;
563             ssr.logAnalyzer = new IXPSimulationLogAnalyser(f, rt);
564             simulations.add(ssr);
565         } catch (Exception e) {
566           problem = true;
567           try {
568             System.err.println("Exception while analysing folder "↩
                    + f.getCanonicalPath());
569           } catch (IOException e1) {
570             e1.printStackTrace();
571           }
572           e.printStackTrace();
573           break;
574         }
575       }
576       // sort simulation list ascending by throughput
577       Collections.sort(simulations, new Comparator<↩
             SingleSimulationRun>() {
578         @Override
579         public int compare(SingleSimulationRun o1, ↩
               SingleSimulationRun o2) {
580           Matcher m1 = simulationRunFolderPattern.matcher(o1.f.↩
                 getName());
581           m1.matches();
582           int t1 = Integer.parseInt(m1.group(1));
583           Matcher m2 = simulationRunFolderPattern.matcher(o2.f.↩
                 getName());
584           m2.matches();
585           int t2 = Integer.parseInt(m2.group(1));
586           return t1 - t2;
587         }
588       });
589
590       // assign simulation id and throughput
591       for (SingleSimulationRun ssr : simulations) {
592         ssr.simulationId = nextSimulationId++;
593         Matcher m = simulationRunFolderPattern.matcher(ssr.f.↩
               getName());
594         m.matches();
595         ssr.incomingThroughput = 3 * Integer.parseInt(m.group(1));
596       }
597
598       // write output files
599       if (!problem) {
```

```java
600         Matcher m = experimentFolderPattern.matcher(←
                experimentFolder.getName());
601         int experimentNo = 0;
602         if (m.matches()) {
603           experimentNo = Integer.parseInt(m.group(1));
604         }
605         File outFolder = experimentFolder.getParentFile();
606         String ePrefix = experimentFolder.getName() + "_";
607
608         // print new packet logs for all simulations in experiment
609         for (SingleSimulationRun ssr : simulations) {
610           @SuppressWarnings("unused")
611           String prefix = ePrefix + ssr.f.getName() + "_";
612           // ssr.logAnalyzer.printCsv(new File(outFolder, prefix +←
                  "packets.csv"));
613         }
614
615         // Generate delay charts
616         List<JFreeChart> charts = new ArrayList<JFreeChart>();
617         for (SingleSimulationRun ssr : simulations) {
618           JFreeChart chart = GraphProducer.←
                  createDelayWithLinearFit(ssr.logAnalyzer.packets,
619             ssr.logAnalyzer.regressions.get("←
                    PacketDelayAsFunctionOfTime"), "Simulation "
620               + ssr.simulationId + ", " + ssr.←
                    incomingThroughput + " Mbps",
621             (experimentNo == 5) ? 5000 : 0);
622           charts.add(chart);
623         }
624         // render SVG graphs in groups of up to 6 per page
625         for (int i = 0; i < charts.size(); i += 6) {
626           GraphProducer.exportChartAsSVG(charts.subList(i, Math.←
                  min(i + 6, charts.size())), new File(outFolder,
627             ePrefix + "delay_" + (i / 6) + ".svg"));
628         }
629
630         // Generate slope chart
631         XYSeries slopeSeries = new XYSeries("slope");
632         for (SingleSimulationRun ssr : simulations) {
633           slopeSeries.add(ssr.incomingThroughput, ssr.logAnalyzer.←
                  regressions.get("PacketDelayAsFunctionOfTime")
634             .getSlope());
635         }
636         JFreeChart slopeChart = GraphProducer.createSlopeOfDelay(←
                slopeSeries);
```

```
637        charts.clear();
638        charts.add(slopeChart);
639        GraphProducer.exportChartAsSVG(charts, new File(outFolder,↩
               ePrefix + "slope.svg"));
640
641        // print experiment cross simulation summary
642        printCrossSummaryTexJitterTable(new File(outFolder, ↩
               experimentFolder.getName() + "_jitter.tex"),
643          simulations);
644      }
645    }
646 }
```

## A.6  utilities.uc

Listing A.6: utilities.uc

```
1
2
3  //
4  //  C = A * B
5  //
6  // 16 instructions
7  #macro imul64x32[c2, c1, c0, a1, a0, b0]
8  .begin
9    .reg hi32
10
11   // C_high64 = A_high32 * B
12   mul_step[a1, b0], 32x32_start
13   mul_step[a1, b0], 32x32_step1
14   mul_step[a1, b0], 32x32_step2
15   mul_step[a1, b0], 32x32_step3
16   mul_step[a1, b0], 32x32_step4
17   mul_step[c1, --], 32x32_last
18   mul_step[c2, --], 32x32_last2
19
20   // C_low64 = A_low32 * B
21   mul_step[a0, b0], 32x32_start
22   mul_step[a0, b0], 32x32_step1
23   mul_step[a0, b0], 32x32_step2
24   mul_step[a0, b0], 32x32_step3
25   mul_step[a0, b0], 32x32_step4
26   mul_step[c0, --], 32x32_last
```

```
27    mul_step[hi32, --], 32x32_last2
28
29    // C_mid32 = (C_low64 >> 32) + (C_high64 & 0xffffffff)
30    alu[c1, c1, +, hi32]
31
32    // C_high32 = (C_high64 >> 32) + carry_out from computing ↩
          C_mid32
33    alu[c2, c2, +carry, 0]
34
35  .end
36  #endm
37
38  //
39  //
40  //   A is a 64 bit integer described by two 32 bit registers,
41  //   a1 is the high 32 bits of A, and a0 is the low 32 bits of A.
42  //
43  //   Q = (A * Reciprocal) >> (32+nbits)
44  //   rem = A - Q * divisor
45  //   while ( Rem >= divisor )
46  //     Rem -= divisor
47  //
48  //
49  #macro remainder64x32[rem, a1, a0, divisor, reciprocal, nbits]
50  .begin
51    .reg q1 q0 c0 kd
52    .reg c2 c1 remh
53
54    //  Q = (A * Reciprocal) >> (32+nbits)
55
56    // Q = A * Reciprocal
57    imul64x32[q1, q0, c0, a1, a0, reciprocal]
58
59    // q0 = q0 >> nbits
60    alu[--, nbits, OR, 0]
61    alu[q0, --, B, q0, >>indirect]
62
63    // q0 = q0 | (q1 << (32-nbits))
64    alu[c0, 32, -, nbits]
65    alu[--, c0, OR, 0]
66    alu[q0, q0, OR, q1, <<indirect]
67
68    // q1 = q1 >> nbits
69    alu[--, nbits, OR, 0]
70    alu[q1, --, B, q1, >>indirect]
```

134

```
71
72      // 23 instructions
73
74
75      // Q is now a 64 bit integer, with the high
76      // 32 bits at q1 and low 32 bits at q0
77
78      // rem = A - (Q * divisor)
79
80      // C = Q * divisor
81
82      imul64x32[c2, c1, c0, q1, q0, divisor]
83
84      // FIXME: fix for the case where the high 32 bits in A and Q *←
                divisor
85      //          are not the same.
86      // we here assume that the high 32 bits in both A and (Q * ←
            divisor)
87      // are the same, hence we only need to look at the low 32 bits←
            .
88      //mul_step[q0, divisor], 32x32_start
89      //mul_step[q0, divisor], 32x32_step1
90      //mul_step[q0, divisor], 32x32_step2
91      //mul_step[q0, divisor], 32x32_step3
92      //mul_step[q0, divisor], 32x32_step4
93      //mul_step[c0, --], 32x32_last
94      //mul_step[c1, --], 32x32_last2 ; discard
95
96      alu[rem, a0, -, c0]
97      alu[remh, a1, -carry, c1] ; remh should be 0
98
99        .if ( c1 != a1 )
100           nop; should never happen but included here to enable ←
                debug breakpoint
101       .endif
102
103     // 30 instructions
104
105     // at this point the remainder might be up to 3*divisor off ←
            from the
106     // actual remainder, so iterate until remainder is less than ←
            divisor.
107
108     // NOTE: for some reason the remainder ends up being k*divisor←
            off from
```

```
109    // the actual remainder, so just subtract k*divisor from rem.
110
111    alu_shf[kd, --, B, divisor, <<5]
112    .while ( rem >= kd )
113      alu[rem, rem, -, kd]
114    .endw
115
116    alu_shf[kd, --, B, divisor, <<4]
117    .if ( rem >= kd )
118      alu[rem, rem, -, kd]
119    .endif
120
121    alu_shf[kd, --, B, divisor, <<3]
122    .if ( rem >= kd )
123      alu[rem, rem, -, kd]
124    .endif
125
126    alu_shf[kd, --, B, divisor, <<2]
127    .if ( rem >= kd )
128      alu[rem, rem, -, kd]
129    .endif
130
131    alu_shf[kd, --, B, divisor, <<1]
132    .if ( rem >= kd )
133      alu[rem, rem, -, kd]
134    .endif
135
136    .if ( rem >= divisor )
137      alu[rem, rem, -, divisor]
138    .endif
139
140  .end
141 #endm
```

## A.7    routing_table.uc

Listing A.7: routing_table.uc

```
1 #include <utilities.uc>
2
3 .sig volatile sig_rtable
4
5 #define CAM_VALID 1
```

```
6  #define CAM_INVALID 0
7
8  #macro hash_to_lwoffset[ lwoffset, in_hash ]
9     .begin
10 #if ( GRIDSTAT_HTABLE_SIZE == 128 )
11    // lwoffset = in_hash & 0x7F
12    alu[lwoffset, in_hash, AND, 0x7F]
13 #elif ( GRIDSTAT_HTABLE_SIZE == 256 )
14    // lwoffset = in_hash & 0xFF
15    ld_field_w_clr[lwoffset, 0001, in_hash]
16 #elif ( GRIDSTAT_HTABLE_SIZE == 512 )
17    // lwoffset = in_hash & 0x1FF
18    ld_field_w_clr[lwoffset, 0011, in_hash]
19    alu_shf[lwoffset, lwoffset, ~AND, 0xFE, <<8]
20 #elif ( GRIDSTAT_HTABLE_SIZE == 1024 )
21    // lwoffset = in_hash & 0x3FF
22    ld_field_w_clr[lwoffset, 0011, in_hash]
23    alu_shf[lwoffset, lwoffset, ~AND, 0xFC, <<8]
24 #else
25 #error "GRIDSTAT_HTABLE_SIZE must be one of (128, 256, 512, ↩
       1024)"
26 #endif
27    .end
28 #endm
29
30
31 #macro getrtaddr[ sram_offset, variableid ]
32 .begin
33    .reg varid1 varid2 varid3
34    move( varid1, 1 )
35    move( varid2, 2 )
36    move( varid3, 3 )
37    // get base offset
38    .if ( variableid == varid1 )
39      move( sram_offset, (GRIDSTAT_RTABLE_BASE+GRIDSTAT_RTENT_SIZE↩
          ) )
40    .elif ( variableid == varid2 )
41      move( sram_offset, (GRIDSTAT_RTABLE_BASE+2*↩
          GRIDSTAT_RTENT_SIZE) )
42    .elif ( variableid == varid3 )
43      move( sram_offset, (GRIDSTAT_RTABLE_BASE+3*↩
          GRIDSTAT_RTENT_SIZE) )
44    .else
45      // default entry if no variableid match
```

```
46      move( sram_offset, (GRIDSTAT_RTABLE_BASE+0) ) // use varid 0↩
            which corresponds to no entry
47    .endif
48  .end
49  #endm
50
51  // output sram_offset will be 0 if no entry matching variableid ↩
       can be found
52  #macro init_rtable_offset[ variableid ]
53  .begin
54    .reg sram_offset lsmask tmp cam_res lru lm_a i
55
56    // lockup variableid in cam
57    cam_lookup[cam_res, variableid]
58
59    // check for hit or miss, if hit, then we don't need to load ↩
         rt data from sram
60    alu_shf[−−, 0x1, AND, cam_res, >>7]
61    bne[rt_entry_in_lm#]
62
63    // cam lookup miss, load routing table entry for variableid ↩
         into local memory
64    alu_shf[lru, 0xf, AND, cam_res, >>3] ; find least recently ↩
         used entry
65    cam_write[lru, variableid, CAM_INVALID] // update cam
66
67      /* The following commented section uses the hashing ↩
            instruction to lookup
68      * routing table entry in a hashtable fashion, it is ↩
            ommitted for simplicity.
69      * The following commented code section hasn't been properly↩
            tested and may
70      * contain bugs.
71      */
72
73      /*
74      xbuf_alloc($hash_data, 2, read_write)
75
76    move[$hash_data0, variableid]
77    move[$hash_data1, 0]
78
79    hash_48[$hash_data0, 1], sig_done[sig_rtable]
80
81    // wait for hash io, there are 2 signals to wait for
82    alu[lsmask, −−, B, 0]
```

138

```
83    alu [tmp, --, B, &sig_rtable]
84    alu[--, tmp, OR, 0]
85    alu_shf [lsmask, lsmask, OR, 1, <<indirect]
86    alu [tmp, tmp, +, 1]
87    alu[--, tmp, OR, 0]
88    alu_shf [lsmask, lsmask, OR, 1, <<indirect]
89    local_csr_wr [active_ctx_wakeup_events, lsmask]
90    ctx_arb[--]
91    .io_completed sig_rtable
92
93      hash_to_lwoffset [sram_offset, $hash_data0]
94    // TODO: include multiple mode functionality
95    add_shf_left (sram_offset, GRIDSTAT_HTABLE_BASE, sram_offset, ←
          2) ; <<2 to convert sram_offset to byte offset
96    sram [read, $hash_data0, sram_offset, 0, 1], ctx_swap[←
          sig_rtable] ; get address of first entry in list
97
98    // The following loop iterates through the single linked list
99    // until variableid matches one of the entries
100   htable_list_iterate#:
101     alu [sram_offset, --, B, $hash_data0]
102     beq[htable_list_iterate_end#] // if there is no rt entry, stop←
              iterating
103     sram [read, $hash_data0, sram_offset, 0, 2], ctx_swap[←
          sig_rtable]
104     alu[--, variableid, -, $hash_data1]
105     bne[htable_list_iterate#]
106   htable_list_iterate_end#:
107
108     xbuf_free ($hash_data)
109        */
110
111     // for now the routing table is hardcoded, so find the actual ←
          sram address
112     getrtaddr [sram_offset, variableid]
113
114     // advance pointer to start of subintervals
115     alu [sram_offset, sram_offset, +, 16]
116
117     xbuf_alloc ($filter_data, 6, read)
118
119     // size of one routing table entry is 36 long words
120     ; lm_a = lru * 36
121     alu_shf [lm_a, --, B, lru, <<5]    ; lm_a = lru * 32
122     add_shf_left (lm_a, lm_a, lru, 2) ; lm_a = lm_a + lru * 4
```

```
123    alu_shf [lm_a , −−, B , lm_a , <<2]   ; move to offset used ←
           ACTIVE_LM_ADDR_0 write
124
125    local_csr_wr [ACTIVE_LM_ADDR_0 , lm_a ]
126
127    alu [ i , −−, B , 0]
128    .while ( i < 6 )
129
130      sram [read , $filter_data0 , sram_offset , 0, 6] , sig_done [←
             sig_rtable ]
131
132      alu [ i , i , +, 1]
133      alu [sram_offset , sram_offset , +, 24]
134
135      // wait for previous sram read from routing table
136      ctx_arb [sig_rtable ]
137
138      alu [∗l$index0++, −−, B , $filter_data0 ]
139      alu [∗l$index0++, −−, B , $filter_data1 ]
140      alu [∗l$index0++, −−, B , $filter_data2 ]
141      alu [∗l$index0++, −−, B , $filter_data3 ]
142      alu [∗l$index0++, −−, B , $filter_data4 ]
143      alu [∗l$index0++, −−, B , $filter_data5 ]
144
145    .endw
146
147    cam_write_state [lru , CAM_VALID ] // set cam entry as valid
148
149    xbuf_free ($filter_data )
150
151 rt_entry_in_lm#:
152
153 .end
154 #endm
155
156 #macro get_stream_info [pubint ]
157 .begin
158   // TODO: download value of pubint from routing table
159   alu [pubint , −−, B , 10]
160 .end
161 #endm
162
163 //
164 //   Assumes that there is an array of FILTER_ENTRY starting
```

```
165  //   at sram_addr. The array is terminated when (element[i].↩
         divisor == 0)
166  //
167  //
168  //   typedef struct s_filter_entry {
169  //      unsigned long divisor;
170  //      unsigned long reciprocal;
171  //   } FILTER_ENTRY;
172  //
173  //   bits 27:31 of divisor contains the number of significant ↩
         bits in
174  //   reciprocal, while bits 0:26 contains the divisor itself.
175  //
176  //
177  //
178  #macro dl_gridstat_filter[ forward, out_port, variableid, ↩
         timestamp1, timestamp0, pubint]
179  .begin
180    .reg divisor reciprocal nbits rem sram_offset txreg_addr ↩
         cam_res cam_e lm_a i
181
182    br[try_cam_valid#]
183
184  wait_cam_load#:
185    ctx_arb[voluntary]
186
187  try_cam_valid#:
188    cam_lookup[cam_res, variableid]
189    alu_shf[tmp, 0x1, AND, cam_res, >>7]
190    sub_shf_right(tmp, CAM_VALID, cam_res, 8) // extract state ↩
           bits
191    bne[wait_cam_load#] ; branch if CAM_VALID is not set
192
193    alu[forward, --, B, 0] // default to dropping packet
194
195    alu_shf[cam_e, 0xf, AND, cam_res, >>3] ; find cam entry
196    ; lm_a = cam_e * 36
197    alu_shf[lm_a, --, B, cam_e, <<5]    ; lm_a = cam_e * 32
198    add_shf_left(lm_a, lm_a, cam_e, 2) ; lm_a = lma_a + cam_e * 4
199
200    .if ( out_port == 1 )
201      alu[lm_a, lm_a, +, 12]
202    .elif ( out_port == 2 )
203      alu[lm_a, lm_a, +, 24]
204    .endif
```

```
205
206     alu_shf[lm_a, --, B, lm_a, <<2]   ; move to offset used ↩
            ACTIVE_LM_ADDR_0 write
207
208       //alu[i, --, B, 1]
209       //.while (i>0)
210     alu[forward, --, B, 0] // default to dropping packet
211     local_csr_wr[ACTIVE_LM_ADDR_0, lm_a]
212     nop
213     nop
214     nop
215
216     // The following loop will terminate if a decision to forward ↩
            the packed
217     // is reached, or when we get to the end of the filter_entry ↩
            array.
218     .while( *l$index0 && !forward )
219
220       alu[tmp, --, B, *l$index0++]
221       alu_shf[nbits, --, B, tmp, >>27]
222       alu_shf[divisor, tmp, AND~, 0x1F, <<27]
223       alu[reciprocal, --, B, *l$index0++]
224
225       // if ( (timestamp + pubint/2) % divisor < pubint ) then ↩
            FORWARD (i.e. don't drop)
226       // Note that the timestamp should already have pubint/2 ↩
            added to it.
227
228       remainder64x32[rem, timestamp1, timestamp0, divisor, ↩
            reciprocal, nbits]
229
230       .if ( rem < pubint )
231         // FORWARD PACKET
232         alu[forward, --, B, 1]
233       .endif
234
235     .endw
236
237       // alu[i, i, -, 1]
238       //.endw
239
240       alu[forward, --, B, 1] // forward packet
241
242   .end
243  #endm
```

```
244
245  #macro rt_next_port []
246  .begin
247    ;alu[sram_addr, sram_addr, +, (8*8)] ; sram_addr += sizeof(←
          FILTER_ENTRY) * MAX_SUBINTS
248  .end
249  #endm
```

## A.8  packet_echo.uc

Listing A.8: packet_echo.uc

```
1  /*
2   * This file is the entry point of the Route microblock and ←
       implements
3   * a simple GridStat status router.
4   *
5   * This file sources packets from the packet_rx[] microblock and←
        forwards them
6   * to the corresponding port transmit queue for the ←
       sphy_mphy4_tx[] microblock.
7   *
8   */
9
10  #include <dl_system.h>
11  #include <stdmac.uc>
12  #include <hardware.h>
13  #include <dl_meta.uc>
14  #include <dl_buf.uc>
15  #include <dl_misc.uc>
16  #include <sig_macros.uc>
17
18  #include <routing_table.uc>
19
20  /*
21   * Signal and register definitions
22   */
23
24  // absolute registers
25  //.reg  @count_p0
26  //.reg  @count_p1
27  //.reg  @count_p2
28  //.reg  @count_p3
```

```
29  //.reg    @count_empty
30  //.reg    @debug_cnt
31
32  // External reply interface from packet copier
33  .reg global volatile read $rx_copy_handle
34  .sig visible volatile sig_copy_done
35  .addr sig_copy_done PACKET_COPY_REPLY_SIGNAL
36
37  // thread ordering signals and registers
38  .sig volatile    sig_prev
39  .addr      sig_prev         DL_SIG_WAKE
40
41  // scratch ring put msgs
42  //xbuf_alloc($tx_pktc, 3, write) ; request msg to packet copier
43  //xbuf_alloc($txreq, 3, write) ; tx request message to ethernet ←
        TX
44  //xbuf_alloc($$gdata, 4, read)
45  .reg global volatile write $tx_pktc0 $tx_pktc1 $tx_pktc2
46  .xfer_order $tx_pktc0 $tx_pktc1 $tx_pktc2
47  .reg global volatile write $txreq0 $txreq1 $txreq2
48  .xfer_order $txreq0 $txreq1 $txreq2
49  .reg global volatile read $$gdata0 $$gdata1 $$gdata2 $$gdata3
50  .xfer_order $$gdata0 $$gdata1 $$gdata2 $$gdata3
51
52  .reg global volatile $rwxf_meta0 $rwxf_meta1 $rwxf_meta2 ←
        $rwxf_meta3
53  .xfer_order $rwxf_meta0 $rwxf_meta1 $rwxf_meta2 $rwxf_meta3
54  .sig  volatile sig_meta
55
56  .reg global volatile sig_mask
57  .reg global volatile dl_buf_handle     // The current Buffer ←
        handle.
58  .reg global volatile dl_buf_copy_handle  // The current Buffer ←
        copy handle.
59  .reg global volatile dl_eop_buf_handle   // For large packets, ←
        this is the last buffer in the chain
60  .reg global volatile dl_next_block     // Next block that ←
        should process the buffer/packet
61  .reg global volatile port_mask forward t1 t0 varid
62
63  // This is a global signal that indicates that
64  // system initialization is done
65  .sig  volatile system_init_done
66  .addr   system_init_done  me_init_signal
67
```

```
68  .sig  volatile scr_get scr_put_copy scr_put0 scr_put1 scr_put2 ←
        sig_dram; signal for scratch put
69
70
71  #define NUM_PORTS 3
72
73  .reg global volatile port[3]  // packet handles or UC_NULL
74
75  /*
76   * The request_packet_copy macro was written to interface
77   * with the packet_copier microblock. That microblock and
78   * this macro is no longer in use, as we now use a different
79   * strategy to enable multicasting.
80   */
81  #macro request_packet_copy[sig_mask]
82  .begin
83    .reg tmp tmp_lw port_out0 port_out1 port_in
84
85    .set dl_meta[1] dl_meta[2] dl_meta[3]
86
87    dl_meta_get_input_port[port_in]
88
89    // Calculate output port (0−>1, 1−>2, 2−>0)
90    alu[port_out0, port_in, +, 1]
91    .if(port_out0 == 3)
92      alu[port_out0, −−, b, 0]
93    .endif
94
95    // compute output port of target, if source/target input port ←
            is 0, the target output port is
96    // either 1 or 2 (the one that's not source output port).
97    alu[tmp_lw, port_in, +, port_out0]
98    alu[port_out1, tmp_lw, XOR, 0x3]  // port_out1 = the remaining ←
            port
99
100   // Prepare copy request for the packet copier
101
102   //dl_meta_get_packet_size[pktc_packet_size]
103
104   alu[$tx_pktc[0], −−, B, dl_buf_handle]
105
106   dl_meta_get_buffer_size[tmp]  // use only the first buffer for ←
            now
107   alu_shf[tmp_lw, −−, B, tmp, <<16]
108   dl_meta_get_offset[tmp]
```

145

```
109    alu[tmp_lw, tmp_lw, OR, tmp]
110    alu[$tx_pktc[1], --, B, tmp_lw]
111
112    local_csr_rd[ACTIVE_CTX_STS]
113    immed[tmp, 0]
114    alu[tmp, tmp, AND, 0x7] // tmp = ctx
115    alu_shf[tmp_lw, port_out1, OR, tmp, <<4] // tmp_lw = port_out1↩
           | (ctx << 4)
116    move[tmp, __UENGINE_ID] // tmp = ME
117    alu_shf[tmp_lw, tmp_lw, OR, tmp, <<7] // tmp_lw = port_out1 | ↩
           (ctx << 4) | (ME << 7)
118    alu_shf[tmp_lw, tmp_lw, OR, &$rx_copy_handle, <<12] // tmp_lw ↩
           = port_out1 | (ctx << 4) | (ME << 7) | ($rx <<12)
119    alu[$tx_pktc[2], --, B, tmp_lw]
120
121    br[pktc_test_ring_full#]
122
123  pktc_ring_full#:
124    nop
125    ctx_arb[voluntary]
126
127  pktc_test_ring_full#:
128    br_inp_state[SCR_RING/**/PACKET_COPIER_SCR_RING/**/_STATUS, ↩
           pktc_ring_full#]
129
130    // Send a copy request to the packet copier
131    alu_shf[tmp, --, B, PACKET_COPIER_SCR_RING, <<2]
132    scratch[put, $tx_pktc[0], 0, tmp, 3], ctx_swap[scr_put_copy]
133
134    // sig_mask |= (1<<sig_copy_done)
135    alu[tmp, --, B, &sig_copy_done]
136    alu[--, tmp, OR, 0]
137    alu_shf[sig_mask, sig_mask, OR, 1, <<indirect]
138
139  .end
140  #endm
141
142  #macro pc_get_copy_handle[_in_handle, _out_handle]
143  .begin
144      .reg tmp_buf_sram_size
145    move(tmp_buf_sram_size, BUF_SRAM_SIZE);
146      alu[_out_handle, _in_handle, +, tmp_buf_sram_size]
147  .end
148  #endm
149
```

```
150  #macro prepare_port_out_handles[]
151  .begin
152    .reg port_out0 port_out1 port_in tmp
153
154    alu[port[0], --, B, UC_NULL]
155    alu[port[1], --, B, UC_NULL]
156    alu[port[2], --, B, UC_NULL]
157
158    dl_meta_get_input_port[port_in]
159
160    // Calculate output port (0->1, 1->2, 2->0)
161    alu[port_out0, port_in, +, 1]
162    .if(port_out0 == 3)
163      alu[port_out0, --, b, 0]
164    .endif
165
166    // compute output port of target, if source/target input port ←↩
                is 0, the target output port is
167    // either 1 or 2 (the one that's not source output port).
168    alu[tmp, port_in, +, port_out0]
169    alu[port_out1, tmp, XOR, 0x3] // port_out1 = the remaining ←↩
                port
170
171
172    // port[port_out0] = dl_buf_handle
173    alu_shf[tmp, --, b, port_out0, <<1]
174    jump[tmp, port_out0_0#], targets[port_out0_0#, port_out0_1#, ←↩
                port_out0_2#]
175    port_out0_0#:
176      alu[port[0], --, B, dl_buf_handle]
177      br[port_out0_end#]
178    port_out0_1#:
179      alu[port[1], --, B, dl_buf_handle]
180      br[port_out0_end#]
181    port_out0_2#:
182      alu[port[2], --, B, dl_buf_handle]
183      //br[port_out0_end#]
184    port_out0_end#:
185
186
187    // inline packet copy
188
189    // Load source packet meta data
190    dl_meta_load_cache[dl_buf_handle, $rwxf_meta, sig_meta, 0, 4]
191
```

```
192     // Get available packet buffer
193       pc_get_copy_handle[dl_buf_handle, dl_buf_copy_handle]
194
195     // Wait for source packet meta data to load
196     ctx_arb[sig_meta]
197
198     // Setup local cache in dl_meta GPRs
199     dl_meta_init_cache[$rwxf_meta[1], $rwxf_meta[1], $rwxf_meta↩
            [2], $rwxf_meta[3], 0, 0, 0, 0]
200
201     // get input and output port of source packet.
202     //dl_meta_get_input_port[port_in]
203     //dl_meta_get_output_port[port_out]
204
205     // Change output port of local meta data cache to reflect that↩
             of the target packet.
206     dl_meta_set_output_port[port_out1]
207
208     // Write a subset of meta data cache to the target packet, ↩
            this is correct as
209     // the target (copy) is to have the same meta data as the ↩
            source except for
210     // output port and next buffer/packet pointers.
211     dl_meta_flush_cache[$rwxf_meta, dl_buf_copy_handle, sig_meta, ↩
            SIG_NONE, 0, 4]
212
213
214
215     // port[port_out1] = copied handle // set it to 0 for now ↩
            which is differenct from UC_NULL (0xFF)
216     alu_shf[tmp, --, b, port_out1, <<1]
217     jump[tmp, port_out1_0#], targets[port_out1_0#, port_out1_1#, ↩
            port_out1_2#]
218     port_out1_0#:
219       alu[port[0], --, B, dl_buf_copy_handle]
220       //alu[port[0], --, B, UC_NULL]
221       br[port_out1_end#]
222     port_out1_1#:
223       alu[port[1], --, B, dl_buf_copy_handle]
224       //alu[port[1], --, B, UC_NULL]
225       br[port_out1_end#]
226     port_out1_2#:
227       alu[port[2], --, B, dl_buf_copy_handle]
228       //alu[port[2], --, B, UC_NULL]
229       //br[port_out1_end#]
```

```
230    port_out1_end#:
231
232  .end
233  #endm
234
235  #macro request_packet_gridstat_data[]
236  .begin
237    .reg sdram_data_base sdram_offset tmp
238
239    dl_buf_get_data[sdram_data_base, dl_buf_handle]
240
241    // set sdram_offset to start of packet data
242    dl_meta_get_offset[tmp]
243    alu[sdram_offset, --, B, tmp]
244
245    // offset  0, size 14 - Ethernet
246    // offset 14, size 20 - IP
247    // offset 34, size  8 - UDP (skipped in simulation)
248    // offset 42, size xx - Gridstat
249    alu[sdram_offset, sdram_offset, +, 34]
250
251    dram[read, $$gdata0, sdram_data_base, sdram_offset, 2], ↩
          sig_done[sig_dram]
252  .end
253  #endm
254
255  // must call request_packet_gridstat_data before calling this ↩
        function
256  #macro get_packet_gridstat_data[t1, t0, varid]
257  .begin
258
259    ctx_arb[sig_dram]
260
261    // NOTE: extracting timestamp and variableid assumes the ↩
          gridstat message
262    // starts at offset 42 and uses network byte order (BIG ENDIAN↩
          )
263
264    // GRIDSTAT message format:
265    // offset  0, size 8 - timestamp
266    // offset  8, size 4 - variable ID
267
268    // offset 42, size 4 - high 32 bits of timestamp
269    ld_field_w_clr[t1, 1100, $$gdata0, <<16]
270    ld_field[t1, 0011, $$gdata1, >>16]
```

149

```
271
272    // offset 46, size 4 − low 32 bits of timestamp
273    ld_field_w_clr[t0, 1100, $$gdata1, <<16]
274    ld_field[t0, 0011, $$gdata2, >>16]
275
276    // offset 50, size 4 − 32 bit variable ID
277    ld_field_w_clr[varid, 1100, $$gdata2, <<16]
278    ld_field[varid, 0011, $$gdata3, >>16]
279
280  .end
281  #endm
282
283  #macro dl_gridstat_route[sig_mask]
284  .begin
285    // Make sure current packet is intended for GRIDSTAT routing
286    alu[−−, dl_next_block, −, BID_GRIDSTAT]
287    bne[end_dl_gridstat_route#]
288
289    request_packet_gridstat_data[]
290
291    //request_packet_copy[sig_mask]
292
293    prepare_port_out_handles[]
294
295    // Lookup what ports might want current packet
296    alu[port_mask, −−, B, 0x7] // all 3 ports
297
298    .reg port_in
299    dl_meta_get_input_port[port_in]
300
301    // except incoming port
302    bits_clr(port_mask, port_in, 1)
303
304
305    // Rate filter per port
306
307    .reg cur_port
308    alu[cur_port, −−, B, 0]
309
310      // 50 instructions
311
312    // wait for dram read to complete then extract timestamp and ↩
            variableid
313    get_packet_gridstat_data[t1, t0, varid]
314
```

```
315    .reg tmp
316
317    alu[tmp, t0, AND, 0xff]
318    .if (tmp == 0x34 || tmp == 0x84 )
319      alu[tmp, --, B, tmp]
320      nop
321    .endif
322
323    init_rtable_offset[varid]
324
325      // 158 instructions if RT entry is not in CAM, 103 ←
              instructions if it is.
326
327    // extract routing info common to all output ports
328    // for now this means publish interval only
329    .reg pubint
330    get_stream_info[pubint]
331
332    // timestamp = timestamp + (pubint/2)
333    alu[tmp, --, B, pubint, >>1] ; tmp = pubint / 2
334    alu[t0, t0, +, tmp]
335    alu[t1, t1, +carry, 0]
336
337      // 164 instructions including initial following while check.
338
339    .while ( cur_port < NUM_PORTS )
340      // for all ports
341
342      alu[--, cur_port, OR, 0]
343
344      alu_shf[tmp, 0x1, AND, port_mask, >>indirect]
345
346      .if ( tmp )
347
348        dl_gridstat_filter[forward, cur_port, varid, t1, t0, ←
              pubint]
349
350        .if ( forward == 0 ) // whether to drop packet
351          // clear bit if packet should not be sent to cur_port
352          bits_clr(port_mask, cur_port, 1)
353
354          //alu_shf[tmp, --, b, cur_port, <<1]
355          //jump[tmp, port0#], targets[port0#, port1#, port2#]
356          //port0#:
357          //  alu[port[0], --, B, UC_NULL]
```

151

```
358          //   br [ port_end #]
359          //port1#:
360          //    alu [ port [1] ,  --, B,  UC_NULL]
361          //   br [ port_end #]
362          //port2#:
363          //    alu [ port [2] ,  --, B,  UC_NULL]
364          //   //br [ port_end #]
365          //port_end#:
366
367        . endif
368      . endif
369
370      rt_next_port []
371      alu [ cur_port ,  cur_port ,  +,  1]
372
373    . endw
374
375    move [ dl_next_block ,  BID_GRIDSTAT_NEXT1 ]
376
377  end_dl_gridstat_route #:
378
379  . end
380  #endm
381
382
383  #define_eval  TXR_VAL0  ( PACKET_TX_SCR_RING_0 <<2)
384  #define_eval  TXR_VAL1  ( PACKET_TX_SCR_RING_1 <<2)
385  #define_eval  TXR_VAL2  ( PACKET_TX_SCR_RING_2 <<2)
386  #define_eval  TXR_VAL3  ( PACKET_TX_SCR_RING_3 <<2)
387
388  #macro  dl_gridstat_write_tx_req [ REQ_PORT ]
389  . begin
390
391    . reg  __tmp_a
392    . if  (  port [ REQ_PORT ]  == 0  )
393      // if  buffer  handle  has  the  value  of  0 ,  then  it 's  the
394      // target  of  packet_copier
395
396          // skip  packet  copier  step  for  now
397          br [ end_write_tx #]
398          nop
399          nop
400
401      alu [ port [ REQ_PORT ] ,  --, B,  $rx_copy_handle ]
402    . endif
```

```
403    alu_shf[__tmp_a, --, b, REQ_PORT, <<24]          ; 27:24 = output↩
           port
404    ld_field[__tmp_a, 0111, port[REQ_PORT]]          ; 23:00 = ptr to↩
           sop buffer desc
405    .io_completed scr_put/**/REQ_PORT/**/
406    alu[$txreq/**/REQ_PORT/**/, __tmp_a, OR, 1, <<31] ; 31    = ↩
           set valid bit (31:28 = reserved )
407
408    // check whether we are supposed to drop packet
409    //alu[__tmp_a, --, B, REQ_PORT]
410    //alu[--, __tmp_a, OR, 0]
411    //alu_shf[__tmp_a, 0x1, AND, port_mask, >>indirect]
412    //.if ( ! __tmp_a && port[REQ_PORT] != UC_NULL )
413      // drop packet
414      // TODO: support chained buffer drops
415    //   dl_buf_free(port[REQ_PORT], BUF_FREE_LIST0)
416    //   br[end_write_tx#]
417    //.endif
418
419    br[write_ring#]
420
421  full_ring#:
422    ctx_arb[voluntary]
423
424  write_ring#:
425
426  #if ( REQ_PORT == 0 )
427  #define_eval SCR_STATUS SCR_RING/**/PACKET_TX_SCR_RING_0/**/↩
         _STATUS
428  #elif ( REQ_PORT == 1 )
429  #define_eval SCR_STATUS SCR_RING/**/PACKET_TX_SCR_RING_1/**/↩
         _STATUS
430  #elif ( REQ_PORT == 2 )
431  #define_eval SCR_STATUS SCR_RING/**/PACKET_TX_SCR_RING_2/**/↩
         _STATUS
432  #else
433  #error "REQ_PORT must be in {0,1,2}"
434  #endif
435    br_inp_state[SCR_STATUS, full_ring#]
436  #undef SCR_STATUS
437
438    //.set $txreq/**/REQ_PORT/**/
439    //.io_completed scr_put/**/REQ_PORT/**/
440
441    .reg zero
```

```
442    alu[zero, --, B, 0]

443

444    // test whether range of buffer handle is valid, used for ←
           debugging
445    ld_field_w_clr[__tmp_a, 0111, port[REQ_PORT]]
446    .if ( __tmp_a >= 0x800 )
447      nop
448    .endif

449

450    scratch[put, $txreq/**/REQ_PORT/**/, zero, TXR_VAL/**/REQ_PORT←
           /**/, 1], ctx_swap[scr_put/**/REQ_PORT/**/]

451

452    //alu[@count_p/**/REQ_PORT/**/, @count_p/**/REQ_PORT/**/, +, ←
           1]
453    // sig_mask |= (1<<scr_put0)

454

455    //alu[tmp, --, B, &scr_put/**/REQ_PORT/**/]
456    //alu[--, tmp, OR, 0]
457    //alu_shf[sig_mask, sig_mask, OR, 1, <<indirect]

458

459  end_write_tx#:

460

461  .end
462  #endm

463

464

465  /*
466   * Macro "dl_gridstat_sink":
467   *
468   * This macro always performs strict thread order.
469   */
470  #macro dl_gridstat_sink[sig_mask, scr_put]
471  .begin
472    .reg tmp

473

474    // Force strict thread order here

475

476    // sig_mask |= (1<<sig_prev)
477    alu[tmp, --, B, &sig_prev]
478    alu[--, tmp, OR, 0]
479    alu_shf[sig_mask, sig_mask, OR, 1, <<indirect]

480

481    local_csr_wr[active_ctx_wakeup_events, sig_mask]  ; tells ←
           which signals to wait on.
482    ctx_arb[--]
```

```
483
484   //. io_completed  sig_copy_done
485   //. io_completed  scr_put_copy
486
487     // clear mask now that we waited on all signals
488     alu[sig_mask, --, B, 0]
489
490
491     // Because we might get a context switch before the scratch[↩
            get, ...]
492     // in dl_gridstat_source, we cannot signal next thread until ↩
            after
493     // scratch[get, ...] in dl_gridstat_source
494
495     // Make sure current packet is intended for Packet TX
496     alu[--, dl_next_block, -, BID_GRIDSTAT_NEXT1]
497     bne[end_dl_gridstat_sink#]
498
499     // TODO: if a packet is to be dropped, make sure to free the ↩
            corresponding buffers.
500     // TODO: use port mask rather than NULL buffer handles to ↩
            determine whether to forward a packet or not
501
502     .if ( port[0] != UC_NULL )
503       dl_gridstat_write_tx_req[0]
504     .endif
505
506     .if ( port[1] != UC_NULL )
507       dl_gridstat_write_tx_req[1]
508     .endif
509
510     .if ( port[2] != UC_NULL )
511       dl_gridstat_write_tx_req[2]
512     .endif
513
514   end_dl_gridstat_sink#:
515
516   .end
517
518   #endm /* dl_gridstat_sink */
519
520
521
522   #macro dl_gridstat_source[THREAD_ORDER, sig_mask]
523   .begin
```

155

```
524    .reg tmp me
525    xbuf_alloc[$rdata, 5, read]
526
527  #if (THREAD_ORDER == DL_THREAD_ORDER)
528    // Wait for signal from previous thread to ensure thread ↩
           ordering
529    ctx_arb[sig_prev]
530  #endif /* DL_THREAD_ORDER */
531
532    // Source packet from scratch ring GRIDSTAT_RING_IN
533    alu_shf[tmp, --, B, GRIDSTAT_RING_IN, <<2]
534    scratch[get, $rdata0, 0, tmp, 5], sig_done[scr_get]
535
536    // Signal that it's safe for next thread to sink and source
537    br=ctx[7, thread_7#]
538    signal_next_ctx[DL_SIG_WAKE]
539  thread_7#:
540
541
542  wait_for_io_signals#:
543
544    // We wait for the above scratch[get...] and the scratch[put↩
           ...] in dl_gridstat_sink.
545
546    // Add the signal to the signal mask on which to wait
547    alu[tmp, --, B, &scr_get]
548    alu[--, tmp, OR, 0]
549    alu_shf[sig_mask, sig_mask, OR, 1, <<indirect]  ; shift 1 by ↩
           scr_get times
550
551    // wait
552    local_csr_wr[active_ctx_wakeup_events, sig_mask]  ; tells ↩
           which signals to wait on.
553    ctx_arb[--]
554
555    // clear mask now that we waited on all signals
556    alu[sig_mask, --, B, 0]
557
558    // Assembler doesn't know that we actually wait on the I/O (↩
           scratch[get..])
559    // and (scratch[put..])
560  .io_completed scr_get
561  .io_completed scr_put0
562  .io_completed scr_put1
563  .io_completed scr_put2
```

```
564
565    br!=ctx[7, thread_7_2#]
566    immed[me, __UENGINE_ID]
567    .if ( me == LAST_ME )
568      signal_me[FIRST_ME, DL_SIG_WAKE]
569    .else
570      signal_next_me[DL_SIG_WAKE]
571    .endif
572  thread_7_2#:
573
574
575    // Check for empty ring
576    alu[--, $rdata0, -, 0]
577    beq[ring_empty#]
578
579
580    // Setup meta data associated with packet
581
582    dl_meta_init_cache[0, $rdata2, $rdata3, $rdata4, 0, 0, 0, 0]
583      //dl_meta_set_nexthop_id[0xffff]  // set nexthop ID to -1
584
585    move[dl_buf_handle, $rdata0]
586    move[dl_eop_buf_handle, $rdata1]
587
588    move[dl_next_block, BID_GRIDSTAT]
589
590    // Exit macro
591    br[end_dl_gridstat_source#]
592
593  ring_empty#:
594
595    //alu[@count_empty, @count_empty, +, 1]
596
597    // Output NULL Packet
598    move[dl_buf_handle, IX_NULL]
599    move[dl_next_block, IX_NULL]
600
601  end_dl_gridstat_source#:
602
603    xbuf_free($rdata)
604
605  .end
606  #endm /* dl_gridstat_source */
607
608
```

```
609
610  #macro dl_gridstat_init[]
611  .begin
612
613    .reg tmp me
614
615    immed[me, __UENGINE_ID]
616
617
618      //  Context common initializations
619
620      // get context num
621      //local_csr_rd[active_ctx_sts]
622      //immed[ctx_num,0]   //immed operand contains local_csr_rd ↩
            data
623      //alu[ctx_num, ctx_num, AND, 0x07]
624
625      move[sig_mask, 0]
626      cam_clear
627
628
629      //  Context specific initializations
630
631      .if( ctx() == 0 )
632        //alu[@count_p0, --, b, 0]
633        //alu[@count_p1, --, b, 0]
634        //alu[@count_p2, --, b, 0]
635        //alu[@count_p3, --, b, 0]
636        //alu[@count_empty, --, b, 0]
637        //alu[@debug_cnt, --, b, 0]
638
639        .reg a1 a0 rec n d rem
640
641        immed32(a1,  0x00000039)
642        immed32(a0,  0x000a0007)
643        immed32(n, 4)
644        immed32(d, 30)
645        immed32(rec, 0x88888888)
646        remainder64x32[rem, a1, a0, d, rec, n]
647
648      .endif
649
650
651      //  All threads wait for system go signal
652
```

```
653        .if( ctx() == 0 )
654
655          // this is the signal that all blocks are to
656          // wait on to indicate that system initialization
657          // is done
658
659          ctx_arb[system_init_done]
660
661          // signal all other contexts to proceed
662
663          alu[tmp, --, B, &sig_prev, <<3]
664          alu[tmp, tmp, +, 1]
665          local_csr_wr[SAME_ME_SIGNAL, tmp] //thread 1
666          alu[tmp, tmp, +, 1]
667          local_csr_wr[SAME_ME_SIGNAL, tmp] //thread 2
668          alu[tmp, tmp, +, 1]
669          local_csr_wr[SAME_ME_SIGNAL, tmp] //thread 3
670          alu[tmp, tmp, +, 1]
671          local_csr_wr[SAME_ME_SIGNAL, tmp] //thread 4
672          alu[tmp, tmp, +, 1]
673          local_csr_wr[SAME_ME_SIGNAL, tmp] //thread 5
674          alu[tmp, tmp, +, 1]
675          local_csr_wr[SAME_ME_SIGNAL, tmp] //thread 6
676          alu[tmp, tmp, +, 1]
677          local_csr_wr[SAME_ME_SIGNAL, tmp] //thread 7
678          alu[tmp, tmp, +, 1]
679
680        .else
681
682          .set_sig sig_prev
683          ctx_arb[sig_prev]
684
685        .endif
686
687      // kick start next thread sequencing by thread 0 of the ←
             first ME signaling itself
688
689      .if( ctx() == 0 && me == FIRST_ME )
690
691          alu[tmp, --, b, &sig_prev, <<3]
692          local_csr_wr[SAME_ME_SIGNAL, tmp]
693
694      .endif
695
696  .end
```

```
697 #endm /* dl_gridstat_init */
698
699
700 /*
701  *
702  * GRIDSTAT DISPATCH LOOP
703  *
704  * This is where the microengine start executing code.
705  *
706  */
707
708 // Initialize
709 dl_gridstat_init[]
710
711 // Get packet from Packet RX
712 dl_gridstat_source[DL_THREAD_ORDER, sig_mask]
713
714 // instruction count 0 here.
715
716 loop#:
717     // Compute meta data associated with this packet.
718     dl_gridstat_route[sig_mask]; // 239 instructions + 88 * ↩
        RTcomplexity.
719
720     // Send packet to Packet TX
721     dl_gridstat_sink[sig_mask, scr_put]; typically 50 ↩
        instructions.
722
723     // Get next packet from Packet RX.
724     // We can call the source macro with no thread order here ↩
        because
725     // strict thread order is already enforced in ↩
        dl_gridstat_sink macro
726     // called immediately before source.
727     dl_gridstat_source[DL_THREAD_NO_ORDER, sig_mask]; 30 ↩
        instructions.
728
729 br[loop#]
730
731 // makes a total of 320 + 88 * RTcomplexity instructions.
732 // e.g. RTcomplexity of 2 makes 496 instructions.
```

160