

UNIVERSITY OF OSLO
Department of Informatics

**Investigating the
distribution of
functionality for
building a video
server hypercube
with IXP2400 cards**

Master thesis

Andreas Petlund

11th November 2005



Acknowledgements

I want to thank my advisors, Pål Halvorsen and Carsten Griwodz for excellent guidance and unlimited patience. I also want to thank my fellow students Tom Anders Dalseng and Øystein Yri Sunde for inspiration, discussions and help. Finally, I want to thank my girlfriend Anne for support and structured help, and my family for never giving up.

Abstract

Streamed multimedia is becoming common on the Internet as the bandwidths increase for end-users. To stream data with high bitrates to many concurrent users, servers are needed that can handle these extreme loads. Implementing servers in clusters to be able to meet demands has proved to be a good strategy, providing scalability and performance. One commercial actor that has done this successfully is nCube. Their n4x solution is a server cluster based on a hypercube interconnection topology, and their reported server performance is promising. However, the use of special hardware for offloading routing increases the cost for deploying this system. Network processing units have many similar properties to the special hardware used by nCube. They are fully programmable and optimized for networking tasks. In this thesis, we start by charting the functionality and capabilities of the IXP2400 network processor by implementing a series of test applications. Using knowledge gained from this exploration, the design and implementation of a video server hypercube prototype is done. We present an evaluation of the IXP2400 hardware platform based on the test applications. Our video server cube prototype is also presented and evaluated.

Contents

1	Introduction	1
1.1	Background and motivation	1
1.2	Thesis domain	2
1.3	Document structure	3
2	Internet exchange architecture	5
2.1	Network Processing Units	5
2.2	Intel Internet Exchange Processors	7
2.2.1	IXP1200	7
2.2.2	ENP2505	7
2.2.3	IXP2400	8
2.2.4	Radisys ENP2611	11
2.2.5	New hardware	12
2.3	IXA software libraries and tools	14
2.3.1	Intel IXA SDK 3.51	14
2.3.2	Radisys ENP SDK 3.5	16
2.3.3	IXP programming paradigms	17
2.3.4	MontaVista Linux	22
2.4	Summary	23

3	IXP2400 Evaluation	24
3.1	Exploring the IXP2400 hardware	24
3.2	Static forward	25
3.3	IP header switch application	26
3.3.1	μ Engines implementation	27
3.3.2	XScale implementation	28
3.3.3	IP header switch application evaluation	29
3.4	Performance gain by offloading	35
3.4.1	Test	36
3.4.2	Discussion	36
3.5	Summary	38
4	Multimedia systems	40
4.1	Multimedia applications	40
4.2	Multimedia requirements	41
4.3	Improving server performance	43
4.4	Multimedia server implementations	43
4.4.1	Single server implementation	43
4.4.2	Server cluster implementation	45
4.4.3	Proxies	46
4.5	Hypercube/n4x server solution	46
4.5.1	Hypercube multicomputer structure	46
4.5.2	The c-cor n4x multimedia server architecture	48
4.6	Protocols	50
4.6.1	RTSP	50
4.6.2	RTP	53
4.7	Summary	56

5	The VS³ Video Server Cube	57
5.1	Hypercube server general design	57
5.2	The legacy design	59
5.2.1	The INF5070 implementation	59
5.3	SDK code base	60
5.4	XScale-Host communication	61
5.5	The raw socket design approach	62
5.5.1	TCP Handling	62
5.6	Egress TCP design	64
5.7	Cube routing	65
5.7.1	Intra-cube header and extension	67
5.7.2	IC packet types	68
5.7.3	Intra-cube session	69
5.7.4	Partitioning of tasks	70
5.7.5	μ Engine tasks	71
5.7.6	XScale tasks	73
5.7.7	Host-XScale data path	75
5.7.8	Host tasks	76
5.8	RTP-session	77
5.9	Microblock designs	78
5.10	Evaluation	80
5.10.1	Packets routed through a switched network	80
5.10.2	IXP hypercube routing tests	80
5.10.3	Comparison of routing times	83
5.10.4	Processing times given increased load	84
5.11	Summary	86
6	Conclusion	87
6.1	Summary	87
6.2	Results	87
6.3	Further work	88

A	VS³ video server cube source	94
A.1	ICrouter microblock source	94
A.1.1	dl_system.excerpt.h	94
A.1.2	ICrouter.h	96
A.1.3	ICrouter.c	97
A.2	XScale source	105
A.2.1	dcube.h	105
A.2.2	dcube.c	109
A.2.3	dcube_utils.h	120
A.2.4	dcube_utils.c	121
A.2.5	bogus_tcp.h	134
A.2.6	bogus_tcp.c	137
A.2.7	bogus_tcp_utils.c	143
A.3	Host source	148
A.3.1	srtsp.h	148
A.3.2	srtsp.c	154
A.3.3	srtsp_utils.h	161
A.3.4	srtsp_utils.c	161
A.3.5	ic.h	170
A.3.6	rtp.h	171
A.3.7	dbg_msg.h	172

List of Figures

2.1	Evolution of technologies [22].	6
2.2	IXP1200 block diagram [31].	8
2.3	ENP2505 block diagram [12].	9
2.4	IXP2400 block diagram [32].	10
2.5	ENP2611 Block diagram [13]	13
2.6	RX and TX block functionality and interface.	17
2.7	IXA Hierachy model [6]	19
2.8	Struct describing packet metadata.	20
2.9	Calculation of metadata and data offsets from buffer handle values.	21
2.10	How buffer handles are passed on scratch rings between processing components.	22
3.1	Radisys static forward application block diagram.	26
3.2	IP header switch application on μ Engines.	27
3.3	IP header switch application using μ Engines and XScale.	29
3.4	Generic test setup.	30
3.5	User space packet echo application test setup.	31
3.6	Kernel space packet echo using iptables.	32
3.7	IXP XScale ICMP echo test setup.	32
3.8	IXP μ Engine ICMP echo test setup.	33
3.9	Comparison of average times for the IP header switch application.	34
3.10	Times spent on tar process with different implementations of IP header switch application.	38
3.11	Average total time spent on tar process.	39

3.12	Average time spent in user space for tar process.	39
3.13	Average time spent in kernel space for tar process.	39
4.1	Simple general server architecture [23]	44
4.2	Interconnection topologies	45
4.3	Diagram of a possible proxy layout [23].	47
4.4	Hypercube scaling	49
4.5	A typical RTSP session [7]	53
4.6	RTP header format.	54
4.7	RTP header extension format.	55
5.1	Basic design of multimedia server cube with IXP cards.	58
5.2	Dataflow inf5070 implementation	60
5.3	Design with forwarding of TCP to machine hosting the file.	63
5.4	Design with egress handling TCP connections.	64
5.5	IC header and header extension format.	67
5.6	IC session struct.	69
5.7	Flow chart for <i>ICrouter</i> microblock when a packet is received.	72
5.8	Flow chart for XScale int_A thread.	74
5.9	Flow chart for Host when receiving data from XScale.	76
5.10	Three microblock structuring alternatives.	79
5.11	Test setup: Routing using switched network.	81
5.12	Test setup: Timing of routing in the IXP hypercube.	82
5.13	Comparison of average routing times.	84
5.14	Development of processing time with increased load.	85

List of Tables

2.1	Main differences between ENP2505 and ENP2611	14
2.2	IXP2400 data terminology [32].	16
3.1	IP header switch application times.	33
3.2	Time used for tar process with different IP header switch implementations.	37
5.1	Routing times for IXP cube and switched network routing.	83
5.2	Routing times on μ Engines given increased load.	85

Chapter 1

Introduction

Streaming data with high bitrates to a potentially high number of concurrent users presents us with great challenges regarding scalability and throughput. The network bandwidth rates needed for multimedia streaming systems today strain the processing power of server solutions, stealing resources that could be applied to other important tasks. This thesis will investigate the possibility of using network processing units (NPUs) to offload networking tasks and, in particular, to implement routing in a distributed server cluster topology.

1.1 Background and motivation

The Internet, since its beginning as an experimental network in 1968, has experienced an astounding development. Especially in the last ten years, it has become common property, and a variety of services are now available to the public. During the last years, the bandwidth available to end users has been steadily increasing in most developed countries. With the high bitrates becoming available, the transmission of audio and video has become the new killer application driving the development of networks and infrastructure.

The most common way to access digital multimedia content is still to download a file and play it locally. The alternative is to transmit the data as the content is consumed. This is known as data streaming. Streamed services are becoming more common every day due to the increased capacity available for end users. For instance, most national radio-stations in Norway are being streamed from the servers of the National Broadcasting Company (NRK) [44]. They also offer the possibility to watch recently transmitted TV programs as a streamed service. The fact that the data is streamed implies that the receiver will require a certain quality to be satisfied. If the stream halts for too long, the presentation of the video will stop, and the customer will probably lose interest. Other multimedia content, like teleconferencing, is even more sensitive to changes in the quality of the data

transmitted. If the different contributions to a teleconference are delivered at the wrong time, the service will be pointless.

One area of multimedia streaming that is demanding in terms of data rates is the streaming of high quality video. The servers that have to deliver such streams to many concurrent end users will be strained to meet demand unless some techniques are applied to enhance server performance. A multimedia streaming server will not only have to fetch the data, and prepare it for streaming, it will have to handle the transmission of a huge number of network packets. There are many strategies that can be applied to handle these challenges. One of these is to build an immensely powerful server with several CPU's, a multitude of network cards and an extensive amount of memory and storage capacity. Such a single server would be very expensive, and would soon have to give up if the demand were to increase.

To solve the scalability problem, a common approach is to distribute the load on several interacting computers. Several of these approaches will be discussed in this thesis. One variant of this strategy is to deploy a multimedia server as a distributed system consisting of many interconnected computers cooperating to bring the media data to the users. Such a solution would be scalable by adding more computers to the system, and dividing the load between the nodes. Data would have to be collected from the correct node in the interconnection topology, and routed to the egress. The successful interconnection would, therefore, demand resources in terms of routing and other network processing. With the servers busy pulling data from the disks, and preparing it for streaming, reserving resources to process the routing will be a challenge.

One way to ensure effective network processing inside such a topology would be to off-load such tasks to a different functional unit. Some commercial actors have done this by implementing routing functionality on custom hardware. The development of hardware, however is an expensive and time-consuming task. On the other hand, network processors optimized for the efficient handling of packets are available today. Although most commonly used in routers and switches, implementations of network processors that can be used in PCs are available. These platforms are highly programmable and configurable to different networking tasks. By building a distributed topology multimedia server, and off-loading internal routing tasks to network processors, it would be possible to achieve high performance in a scalable system without having the disadvantages of having to build custom hardware.

1.2 Thesis domain

One commercially available server that has been successful at implementing a distributed server topology is the nCube [42] n4x solution. This is, however, a system that depends

heavily on custom hardware for internal routing, making it expensive. The design of special hardware is a task that is time consuming. This makes the time-to-market for such products longer than for programmable solutions. Finding a way to implement a similar server topology without having to use custom hardware would enable much cheaper server solutions with the kind of scalability and performance that is needed to meet competition.

This thesis aims at answering the question of whether a distributed multimedia server cube, similar to the n4x solution, can be implemented using IXP2400 [32] network processors. In pursuing this goal, different areas of the network processor programmability have been explored to find the most effective ways of implementing various applications. We also seek to explore whether offloading network tasks would enhance the performance of a host machine. An implementation of a multimedia server hypercube was made in order to identify strengths and weaknesses of such solutions, and to demonstrate that an NPU offloading of the routing framework in such a solution would work.

The thesis shows that network processors are more efficient at routing tasks than an application running on a Linux host. This proved also to be true for the host application when implemented in the kernel. Moreover, it shows that an interconnected network topology, a hypercube, using IXP2400 cards for routing was far superior to a switched network topology for delivering data packets. Tests showed that there is a noticeable performance gain by offloading network tasks from the host to the IXP2400 card. The video server cube (VS^3) system was implemented, and proved to work with the proposed routing framework. In conclusion, our prototype shows that the IXP2400 NPU efficiently can offload the host machine and provide a favourable way of implementing closely interconnected servers in a cube topology.

1.3 Document structure

In the following chapters, the requirements for successfully implementing a video server cube using IXP2400 cards will be investigated. This will be done by analyzing hardware and software capabilities on a general scale, and with regard to our implementation goal. The requirements and implementation steps taken are presented, and evaluations of the different results are given.

The focus in chapter 2 and 3 is mainly on general aspects of the network processors and the applications and tests made to gain the necessary understanding for programming the IXP2400 platform. Chapter 4 and 5 focus on the requirements for multimedia applications, and the building of the VS^3 server cube application in particular. An evaluation of the results is given for all stages of implementation.

Chapter 2 describes the concept of network processors with the main focus on the Intel internet exchange architecture (IXA). An introduction is given to the hardware boards (Radisys ENP2611 [13]) used in the thesis, and the network processing units that resides on them. The software tools provided by the manufacturers is described, as well as the general programming techniques and strategies.

Chapter 3 presents the applications developed to explore and evaluate the ENP2611 network processor card based on the IXP2400 chipset. We describe the different software and hardware building blocks, and how they perform given different tasks. The results of tests performed on these applications is also presented, and some general conclusions are drawn on how to use these NPUs in general.

Chapter 4 gives an introduction to multimedia systems. Different multimedia applications are presented. The requirements that have to be met when dealing with multimedia applications is discussed. Different solutions for meeting the requirements is also presented. The nCube n4x solution to these challenges is described, as well as some of the protocols that can be used for multimedia streaming.

Chapter 5 describes all the stages in implementing the VS^3 video server cube solution. Different design approaches are presented, and the routing framework is described in detail. The main bottlenecks of this system are located, and the results of tests performed on the system are shown and discussed.

Chapter 6 concludes this thesis, and outlines further work.

The source code for the VS^3 server cube system is included as an appendix. The appendix is divided into subsections based on whether the code is written for μ Engines, XScale or Linux host.

Chapter 2

Internet exchange architecture

Network processing Units (NPUs) are special processor architectures used for demanding networking tasks such as backbone routing and switching. This chapter will give a short introduction to NPUs, and especially the internet exchange architecture (IXA) platforms. We will focus on the two network boards that have been available for testing, and especially the hardware and software for the ENP2611 [13], based on the IXP2400 [32] chipset, upon which this thesis bases its studies. The reason why the ENP2505 [12], based on the IXP1200 [31] chipset, is described as well, is that this was the NPU board we used before receiving the ENP2611. Some of the areas which have been explored have also been tested on this platform. The comparison of performance on different levels of the architectures, and the challenges of implementing on each platform have also been important issues.

2.1 Network Processing Units

After 40 years of Moore's law [18], we still see processing power expand at an exponential rate. The successful use of multithreading in the last generations of processors have made the architectures more efficient. This increases the heat produced by the logic units [41]. This has made the silicon manufacturers think about new ways of continuing the current progress. The picture we have been shown of future architectures outlines multiple core processors with simpler structure, thus using more parallel computations.

In a similar manner to how Moore's law applies to microprocessors, we have had an exponential growth in network bandwidth capacity. Figure 2.1 shows how networking bandwidth has increased compared to processing power over the last years. There is no indication of the development rate slowing down, but it will eventually have to halt due to limitations of the transmission medium (if new technologies are not introduced). As of

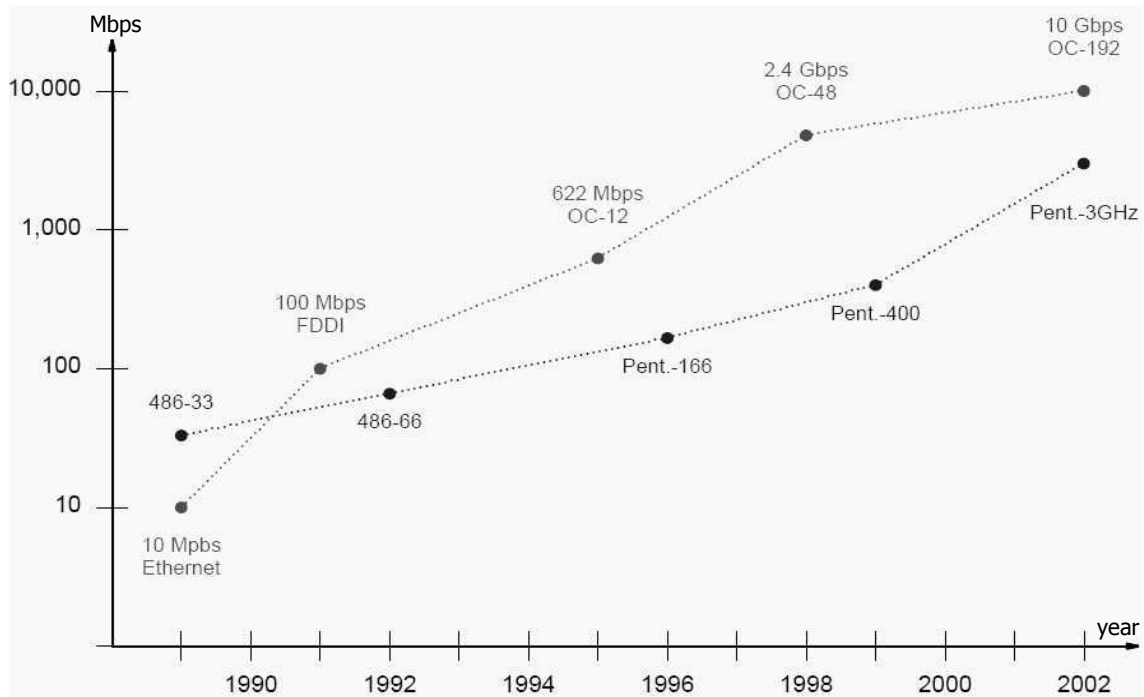


Figure 2.1: Evolution of technologies [22].

yet, the network speed growth is faster than the CPU power growth. Thus today's possible bandwidths sustain an increasing demand for processing power. With the decreasing inter-arrival time of network packets, computers will be strained to keep up with protocol handling, checksumming and other necessary tasks without sacrificing capacity that other processes need.

One solution to this challenge is to offload network tasks to separate functional units. These units will be optimized for efficient packet handling and throughput. Since NPUs are to be optimized for packet handling, the design will differ from traditional computing chipsets. A typical design is a series of several small, symmetric processing units working in parallel. The parallel structure of the packet processing enables the tasks to be performed in a pipeline, with each functional unit performing a special task. This removes the bottleneck of single CPU processing that has to schedule the processing capacity between the different tasks.

A wide range of different companies are manufacturing NPUs for different platforms and purposes, these include Agere [1], AMCC [3], IBM [25], Intel [27], Internet Machines [26], Motorola [38], PMC-Sierra [48], and Vitesse [61]. All these are more or less based on the same offloading ideas, but their implementations vary greatly.

In the past few years, in particular, Intel has focused on network processing, creating a series of network processor platforms called the Internet Exchange Architecture (IXA). This platform group includes three main series of NPUs, namely the IXP12xx, IXP2xxx and IXP4xx series of network processors [43].

2.2 Intel Internet Exchange Processors

The IXA platform from Intel includes, as of now, three main series: IXP12xx, IXP2xxx and IXP4xx. The IXP12xx series was in common use up to last year, when Intel decommissioned it in favour of the IXP2xxx series. The IXP4xx series is targeted at the home and small business market, while the IXP2xxx series aims at the heavier applications. In this section, the main features of the IXP1200 will be outlined. A more detailed description will be given of the IXP2400 as this is the chipset on which most work has been done in this thesis.

2.2.1 IXP1200

The Intel IXP1200 [31] network processor chipset is designed to meet the wide requirements placed on network equipment in high performance systems and consists of six main functional units as shown in figure 2.2. The core processing unit is a 32 bit StrongARM processor running at 232MHz. There are six special purpose microprocessors called microengines (μ Engines) running in parallel, also at 232MHz. Each μ Engine can accommodate a maximum of four contexts. The DRAM unit provides access to memory for storing packet data, and the SRAM unit gives access to memory for storing shared variables and metadata. The IX bus connects the internal IXP chipset devices, and the PCI bus enables the card to interface with other PCI devices. In addition to these six main components there are 4KB of scratchpad memory used for high-speed communication between functional units. There are also special registers for inter-unit communication.

2.2.2 ENP2505

The Radisys ENP2505 (see figure 2.3) integrates the IXP1200 chipset on a network board. It has four 10/100 Mbit Ethernet ports for communication. For boot code and other stored procedures, there are 8MB of Flash memory. The card has a total of 8MB of SRAM and 256MB of DRAM. In addition there are a serial port for interfacing and debugging, and a PCI connector.

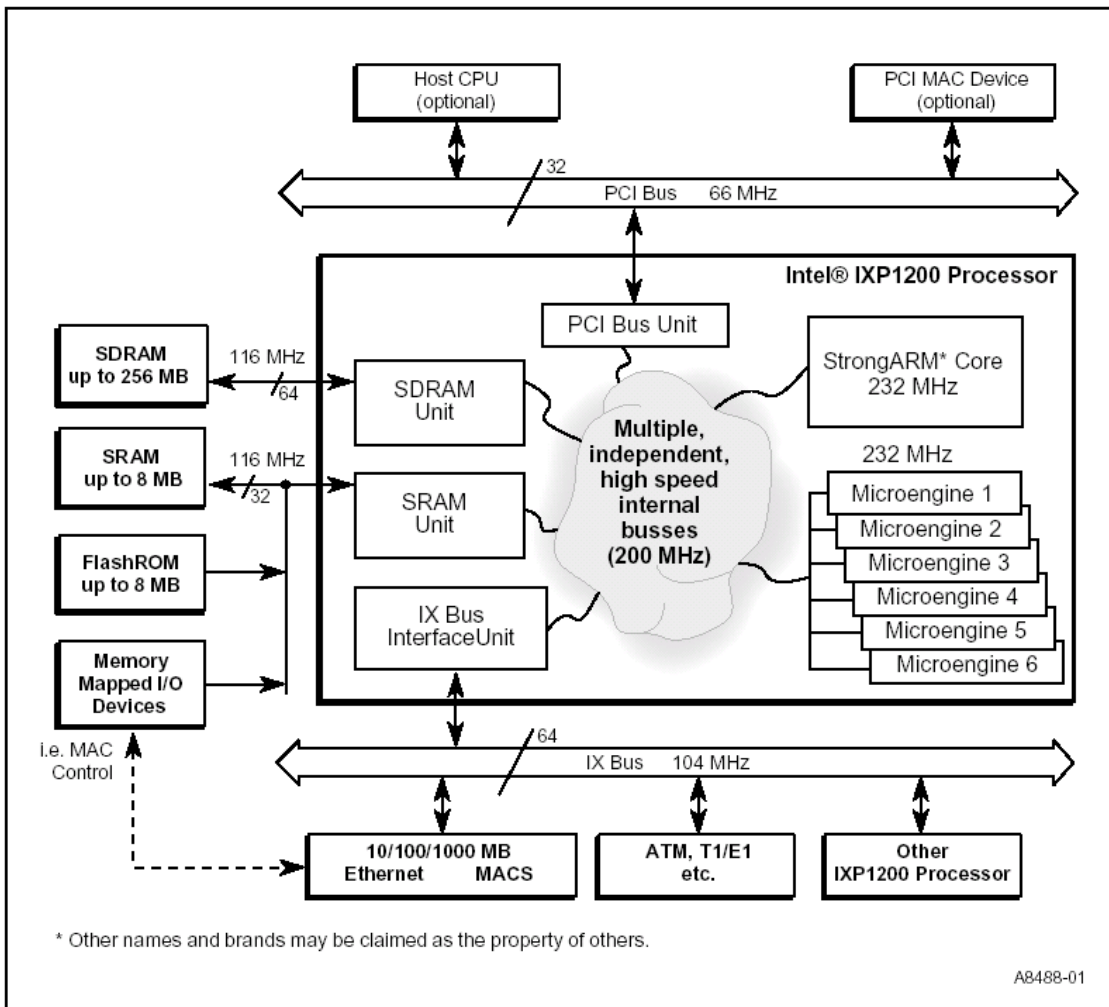


Figure 2.2: IXP1200 block diagram [31].

2.2.3 IXP2400

The IXP2400 chipset [32] is a second generation NPU platform from Intel. It is designed to handle a wide range of access, edge and core applications. It has a more powerful CPU and microengines than the IXP1200 and is better suited for heavy networking tasks. The physical interfaces are customizable and can be chosen by the manufacturer of the device on which the IXP chipset is integrated. The number of network ports and the network port type are also customizable. The major functional blocks of the IXP2400 chipset are shown in figure 2.4:

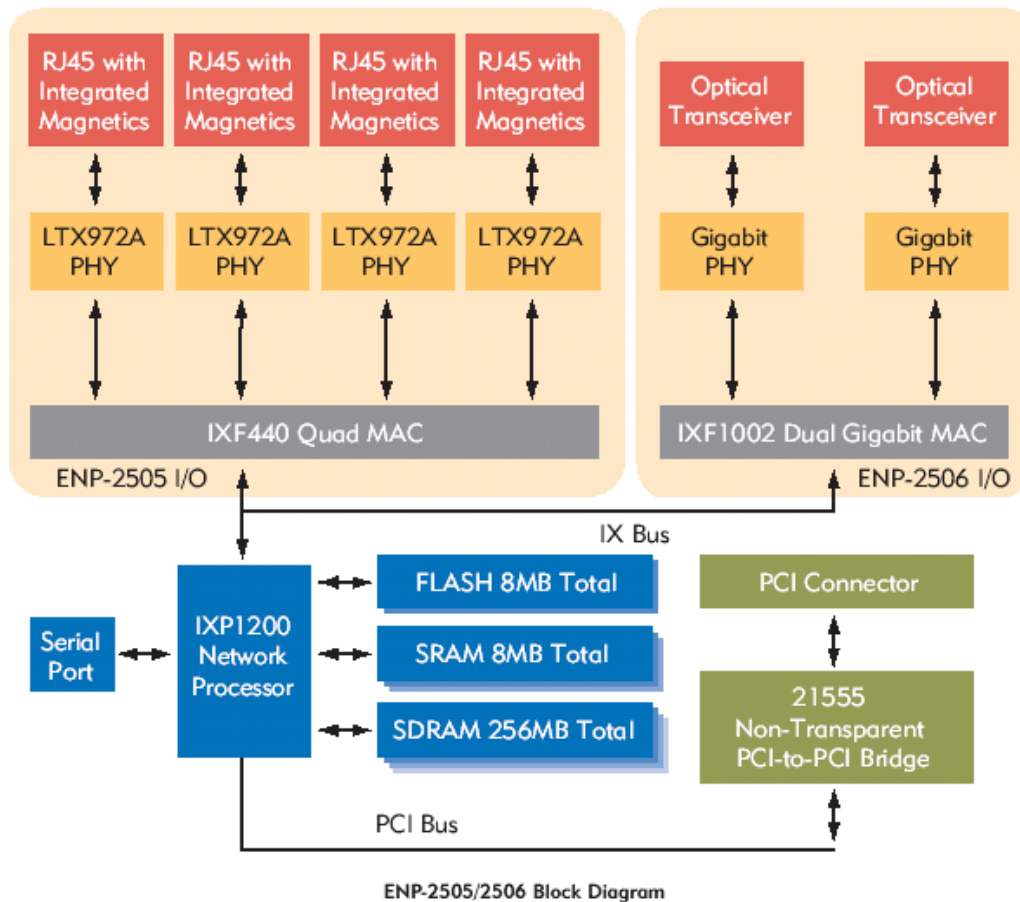


Figure 2.3: ENP2505 block diagram [12].

- **An Intel 600MHz XScale core:** The XScale core is capable of running an independent operating system (in our case, MontaVista Linux for embedded platforms [37]). This unit is used to initialize and manage the chip. In a network application, it typically controls some of the higher layer network processing tasks like updating IP routing tables.
- **8 600MHz μ Engines:** For efficient handling of general packet processing, the eight μ Engines can be used. These are separate 32 bit programmable units that are specialized for network processing. Each μ Engine has a maximum of eight threads (contexts) that can enhance performance further.
- **Two independant SRAM controllers:** The two SRAM controllers can independently access one SRAM channel each. This type of memory is generally used for

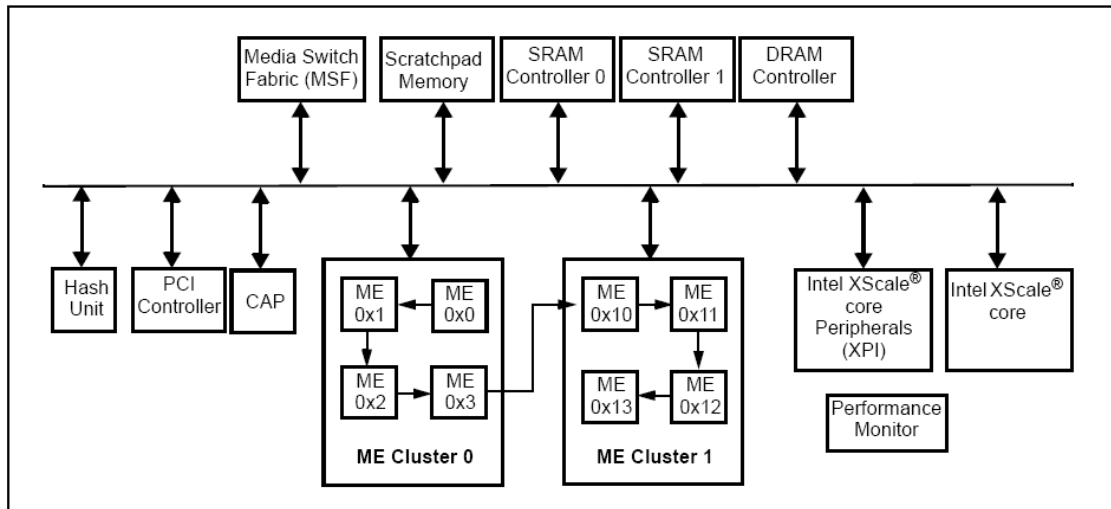


Figure 2.4: IXP2400 block diagram [32].

packet metadata, control data, shared counters and variables. Atomic control operations are supported for synchronization purposes.

- **One DDR DRAM controller:** DRAM memory is generally used to store packet data. When a packet is received, it can be stored here, and all steps of the packet handling pipeline can modify the packet from the same memory area.
- **Media switch fabric (MSF):** The MSF is the IXP interface to the framing hardware. This is where connections with the physical layer is configured. It contains transmit and receive buffers, and packets that are about to be transmitted are divided into smaller MPackets¹ upon transmission to make the IXP chipset compatible with as many different types of media hardware as possible.
- **PCI 2.2 compliant controller:** In order to connect with other components, like a host machine or PCI compliant peripheral devices, the IXP has a PCI controller.
- **Scratchpad, hash and CAP (SHaC) unit:** Three of the most useful functions of the IXP chipset resides in the SHaC unit:
 - **Scratchpad memory - 16KB:** The scratchpad memory is a 16 KB storage for general purpose use with atomic operations and ring support. This is widely used to convey packet buffer handles between different processing units.

¹To make the IXP system compatible with different physical interface standards, packets are divided into a basic unit called mpackets before they are sent to the interface controllers via the MSF. On the IXP2400 mpackets can be 64, 128 or 256 bytes, but once configured, the mpacket size must remain constant [6].

- **Hash unit:** The hash unit can be used to achieve hardware supported hash calculations, thus giving increased performance when making lookups in tables etc.
- **Chip-wide control and status registers (CAP):** The CAP unit is used to handle inter-processor communications.
- **XScale Core peripherals:** These include an interrupt controller, four timers, one serial UART port, 8 general purpose input/output circuits (GPIOs), and an interface for low-speed off-chip peripherals. In addition, the IXP chipset has a performance monitor with registers for analyzing and tuning performance.

The sum of these components is the IXP2400 chipset. In the hardware board that is available to us, this chipset resides on a network card made by Radisys [51].

2.2.4 Radisys ENP2611

The Radisys ENP2611 [13] is a network board that integrates the IXP2400 chipset, peripherals, memory and physical interfaces. Figure 2.5 gives a schematical representation of the hardware. The following main components are present:

- **256MB (DDR) DRAM:** The DRAM is, as mentioned in section 2.2.3, primarily used to store packet data.
- **8MB SRAM:** The SRAM is accessible via two different channels to optimize performance. It is generally used for metadata and shared variables.
- **16MB StrataFlash memory:** As on the ENP2505 this memory keeps the boot code and utilities for the board. The increased size compared to the ENP2505 allows for a Linux kernel to be loaded into the flash memory.
- **Three gigabit ethernet interfaces:** The physical communication with the network is made possible by the optical transceivers that are controlled by the PM3386 (controls two optical interfaces) and PM3387 (controls one optical interface) Gigabit Ethernet controllers. These give the board a total of 3 optical interfaces that can be used freely by the IXP hardware.
- **SCSI parallel interface v3 (SPI-3) bridge FPGA:** The SPI-3 is used for connection between the physical interfaces and the IXP2400 MSF. It is the link between the PM3386 and PM3387 controllers and the IXP2400.

- **Two PCI-PCI bridges:** The first is a non-transparent Intel 21555 PCI-bridge chip which connects the internal 64-bit PCI bus to the backplane host 32 or 64 bit PCI bus. This gives the hardware the possibility to provide inter-processor communication and interrupts. The second is a TI PCI2150 transparent PCI bridge which links the internal 64-bit PCI bus to a downstream 32-bit PCI bus. This is used to connect the debug 10/100 interface to the rest of the chipset.
- **Ethernet controller:** An Intel 82559 PCI Ethernet controller is used to control the 10/100 ethernet debug interface.
- **SPI-3 Option Board connector:** This can be used to connect the chipset to a service specific NPU co-processor that can be used to provide further hardware support for specific tasks.
- **10/100 Ethernet interface:** This debug port can not be used as a part of the IXP topology, but is rather used as a tool for loading images into flash ROM and mounting NFS filesystems by the Linux kernel running on the XScale core.
- **Clock generation circuitry:** The clock generation circuitry includes the general IXP2400 system clock, and the interface clocks for the IXP2400 MSF/FPGA and the FPGA/PM338x interfaces.
- **Reset and initialization circuitry:** This circuitry connects to a switch located on the board to enable a manual reset. Software resets can be initialized both from the board and via the PCI bridge.
- **Power:** Power supply circuitry is needed for the different logical parts of the chipset.

These components are the external hardware needed to complete a functional IXP chipset environment. The hardware manufacturers also provide some software tools to make developing applications for the platform easier.

2.2.5 New hardware

After some time working on the ENP2505 [12], it became clear that Intel intended to discontinue support of the IXP1200 [31] chipset, and also the IXA SDK 2.01 [29], which was the latest SDK version available for the IXP1200 chipset. At about the same time, we received the ENP2611 [13] cards. The decision was then made to switch to the ENP2611, in the hope that the new Intel SDK 3.51 [30] would be a more stable and thorough release than the 2.01 which had proved to be especially challenging to make work smoothly. This,

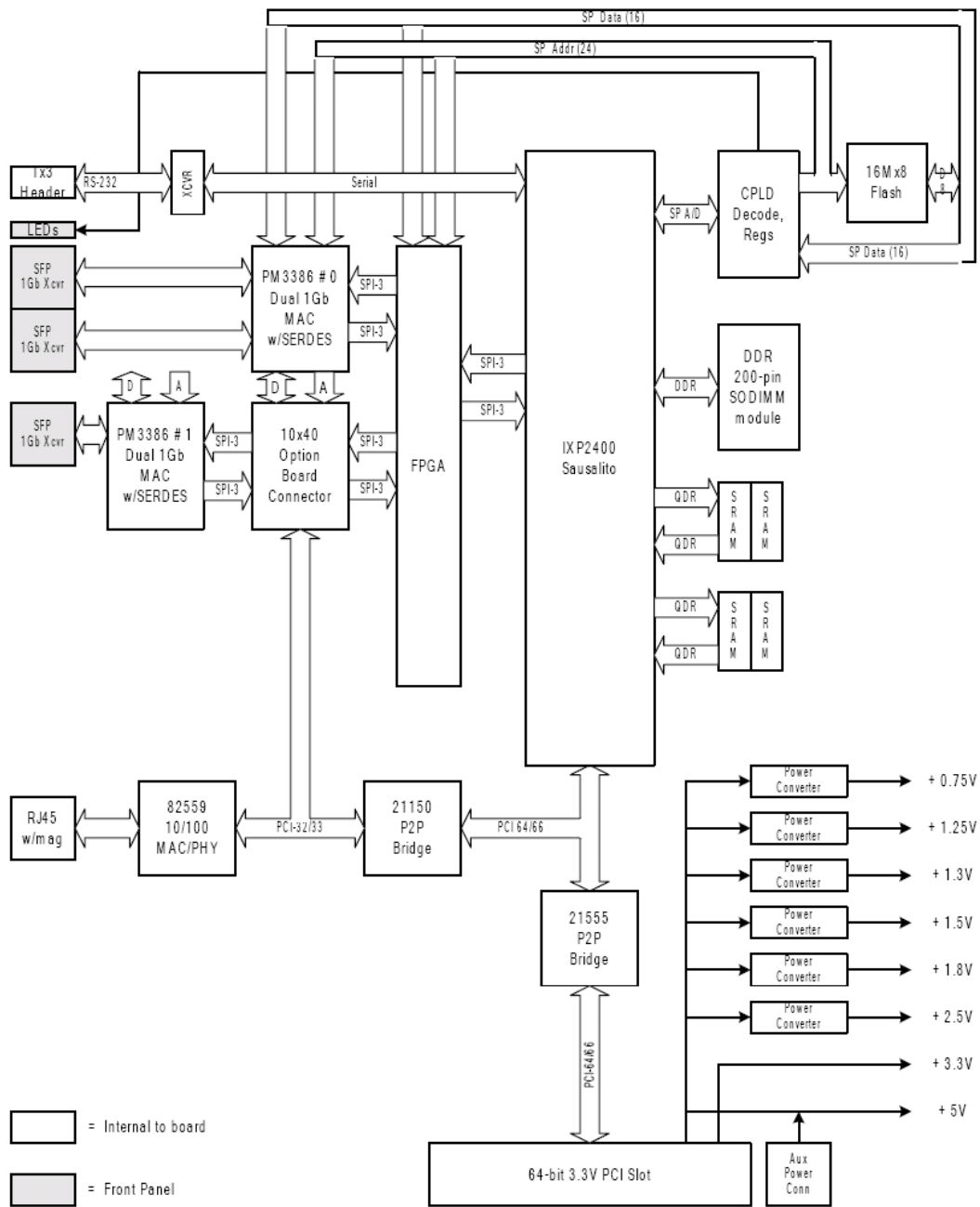


Figure 2.5: ENP2611 Block diagram [13]

combined with the fact that developing applications for a platform that is no longer in use would be of little future value, made it feasible to commence with the hardware change².

Component	ENP2505	ENP2611
Core	StrongARM	XScale
Core clock frequency	232MHz	600MHz
Number of μ Engines	6	8
μ Engine clock frequency	232MHz	600MHz
Threads per μ Engine	4	8
Scratch memory	4KB	16KB
Number of network interfaces	4	3
Type of network interfaces	10/100 Base-T	Gigabit Optical

Table 2.1: Main differences between ENP2505 and ENP2611

The most significant changes in hardware between the two platforms are shown in table 2.1. The structural changes of the Intel IXA SDK versions 2.01 and 3.51 are so significant that porting IXP1200 applications to IXP2400 would have taken excessive work. This task would be further complicated by the fact that the different hardware peripherals would demand different drivers and configuration. The consequence of this was that the applications had to be implemented from scratch on the new platform to make it conform to the new SDK and hardware.

2.3 IXA software libraries and tools

The Intel IXA platform comes with software and libraries to support the programming task. The main bundle is the Intel IXA software development kit (IXA SDK) that provides debug tools, compilers and support libraries. There is also an SDK for the board implementation that contains sample code and board specific drivers. This section will give an introduction to these libraries.

2.3.1 Intel IXA SDK 3.51

The Intel IXA SDK [30] is a set of libraries and reference designs provided by Intel to support programming the IXP chipsets. It also includes some debug tools that interface

²This was considered to be profitable even though it meant starting almost from scratch with new hardware and a new microengine programming language (micro-C). The setup and configuration process would also have to be figured out.

with the card, and ease debugging of μ Engine applications³. The main components of the SDK are as follows:

- **Development Workbench:** The development workbench is a Windows application that can be used to simulate the running of μ Engine programs. It can also be used to debug the program directly on the hardware by means of the 10/100 network debug port. Since μ Engines cannot directly provide output to screen, this is a tool that is invaluable when it comes to verifying that your program performs as it should. It also has the possibility of simulating a packet stream, so as to see how the application responds to different input. It allows for analysis of the amount of cycles different operations consumes, and to use this to eliminate performance bottlenecks.
- **Software libraries for μ Engines and XScale:** Intel's software libraries, both for μ Engines and XScale, provides means of executing operations and intrinsics, saving implementation time on several tasks by reducing the necessary number of code lines. There are also libraries for debugging via the XScale, and hardware abstraction libraries for accessing μ Engine functionality from the XScale.
- **Compilers, assemblers and linkers for μ Engines:** In the IXA SDK 3.5, Intel has added Linux support for the μ Engine C language, and provided enhanced compilers and assemblers. This makes development on an all-Linux environment easier than in was using previous editions like the IXP1200 platform [31]⁴.
- **Software framework and sample applications:** The software framework is a collection of libraries and sample applications that gives an insight into how different networking functionality can be implemented using standard Intel tools.

There are several layers of abstraction provided with the SDK. On the bottom level microcode instructions can be used to manually set all needed registers. Common instruction combinations are gathered in the intrinsics library for the μ Engines. The micro-C language provide an even greater level of abstraction by eliminating the need for assembly-style programming.

When programming for both XScale and μ Engines, the hardware abstraction layer (HAL) and operating system services layer (OSSL) libraries provide access to microengine functionality from the XScale. For larger scale programming there are a software framework that manages resources and enables structures and methods for acheiving common networking tasks.

³Debug printouts and similar techniques are difficult due to the fact that the μ Engines do not allow any screen output. All debugging of this kind will have to go via the XScale.

⁴This was previously only supported for Windows.

We have experienced that the abstractions provided by the software framework often are small, and therefore it may, for some tasks, be better to use the hardware abstraction libraries upon which the framework is built. This is particularly relevant for custom tasks that does not conform with standard network application procedures.

The datatypes used when programming the platform are common for both μ Engines and XScale. The byte order of the datatypes is big-endian. The most common datatypes are shown in table 2.2.

Term	Words	Bytes	Bits
Byte	$\frac{1}{2}$	1	8
Word	1	2	16
Longword	2	4	32
Quadword	4	8	64

Table 2.2: IXP2400 data terminology [32].

The IXA SDK is the main tool of development when writing software applications for the IXP2400 platform. Though the software framework is superfluous in many cases, the basic libraries are very useful for development on all layers of the IXP2400 hardware.

2.3.2 Radisys ENP SDK 3.5

Radisys supplies an SDK [33] which provides a range of services that are specific to the ENP2611 board. These services include:

- Drivers for the PM3386 and PM3387 optical interface controllers and for the SPI-3 bridge.
- Sample application code with microblocks for transmitting and receiving packets, that are adapted to the current physical interfaces.
- Additional drivers that provide connection to the host over PCI.
- Kernel image of the Monta Vista Linux [37] operating system that can be used on the XScale.

The drivers for the PM3386m, PM3387 and SPI-3 give us the opportunity of tuning parameters on the hardware devices (like turning on and off hardware ethernet checksum calculation). Prior to running an IXP application, these drivers have to be initialized.

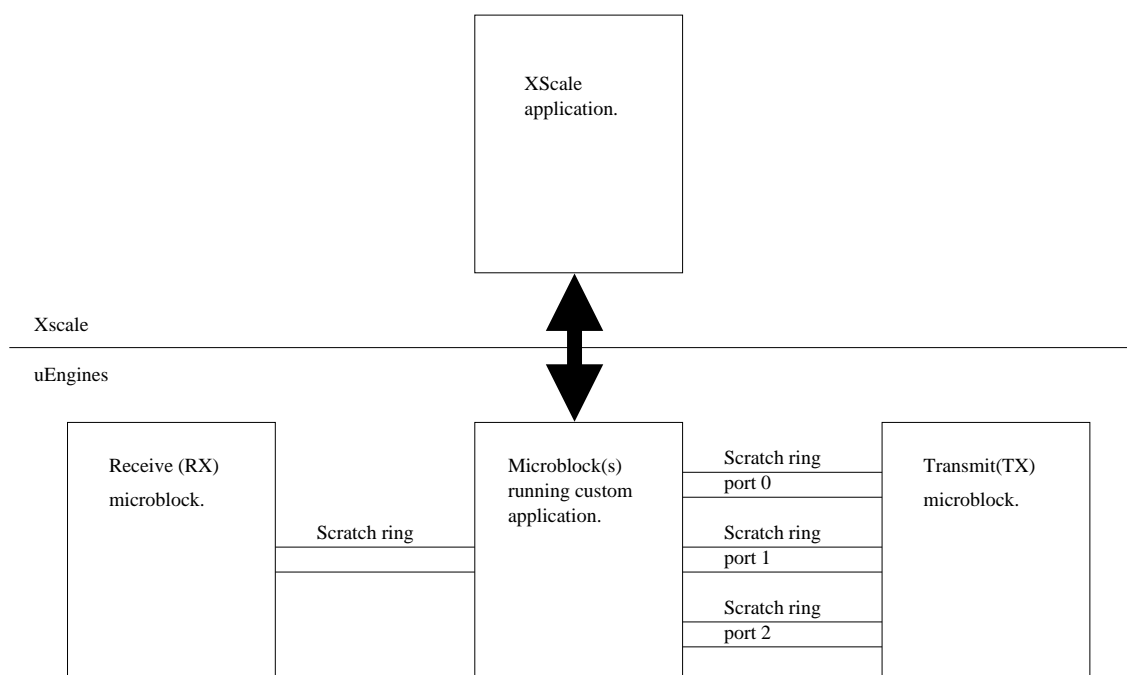


Figure 2.6: RX and TX block functionality and interface.

The microblocks provided have been an important tool. They are in principle the same as the Intel IXA SDK transmit and receive microblock reference designs, but the parameters have been adapted to suit the ENP2611 interfaces. Thus, we can integrate parts of this code to provide us with suitable transmit and receive blocks for our purposes. Figure 2.6 shows a typical application setup. The Receive (RX) microblock receives data from all ports, and passes a reference to the data on one scratch ring (explained later in section 2.3.3). An application on μ Engines or XScale can then process the data. To send data, a reference has to be put on a scratch ring. The Transmit (TX) block is designed with three input scratch rings. Which port the packet is to be transmitted on is decided by which scratch ring the reference is passed on. The drivers and board-specific SDK provided by Radisys is an extension of the more extensive SDK provided by Intel to support programming the IXP chipset itself.

2.3.3 IXP programming paradigms

To handle the large throughput that is expected from a networking application, a method that has been proven effective is to give the application a hierarchical structure [6]. Figure 2.7 shows a diagram of the levels of structure. In practice, this means that most of the

packets arriving have to be handled by the lower processing levels (i.e., the μ Engines). Some packets can be sent to higher levels (XScale) for further processing, and a few packets can be forwarded further (for instance to the host machine). The reason for this is that the performance of the lowest level is highly optimized for effective forwarding. As you look at higher levels, you will find more general purpose processing units that can handle more diverse tasks at a penalty of throughput.

The operational layers are commonly divided in two. The “data plane” handles the high speed processing and forwarding of the majority of network packets. The key elements are real time forwarding and efficiency [41]. The “control plane” involves non-wire-speed general purpose processing that can include table creation and updating, data plane exceptions or computationally intensive tasks [41]. In an IXP2400 setting the data plane typically resides on the μ Engines, and the control plane is represented by the XScale.

Both hardware and software design reflects this design philosophy, and as a result, trying to move large datarates between XScale and μ Engines will yield poor results⁵. To ensure efficient packet handling, steps must be taken upon implementing applications, so that the majority of tasks are performed on the data-plane (if the tasks are not so processor-intensive that they hamper the data pipeline throughput).

Packet handling As the goal of the lower layers of the network processor is to handle the packets as swiftly and efficiently as possible, it is essential to avoid unnecessary copy operations. Once the packet has been written to DRAM, it resides in the same memory buffer until packet processing is finished, and the packet is transmitted or discarded. What passes between the different processing units is a longword with data (buffer handle) from which the packet location and metadata can be derived.

In the reference design from Intel, the RX and TX microblocks use buffer handles as the means to convey the reference to a packet between processing units. From each buffer handle, the location of two related buffers can be extracted:

- **Packet metadata:** The metadata contains information on the packet size, the port that received the packet, the port that is to transmit the packet, the linking of buffers into larger packets etc (see figure 2.8). This data resides in SRAM. (More details can be found in the IXA framework reference manual [28].)
- **Packet data:** Located in DRAM, this is the buffer where the packet data is actually stored. The RX block reserves some space in front of the packet and after the packet ends in case additional headers have to be appended or prepended.

⁵We have achieved datarates of up to 100 Mbps between XScale and μ Engines in experiments using batch queuing of packets on a scratch ring.

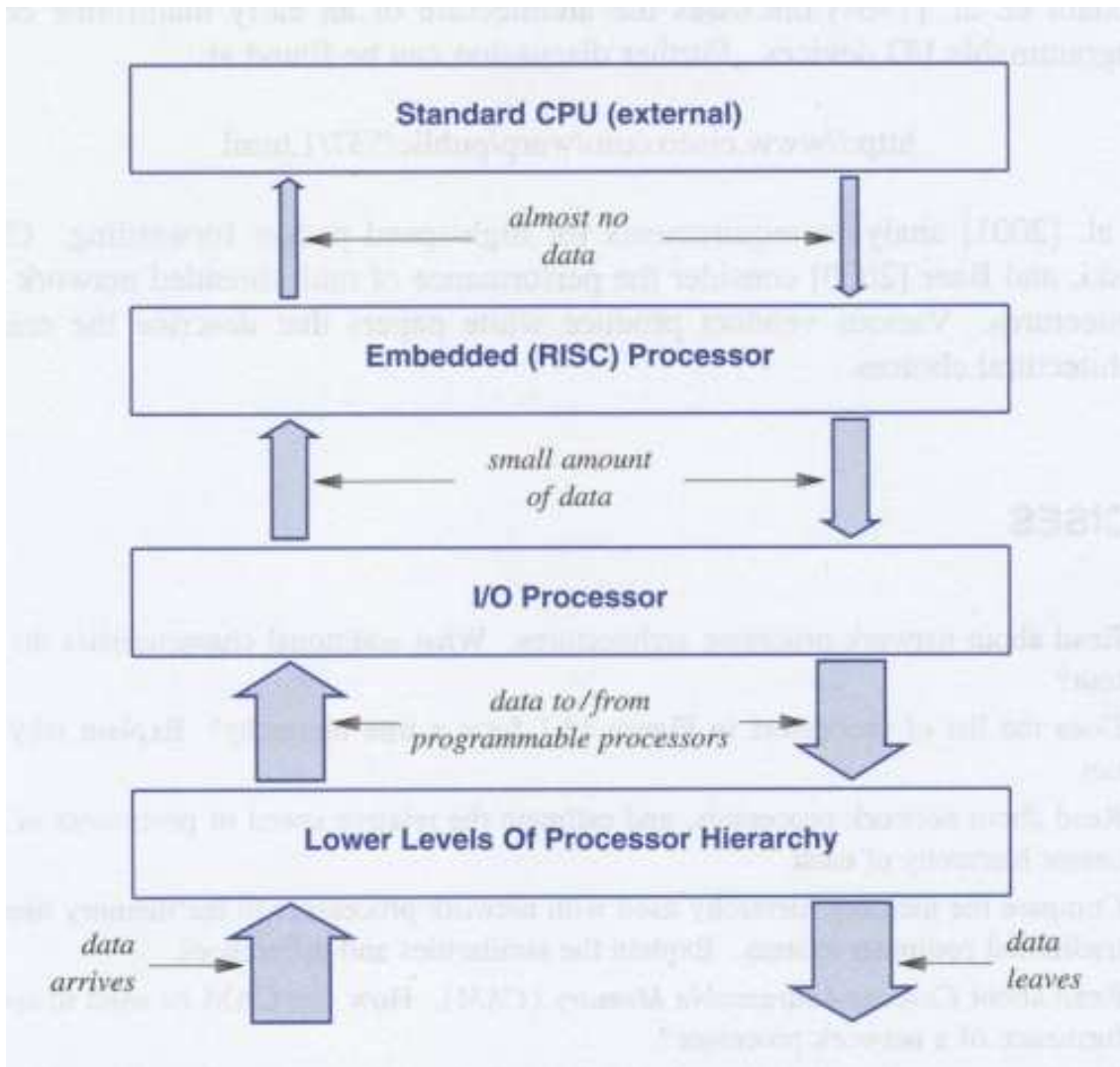


Figure 2.7: IXA Hierarchy model [6]

The metadata buffers, although customizable in size, are typically 32 bytes long. The data buffers are defined by the need for space (2048 bytes would be an appropriate size if you want to store ethernet packets with an maximum transmission unit (MTU) of 1500 bytes). To be able to use the current buffer handle structure, the size of both the metadata buffer and the data buffer have to be power of two. Figure 2.9 give an example of how metadata and data offsets are extracted from the buffer handle in a typical data structure. The SRAM and DRAM offsets calculated from the buffer handle give the starting position of the respective buffer. The packet data in DRAM, however, starts a number of bytes into

```

typedef __declspec(packed) union {
    struct {
        dl_buf_handle_t buffer_next;

        uint16_t buffer_size;
        uint16_t offset;

        uint32_t packet_size : 16;
        uint32_t free_list_id : 4;
        uint32_t rx_stat : 4;
        uint32_t header_type : 8;

        uint16_t input_port;
        uint16_t output_port;

        uint32_t next_hop_id : 16;
        uint32_t fabric_port : 8;
        uint32_t reserved : 4;
        uint32_t nhid_type : 4;

        uint32_t color_id : 4;
        uint32_t reserved_1 : 4;
        uint32_t flow_id : 24;

        uint16_t class_id;
        uint16_t reserved_2;

        uint32_t packet_next;
    };
    uint32_t value[8]; /* aggregate for the above fields */
} dl_meta_t;

```

Figure 2.8: Struct describing packet metadata.

the buffer. This is to leave room in front of the packet in case we need to prepend the packet data with a new header. The offset where the packet data begins can be found in the metadata (see the “offset” field in figure 2.8). The 24 least significant bits provides us with the offset of metadata in SRAM and to the data in DRAM. The two most significant bits indicate if the buffer is start of packet (SOP), end of packet (EOP) or both. This is to make sure that packet sizes larger than the internal buffer size can be accomodated. Bits

SRAM metadata buffer size: 32 Bytes.
 DRAM data buffer size: 2048 Bytes.

New buffer handle = old buffer handle + 8.

Buffer handle 1: 0x2000

SRAM offset = $0x2000 \ll 2 = 0x8000$
 DRAM offset = $0x2000 \ll 8 = 0x200000$

Buffer handle 2: 0x2008

SRAM offset = $0x2008 \ll 2 = 0x8020$
 DRAM offset = $0x2008 \ll 8 = 0x200800$

Buffer handle 3: 0x2010

SRAM offset = $0x2010 \ll 2 = 0x8040$
 DRAM offset = $0x2010 \ll 8 = 0x201000$

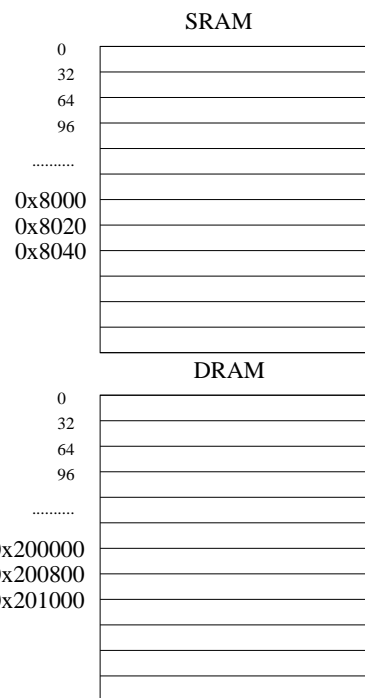


Figure 2.9: Calculation of metadata and data offsets from buffer handle values.

24 through 30 state the cell count, indicating how many cells the buffer contains⁶.

Scrath ring mechanics The SHaC unit of the IXP chipset contains 16KB scratchpad memory that is accessible to the XScale and all μ Engines. Listed below are its most important properties:

- Normal read and write: The memory is accessible on 32 bit boundaries. That implies that you cannot read or write less than 32 bits at a time. You can read or write up to 16 longwords with a single command.
- Atomic read-modify-write operations: You can set or clear bits, increment or decrement, add or subtract in an atomic operation. These operations can also return the pre-modified value of the written data.
- Sixteen hardware assisted rings for interprocess communication: These rings are implemented as FIFO-queues with a head and tail pointer.

⁶Cells are used for Asynchronous Transfer Mode (ATM) processing. When not in ATM mode, bits 24 through 30 can be used for other purposes.

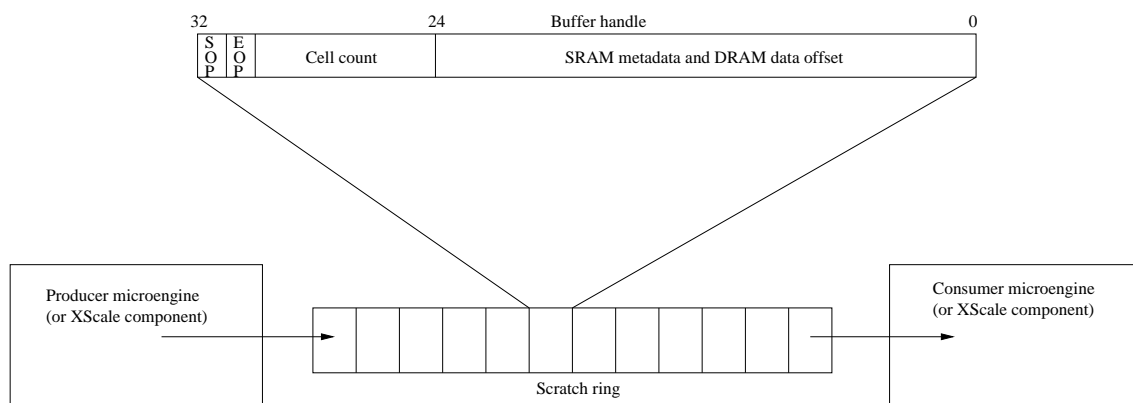


Figure 2.10: How buffer handles are passed on scratch rings between processing components.

The most important generic use of scratch memory is as scratch rings. Given the buffer handle structure described above, you can assign a scratch ring as “egress” for one μ Engine and as “ingress” for another one (or the XScale core). It is, in other words, used to facilitate message passing between different structural blocks. Since the location of the packet metadata can be derived from the buffer handle, all necessary information can be supplied by passing this one buffer handle on a scratch ring. Figure 2.10 shows the use of a scratch ring to move buffer handles from a producer component to a consumer component. For optimization purposes, however, it is possible to pass more information about a packet or buffer on the ring. This can save SRAM accesses in cases where every cycle is valuable.

2.3.4 MontaVista Linux

The final piece of software needed to use the ENP2611 as an operational network platform is an operating system running on the XScale core. For the standard ENP SDK, this OS is a small Linux distribution optimized for embedded platforms. The kernel image is developed by MontaVista [37], and provided with the Radisys SDK. The MontaVista preview kit for ENP2611 is, as the name suggests, a limited edition for evaluation purposes. It is based on the 2.4.14 Linux kernel. The hardware drivers provided by Radisys and Intel extend the distribution functionality, and give the programmer full control of the IXP hardware. There is, however, an initiative in the open source community to provide up to date Linux kernels for this platform [34].

2.4 Summary

There is a wide range of possible uses for NPUs. Though their use, up to now, mostly has been in free-standing embedded devices like routers and managed switches, network boards are being developed that can add NPUs to conventional machines. In environments that have need for heavy network processing, these can provide valuable offloading functionality.

The parallel processing done by the IXP chipsets provides effective packet processing pipelines. More complicated tasks, like managing data structures and handling exceptions, can be forwarded to the core CPU, though most packets should be handled on the data plane. The extensive programmability of the IXP platform give us the opportunity of building custom functionality into μ Engines and core CPU and adapting the NPU functionality to our needs.

There are many aspects of NPU functionality that have to be explored in order to implement an application which uses the host, the NPU core and NPU data plane. In the next chapter, we will explore some of these key features on the ENP2611 to try to evaluate the performance of the hardware, and how such an application best can be implemented.

Chapter 3

IXP2400 Evaluation

We have, for some time now, been implementing applications on the ENP2505. The ENP2611 was radically different, both in hardware functionality and in SDK structure and use. To be able to efficiently program the IXP2400 hardware platform, and learn the basics of configuring and using the hardware, it was necessary to start with small applications and expand the use as the needed knowledge was acquired. With a hardware platform so different from the programming environments that is common knowledge, acquiring the information needed and successfully applying it was a great challenge. This chapter describes the most important applications developed in this process, and some measurements that have been valuable, both for designing the video cube solution, and for planning IXP2400 application development in general.

3.1 Exploring the IXP2400 hardware

This section describes the applications developed to explore the IXP2400 platform. The first challenge was to be able to successfully transmit and receive packets. The next goal was to modify packet data, and, finally, to do some packet processing on the XScale. The ability to move data efficiently between the host and the IXP was also an area in question (this is discussed in more detail in section 5.4). One of the main challenges was to find methods of implementation that are as simple as possible, but still effective enough to solve the task at hand. This was emphasized by the multitude of different approaches that, judging by the sample code applications provided in the IXA SDK, all would achieve the same goal.

After reviewing the domain of the server cube application, it became clear that the implementation would need hardware resources on three different levels (see section 2.3.3).

The media stream was to be delivered from a server application running on the host machine. Some of the setup and routing structure maintenance had to be done on the XScale, but most of the large-scale packet processing had to be done on μ Engines. To ensure that an expedient implementation could be made on the ENP2611 with IXP2400 chipset, the workings of this hardware would have to be explored and some central questions would have to be answered:

- How can operations be implemented efficiently on the μ Engines?
- Should we use microcode or micro-C on μ Engines?
- How can data be transported to and from the XScale level?
- How large bitrates can be moved to the XScale?
- How can control-plane program functionality be implemented on the XScale?
- How can data be transported to and from the host machine?

The questions regarding the IXP platform, when not exchanging data with the host, is answered in this section, the questions that imply communication with the host machine will be further discussed in section 5.

3.2 Static forward

The key task when designing a network application is the ability to receive and transmit packets efficiently. With the SDK supplied by Radisys, a small application was included that had microblocks designed for these tasks and configured to work with the ENP2611 hardware. This is an application with a three microblock structure as shown in figure 3.1. The microblocks perform the following tasks:

- Microblock 1 - RX microblock: This microblock receives data from the medium. It then reassembles the MPackets (see section 2.2.3) into the original ethernet packets. The packet is written to DRAM, metadata is created in SRAM and the buffer handle (see section 2.3.3) is enqueued on a scratch ring.
- Microblock 2 - Packet echo: In this microblock, 5 longwords, including buffer handle and input port, are read from the RX scratch ring. The input port number is checked and the buffer handle is enqueued for sending on an outbound scratch ring (TX) based on which input port it arrived on: $0 \rightarrow 1$, $1 \rightarrow 2$, $2 \rightarrow 0$.

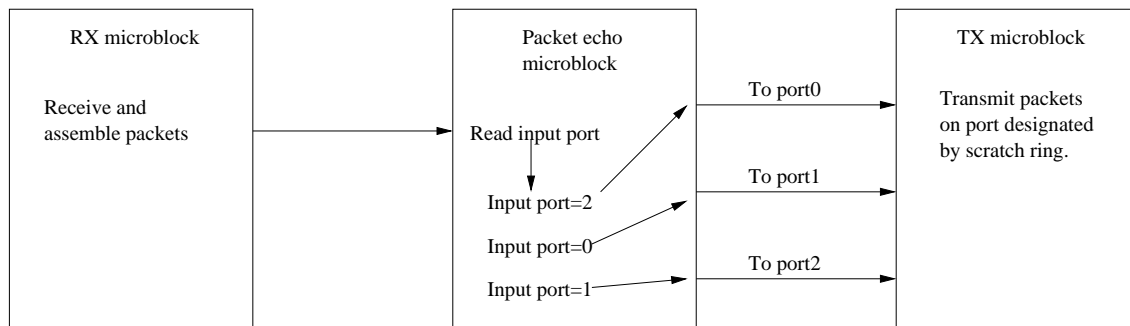


Figure 3.1: Radisys static forward application block diagram.

- **Microblock 3 - TX microblock:** The buffer handles are dequeued from three different scratch rings by this microblock. An output port is then selected based on which scratch ring the handle arrived on. The packet is split into MPackets, and transmitted on the chosen port.

The RX and TX microblocks maintain statistics on the number of received and transmitted bytes and packets on each port. These variables are shared, and can be read by the XScale application. The XScale application is also responsible for initializing the spi3br and pm3386 (see section 2.3.2) drivers, load μ Engine object file images into the μ Engines and start the μ Engines.

It is worth noting that the RX microblock puts not one, but 5 longwords on the scratch ring bound for the packet echo microblock. The first longword is the buffer handle, the four next are relevant metadata. The reason for this is probably that the receiving microblock will save at least one SRAM access by having received the packet offset on the scratch ring.

After some attempts of trying to design RX and TX code from scratch, it was decided to try to integrate the microblocks from the static forward application into a custom application by replacing the “packet echo” microblock. This proved to be a successful strategy, and enabled the building of the next application.

3.3 IP header switch application

In order to find out more about the performance of μ Engines and XScale, a simple application that had to modify packet data was implemented. The first implementation performs all task on μ Engines, the second one forwards all packets to the XScale for processing. The designs are described in the following sections.

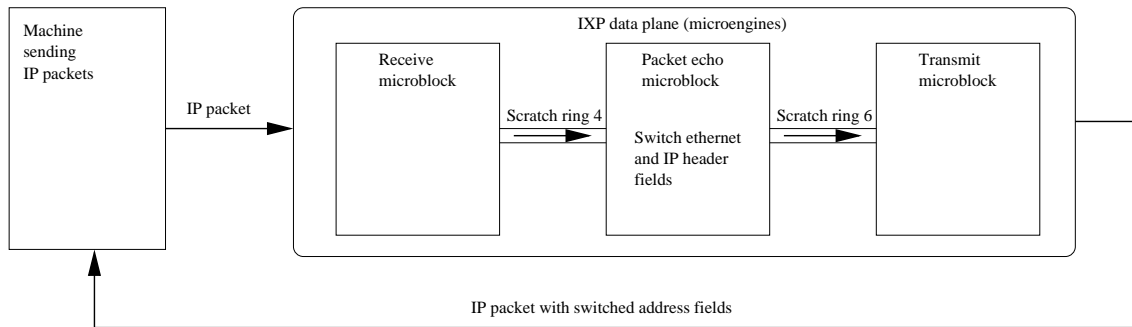


Figure 3.2: IP header switch application on μ Engines.

3.3.1 μ Engines implementation

To try to measure the time consumed by a simple application, and to answer the question of micro-C or microcode, it was decided to make a simple modification of the static forwarding application. This application also enabled us to test the mechanisms needed to change the packet data in DRAM from the μ Engine level. The following changes were implemented:

- The packet echo microblock (originally written in microcode) was replaced with a version written in micro-C.
- The packet was to be transmitted on the same port it arrived.
- The source and destination ethernet addresses were switched.
- The source and destination IP addresses were switched.

The use of micro-C proved to make the programming of μ Engines easier. The code was easier to read and survey, and a breakdown of the microcode produced by the compiler showed that the amount of cycles used was not far from the microcode version.

The greatest change from the static forward application was that this application would actually modify the incoming packets before transmitting them. The RX and TX microblocks from the static forward applications were reused. To avoid having to calculate ethernet checksums in software, ethernet checksumming in hardware was enabled for the interfaces. This was done in the driver initialization of the XScale code. The packet echo block was given one context on one μ Engine, running in a loop receiving packets, switching headers, and transmitting packets. Figure 3.2 shows a diagram of the microblock usage and packet flow.

The next thing that had to be explored was the passing of buffer handles to the XScale, and how packet processing could be performed there. From the example applications in the IXA SDK, there seemed to be several feasible approaches.

3.3.2 XScale implementation

To test the communication capabilities between μ Engines and the XScale, it was decided to implement the IP address switch functionality on the XScale. In achieving this, a choice had to be made between using the kernel-mode IXA Software framework structure, or to manage with the hardware abstraction layer (HAL) libraries.

The software framework was first tried, in the hope that it would provide greater flexibility and a larger degree of abstraction. The experience from this experiment was that most of the methods introduced a greater aspect of complexity without enhancing functionality, at least not for such simple tasks. One more complicating factor was that the software framework required the application to be implemented as a kernel module, thus excluding the possibility of making use of general user-level libraries in the development. The main argument for running the XScale part of a network application in kernel mode, was to avoid extra context switches (making it more efficient). The nature of IXP hardware (see section 2.3.3), however, restricts the data flow between μ Engines and XScale. This means that the bottleneck will not be XScale processing, but the moving of data between XScale and μ Engines.

Another consideration that had to be made was that, in the server cube application, data had to be moved between the host and the XScale. One way of implementing this, that looked promising was to use a socket on the ENP2611 debug port (see section 5.4). Implementing in kernel mode would make programming on a socket interface much more difficult. The sum of these arguments led to a decision to redesign the application using the HAL and operating system services layer (OSSL) libraries in user mode. Figure 3.3 shows the basic structure of this application.

The mechanism that seems to be the best for getting the buffer handles to the XScale is to spawn a thread that is dedicated to listening for an interrupt triggered by the μ Engine code. There is functionality for this in the HAL. When the interrupt is generated, a given procedure is called. This procedure processes the packet, and if a packet has to be sent to another stage in the processing pipeline, it can be enqueued on the appropriate scratch ring. The μ Engine that wants to send a packet to the XScale simply inserts the buffer handle on the scratch ring designated for communicating with the XScale and generates the appropriate interrupt.

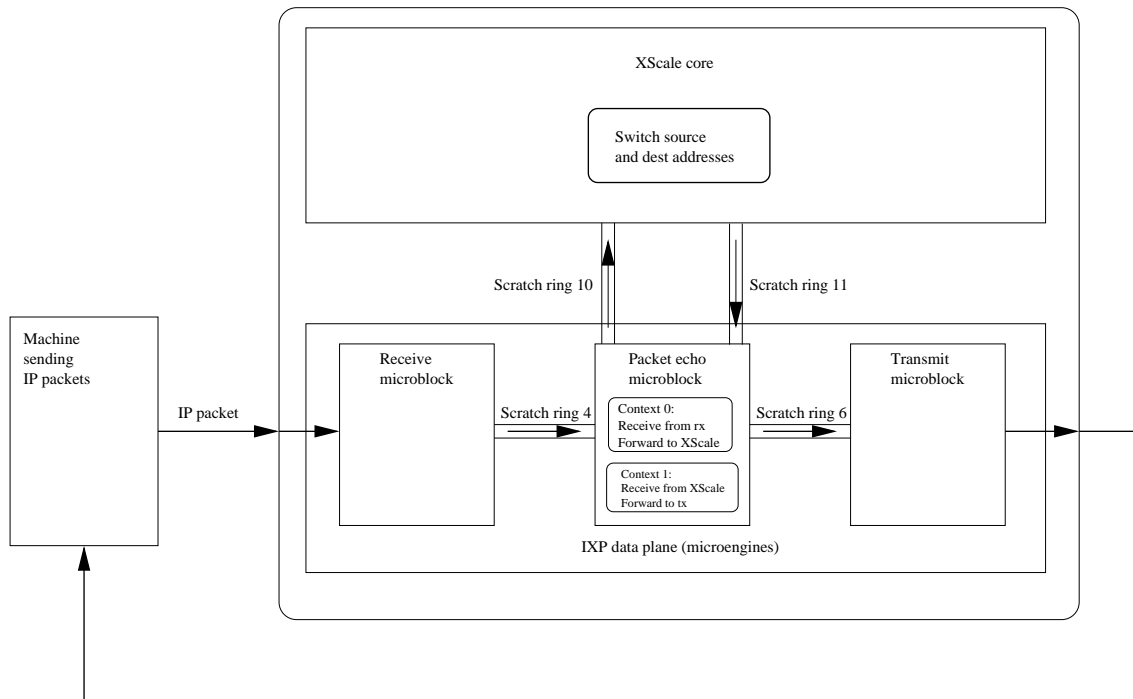


Figure 3.3: IP header switch application using μ Engines and XScale.

3.3.3 IP header switch application evaluation

In order to determine the general abilities of the different networking components of the IXP2400, the IP header switch application was created. This application receives an IP packet, switches the IP source and destination addresses, switches the source and destination ethernet addresses, and transmits the packet on the same interface that it arrived on. The nature of IP checksumming (a ones compliment sum), makes it unnecessary to recalculate the IP checksum.

In the following sections, the result of tests performed on this functionality implemented on different areas of the IXP platform and on a Linux host machine is presented. The measurements are done by timing a packet on the way out from one machine. The packet is processed on another machine, and returned. The difference between the send time and the receive time is then taken. The reason why the timing is done on the sending machine, not the processing machine, is to ensure that the measurements will be comparable. The difference in platforms and implementations would lead to results that could not be compared if measured on the processing machine. When looking at the results (in μ s) we have to take into account that the times are not only processing time for the machine performing the IP header field switch, but the processing time of the sender upon send and receive

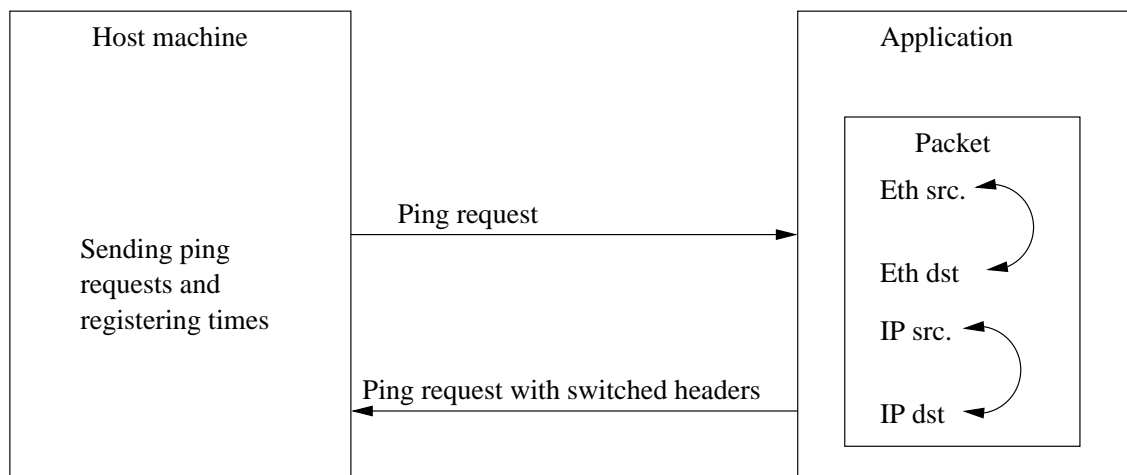


Figure 3.4: Generic test setup.

as well as the transmission times over network link and busses.

Test setup The purpose of the test is to measure the time used by an application to switch the source and destination IP header fields and the source and destination ethernet header fields. The application that performs the task is implemented in four different ways: On a Linux host in user space, on a Linux host in kernel space, on an IXP card on XScale and on an IXP card on μ Engines.

The test itself is done by pinging the interface that the IP header switch application is assigned to. The ping packet will be echoed back to the sender, which will generate a reply. The reply is also echoed to the sender. The measurement of elapsed time is done by having tcpdump listen to the interface on the sending host. Each ping request (or reply) has a sequence number. By calculating the time between ICMP packets with matching sequence numbers, the processing time for the packet pair can be acquired. In this test set, two packet sizes was used: 98 and 1497 bytes packets. This can help us to determine how big part of the whole operation copy operations related to the larger packet size represents. Figure 3.4 shows the generic test setup for the four tests.

The *first* implementation was done on a standard (SUSE 9.3) Linux host with a user space application. To be able to manipulate the IP header freely, a packet socket was created. This socket makes a copy of received packets that complies with the packet socket configuration¹. The packet is then copied to a local application buffer, and the source and destination IP addresses are switched. The packet is, finally, returned to the

¹The implemented packet socket was configured to receive all IP packets on eth0 (gigabit interface).

sender via the packet socket. The kernel ethernet layer takes care of prepending the packet with the correct ethernet header, and the packet is then transmitted on the interface. In order of preventing the kernel from generating replies to the received packets (the user space application should have total control), IP packets on the given interface is firewalled. Figure 3.5 shows a diagram of the test setup for this implementation.

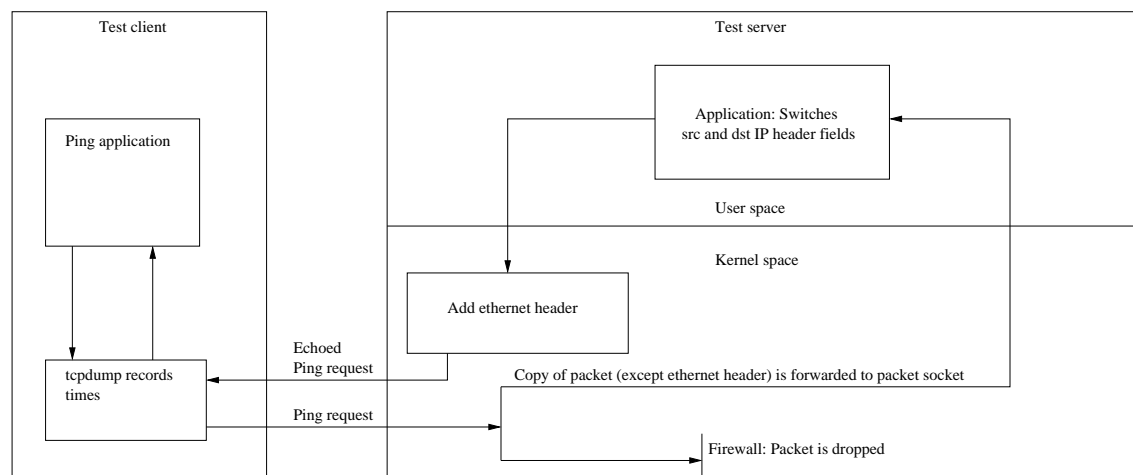


Figure 3.5: User space packet echo application test setup.

The *second* implementation uses Linux iptables to perform the header switch in kernel space. This eliminates the need of context switches and copy operations from kernel to user space. There is, however, a possibility that the complexity of iptables state matching and routing operations can slow the process down a little compared to what an application dedicated only to doing this one task can perform. The test was done on the same host as the user space implementation test. In figure 3.6, the test setup for this implementation is shown. The iptables NAT PREROUTING table rule rewrites the IP header destination field to the sending machine's address. The packet is then sent to routing. The Linux routing table directs the packet back to the interface on which it was received (based on the rewritten IP header dst field). Finally, the iptables NAT POSTROUTING table rule rewrites the source IP header field with the IP address of the local machine, and the packet is transmitted.

The *third* test implementation uses the IXP card, but performs the header switching task on the XScale core. The task is completely independent of the Linux Host, and all steps are performed on the IXP platform. Figure 3.7 describes the test setup for this implementation. The application receives the packets in the RX block and forwards them to the packet_echo block. They are then sent to the XScale core by calling the "A" interrupt. The XScale core takes care of switching source and destination headers, and sends the

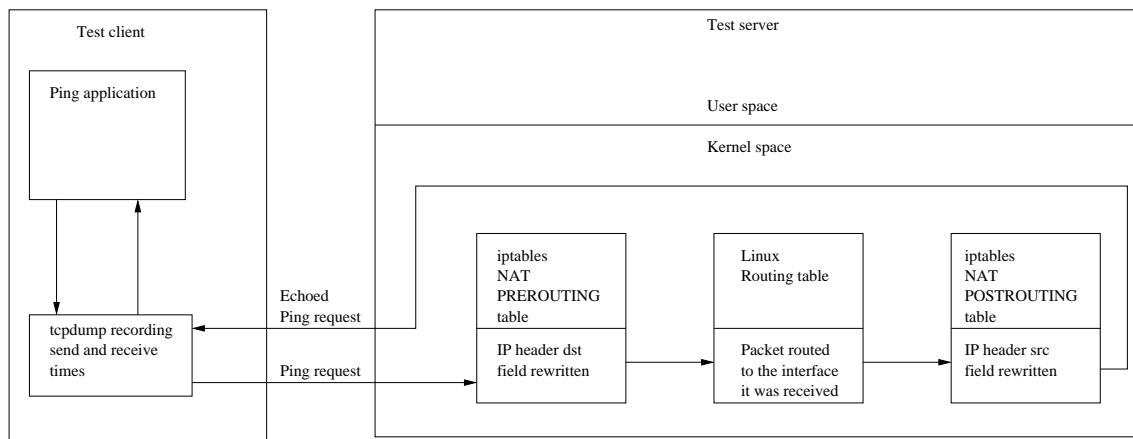


Figure 3.6: Kernel space packet echo using iptables.

packet back to the packet_echo block. The packet is finally passed to the TX block and transmitted on the same port it was received.

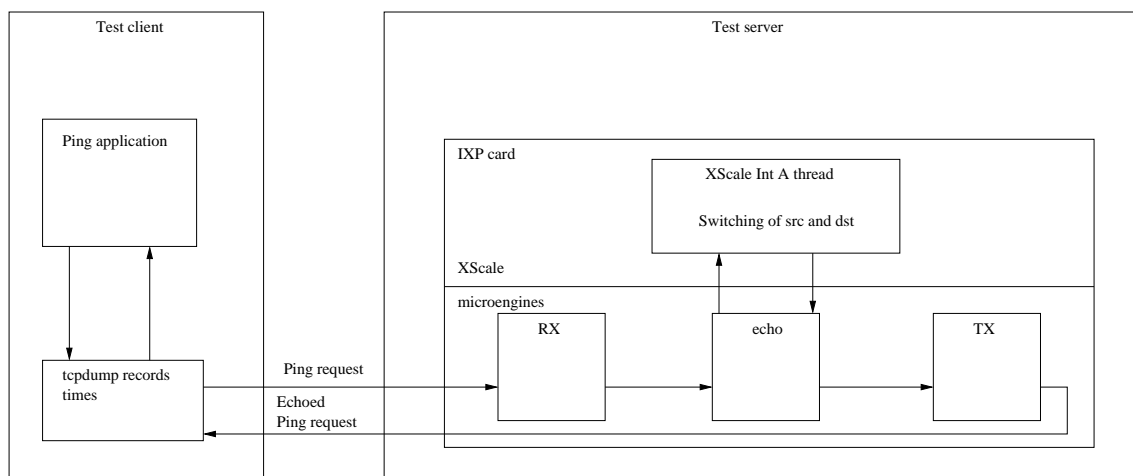


Figure 3.7: IXP XScale ICMP echo test setup.

In the *fourth* IP header switch implementation to be tested, all packet processing is done on the IXP μ Engines. In an application such as this, with limited processing power needs, we can expect a fairly efficient processing due to the parallel processing capabilities on the μ Engines. If the need for processing power should increase, it is possible to distribute processing further. In this implementation, however, the task of switching IP source and destination addresses (and mac addresses) is done on a single microengine. A diagram of the test setup is shown in figure 3.8. The ICMP packets are received by the TX block, and

forwarded to the packet_echo microblock. The ethernet source and destination fields are switched, then the IP source and destination fields are switched. The packet is then sent to the TX block to be transmitted back to the host who sent it.

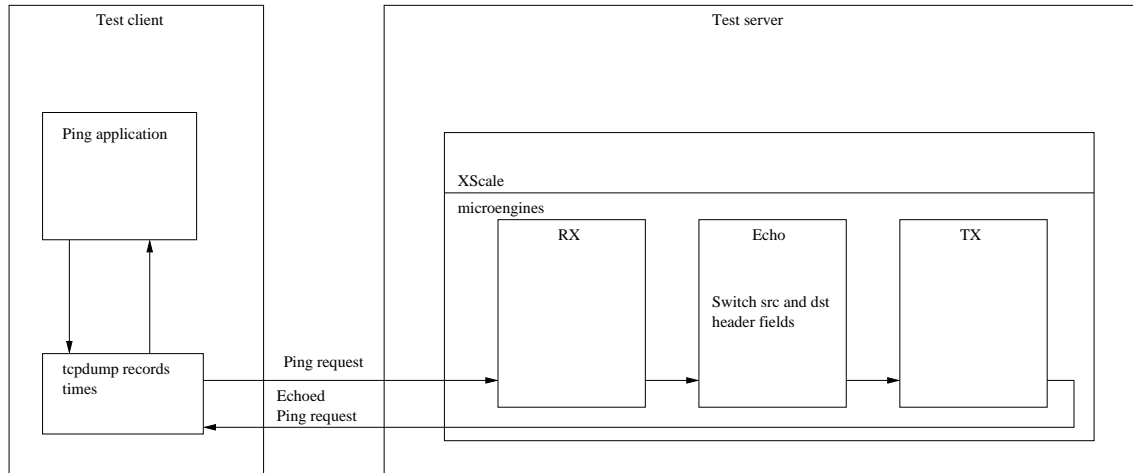


Figure 3.8: IXP μ Engine ICMP echo test setup.

Packet size	Max time	Min time	Avg time	Median	Std. dev.	Count
Linux host user space implementation						
98 Bytes	19331 μ s	11 μ s	111 μ s	105 μ s	310 μ s	99380
1497 Bytes	16216 μ s	16 μ s	172 μ s	174 μ s	102 μ s	99960
Linux host kernel space implementation						
98 Bytes	19723 μ s	11 μ s	111 μ s	101 μ s	378 μ s	98902
1497 Bytes	19092 μ s	16 μ s	168 μ s	178 μ s	182 μ s	99806
IXP XScale implementation						
98 Bytes	1145 μ s	54 μ s	123 μ s	109 μ s	44 μ s	100000
1497 Bytes	1146 μ s	110 μ s	171 μ s	173 μ s	45 μ s	100000
IXP μEngines implementation						
98 Bytes	1169 μ s	28 μ s	98 μ s	99 μ s	38 μ s	100000
1497 Bytes	1150 μ s	84 μ s	151 μ s	162 μ s	50 μ s	100000

Table 3.1: IP header switch application times.

Evaluation The purpose of the above tests was to compare the processing time of the same functionality implemented on the IXP platform and on a host machine. Table 3.1

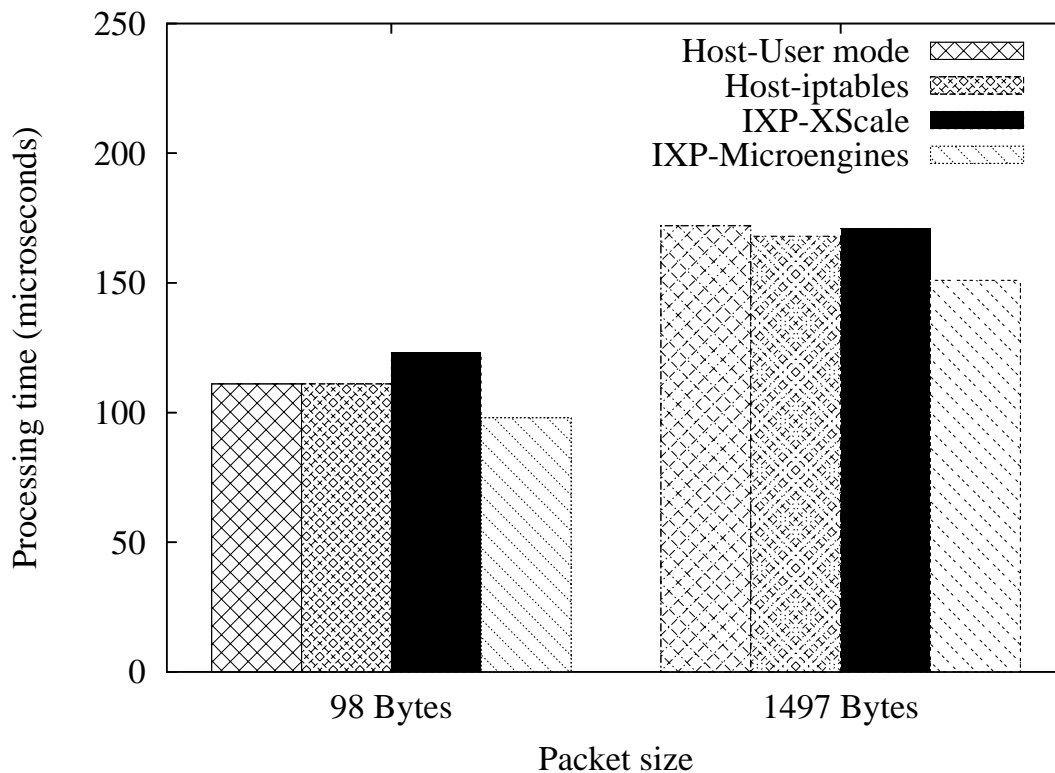


Figure 3.9: Comparison of average times for the IP header switch application.

shows the times and statistical data from the four IP header switch tests. A comparison of the average times from the tests can be found in figure 3.9.

The fact that there are similar results on the iptables implementation and user space implementation for the 98B/packet tests on the host can possibly be ascribed to the string matching support of iptables that makes it necessary to do the packet inspection on the application layer [5]. The routing process will also consume some cycles. In the user mode application, no other logic than switching the headers is applied. The packet is copied directly from the network layer, and immediately returned, thus bypassing several costly operations.

The method that produced the highest processing time for the 98B/packet tests, was the XScale implementation. This is not unexpected as the architecture is not supposed to do processing of high-performance tasks on this level, but rather handling of exceptions and data structure updating (see section 2.3.3).

It is reasonable to assume that the applications will have to spend more time on copy operations for the larger packet size, and that this is what evens the numbers somewhat.

We can see that the overall times have increased for all implementations when the packet size is 1497B. The average processing time is very similar for the user space implementation, the iptables implementation and the XScale implementation when the packet size is 1497B. The reason why the XScale results have improved in comparison to the other two may be attributed to the fact that no copying of the packet occurs on the IXP platform after it is received. Thus, more resources should be saved when the packet size is larger. The margin between the μ Engines implementation and the other three has also increased slightly. The difference between the XScale and μ Engines results for 1497B is, however, very close to the difference between the same implementations with 98B. This makes sense when considering that the cost of copy operations is the main thing separating the two experiments, and the fact that no excessive copying takes place on either of the IXP implementations.

The IXP μ Engine implementation has the shortest average time of the four by a reasonable margin for both packet sizes. This is probably due to the fact that none of the components in the processing pipeline has to wait for any system resources other than the read and write operations. The system is entirely and solely dedicated to the one task.

Although the processing speed of each functional unit of the IXP is slower than on the generic computer used in the test, it seems that the fact that the platform is dedicated to this task gives an advantage in comparison to the host implementation. For the implementation that had to pass all the packets to the XScale, however, the results were poorer. This matches the expectations outlined in section 2.3.3 about IXP programming paradigms. For the implementation of other applications, this implies that packet altering microblock applications does not stand back in performance to similar implementations on a generic host, but that XScale processing should be used for exception packets. The XScale was, however, not so slow as to exclude using it to handle situations that could arise quite frequently.

Offloading network functionality to μ Engines will, in other words, not only free host resources, but also improve on the processing speed as compared to the what the host would be able to handle.

3.4 Performance gain by offloading

In this section, the possible system performance gain by offloading network services will be tested and discussed. As presented in section 2.1, the offloading of network tasks is a sensible step to take in order to free more system resources on the host.

As the host is given extra network load, an increase in the elapsed time spent on other processes is expected due to the fact that both the kernel and, in most cases, user space

applications will have to process network data. When offloading the network task to the IXP, there should be no noticeable increase in the processing time.

The following set of tests will try to indicate what the expected gain from offloading a “low-cost” task to a network processor can be. “Low-cost” in this context is that the process of switching IP source and destination header fields of ICMP packets requires little computational power. The fact that transport layer protocols are not involved is also a simplifying factor.

3.4.1 Test

The way the tests were conducted was to measure the time spent in the process of compressing a set of folders containing about 870MB of data with the “tar” command. While this is done, the system is subject to flood pinging with packets of 1497 bytes. The ping packets are processed with the IP header switch applications described in section 3.3.3. The time spent was measured by running the “tar” command through the “time” command. This gave three results for each experiment:

- Real time: The total elapsed time for the command to run. This includes user space time, kernel space time and the time the process is swapped away, and not running.
- User time: The active time the process has used in user space.
- Sys time: The active time the process has spent in kernel space.

A total of four tests was done, all with different types of IP header switch applications running on the tested host. Table 3.2 shows the data for the performed tests. The statistical data for the first test shows how much time was spent on the process with no extraordinary network load. This should be used as a reference. In the second test, a user space application is used to switch IP headers on the received network packets while the tar process is running. The third test’s results display the times for the tar process while the IP header source and destination field switch was handled by iptables in the kernel. In the fourth test the switching is offloaded to the IXP card, and performed on the μ Engines. Any interaction with the host system that may occur is in order to read from the NFS file system.

3.4.2 Discussion

Figure 3.10 shows the average load times for this test divided in real time, user time and sys time. Figures 3.11, 3.12 and 3.13 show a zoom in of the top of each group from figure 3.10.

Time class	Max time	Min time	Avg time	Median	Std. dev.	Count
Reference time - No extraordinary load.						
Real time	153.07s	132.57s	142.42s	142.35s	4.0946s	201
User time	74.5s	73.16s	73.41s	73.39s	0.1367s	201
System time	6.94s	6.32s	6.6105s	6.61s	0.1152s	201
User space IP header switch.						
Real time	165.53s	142.12s	153.43s	153.61s	4.7797s	201
User time	80.17s	77.18s	77.62s	77.57s	0.3219s	201
System time	10.23s	7.57s	9.3971s	9.44s	0.4537s	201
Iptables IP header switch						
Real time	162.21s	136.06s	146.56s	146.46s	4.6399s	201
User time	76.89s	75.85s	76.2s	76.18s	0.1639s	201
System time	8.71s	7.36s	7.8907s	7.88s	0.2067s	201
IXP μEngines IP header switch						
Real time	156.13s	130.68s	141.12s	141.13s	4.2010s	201
User time	73.88s	73.05s	73.3687s	73.36s	0.134s	201
System time	6.83s	6.14s	6.4602s	6.46s	0.127s	201

Table 3.2: Time used for tar process with different IP header switch implementations.

Because offloading the IP header switch application to the IXP card should leave the host undisturbed to attend to other tasks, the “no-load” tests should yield the same results as the IXP-implementation. There is, however, a slight difference in the numbers. It seems that the times for the IXP implementation tests are slightly lower than for the “no-load” tests. This is probably due to the difference in network setup that is necessary to perform the tests. This can lead to system services slightly changing behaviour, and the results may therefore be somewhat affected.

The trend for the user space IP header switch implementation test is as expected. The numbers tells us that when subjected to this network load, the time used by the tar process increases by 8,73 percent compared to the IXP implementation test. The user space load times and kernel space load times confirm this tendency.

The iptables implementation load test times show a significant improvement over the user space implementation tests. It is still somewhat higher than the IXP implementation test times. The elimination of costly context switches, and also copy operations can probably explain most of the improvement. The elapsed time for this test is still 3,86 percent slower than for the IXP implementation test.

The IP header switch application applied to ICMP packets is a very simple task. For more complex protocol handling tasks, the difference in times would be bigger since the kernel tasks would have priority over the tar process. The potential of the IXP card to handle the

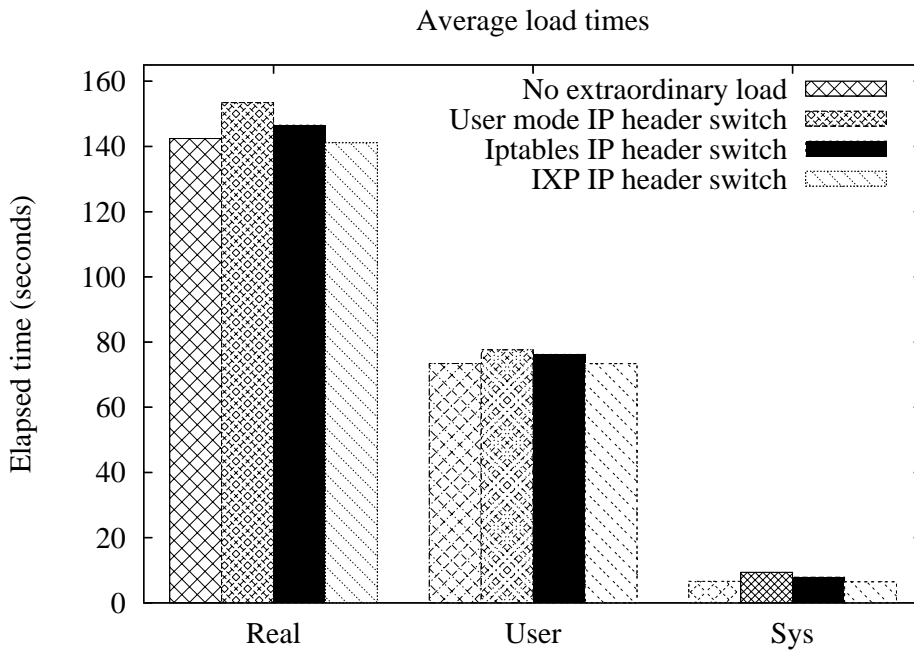


Figure 3.10: Times spent on tar process with different implementations of IP header switch application.

more complex tasks with equal efficiency (due to pipelining) increases the probability of even higher benefits from offloading tasks.

3.5 Summary

In this chapter we have evaluated the ENP2611 NPU board. The evaluation have consisted of tests related to packet handling on μ Engines and on the XScale. We have also examined the possible performance gain by offloading simple network functionality.

The IP header switch tests confirm the control plane / data plane implementation model, but also indicates that reasonable processing speed can be achieved on the XScale. This has to be taken into consideration when partitioning tasks between IXP system layers. The load tests shows that offloading network tasks to a NPU will free considerable resources. When the bandwidth increases, this gain will probably be crucial to high-throughput applications like a multimedia server.

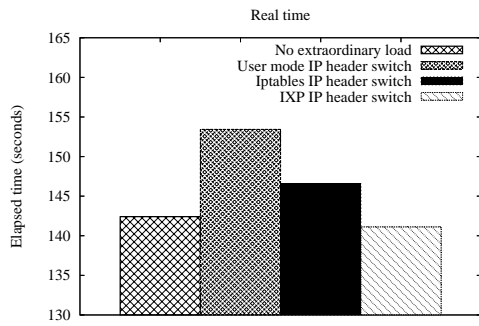


Figure 3.11: Average total time spent on tar process.

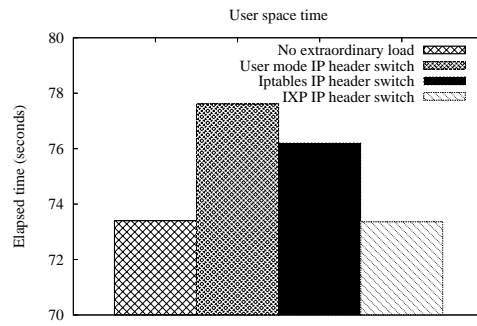


Figure 3.12: Average time spent in user space for tar process.

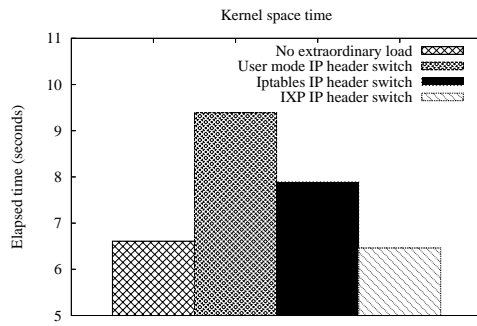


Figure 3.13: Average time spent in kernel space for tar process.

Chapter 4

Multimedia systems

In order to understand the mechanisms needed to implement a solution in the area of multimedia systems, we have to take a look at the properties of multimedia systems and the challenges they present. This chapter will present some multimedia applications and the requirements needed by such systems. An overview of some existing solutions and related protocols will also be presented.

4.1 Multimedia applications

One attempt to define the term multimedia tells us that it is more than one concurrent presentation medium [40]. In this thesis, we will refer to multimedia as one or more presentation media delivered from one computer to another over a network connection. Note that this does not exclude the cases where only one medium (i.e., only video) is delivered. The main point is that the systems that is described has the functionality required to deliver more than one type of media.

There exists many different multimedia applications today, and new, inventive uses are steadily being developed. Some common areas of use today can be:

- **Teleconferencing:** Giving people located physically far apart the possibility of speaking to each other and see the other participants of the meeting. Previously used with only sound over regular phone lines, teleconferencing over data links opens up new possibilities like sharing presentations, documents or mark up a common whiteboard.
- **Video on demand (VoD):** Video (and audio) delivered to a client upon request. This gives the customer the possibility of requesting a video whenever he/she wants.

The greatest challenge in this area has been the large datarates involved. However, the improvement of the network infrastructure makes these solutions more common. An example of an operative VoD site is SF-anytime [58].

- **Audio on demand:** Due to effective compression techniques like RealAudio [53], mp3 [39] and Ogg Vorbis [45], the streaming of audio has been a familiarity on the internet for several years. Common streamed audio services can be radio channels broadcasted over Internet, or sound used to add effect to a website. The data stream can also provide metadata for the audio transmitted (like telling you which radio program you are listening to).
- **Gaming applications:** The number of people using online games and virtual communities have exploded the last years. As an example, there are more than 2 million users of the role playing game “World of Warcraft” [62]. The data that has to be transmitted is mainly information about player actions and environment, although pictures, sound and executable code are also common to exchange. In this field there are a lot of challenges related to keeping the game world consistent to all players, and to minimize the data that has to be passed to each player.
- **Education:** There are many ways to use multimedia services for educational purposes. Students can be shown presentations, ask questions (written or by audio transfer), solve graphical tasks or participate in discussions. Lectures can be transmitted, both live and on demand, allowing students in remote locations the same possibilities as those able to travel to the lecture location.

These are only some of the myriad of multimedia applications that are in use today. With the ascent of common broadband networking, different multimedia content delivered via network to end users are increasingly common. This presents us with some challenges, given the fact that multimedia usually have some properties that are different from the kind of data that has been the most common to access over the Internet until recently.

4.2 Multimedia requirements

When a server has the task of delivering multimedia content, there are some differences in requirements compared to most other networking applications. The most important differences are:

- **High data rates:** Each second of a movie combining audio and video will require a relatively large amount of data to be transferred per time unit over the network in order to yield satisfying results.

- Large amount of data: The data that have to be delivered occupy much storage space, and the internal server operations needed to handle the data will consequently be costly (due to copy operations, buss transfers etc.).
- Many concurrent streams: If several users request content at the same time, the server will have to provide output of very high bit rates to many streams.
- Time-sensitivity: The frames delivered will have to reach the client within certain time-limits to ensure that the presentation will be satisfactory. In applications with several interacting participants, like teleconferencing or online games, the data has to reach all participants in due time in order to make the application work.

In addition, there are challenges related to the combination of different types of media into streams. As an example, there can be a video conference that includes a presentation, video and audio. All of these components have to be delivered, in time, to all of the different participants.

Two very common media types are audio and video. An example of audio data is a 16 bit, 44KHz compact disk (CD), with a play time of 77 minutes and 20 seconds. The data on this CD will occupy 780.56MB [14]. If we were to stream this data uncompressed, the bitrate would be 1.3Mbps. This, however, is vastly inefficient. To solve this problem it is usual to compress the data before transmitting. Audio data like this can be compressed without losing any of the original sound data. Using free lossless audio compression (FLAC) [14], this audio file can be reduced to the size of 413.46MB, giving us a bitrate of 713Kbps. Applying a lossy compression we can further reduce the needed bitrate. With Moving Pictures Expert Group (MPEG) 1/2 Layer 3 (MP3) compression, the bitrate can be reduced to 128Kbps without noticeable loss of audio quality (as subjectively estimated by the listener) [16]¹.

Video streams will require even higher bitrates than audio. A phase-alternating line (PAL) DVD standard can have a resolution of 720x576 pixels, and a framerate of 25 frames per second. This will give an uncompressed bitrate of about 249 Mbps (given a color depth of 24bit/pixel). This would be almost impossible both to store and send efficiently with todays technology. On regular DVD records, however, the data is compressed with MPEG2. This gives a maximum bitrate of about 9 Mbps for the video stream [11]. With the help of MPEG 4 encoding, the bitrate can be brought down to 1-2 Mbps without apparent loss in visual quality [9].

Even with the help of advanced compression techniques, we can see of the above examples that the datarates involved are still formidable. This spurs us to find ways of implementing server solutions that are able to meet the demands that multimedia applications presents.

¹When lossy compression is applied, it is impossible to restore the data to its original form.

4.3 Improving server performance

There are several ways to improve operating systems and traditional server models to enhance performance when applied to multimedia content. Some of the approaches are:

- Enhance server internals and software to allow greater efficiency (on one server). Such enhancements can be to optimize disk performance by placing data on the disk in ways that is ideal for the type of media [24], or to eliminate copy operations when moving data from disk to transmission medium [21]. Enhancements can also be made to how different protocols are processed, saving precious cycles.
- Build server clusters that behave like one server as seen from the outside. These can be directly interconnected, or connected by a switched network topology. Examples of existing solutions are nCube n4x [42], IBM VideoCharger [25] and Oracle inter-Media [46].
- Distribute delivery through proxies. Examples of systems with support for this are Apple Darwin streaming server [8] and Komssys [36]. This kind of approach can be further enhanced by applying multicast techniques like gleaning [23] and patching [20].

In addition, there are different ways to combine one or more of the above components to accommodate the needs of the system that is to be designed.

Next, we will present some server implementation strategies, and show how these deals with the challenges that comes from multimedia requirements..

4.4 Multimedia server implementations

There are many strategies that are applicable when it comes to implementing adapted multimedia servers. In this section, three implementation strategies will be described, namely single server, server clusters and proxies. It has to be said that these strategies are not mutually exclusive. A good single server implementation can be multiplied to build a server cluster structure, which again can be enhanced using a proxy strategy.

4.4.1 Single server implementation

There are many multimedia single-server implementations on the market today, although most of them can be combined to form different types of clusters. Server implementations that are in production today are amongst others: “Quicktime” [50], “Real Helix Server” [54], “Alex Arachnid” [2], “Apple Darwin” [8] and “IBM VideoCharger” [25].

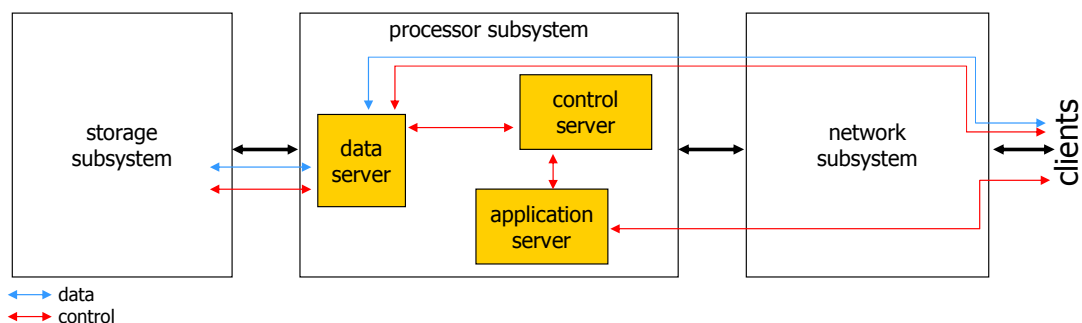


Figure 4.1: Simple general server architecture [23]

When deploying a single server to serve a (potentially) large amount of users, it is common to spend large amounts on hardware (e.g., disks, memory and CPU power) to allow the server to handle periods of great traffic. Most of these resources will then be idle most of the time. Another aspect of the single server solution is that it scales poorly. That implies that hardware upgrades, or even buying a new machine, will be required to meet demand.

Figure 4.1 depicts a common layout for a multimedia server. The diagram divides the server components in three main parts, namely storage, processor and network subsystems. The storage subsystem keeps the data that is to be transmitted, and the networking subsystem handles the actual transfer of data to the client. The processor subsystem is responsible for all the tasks required to control and synchronize the operations. The diagram also divides the typical tasks of the processor subsystem in three; Data server, application server, and control server. The application server functionality is typically the interaction with the user. This part can give a list of available material, organize billing and keep track of users. The data stream itself is delivered by the data server. The control server can be used to guide the whole process, synchronize the operations, and make sure only valid, authorized transactions are committed.

In order to make the servers better adapted to multimedia transactions, there are many alterations that can be made to architecture and operating systems. File systems that specialize on continuous data and large files can be implemented. Examples of specialized file systems are Minorca, Fellini and Presto [23]. We can use multi-processor (and multi-core) architectures to make data processing more effective. Operating systems can be modified to make the data-path more efficient (i.e., DROPS [10] and INSTANCE [47]). These kind of alterations, especially the ones that imply custom built hardware or over-provisioning, are expensive. In order to make cost-efficient server solutions, it is common to use other strategies to meet the scaling obstacles.

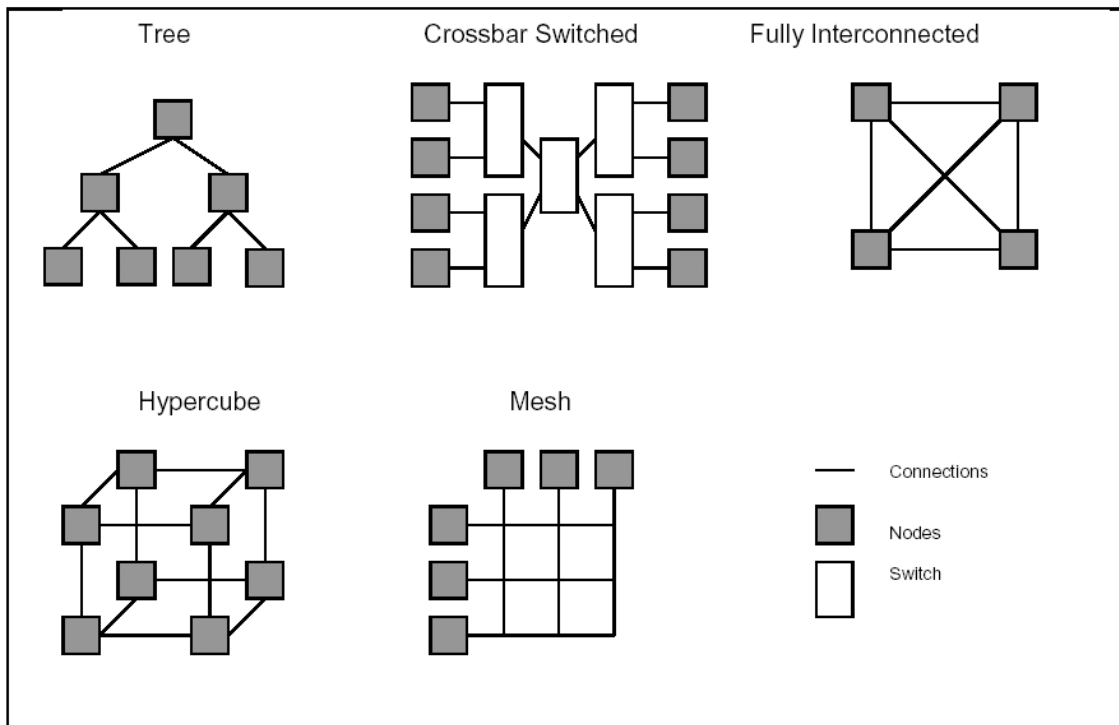


Figure 4.2: Interconnection topologies

4.4.2 Server cluster implementation

To meet the demand for scalability, it is possible to distribute the multimedia server over several interconnected computers. These systems can behave as one server as seen from the outside, but will internally distribute the tasks between several machines. It is also possible to distribute the content between several physically separated servers.

Different topologies offer different benefits to the cluster. This, and the cost of implementing the system (number of connections per node, cost of switches, etc.), must be taken into consideration when deciding upon a topology. Figure 4.2 shows some popular interconnection topologies [19]. The fully interconnected solution, where every node has a direct connection to every other node, is the one with the greatest capabilities in terms of communication, but it is also the most expensive to implement in most cases [19]. From this extreme there are interconnect topology variants which have to pass some data through other nodes to reach its target. For the tree topology, the root is the access point of all data, and therefore a possible bottleneck. Mesh or hypercube topologies will have more evenly distributed data, but will require more complex search and routing algorithms. Switched networking topologies are well-known, and scales well. The switching mechanisms, how-

ever, increases the inter-node latency.

In order to be scalable, most video server solutions have implemented support for clustering in one way or another. This kind of support can range from simple load-balancing algorithms which divide the load between different free-standing servers to more complex interconnection topologies and strategies. Examples of server solutions with support for clustering is “Real Helix” [54] and “IBM VideoCharger” [25]. There are also server solutions that are not based on switched network clustering, but on direct interconnection topologies. An example of this is the “nCube n4x” [42] server solution [4]. This cluster solution is based on a hypercube topology, with directly interconnected nodes.

4.4.3 Proxies

The main idea of proxies is to move the most popular content closer to the end users, and in this way reduce traffic on the main content server(s) and backbone network. A statistical rule of thumb is that a great majority of the requested content is represented by just a few of the most popular titles [63]. This empirical observation makes it possible to design caching strategies that make proxies keep the most wanted media, and forward requests for less popular media upward in the media server hierarchy. Another effect of this strategy is that the proxy servers can actually be moved closer to the end user physically, thus reducing network traffic.

Figure 4.3 shows a possible proxy server layout. In this diagram, the master servers have available all the offered content. When requests are made from the end users, the requests are sent to the master servers which provide the content. The in-between proxy servers can then begin to cache the most requested titles of their region. Multimedia server solutions that support proxying include “Real Helix” [53] and “Komssys” [36],

4.5 Hypercube/n4x server solution

In the myriad of different architectures and server topologies, one system that possesses many of the wanted characteristics is the c-cor n4x [4] (previously nCube n4x [42]). The solution is, however, based on custom hardware, and is quite expensive to deploy. This section will try to describe some of the key aspects of the n4x solution.

4.5.1 Hypercube multicomputer structure

This subsection will describe why a hypercube is a good base for an interconnection topology. To show this, we first have to explain the geometrical properties of a hypercube

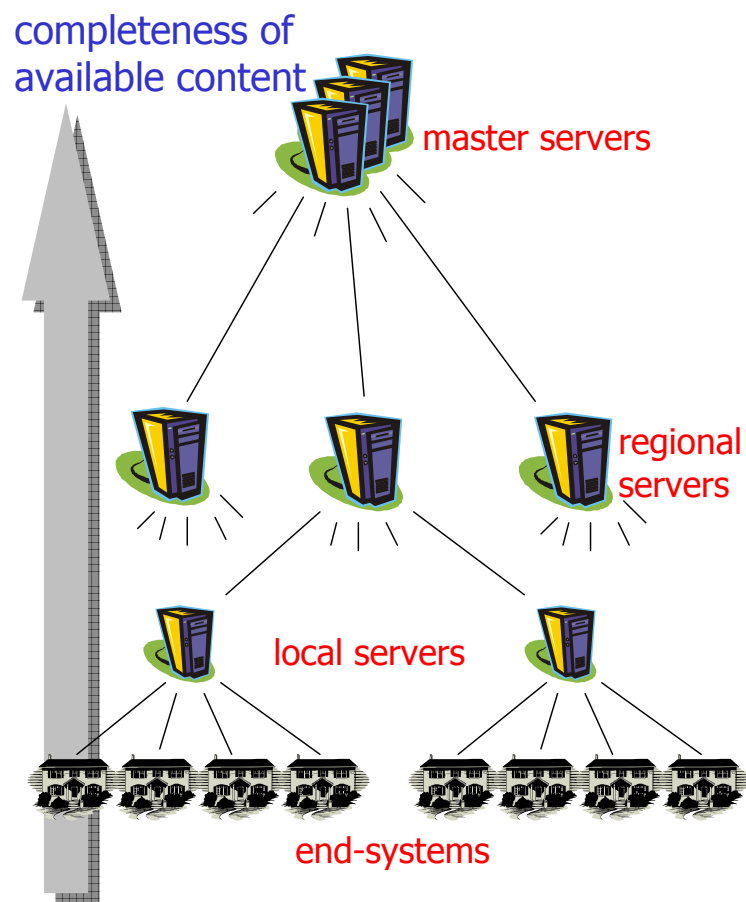


Figure 4.3: Diagram of a possible proxy layout [23].

structure. This is the “Free On-line Dictionary of Computing” (foldoc) [15] definition of a hypercube:

A cube of more than three dimensions. A single ($2^0 = 1$) point (or "node") can be considered as a zero dimensional cube, two (2^1) nodes joined by a line (or "edge") are a one dimensional cube, four (2^2) nodes arranged in a square are a two dimensional cube and eight (2^3) nodes are an ordinary three dimensional cube. Continuing this geometric progression, the first hypercube has $2^4 = 16$ nodes and is a four dimensional shape (a "four-cube") and an N dimensional cube has 2^N nodes (an "N-cube"). To make an N+1 dimensional cube, take two N dimensional cubes and join each node on one cube to the corresponding node on the other. A four-cube can be visualised as a three-cube with a smaller three-cube centred inside it with edges radiating diagonally out (in the fourth dimension) from each node on the inner cube to

the corresponding node on the outer cube.

Each node in an N dimensional cube is directly connected to N other nodes. We can identify each node by a set of N Cartesian coordinates where each coordinate is either zero or one. Two nodes will be directly connected if they differ in only one coordinate.

The simple, regular geometrical structure and the close relationship between the coordinate system and binary numbers make the hypercube an appropriate topology for a parallel computer interconnection network. The fact that the number of directly connected, "nearest neighbour", nodes increases with the total size of the network is also highly desirable for a parallel computer.

There are several advantages to an interconnected topology in comparison with a switched network topology. The direct connections reduce latency, which switched connections would introduce. Powerful switches are expensive, so a switched network doesn't scale well in terms of cost [19]. It must also be considered that when implementing an application for an interconnected network, you are free to make optimized protocols that can enhance performance for your task.

In principle, the hypercube topology can be used to implement different sorts of parallel multicomputers. The topology has a balance between the fully interconnected network, where a connection point to every other node must be supplied, which would be expensive, and topologies where the data have to travel a long distance from one node to another. A multicomputer is a cluster of interconnected cooperating computer nodes, behaving as one computer. It distributes tasks and information by message-passing (as opposed to a *multiprocessor*, where several processors share a common memory area) [52]. Seen from the outside, it will behave as one server instance.

The qualities of hypercube topology ensure that when the system size increases, no new bottlenecks are introduced. Figure 4.4 shows how nodes are added in the cube. The amount of possible routes for information to take in the cube makes an even distribution of message passing possible.

4.5.2 The c-cor n4x multimedia server architecture

The n4x solution [4] makes use of the Hypercube topology to implement a multimedia server cluster that allows for a great bandwidth potential. This is made efficient by high-speed message routing hardware [4]. The solution is scalable because you can expand it exponentially (see section 4.5.1) There is no single point of failure; routing is dynamically processed. The routing is achieved through hardware especially built for this task. This also helps performance when data can be routed by the least trafficked route, or avoid broken routes. Since there are only one copy per content-item, it maximizes storage space. The

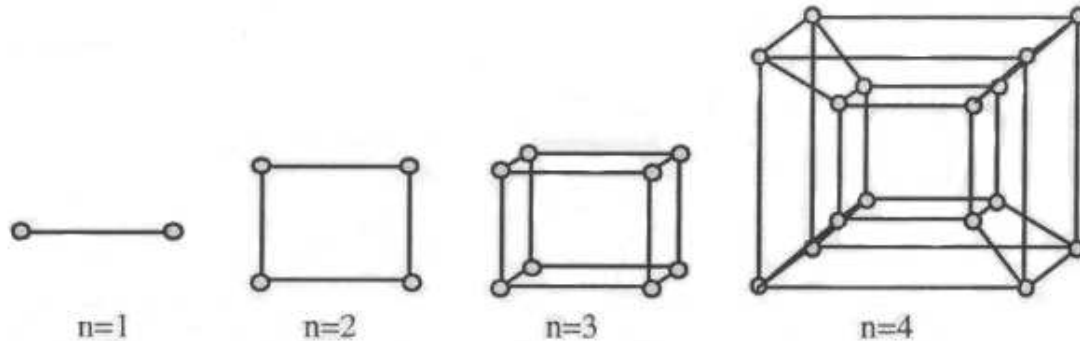


Figure 4.4: Hypercube scaling

system boasts ability to serve over 60000 concurrent streams with rates of 16Gbps per node and over 2 Tbps systemwide [42]. This is achieved by using the topology to evenly distribute data load internally in the cube. Each node has 8 full-duplex, high-speed ports for intra-cube communication. Internal load-balancing is achieved by striping² the data over all disks in the cube. The redundant array of inexpensive disks (RAID) striping also ensures that if a whole node (and all its disks) fails, the system can continue operation uninterrupted.

The n4x MediaHUBs are where the data is stored. These nodes are interconnected with eight hypercube connectors. This can accommodate a hypercube with $n = 8$ or a total of 256 nodes. These require no common memory access, they use message-passing to move data. The messages are controlled by an adaptive routing system. This system will find the most efficient route from point to point, and will avoid heavily trafficked routes or broken links.

The interconnection is supported by special hypercube connector hardware. This hardware supports routing logic by using a built-in vector processor unit [42]. This ensures that sufficient resources are available for video data retrieval and streaming on the server.

To transmit the data delivered by the cube, n4x uses hardware implemented interface modules customized for the needed transmission type. Examples of such modules are: QAM Cable, ATM, DVB-ASI, Ethernet and Gigabit Ethernet. These interface modules have hardware support for i.e. multiplexing, encoding or forward error correction (FEC) operations needed to support its interface type/model. This enables high output data rates from each hypercube node (MediaHUB).

²The data for each media item is evenly distributed among all nodes, requiring a smaller amount of data to be fetched from each node, thus reducing the possibility of bottlenecks.

The custom built components are assembled on a motherboard with Intel [27] 860 chipset and an Intel Xeon processor. The disk array are controlled by Qlogic 12160 [49] SCSI controllers. Put together, these components form a customizable, scalable server solution.

The next section will present some protocols that are instrumental in achieving better multimedia streaming performance. The protocols presented are used for setting up and performing streaming of real-time data.

4.6 Protocols

To set up a multimedia stream, we have to make use of some tools to control the stream (start, stop, pause the stream etc.), to convey information about how the stream progresses, and to be able to ship the data with the necessary information. This can for example be accomplished with the following protocols:

- Real-Time Streaming Protocol (RTSP) [56]: Used to issue control commands to a multimedia server. Enables the client to for instance play, stop or pause a media stream. Also used to setup the stream, and convey necessary information about port numbers for server and client.
- Real-Time Transport Protocol (RTP) [55]: Used for sending media. Provides services like payload type identification, timestamps, sequence numbering and delivery monitoring.

RTP Control Protocol (RTCP) is used in conjunction with RTP to allow the service to optimize streaming performance based on feedback from the client. This protocol is used to send information between the client(s) and server(s) to keep track of the progress and quality of the stream(s). The server(s) can then make adjustments to compensate for any events that may occur during transmission.

4.6.1 RTSP

The RTSP protocol [56] is an application-layer protocol for controlling media streams. An RTSP session is identified by its ID on the server and is completely independent of an eventual TCP connection. Consequently, the client can open and close several TCP connections, and still control the same RTSP session. RTSP can also be conveyed by user datagram protocol (UDP). In most aspects, RTSP is very similar to HTTP/1.1 [17]. The protocol supports the following operations:

- Retrieval of media from media server.

- Invitation of a media server to a conference.
- Addition of media to an existing presentation.

To be able to send a request for media, the client has to have some info about it. This can be placed in a description of the media, and can be made available in different ways, like publishing it on a web server, or making it available through the media server itself. The description will also give information about the transport methods the server is capable of. The alternatives are:

- Unicast with the client determining the port.
- Multicast with the server determining the address and port.
- Multicast with the client determining the address and port.

The stream controlled by RTSP may be sent via a separate protocol, independent of the control channel. A common protocol combination to use is RTP/UDP. The stream is totally independent of RTSP, and once started, it will continue until the media is exhausted if no more commands are received. The RTSP session on the other hand will have to maintain a state depending on which commands it has received. The specification states that the RTSP session identifier has to be randomly generated, and has to be at least eight octets long to make guessing the ID more difficult.

The following RTSP methods will affect the sessions state:

- **SETUP**: Initiates an RTSP session and allocates resources for a given stream.
- **PLAY** and **RECORD**: Start transmission of an already configured stream.
- **PAUSE**: Halts a stream. Does not free server resources. The stream can be restarted if wanted.
- **TEARDOWN**: Releases the RTSP session, and frees the resources allocated for the stream.

The location of the media controlled by RTSP is defined by the RTSP uniform resource locator (URL). It is equivalent to a HTTP URL, and describes the server address, eventual port number and the absolute path of the media on the server [56]:

```
rtsp_URL = ( "rtsp:" | "rtspu:" )
           "/" host [ ":" port ] [ abs_path ]
host      = <A legal Internet host domain name of
           IP address (in dotted decimal form)>
```

The RTSP message is text-based, and each line is terminated with CRLF. The character set used is ISO 10646.

When a request is received and interpreted a response will be sent. The response can have the following format: [56]

```
Response      =      Status-Line
                *(
                |      general-header
                |      response-header
                |      entity-header )
                CRLF
                [ message-body ]
```

The status line consists of the RTSP version type, a status code and a reason-phrase. The reason-phrase is intended to give a short textual explanation of the status code. The status codes can be classified as follows:

- 1xx: Informational - Request received, continuing process
- 2xx: Success - The action was successfully received, understood, and accepted
- 3xx: Redirection - Further action must be taken in order to complete the request
- 4xx: Client Error - The request contains bad syntax or cannot be fulfilled
- 5xx: Server Error - The server failed to fulfill an apparently valid request

The response header gives additional information that could not be included in the status field. This includes additional information about the resource requested like for instance server and client port numbers for the media stream.

Each RTSP request carries a sequence number labeled “CSeq”. The sequence number is incremented by one for each new request. If a given request is not acknowledged, the eventual retransmission of the request have to carry the same sequence number as the original request.

Figure 4.5 shows a typical RTSP session. The first step for the client is to retrieve information about the stream. In this example, it sends a regular HTTP request to a web server, and gets the description file in response. This file could have been made available from the media server too, accessed with an RTSP DESCRIBE request. The client now sends a SETUP request, asking the server to reserve resources for the new stream. If all is well, the server responds with an OK message containing the information needed by the client to receive the stream. When the client is ready, a PLAY request is sent to the server. This tells the server to commence streaming the data. The server sends an RTSP

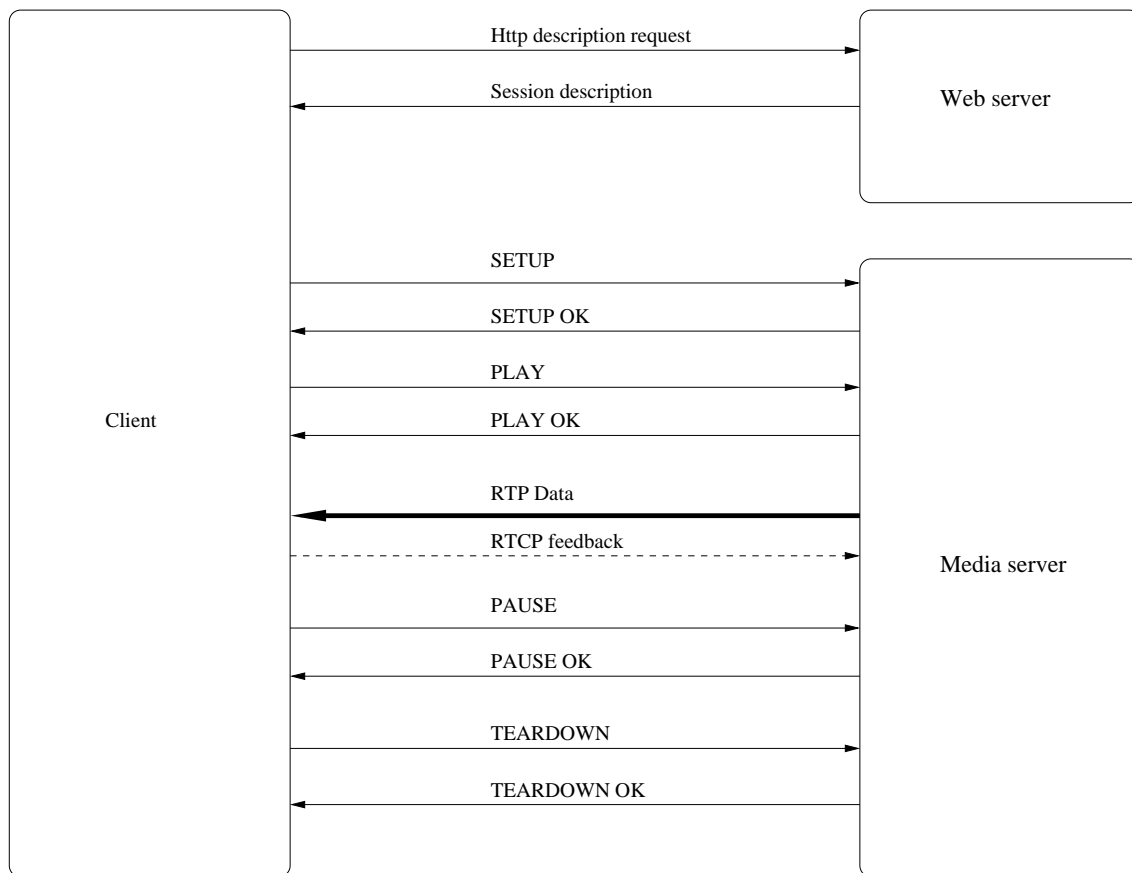


Figure 4.5: A typical RTSP session [7]

OK to the client, and starts to send RTP packets with the requested data. If supported, the client occasionally sends RTCP feedback messages to the server to report on the state of the stream. If a PAUSE or STOP message is received, the server will have to take the appropriate action. To end the RTSP session, and free all related resources, a TEARDOWN message is sent from the client.

4.6.2 RTP

RTP [55] is a protocol which provides a set of services for real-time media streaming purposes. These services include payload type identification, sequence numbering, timestamping and delivery monitoring. RTP is usually being used on top of UDP, but it is also possible to use it with other suitable underlying network protocols. The protocol was primarily designed to facilitate multimedia conferences, but is in common use as a means

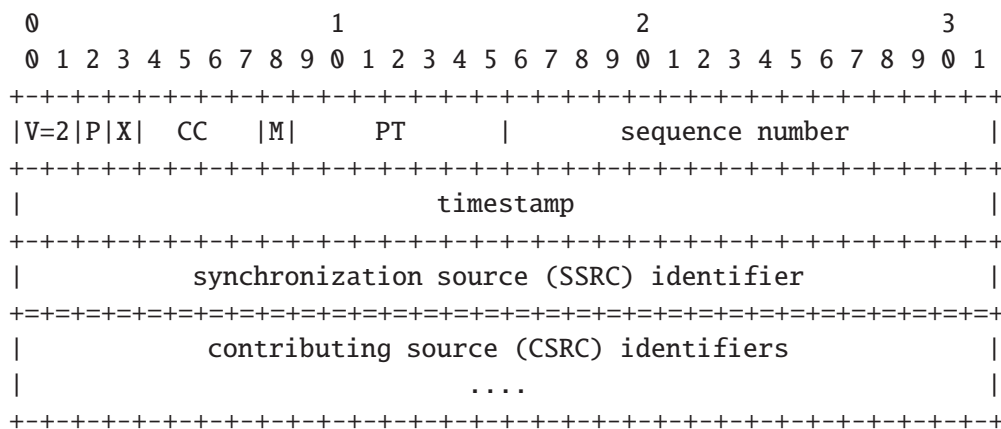


Figure 4.6: RTP header format.

to transport any kind of real-time media stream.

The RTP header is shown in figure 4.6. The version (V) field of the RTP header states which version is used in the implementation. The version designated by RFC3550 is 2.

The padding bit (P) states that at least one octet at the end of the payload is padding. How many octets are padding can be found by reading the last octet. This padding is necessary for some encryption algorithms.

The extension bit (X) states that exactly one header extension follows the RTP header. The potential header extension has a field that states its own length. It is recommended to try to manage challenges using conventional means, and try to avoid using the header extension.

The source of an RTP session is identified by a 32 bit synchronization source (SSRC) identifier (figure 4.6.2). This makes the source independent of network addresses for identification. Packets from a given synchronization source share the same timer and sequence number space. In that way the client can group the received packets by SSRC for playback (for instance where several different RTP sessions are received on the same port). The SSRC is to be chosen randomly, and is meant to be globally unique inside a specific RTP session.

If there are contributors to an RTP session (for instance if several media streams are mixed into one), the RTP header is appended with the SSRC's for the contributing streams. This could for instance happen in an audio conference, where the mixer could identify the speakers contributing to the audio by appending their contributing sources CSRCs to the RTP header.

The CSRC count (CC) field contains a number indicating how many CSRC's are appended to the header.

The usage of the marker bit (M) is defined by a profile. It can be used to indicate significant points in a stream, like frame-boundaries.

The payload type field (PT) contains a code identifying the format of the payload media. Some standard types are defined in RFC3551. If a receiver does not recognize the payload type, the received packets will be discarded.

The sequence number is incremented by one for each packet sent. The initial value is to be chosen randomly to make attacks on encryption more difficult.

The timestamp field represents the sampling time for the payload data. It is derived from a time reading, or from the nominal sample instance. The progression rate of the timestamp can therefore be different depending on the sample rate of the payload media and format. The initial value should be chosen randomly. If, for instance, several packets belong to the same video frame, they should have the same timestamp. They would, however, have consecutive sequence numbers.

An RTP session can be monitored and enhanced using RTCP in tandem with the session. This is done by sending RTCP sender and receiver reports that indicate how all of the parts involved in an RTP session are doing.

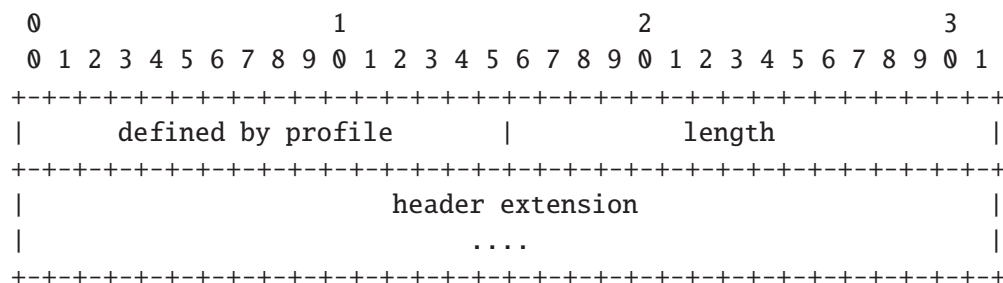


Figure 4.7: RTP header extension format.

RTP header extension (figure 4.7) is a mechanism that allows for custom extensions to be added to the RTP header. It is possible to implement servers and clients that support these extensions, but systems that do not support the custom add-on must be able to ignore the extension without losing any of the original RTP functionality. It is recommended in the specification that other methods should be used to overcome any challenges that an extension would help solving.

4.7 Summary

What is wanted from a multimedia server solution is high performance to serve the clients during eventual peaks, scalability to meet an increasing demand, and reliability to make the users trust the supplier [59]. In this chapter, we have seen that multimedia applications have a range of stringent requirements that must be met. The methods for achieving this are to build optimized servers, to group the servers in special topologies and to distribute the content to proxies, depending on demand. To support such services, enhance transmission and optimize streaming, we can make use of special protocols like RTSP, RTP/RTCP that improve the streaming performance and adaptability.

The c-cor n4x addresses many of the multimedia challenges mentioned. The n4x approach achieves better multimedia streaming performance and is a flexible, freely expandable topology, but uses special hardware that offloads the main processing platform. There are custom hardware both for routing inside the server topology and for multiplexing, encoding and transmitting data out of the cube with different network standards. Although effective, this kind of custom hardware is expensive to develop.

Today, more and more parts of standard computer components are given their own processing power in order to offload the central systems. In the wake of this development, there have surfaced new platforms of programmable hardware with special properties and uses. One common example of this is the graphics processing unit (GPU) commonly used to offload demanding graphical calculations. Another platform that is maturing is the Network Processing Unit (NPU). These are programmable units optimized for network functionality. The NPUs are mass-produced, and relatively cheap. The possibility then arises for programming these units to perform tasks that previously would have to be hardware-implemented, like hypercube routing mechanisms in a server solution.

In the next chapter, the implementation of the basic functionality of a video server cube (VS^3) will be described. The server cube will have a hypercube design topology, and IXP2400 NPUs will be used to handle the routing process.

Chapter 5

The VS³ Video Server Cube

This section describes the steps taken to implement a limited working multimedia server cluster based on a hypercube topology. Small parts of a similar system are first implemented on IXP1200[31] boards as an assignment in the inf5070[23] course, and this system is also described. From there, each step in testing the hardware, and deciding on implementation strategies up to the final implementation is discussed. A major part of this work has been to get knowledge of the IXP2400[32] hardware, the Intel IXA SDK 3.51[30] and the possibilities and disadvantages of different approaches. This has led to several separate experiments to measure how different implementations and components perform given specific tasks. The results of these, and the increasing understanding of the hardware and programming platform have resulted in several incremental design steps that has culminated in the current server solution. The implementation process resulted in several intermediate designs that were rejected. The purpose of these designs, and the reason why they finally were discarded will be discussed. Finally, approaches on how to further improve the system, and remove the current bottlenecks will be presented.

5.1 Hypercube server general design

The system we want to make is an implementation of a multimedia server cluster with hypercube topology similar to the nCube n4x system. The system should use message-passing to communicate internally, and have no shared memory areas between nodes. To make routing efficient, and to offload the host, we want to use IXP2400 network processor cards to handle the routing operations. In figure 5.1, the basic design is outlined. We want to use the Linux hosts as data and streaming servers, and the IXP cards to take care of operations related to locating the media and routing packets to the correct egress node. When a data packet is to be sent to an outside client, the streaming application has to

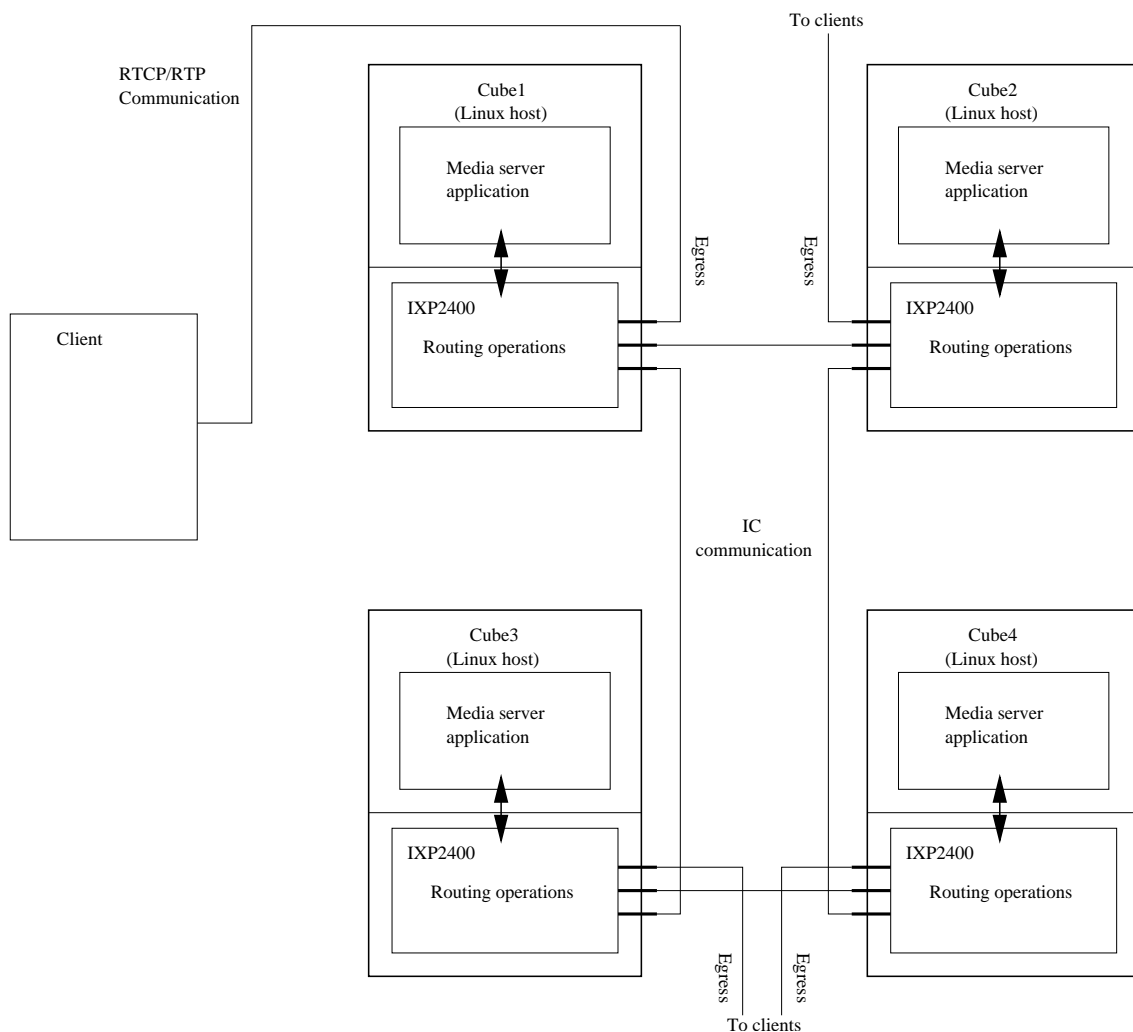


Figure 5.1: Basic design of multimedia server cube with IXP cards.

send it to the IXP card, it will then be routed to the egress node and sent to the client. An adaptive Domain Name Service (DNS)¹ will make sure that requests are distributed evenly between the egress nodes. The cluster will look like one server as seen from the client.

RTSP is the protocol we want to use for setting up and controlling the media stream. The streamed data itself shall be transported using RTP. For intra-cube messages and routing, a new Intra-Cube (IC) protocol is designed.

¹This service could either communicate with the server cluster to find the egress node with the most free resources, or could statically select a new node each time, in a round robin fashion.

As discussed in chapter 4, there are many optimizations that can be applied to a multi-media server cluster. In this assignment, we have implemented the basic functionality needed by a server cluster capable of delivering a video stream. Optimizations like striping the data over all nodes, and dynamically adaptive routing will have to be built into the system at a later stage due to the limited timespan of this master thesis. The n -value of the hypercube topology is limited to 2.² The reason for this is that the ENP2611 network cards have 3 optical interfaces, and one of these interfaces is used for the egress connection. There is the possibility of inserting more than one card into each host to increase the possible n -value of the topology, but this option has not been thoroughly explored due to the increase in complexity this would lead to, and the number of available cards.

5.2 The legacy design

The design of the server system is based on the nCube n4x system. nCube claims that their solution is able to serve 2000 concurrent users with a bandwidth of 3,75Mbps per user with a rack of 8 MediaHUBs (cube with $n = 3$) [42]. Their solution is, however, heavily based on custom hardware, and as of such expensive to develop and produce. The IXP cards are optimized for packet handling and have several network interfaces, in a manner similar to the custom hardware of the n4x system. This makes it possible to envision a similar design using programmable IXP cards for cube interconnection.

The work began with some basic attempts to implement parts of a routing functionality on the IXP1200 card as an assignment for the course INF5070 [23]. We describe the main features of this design in the following section.

5.2.1 The INF5070 implementation

The assignment for INF5070 [57] was to implement an n4x-like solution with RTP based transport and RTSP based control utilizing IXP1200 [31] cards and a given programming framework based on the Intel IXA SDK 2.01 [29]. We chose to implement the routing part, and postpone the RTP/RTSP part. The software framework that was provided had a structure that delivered all packets to the StrongARM core (Figure 5.2). The choice was made to use a source routing algorithm with information conveyed through an intra-cube (IC) protocol. The purpose of the IC-protocol was to handle:

- Routing between cluster nodes.
- Transportation of video data to the egress-node.

²A hypercube topology with $n = 2$ consists of a total of $2^2 = 4$ nodes.

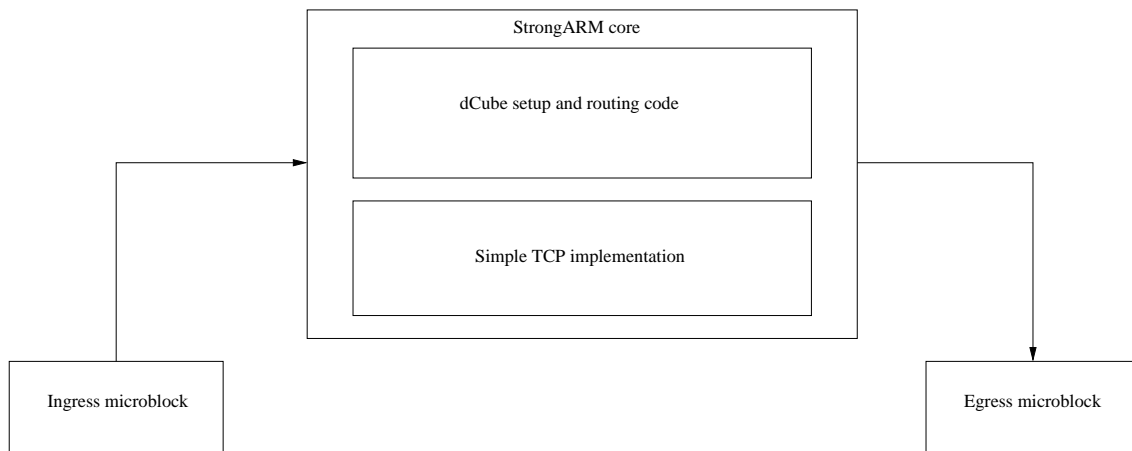


Figure 5.2: Dataflow inf5070 implementation

- Forwarding of video-playback-control to the node serving the stream.
- Handling requests for specific video files, locating the file and setting up the stream.

The final result of the project was an application that managed to stream the requested data from one node to another using intra-cube routing. The bitrates was, however, very limited because of the fact that all traffic, routing, setup information and data packets, had to pass through the StrongARM core on each and every node it passed through. The performance was also reduced by the fact that the data itself was read and streamed by an application running on the StrongARM. This demanded too many resources to achieve high performance.

5.3 SDK code base

The SDKs from Radisys and Intel provide an extensive code base from which applications can be adapted. The majority of these examples, however, focuses on standard networking tasks, and could not easily be used as a basis for the VS^3 application. In the Radisys SDK static forward application, there were components that could be used successfully. To build any network application, functionality has to exist for receiving and transmitting data. This functionality was provided by the RX and TX microblocks from the static forward application. With an asynchronous interface of passing buffer handles on scratch rings, new components could be fit seamlessly in between of these microblocks.

With RX and TX functionality being the only “recyclable” component, the following stages had to be built from the bottom: Microblocks for data plane processing, XScale

components for control plane processing and a streaming application located on the host machine. There were made designs that aimed to be able to use existing host server applications, thus eliminating the need to implement this. The next sections will discuss different designs for implementing the VS^3 system.

5.4 XScale-Host communication

In order to pull data from the host to XScale and μ Engines, the ideal solution would be one that could accommodate very high data rates. Due to the large scope of the assignment, the task of implementing this efficiently was temporarily postponed. A strategy for achieving this can be found in the further work section (see section 6.3). The method that was ultimately chosen was to transmit the data over the 10/100 debug Ethernet port. Based on this decision, three different design choices were outlined and tested:

- **Use network file system (NFS) to read the file directly from the XScale application:** Using NFS would provide an abstraction of all operations that are needed to get the file from the host to the XScale. The performance of the transfer itself, however, would be reduced due to the nature of NFS services. NFS, based on remote procedure calls (RPC) will, in addition to moving the file data, also use much resources on consistency checking and synchronization. The NFS approach would also imply that all tasks related to reading data, generating packets and sending these packets have to be handled by the XScale core. For a scenario with many streams, the load would probably be too big for the XScale core, and lead to performance reduction.
- **Create a raw socket and forward packets directly to the host:** Another approach was to create a raw socket on the XScale. This would give the opportunity of creating tailor-made packets and transmitting them to the host without creating a two-way socket connection. Any available multimedia streaming server software could then be used on the host machine, but some manipulation of packets would be needed to make the server software work as wanted within the cube structure. This concept will be described in more detail in section 5.5.
- **Create a regular TCP socket to communicate with an application running on the host:** The regular socket alternative would demand a server application on the host that could receive a connection from the XScale. The performance of the data transfer would be limited by the TCP/IP stack processing both on XScale and on the host machine. It would, however, provide the service of a byte stream between the applications. The most attractive element of this solution is that it would simplify

some aspects of implementing the streaming (host) application. More information about this design approach can be found in section 5.6.

The implementation work began using the raw socket strategy. This allowed the use of a server system that already existed (i.e. komssys [36]) on top of the routing framework. The next section provides a description of the implementation strategy, and explains why it eventually was discarded in favour of the regular TCP socket design.

5.5 The raw socket design approach

After exploring the hardware possibilities in section 3.1, the first structural design for the whole server cube application was made. If stream setup was to be done by RTSP, a TCP connection would have to be set up. A solution where the Intra-Cube (IC) routing system would forward all TCP packets to the host machine containing the requested media was outlined. The task of delivering packets to the correct egress node would rest on the routing layer (IXP cards), and each host machine would behave like a freestanding multimedia server, oblivious of the routing layer below.

5.5.1 TCP Handling

For this approach to work, the TCP connection from the client would have to address the machine that hosts the wanted media, not the egress machine. In order to do this, the egress node (the node receiving the initial TCP SYN request) would have to do a three-way handshake, then receive the RTSP SETUP message. This is because we don't know where the media file is located upon receiving the initial request, and the client expects a complete TCP connection before sending the RTSP SETUP request. The machine that hosts the file will have to be found (the mechanisms for this are presented in section 5.6). When the route between the egress node and the machine delivering the media has been found, the XScale application on the machine will have to perform a three-way handshake to connect to the application on the host. It can then forward the RTSP SETUP packet. In this way, the host application will never know that the TCP connection has been routed through several other machines, and will believe that it is communicating directly with the client. There are, however, some more complications to doing this. The sequence (seq) and acknowledgment (ack) numbers of the different TCP connections will be different, and will have to be simultaneously translated to match. The content of the RTSP packets will have to be modified to make the server stream to the XScale (which will take care of routing and forwarding). All these extra steps will have to be done on the IXP platform as a part of the routing framework.

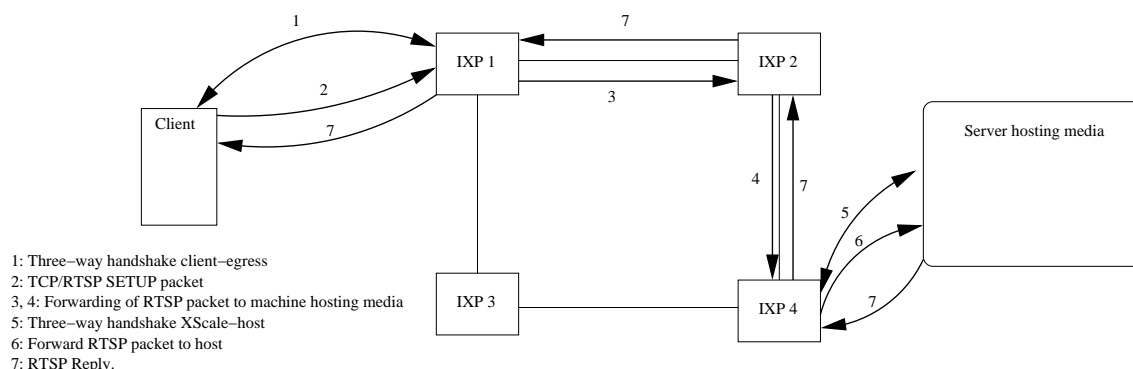


Figure 5.3: Design with forwarding of TCP to machine hosting the file.

In figure 5.3³, the set up of a stream with RTSP is shown. At the egress (IXP 1) the TCP three-way handshake is done (1). The RTSP SETUP packet is then received (2). The packet has to be inspected, and the file must be found in the cube. Info about the TCP session and the IC-session has to be kept both at the egress and at the machine that has the file. This is to ensure that packets from the outside can be routed to the correct host (egress IC information) and that correct headers and routing information can be supplied to build the packets that are outbound (media server IC information). After receiving the RTSP SETUP message, it is forwarded to the correct machine (3, 4). A three-way handshake is then done between the XScale and the host machine (5). The RTSP/TCP packet can be modified to conform with this TCP session, and forwarded to the host machine(6). Finally, RTSP replies can be processed the same way, and the server can start streaming RTP packets to the client (7).

When using raw sockets, received packets are forwarded both to the kernel and to your socket. This means that kernel-level filtering has to be introduced in order to stop all packets bound for your application from being answered by the kernel processes. The Monta Vista preview kit Linux kernel did not have support for iptables in the distributed version, so in order to make this work, the kernel had to be reconfigured and recompiled. Since no extra software was included with this distribution, a cross-compile of iptables had to be made for the XScale. These steps made the necessary filtering possible.

The advantages of this design approach would be that the media streaming server software would not have to be aware of the routing framework below, thus making it easy to replace the streaming server software. It became evident, however, that the process of maintaining several parallel TCP sessions for every TCP stream, together with the process of translating IP addresses and sequence numbers for the streamed packets would introduce a degree of complexity that would make the system error-prone and difficult to

³The process of locating the node that has the media file is not represented in the figure

implement. The conclusion was that the disadvantages of complexity and structure were greater than the advantages gained by this approach. The work with this design-option was discontinued before the routing framework was completed. This was due to the estimation that redesigning the software structure would result in an implementation that was simpler and more likely to be completed inside the time-limits of the thesis.

It was obvious that another approach had to be taken to make a design that was expandable, yet simple enough to implement given the time limit. The new design structure would have to give the egress node responsibility for handling TCP sessions. The next section presents the final design.

5.6 Egress TCP design

The complexity problems introduced by the TCP-translation solutions made it necessary to investigate alternative approaches. An obvious solution was to let the egress node handle the TCP connection, and let the rest of the control mechanism be implemented by the cube application. A consequence of this is that the media streaming server application on the host machine would have to be aware of the underlying mechanisms.

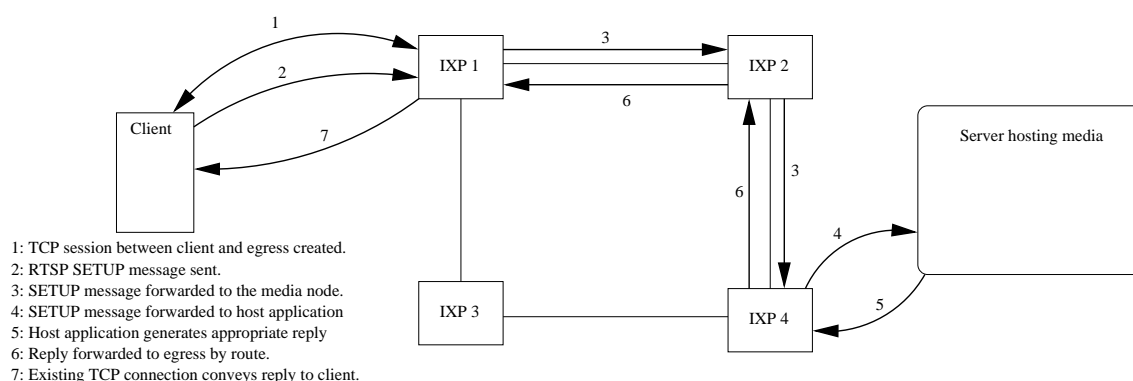


Figure 5.4: Design with egress handling TCP connections.

Figure 5.4 shows the setup process for a stream using the egress TCP design. The main differences from the raw socket design is that the node receiving the request handles all TCP operations. The first step, when a request is made, is to setup a TCP session, and give the correct replies to the client (1). When the client sends the RTSP SETUP request, the receiving node has to acknowledge the packet, and forward the request to the node that has the media that was requested⁴. (2, 3 and 4). The appropriate reply is generated

⁴The process of locating the media is not represented in the figure, but explained in more detail in section 5.7.4.

by the host application (5 and 6) and the egress node wraps the reply in the appropriate TCP header, and transmits it to the client (7).

Not only would this structure make several aspects of routing and TCP handling easier to implement, it would also make a better base for expanding the application when the primary bottleneck, namely host- μ Engine dataflow, could be overcome. Using direct memory access (DMA) to move data between the host and the μ Engines would be a natural way to improve on this bottleneck. This option will be further discussed in section 6.3. In order to achieve this, the host application (or the operating system) would have to be aware of the underlying structure. Another consideration is that the RTSP standard states that it should be possible to create several RTSP connections on one single TCP connection. This would be difficult to implement using the raw socket design.

One of the moments that were crucial in deciding on implementation strategies was how the routing itself would be handled in the cube. The next section will describe how this was designed and implemented.

5.7 Cube routing

Several advantages are obtained by implementing a solution in a distributed hypercube. In a cube with $n > 1$, there is more than one possible route between any two nodes. This gives us the opportunity to avoid congested routes, or to bypass possible broken links.

There are various routing strategies that can be applied to a network topology like the hypercube. One concern in a setting like this must be that the strategy will allow for a certain degree of adaptability. This is to be able to cope with eventual broken links and avoid heavily trafficked routes.

The choice fell on source routing as the routing strategy to implement. Using this strategy, the sender determines the route that has to be taken. The routing information is included in the packet header when the packet is transmitted. This eliminates the need for routing tables at every node. There were two main reasons for this choice. The solution would be fairly simple to implement, and the processing speed for each packet would be high.

There are also possibilities for extending the routing strategy with features that will add adaptivity. There is, amongst others, an internet-draft proposing a dynamic source routing strategy for ad-hoc networks [35]. Some strategies on how to improve adaptability for our source routing strategy will be found in the section 6.3.

Below is a description of how the source routing strategy is implemented in the current video cube routing framework.

- When a new setup request is received, the message is broadcasted through the cube.

- For each node the request visits, the port where it arrived is pushed on the routing field in the IC header.
- When the request arrives at the node that has the media file, it saves the route carried in the packet. It then sends an IC packet via the route to the egress node, instructing it to setup an IC session to be able to route new control packets directly to the node that has the media file.
- The stream then appends all RTP packets with an IC header that routes the packet to the egress node.
- The port number popped from the route field of the IC header can have one of the following values:
 - 0x0: Transmit the packet on port 0.
 - 0x1: Transmit the packet on port 1.
 - 0x2: Transmit the packet on port 2.
 - 0xff: Forward the packet to the XScale and pass it to `tcp_send`.

There are several advantages to this routing strategy. It is fairly simple to implement. It requires no routing information on intermediate nodes, only on the egress and the node that has the media file. It is very efficient to execute at each intermediate node, the microblock code only have to pop the port number and transmit the packet.

The disadvantages, however, touch some of the central concepts of implementing the system as a Hypercube. The current implementation does not take into consideration that the route have to be redirected if traffic patterns changes and transmission slows down. The option of redirecting the route if one or more links are broken also has to be added. It would also be profitable to be able to reserve resources along the route to guarantee that the stream will perform as wanted.

These disadvantages are subjects that should be explored further. Our source routing algorithm could probably be kept without losing the possibility of rerouting in the case of a broken link. A way of doing this could be to keep all routing messages that arrive at the node containing the media⁵ in the IC struct. If a broken link is detected, we could generate a message that tells the IC session to switch to the next route that has not been tried. For large cubes (big n-value) it would probably be more efficient to map possible routes on system startup (and maybe with a predefined interval), and store an index to a route in the IC session struct. It is also possible to transmit “ping” packets along different routes, to try to measure any difference in latency. This, however, can not be done too

⁵All possible paths from the egress to the node inside the limits of the defined TTL.

frequently without affecting the streaming performance that should be prioritized. An approach with reserving resources along a route would probably be a much better alternative than switching routes “on the fly”, given that this approach easily could lead to streams jumping back and forth between route alternatives. The ability to reserve resources along a route is also possible. Each node could keep a record of how many pass-through streams they support. When the packet locating the file is transmitted through the system, it could sum the pass-through number for each node it visits. The media server could then choose the route with the lowest sum.

5.7.1 Intra-cube header and extension

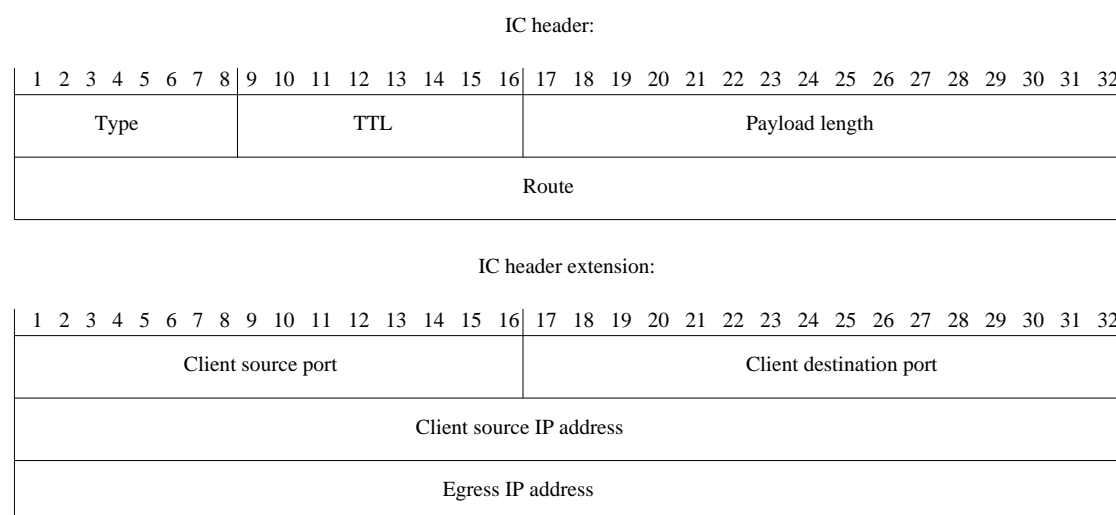


Figure 5.5: IC header and header extension format.

Figure 5.5 shows the format of the IC header. The basic IC header is designed to fulfill the following tasks:

- Separate setup and control packets from data packets.
- Limit the number of jumps a packet can circulate in the cube.
- Give information about payload length (when pushing data through the host-XScale socket).
- Provide routing information for the packet.

The purpose of the IC header extension is to carry information necessary to create IC session structs and to carry information about IC sessions to the egress when packets must be sent through `tcp_send` (port numbers and IP addresses is needed to look up the correct TCP struct).

5.7.2 IC packet types

In the system, as it is implemented in this thesis, there are five different IC packet types: `ICH_FIND_FILE`, `ICH_ROUTE_FB`, `ICH_CTRL_FB`, `ICH_CTRL_MSG`, and `ICH_RTP`. The first three of these packet types implicates that a header extension will follow immediately after the basic IC header.

The `ICH_FIND_FILE` is first created at the egress when a new `SETUP` request arrives. The IC header is appended with the IC header extension because the machine hosting the wanted media has to set up an IC session to be able to stream the content. The `ICH_FIND_FILE` packet is broadcasted throughout the cube until the node hosting the wanted media file is located.

When the file has been located, the node hosting the wanted media generates a packet of the type `ICH_ROUTE_FB`. This packet is transmitted back to the egress through the route found by the `ICH_FIND_FILE` packet. This packet also carries the header extension. This is because when the packet arrives at the egress, an egress IC session has to be created. This session is used to be able to route following RTSP control packets directly to the correct machine in the cube.

The `ICH_CTRL_MSG` package is created when the egress receives an RTSP control packet other than `SETUP`. The IC session that matches the request is found, and the packet is forwarded to the correct host through the cube.

When an RTSP control message is processed by the media streaming server, a reply is generated, and has to be sent to the egress through the cube. These packets are wrapped in IC headers of the `ICH_CTRL_FB` type. A header extension is used with these packets to look up the correct TCP session at the egress. This could have been avoided by adding a IC session index field to the regular IC header, and use this to look up the egress IC session. There are two reasons why this solution was implemented:

1. The main bulk of data packets does not need the extra information. Transmitting it through the cube with each RTP packet would be a waste of resources.
2. The RTSP traffic load in the cube is negligible compared to the RTP data, so the extra 24 bytes on these packets will not make a noticeable impact on cube performance.

The last IC packet type is the ICH_RTP. This is by far the most common in the cube since it wraps all RTP data packets. The purpose of this header type is to bring the payload to the egress, and subsequently to the client, as fast as possible.

5.7.3 Intra-cube session

The purpose of the IC session is to keep the necessary data to handle the functionality needed by the cube infrastructure. As of now, the following data is kept by the IC session struct as shown in figure 5.6:

- **Status:** Information concerning the state of the session. Used to identify active sessions, sessions about to be set up and closed sessions.
- **Route:** The current route the packets have to travel. For the egress node this is the route to the machine serving the media, for the media server, it is the route to the egress.
- **Client source port and server port:** Used to generate UDP headers.
- **Source and destination IP addresses:** Used to build IP headers.
- **Source and destination MAC addresses:** Used to build ethernet headers.

```
typedef struct{
    uint8_t    status;
    uint32_t   route; /* Route */
    uint16_t   sport; /* Client source port */
    uint16_t   dport; /* Server port */

    uint32_t   saddr; /* Client source IP addr. */
    uint32_t   daddr; /* Egress IP addr */

    char       eth_src[6]; /* Ethernet src address */
    char       eth_dst[6]; /* Ethernet src address */
} IC_session_t;
```

Figure 5.6: IC session struct.

The IC session is first created at the server hosting the media when the requested media file is found. The egress saves the session data when the route is set up, and the host is updated

with the session data when the SETUP message is forwarded there. The way the program is implemented now, the media server on the host builds RTP, UDP, IP, ethernet and IC headers for the data packets. This means that the XScale IC session on the media server machine is superfluous. The reason it was implemented in this way was to facilitate the transition from a solution where the server handles header generation, to a future design where this functionality is handled by μ Engines (aided by the XScale). The different aspects of this are discussed in more detail in section 6.3.

Regarding the general design structure and the routing framework that had been decided upon, tasks had to be divided between the different hardware layers available (host, XScale or μ Engines). The next section will describe where the different functionality was placed, and the advantages and tradeoffs this led to.

5.7.4 Partitioning of tasks

The decisions on where to locate the different application parts was based on the general programming paradigms for networking applications (see sections 2.3.3) and on the results of the preliminary tests done on the IXP2400 hardware (see section 3). The way tasks are divided between the different subcomponents in this design is as follows:

- Host machine: The host machine runs the streaming server application. It receives RTSP requests, generates replies, creates streaming threads and forwards data to the XScale. The media data is stored on the host disk, and the host server is responsible for retrieving it.
- XScale core: The XScale core handles the TCP connections (if egress). It receives data from the host, and forwards it to the μ Engines for routing to the egress node. It also receives control packets from μ Engines, and forwards them to the host.
- μ Engines: The μ Engines receives control packets on the egress port. They filter all packets not TCP/IP bound for port 9070⁶ and forwards control packets to the XScale for TCP handling. The μ Engines also receive IC packets and route them according to the port number given in the route header field.

The IP header switch tests (see section 3.1) showed that the μ Engines was very efficient on a simple task that implicated manipulating packet headers. Consequently, it was decided to implement the bulk of the routing functionality on the μ Engines. Filtering of unwanted data that could arrive on the egress port was also to be implemented here.

⁶This port is used because Komssys [36] at first was intended to be the streaming server application run on the host. The default RTSP listening port for this system is TCP/9070.

Since TCP related traffic would represent a small part of the total, it was acceptable to place the main functionality on the XScale. This decision was also supported by the fact that maintaining data structures is a less complicated task on the XScale.

The media streaming application, performing a processor-intensive task, was to be implemented on the host. It was decided that the streaming application would also generate all the necessary headers to leave all but the routing to the μ Engines. When direct transfer of data from the host to the μ Engines is implemented, it will be natural to let the microblocks generate some of the headers. As it was, however, the data was moved to the μ Engines via the XScale, consuming a lot of processing power from both μ Engines and XScale alike. To counter this effect, as much of the packet generating functionality as possible was placed on the Host.

Next, descriptions of how the aforementioned components were implemented on the different hardware layers follow. The structure of the components and the main challenges regarding each step will also be discussed.

5.7.5 μ Engine tasks

The *ICrouter* microblock used in the egress TCP design consisted of two active contexts. This was to better exploit the available resources by switching contexts when memory-accesses or other operations that leave the code waiting for results are performed.

First context The packet handling mechanisms used for the first context of the *ICrouter* microblock are shown in figure 5.7. When a packet arrives at any port, the first thing to be checked is whether the packet is received on the egress port (port 0) or any of the intra-cube ports (port 1 or 2). If the packet is received on the egress port, it is sent through a filter dropping all packets not TCP/IP bound for port 9070, thus eliminating all traffic except the packets bound for this server solution⁷. If the given criteria is met, the packet is forwarded to the XScale for TCP processing. If the packet received arrived at an intra-cube port (port 1 or 2), it is by definition an IC packet with an IC header (see section 5.7.1) and must be routed according to the given routing algorithm. If the packet is of the type ICH_FIND_FILE, it is sent to locate the host containing the wanted media. Such packets are forwarded to the XScale which performs the rest of this processing. All other packets are to be routed to the port given by the route field in the IC-header. The Time To Live (TTL) is decremented. If the TTL is zero, the packet is dropped. If the port number popped from the IC header route-field is 0 (egress port), the IC-header is stripped and the packet is transmitted. Packets bound for port one and two is transmitted on the popped port.

⁷To support RTCP reports, the filtering mechanism has to be able to open up the corresponding port for the connections, depending on the RTP transport method. This is as of yet not implemented.

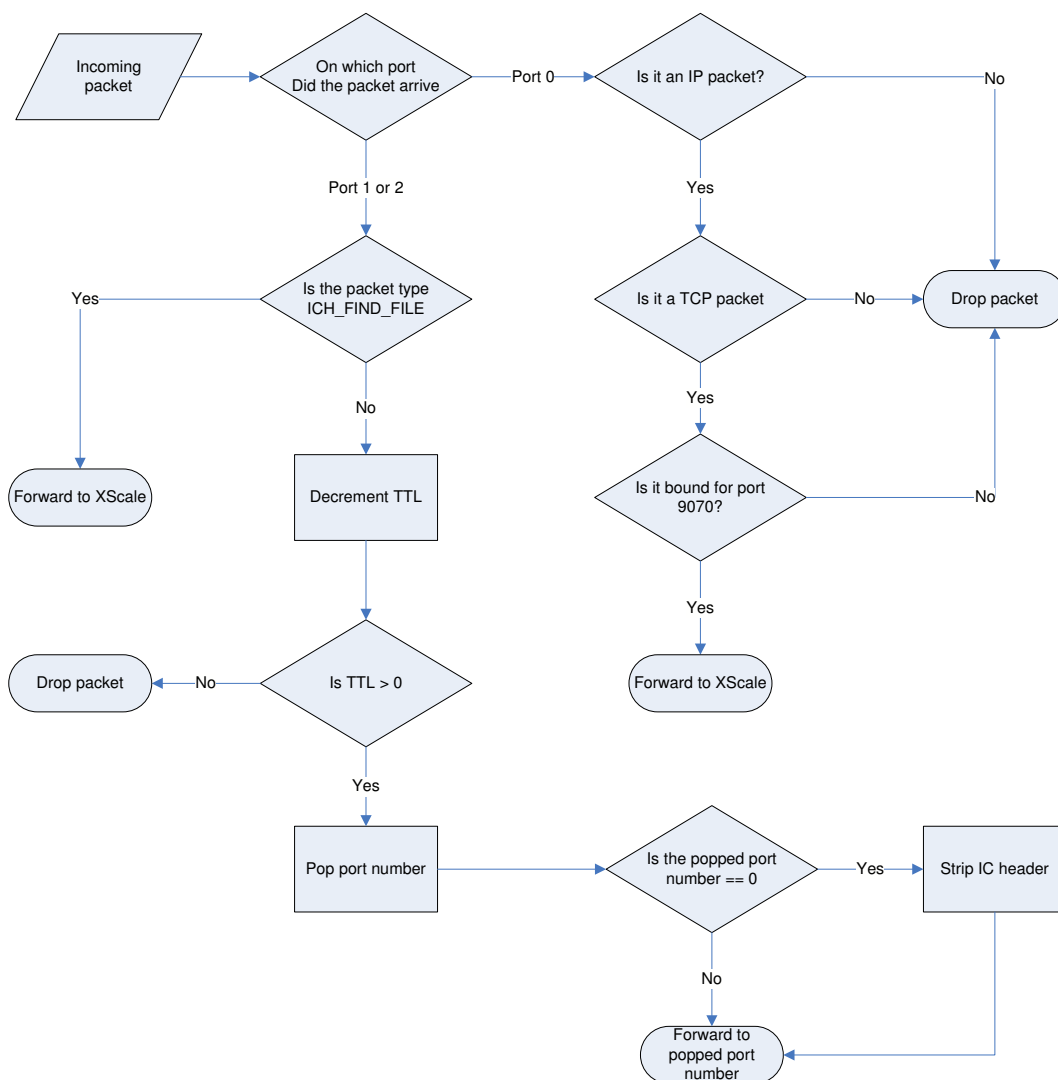


Figure 5.7: Flow chart for *ICrouter* microblock when a packet is received.

Second context The second context of the *ICrouter* μ Engine receives buffer handles from the XScale, and transmits the packets referred to on the port specified by the packet metadata. All setup of packet and header is done by the media server on the host, or by the XScale. Functionality for batch processing of packets⁸ is ready to use, but not yet supported by the XScale code.

⁸In this case batch processing means caching packets and sending them to the μ Engines collectively to achieve higher transfer rates and more efficient resource utilization.

5.7.6 XScale tasks

The IXA SDK has support for spawning specialized threads that will call a specified method upon receiving one of two interrupts branded “INT_A” and “INT_B”. In this implementation, only one such thread for handling data sent by the *ICrouter* microblock to the XScale is used. The procedure for handling packets forwarded from the XScale is described schematically in figure 5.8. The first thing that is done, as in the *ICrouter* microblock, is to check which port the packet was received on.

If a packet is received on the egress port (port 0), we know that it is a TCP packet, and it will be passed to the `tcp_recv` method. This method returns one of three values:

- **B_TCP_NO_DATA:** The received packet is successfully handled, but contained no data. This is for example the case for SYN, ACK or FIN packets.
- **B_TCP_ERROR:** The packet is not successfully processed. This can happen if the TCP session is not found, or an error is detected with respect to sequence numbering.
- **B_TCP_DATA:** There is a data payload in the TCP packet. The method also sets references to the payload and the payload size in order to enable data processing.

Int_A thread processing When a packet is sent to the XScale, the `Int_A` procedure is invoked. If the input port is 0, we call `tcp_recv`. If `tcp_recv` returns `B_TCP_ERROR` or `B_TCP_NO_DATA`, the TCP implementation have processed all required replies, and the call returns without providing any data payload. If data is received, we know that it is an RTSP control packet. The RTSP command will then have to be identified. If the packet is a SETUP request, the program will try to find the file on the local machine. If the file is not found on the machine, an IC header and an IC header extension (containing information that is required to setup a new IC session) are prepended to the RTSP packet (for details, see section 5.7.1). The port number that the packet is received on is pushed on the route field of the IC header, and the packet is sent on both IC ports (port 1 and 2). An egress IC session is not created yet, because the route to be used is determined by the packet that arrives first at the host machine that has the requested file. If we assume that the packet travelling the path with the least load is the first to arrive, this will help the system give the stream an effective route during the current conditions. If the file is found on the machine, however, both a server IC session, and an egress IC session is created (given that this machine is both server and egress). The RTSP packet is prepended with an IC header and header extension, and forwarded to the host on the socket interface.

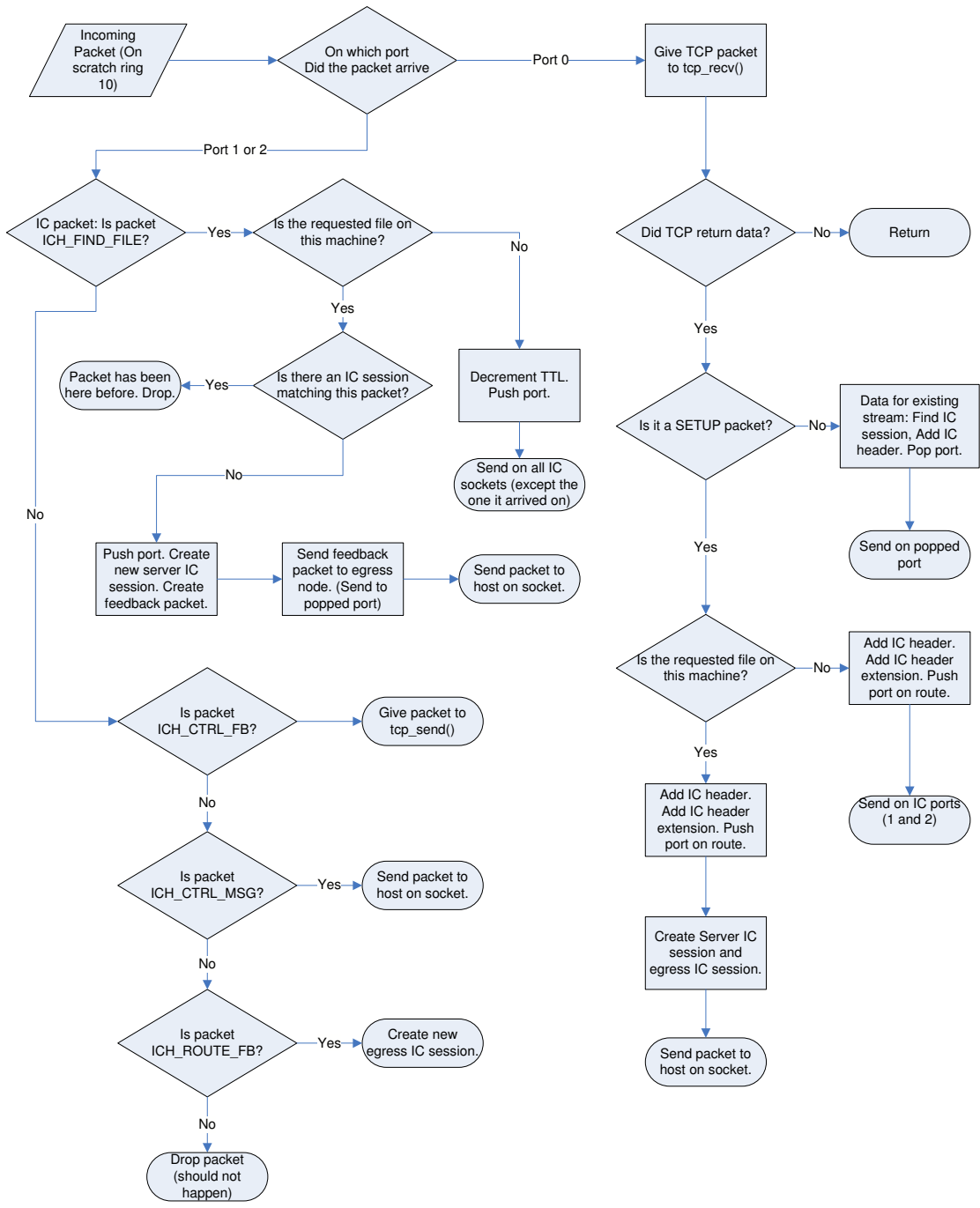


Figure 5.8: Flow chart for XScale int_A thread.

If a packet is received on one of the IC interfaces (port 1 or 2), we know that the packet is an intra-cube packet and has an IC header. The IC packet received can be one of the five types mentioned in section 5.7.1.

As shown in figure 5.8, the handling of each packet depends upon which IC header type the packet has. If the packet is tagged “ICH_FIND_FILE”, the next step is to check if the file resides on this machine. If the file is not found, the TTL is decremented, the port the packet is received on is pushed on the route, and the packet is transmitted on all IC ports except the one it arrives on. If the file is found on this machine, a search is done to try to find an IC session (see section 5.7.3) that matches. If a matching IC session is found, the same setup packet has arrived before through another route, and the packet is dropped. If no matching IC session is found, a new server IC session is created. This session keeps the return route that this stream is to use. To make sure that control messages arrive at the server by the correct route, a feedback IC message is created (“ ICH_ROUTE_FB”) and sent through the route to the egress. The original SETUP packet is then passed to the host on the XScale-host socket.

If the packet received is of the type “ICH_CTRL_FB”, the payload is an RTSP reply message. The payload is passed to `tcp_send` along with IP addresses and TCP ports (to be able to find the correct TCP session), and transmitted to the client.

When a packet of the “ICH_CTRL_MSG” type arrives, the payload is an RTSP request bound for the media server. This packet is directly forwarded to the host on the XScale-host socket.

The last IC packet type handled by the “int A” thread is the “ICH_ROUTE_FB”. This packet type is transmitted to the egress machine after a route has been found and the streaming server has created its IC session. An egress IC session used to guide incoming RTSP packets to the correct machine is created, and the packet is discarded.

If a packet arrives that matches none of the above categories, the packet is dropped, but this should not happen if the system works correctly. The reason “ICH_RTP” never reaches this stage of processing is that once sent into the routing framework, these packets never surface to be examined by the XScale, but are forwarded straight to the egress node.

5.7.7 Host-XScale data path

In addition to the “int A” thread, there is a program loop running on the XScale that is polling the host-XScale socket for incoming data. This data is formatted in such a way that an IC header always arrives first. This header contains a byte count for the packet payload, so that the packets can be distinguished from one another. If the packet is of the type ICH_ROUTE_FB or ICH_CTRL_FB, we know that an IC header extension comes after the header (see section 5.7.1). If the destination port (popped from route) is 0xff and

the packet is of the type ICH_CTRL_FB, the packet is passed to `tcp_send`. In all other cases, the packet is routed as is on the port popped from the route field of the IC header.

5.7.8 Host tasks

The host application is implemented in userspace with a TCP socket providing the data stream to the XScale. For each RTP stream, a processing thread is generated. It opens the file, prepends each data chunk with the needed headers, and transmits the data on the host-XScale socket.

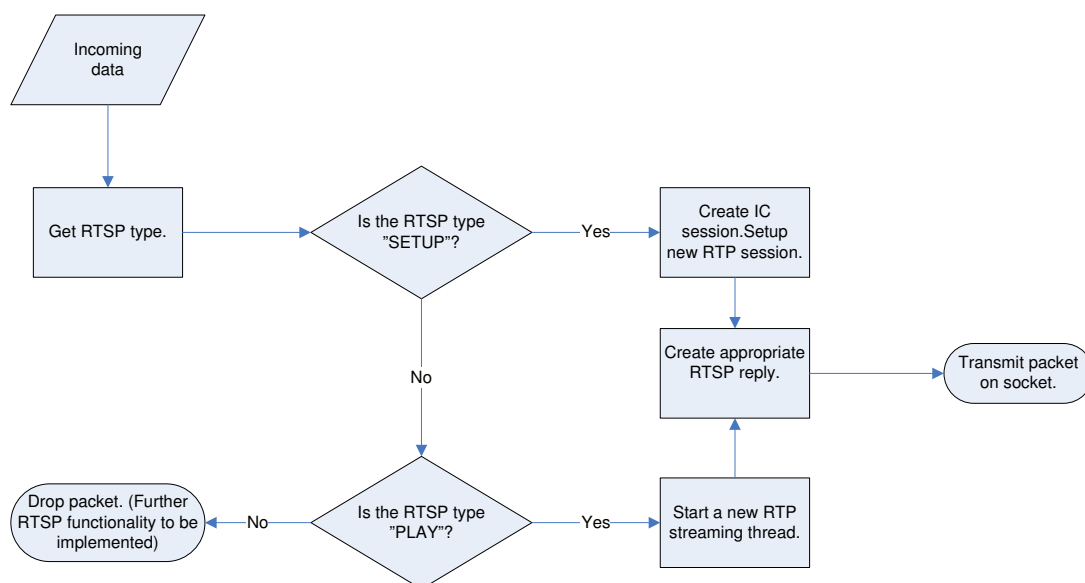


Figure 5.9: Flow chart for Host when receiving data from XScale.

Figure 5.9 shows the data flow for the host application when an incoming packet is received. When data is received from the XScale, the first step is to get the payload size and the IC type. Any IC header extension is also read. The RTSP packet type is then read. If it is a SETUP packet, a new IC session is created, and a new RTP session (see section 5.8) is prepared. We already know that the file is located on this machine. This was checked on the XScale. The RTSP reply to the host is then generated stating the server port numbers, the reply is wrapped in an IC header, and the IC packet is sent on the socket. If a PLAY message arrives, an RTP streaming thread is created. The program the generates an RTSP reply message, and transmits the reply on the socket to the XScale.

5.8 RTP-session

Though the current version of the cube does not implement genuine RTP functionality, it keeps some information about each RTP stream. This information is kept in the RTP session struct, which is comprised of the following components:

- **Thread ID:** An identifier for the streaming thread. Used to send signals or forcefully shut down the stream.
- **Client RTP Port:** The port negotiated by RTSP to deliver data to the client.
- **Client RTCP port:** The port negotiated by RTSP to deliver RTCP data. Not in use as of this version.
- **Server RTP port:** The server port delivering the stream. In this implementation, this port is a “phoney”, because a regular socket mechanism is not used.
- **Server RTCP port:** The port used by the server to send RTCP packet. Not in use as of this version.
- **Index of IC session:** The index of the IC session corresponding to this RTP session. Used to look up IP and MAC addresses.
- **Filename:** The name of the wanted media file. Used to open the file for streaming at the server.

As the description shows, there are many crucial RTP components that is not yet implemented. These include RTCP feedback and rewinding or forwarding the stream. Eventual RTCP packets will be discarded by the μ Engine code (see section 5.7.5). The current implementation supports the streaming of a file from beginning to end, as is, and with a fixed bitrate. The reason for this is that the focus has been on implementing a working routing framework for the cube and exploring the benefits of offloading this kind of routing functionality on the IXP μ Engines.

When an RTSP PLAY message is received, the media server on the host creates a new media streaming thread. The first operation that is performed when a new thread is started is to build the RTP, UDP, IP, ethernet and IC headers that will wrap the media data. Since only a few of the header fields are changed from packet to packet, these headers are saved, rather than creating the entire header again for each data packet. This is done by having only one transmission buffer of the current MTU size. The headers are kept statically in the beginning of the buffer, and the data from disk is copied into the payload part. For each of the data chunks read from disk, the length field is updated in the IC, IP and UDP header. The sequence number and timestamp are updated in the RTP header. Then, checksums are calculated for UDP and IP, and the packet is sent to the XScale on the socket. The thread then goes to sleep until it is time for another packet to be sent.

5.9 Microblock designs

How the different tasks should be prioritized between the microblocks was an issue that strongly depended on the choice of routing algorithm. Another deciding factor was the method of moving data from the host to the μ Engines. Given the possibility of using 8 μ Engines in parallel, there were many tasks that could be distributed and optimized by distributing them on several μ Engines. The fact that each μ Engine also supports eight program threads, further enhanced this possibility.

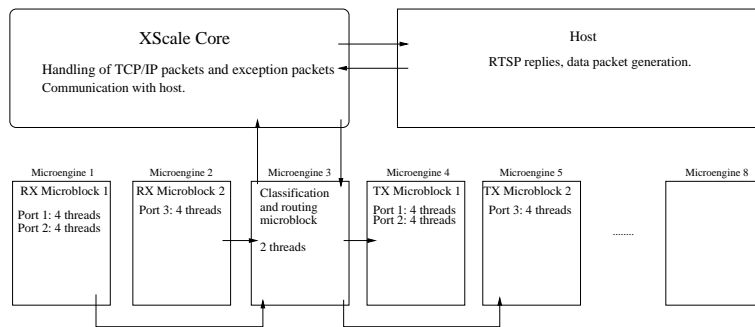
If the network processing application that is to be written expects a very heavy load, steps can be taken to enhance RX and TX performance. The reference designs of RX and TX microblocks provided by Intel give the programmer the possibility of distributing the task (RX or TX) on two microblocks (figure 5.10(a)). Four threads are assigned to handling each of the ports for receive and transmit, giving the performance of the microblocks a potential efficiency increase in comparison to the one-microblock RX and TX designs. The actual effects of this enhancement depends on the transmission medium. Implementing this solution reduces the chance of queues that in the worst case can lead to packets being dropped. This measure would, however, be in vain if the packet processing between RX and TX is slower than the receive and transmit process. In order to meet with this problem, steps should be taken to optimize other packet processing tasks. Figure 5.10(b) shows the same application with the two main tasks split between two μ Engines.

As long as each step in the pipeline between RX and TX is kept at a level where fewer cycles is consumed for each packet than would the RX and TX block, the packet flow would be optimal. From this a rule of thumb for microblock design can be derived:

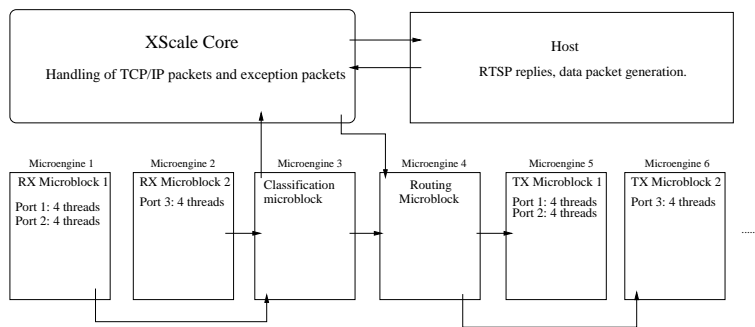
$$C_{PP} \leq C_{TX} \leq C_{RX} \quad (5.1)$$

where C represents cycles used, PP is each individual microblock used in packet processing, TX is transmit microblock(s) and RX is receive microblock(s). In practice this means that the system should be able to transmit packets at least as quick as they can be received, and every individual stage of the packet processing pipeline in between should spend less time on processing one packet than the transmit block(s).

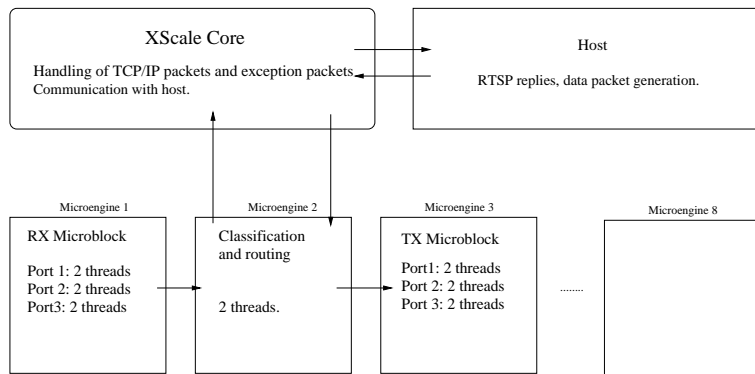
The handling of exception packets and other traffic that has to be forwarded to the XScale will of course impede the packet processing pipeline severely. It is therefore crucial to make sure that the bulk of the packet traffic is handled at the μ Engines. In the implemented system, however, the XScale-host transfer bottleneck slows the system down to a point where the μ Engines can easily handle the traffic load it is presented with. The packets are, as of now, being generated on the host, then moved by TCP socket to the XScale for then to be sent to the μ Engines. Until a more efficient way of moving data from the host to the μ Engines can be implemented (see section 6.3), it is sufficient to keep the RX and



(a) RX and TX blocks split on two μ Engines.



(b) Two microblocks for RX, TX and packet processing.



(c) The current microblock design.

Figure 5.10: Three microblock structuring alternatives.

TX functionality on one microblock each, and the packet processing on one microblock utilizing two program threads. This design is visualized in figure 5.10(c).

When data can be moved directly from the server to the μ Engines, more of the functionality of generating packet headers can be delegated to the μ Engines. This will be discussed in more detail in section 6.3.

5.10 Evaluation

This set of tests is done with the purpose of measuring the forwarding time of packets routed in the cube server implementation. Since the performance of data transport from the host down to the μ Engines yields bitrates of less than 10Mbps due to the lack of efficient host- μ Engine communications, a test of the system as a whole would give little or no information about how packets would move through a heavily loaded system once this bottleneck is removed. The test focus is, therefore, on the underlying routing mechanisms. The results from the μ Engine measurements are then compared to packets routed through a conventional switched network to try to determine the possible performance gain by implementing the routing in this network topology using IXP cards.

5.10.1 Packets routed through a switched network

The first test emulates a server cluster interconnected by a dedicated conventional switched network. Figure 5.11 shows the test setup for this configuration. The test is conducted by transmitting ICMP ping requests from the “Cube1” host destined for “Cube2”. The forwarding machines (“Cube2”, “Cube3” and “Cube4”) have iptables DNAT and SNAT rules that forwards the packet to the next host in the chain, thus sending the packets in a complete circle through the hosts and back to the sender (“Cube1”). There is an iptables rule firewalling ICMP requests in order to prevent the kernel from replying to the messages.

The packets that are sent from “Cube1” are timed on the way out by tcpdump [60]. When they arrive again from “Cube4”, the arrival time is recorded. The difference between the send and arrival time give us the time used for routing the packets through the four hosts.

5.10.2 IXP hypercube routing tests

To be able to measure the amount of time spent on the routing process in the IXP cube implementation, a special timestamping mechanism had to be developed. Since every host relates its timestams to its own time domain, the packets had to be timestamped on

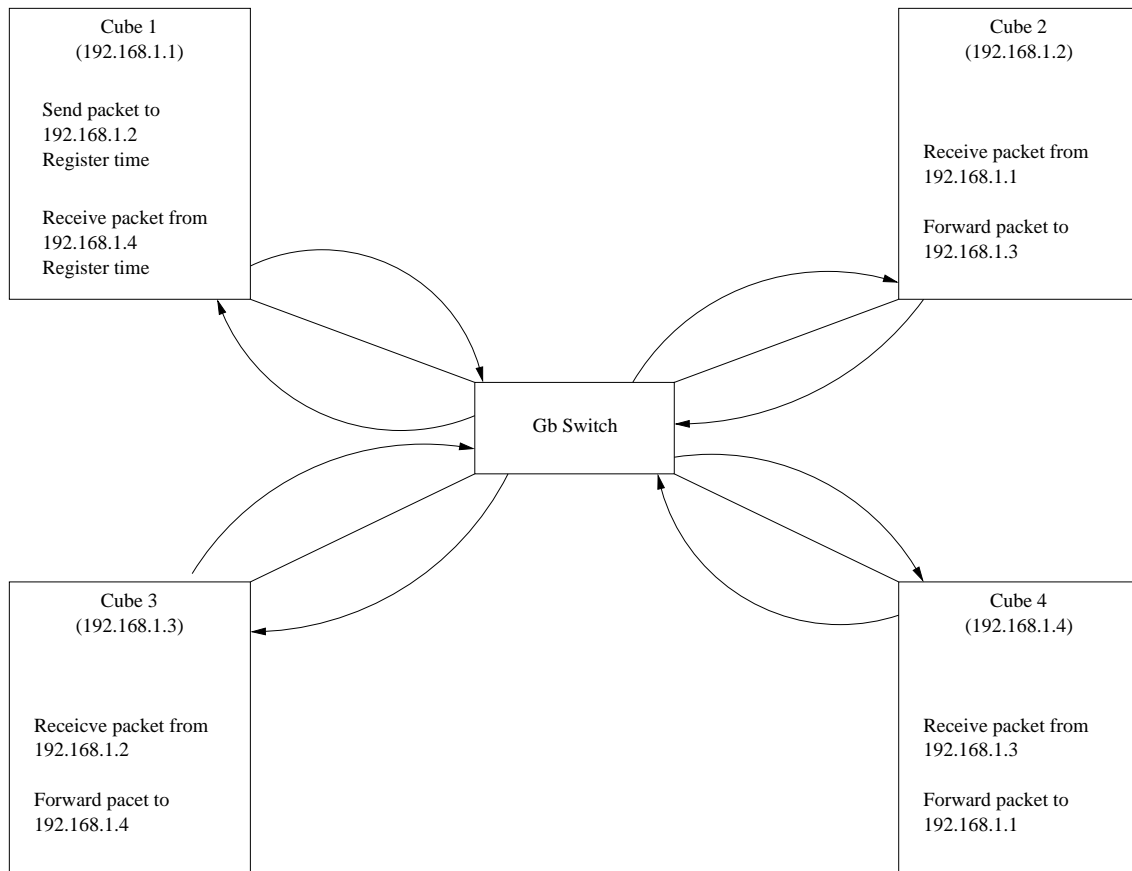


Figure 5.11: Test setup: Routing using switched network.

the same host. To accurately be able to measure the routing mechanisms, the timestamps should be done on μ Engines. Tests have shown that even when synchronized at a designated point of code execution, there are still differences in the clocks of each microengine. This implies that a timestamp has to be made by the same IXP μ Engine on the same host.

This challenge was solved by creating a new IC header type called "ICH_TIMER". Packets given this designation is to be timestamped, traverse the entire "cube" and return to the sending host. Upon return, the packet is timestamped again on the same μ Engine as on the way out. This give us an accurate time for the packet to traverse the entire cube (4 hosts) in a manner similar to the switched network.

The route given to the "ICH_TIMER" packet is processed in the same way as every other IC packet to ensure a realistic estimate of the processing time.

Figure 5.12 shows the test setup for this process. The "ICH_TIMER" packet is created on the host, and passed to the XScale and consequently to the microengines in the same

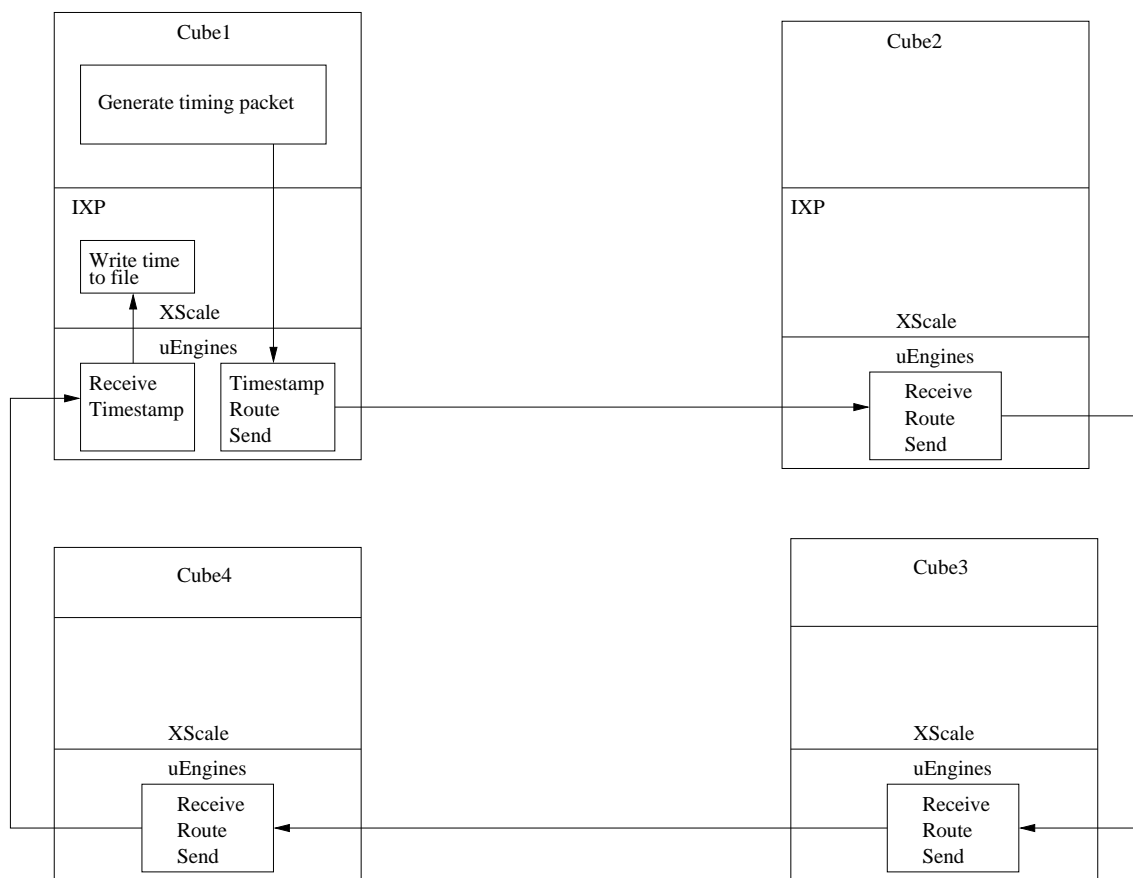


Figure 5.12: Test setup: Timing of routing in the IXP hypercube.

manner as for instance an RTP packet in the *VS³* server solution. The packet is then forwarded to the μ Engine where it is timestamped. The send port is then popped from the IC header route field, and the packet is sent to the next host.

On the “Cube2”, “Cube3” and “Cube4” machines the μ Engines receive the packet, checks the packet type to confirm that it is to be routed, routes the packet to the port given by the IC header route field, and transmits the packet.

When the packet arrives back at the “Cube1” machine, the route field of the IC header has no valid port number, telling the system to forward it to the XScale, and the fact that the packet is of the type “ICH_TIMER” tells the μ Engine to timestamp the packet before forwarding it to the XScale. The timestamp is made to not overwrite the timestamp made on the way out.

On the XScale the two timestamps are read. The difference is then calculated, and printed to a file on the mounted NFS filesystem used by the IXP card. The packet sizes used

for the test was the same as for the switched network test. The results from this test are presented in table 5.1 with the results from the switched network test.

5.10.3 Comparison of routing times

Table 5.1 shows the statistical data of this test and the test discussed in section 5.10.2. Both tests are performed with ICMP packets of 98 bytes and 1500 bytes.

The average times for completing one “circuit” in the machine topology is graphically displayed in figure 5.13. As expected, the total processing time is radically larger for the larger packet size (1500 bytes). We can also see that the benefit of implementing routing in a dedicated interconnected topology (like the hypercube), compared to using a switched network structure, is significant with respect to inter-server communication. This can be attributed not only to the extra transmission times introduced by having the packet move through the switch, but also to the overhead that the layered handling structure of the kernel requires. Since this does not have to be supported by the intra-cube routing mechanisms of the IXP hypercube implementation, the routing costs diminishes. This effect is magnified by the simple nature of the routing mechanism.

Packet size	Max time	Min time	Avg time	Median	Std. dev.	Count
Switched network times						
98 Bytes	9022 μ s	9 μ s	135,98 μ s	129 μ s	71,481 μ s	99928
1500 Bytes	2495 μ s	261 μ s	265.6 μ s	265 μ s	8.2152 μ s	99991
IXP cube routing times						
98 Bytes	30 μ s	27 μ s	28.138 μ s	28 μ s	0.3544 μ s	63965
1500 Bytes	109 μ s	106 μ s	107.0034 μ s	107 μ s	0.1367 μ s	63965

Table 5.1: Routing times for IXP cube and switched network routing.

For 98 bytes packets, the switched network implementation is 483% more time consuming than the IXP hypercube solution (a difference of 108 μ s in average). For 1500 bytes packets, the switched network implementation time is 248% higher than the IXP hypercube time (a difference of 159 μ s in average). This can be explained by the copy operations occupying a greater share of the total time consumed.

Based on these numbers we can say that a directly interconnected network topology will have a vast advantage in comparison to a switched network solution. An increase in routing algorithm complexity will probably make the difference smaller, but still large enough to make a significant difference.

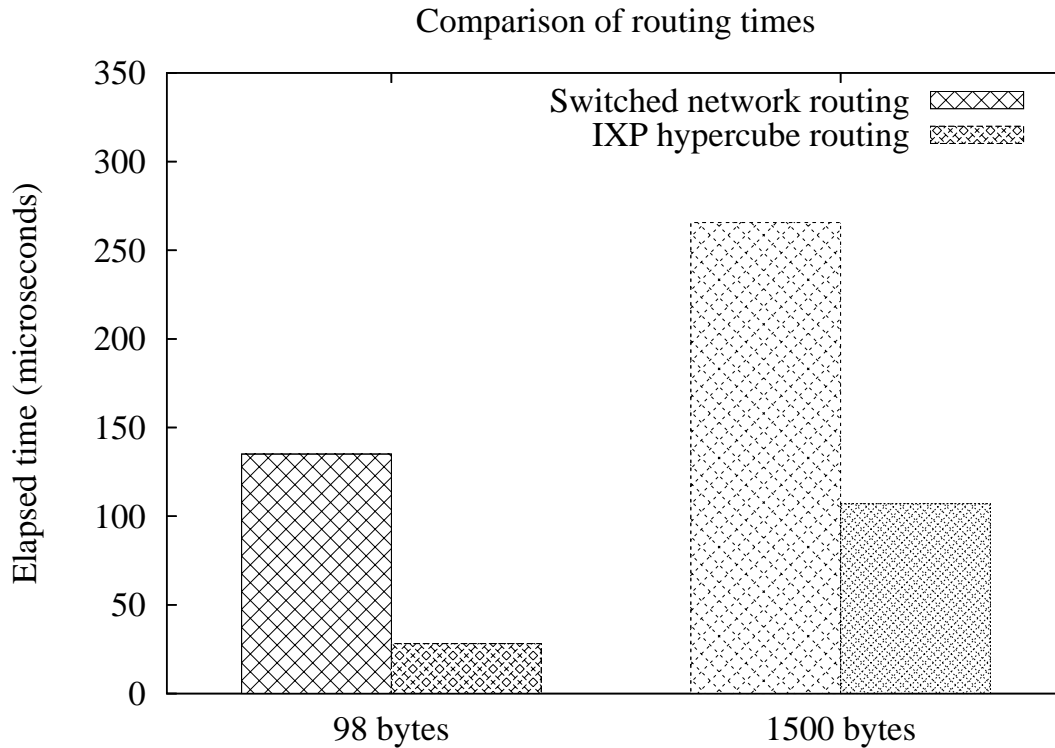


Figure 5.13: Comparison of average routing times.

5.10.4 Processing times given increased load

To try to find out what level of impact an increase in load i.e., will have on the IXP hypercube system, the loop routing test was done again, this time with packets being generated on more than one of the hosts. For each new host generating packets, the load on all hosts increases because all packets have to be forwarded through all hosts. Although the bit rate of each packet stream is not high (about 2.4Mb/s) it will give an indication of how load times will behave.

Table 5.2 shows the statistical data from this test transmitting 1500 bytes packets. To maximize the bitrate given the sending intervals permitted by the host-ixp implementation, a packet size of 1500 bytes was used. The reason why the number of measurements (count) increases by such a degree is that more transmitting nodes give more valid samples. The number of measurements for four nodes should be approximately the same as four times the number of measurements for one node.

Figure 5.14 plots the time spent on traversing the loop given the increased load. The

Transmitting nodes	Max time	Min time	Avg time	Median	Std. dev.	Count
1	109 μ s	106 μ s	107.0034 μ s	107 μ s	0.1367 μ s	63965
2	141 μ s	106 μ s	107.36 μ s	107 μ s	2.6691 μ s	121395
3	142 μ s	106 μ s	108.47 μ s	108 μ s	3.532 μ s	191846
4	202 μ s	106 μ s	108.92 μ s	108 μ s	4.8546 μ s	255709

Table 5.2: Routing times on μ Engines given increased load.

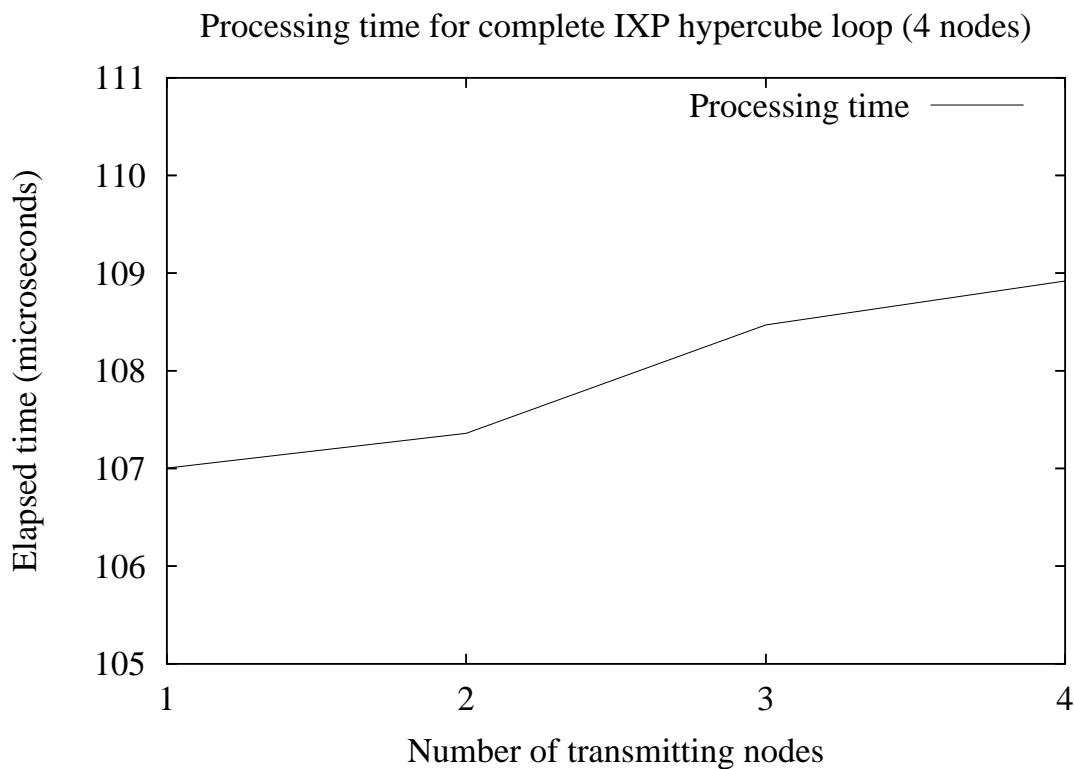


Figure 5.14: Development of processing time with increased load.

graph shows an increase in load time, but not a significant one⁹. This indicates that the system should be able to handle much higher loads, before having to drop packets. The chipset is built to be able to handle 1Gb/s on the three interfaces, so given that an efficient processing pipeline can be implemented, the platform should be able to accommodate high loads.

⁹Note that the y-axis starts on 105, and that the actual difference is below 2 μ s

5.11 Summary

This chapter has described the process of implementing our VS^3 video server cube with routing functionality based on the capabilities of IXP2400 cards. Several designs were proposed to achieve this. The first design aimed at using an existing server application on the host. This was discarded in our prototype, due to the complexity that this would introduce to packet routing and translation, but should be considered in a later stage. The design that was finally chosen reused the existing RX and TX components from the static forward application. The rest of the components had to be built from scratch. The streaming server application running on the host had to be aware of the underlying routing framework. This was to ensure that a future implementation using DMA to effectively move the data to the μ Engines could reuse the code (see section 6.3), and to avoid unnecessary complexity.

A video server hypercube was implemented with routing functionality located on the μ Engines, TCP support and uplink to host on the XScale, and a streaming server on the host. The server solution was tested with several streams delivered from the different content servers through the same cube egress. Setup and routing worked as expected, and the video data was delivered as intended. The bitrates, however, were hampered by the host-XScale bottleneck, thus impeding the performance of the server.

Two tests have been performed on the routing framework. The first was a measurement of the time used by a data packet to traverse a loop of four machines in the system. A similar test has been run on a group of Linux hosts linked in a switched network. The second test measured the increased processing time when doubling and quadrupling the load.

These test show, as expected, that the directly connected hypercube topology delivered the routed packets much faster than the switched network could manage. This is a strong argument for implementing data-intensive servers like a multimedia streaming server in such a topology. When increasing the load, there was a slight increase in processing time. To be able to put serious pressure on the routing system, however, packet generators have to be applied, providing heavier loads.

Chapter 6

Conclusion

This chapter concludes the master thesis. We will present a short summary of what has been done, then the most important results of this work, and finally some key issues that should be pursued in the future.

6.1 Summary

In this thesis, we have investigated the behaviour of Intel IXP network processors, with focus on the IXP2400 and the IXA SDK. A series of test implementations have been done to find the strengths and weaknesses of both the hardware and software platform. Tests have been run on the different applications to find out whether offloading network tasks to this platform can be valuable. The knowledge gained from the experiments has been used to construct a multimedia server cube solution, similar to the nCube n4x, using network processors for offloading routing tasks.

6.2 Results

From the work done on the ENP2611 and our *VS*³ video server cube we have gained valuable knowledge of hardware and software issues related to programming the IXP2400 NPU. We have also run a number of tests that were done to evaluate performance of the implemented systems. The tests performed on the applications showed that offloading network operations on a Linux host frees a noticeable amount of resources, even for simple networking tasks. We could also observe that the time needed to perform a simple task of editing some protocol header fields was performed more efficiently on the NPU than on a Linux host. We conclude that wire-speed processing should, to as large degree as

possible, be kept on the μ Engines, while processor-intensive tasks and control operations should be put on the XScale. Measurements showed that when using the interrupt “int A” method of passing packet references to the XScale, the processing times were only slightly higher than the results for the Linux host and the μ Engines. This shows that the XScale can be used quite efficiently without impeding performance significantly.

From the VS^3 implementation we found that the greatest challenge was to move data efficiently from the host to the IXP as there is no support for DMA transfers yet. However, focusing on the internal cube communication, tests on the implemented routing framework showed that this method of intra-cube routing was efficient. Tests showed that, compared to implementing the server cluster in a switched network, the hypercube topology with direct interconnection was far more efficient.

In conclusion, our prototype shows that the IXP2400 NPU can efficiently offload the host machine and provide a favourable way of implementing closely interconnected servers in a cube topology.

6.3 Further work

In this section, we present a short description of what we deem to be the most important issues to address in future work on this subject.

During the course of implementing the VS^3 , it soon became clear that the main bottleneck of the system would be the host-IXP data transfer. Hardware documentation and experiences from research projects elsewhere using the same hardware tells us that DMA can be used to move data directly from the host to the μ Engines over the PCI bus. This can be done without having to go through the XScale as of the present implementation. This would help the system in two ways:

1. It would dramatically enhance transfer speeds. If a 64 bit bus could be used effectively, it would be able to feed enough data to fully utilize the optical interfaces of the IXP card.
2. It would enable the host to transmit the data only, leaving some or all of the header generation operations to be done on μ Engines. Allocating one or more μ Engines to the task of packet preparation could be done without stealing resources from neither host nor existing IXP functions.

Work is currently in progress in our research group to find methods of doing this efficiently on the IXP2400.

To be able to reap the full benefits of RTP/RTCP and RTSP on the VS^3 , many improvements have to be done. As of this implementation, the system only has support for RTSP

SETUP and PLAY, and basic RTP sending. To enable further control of the streaming, we should complete the RTSP and RTP implementations. To be able to receive RTCP receiver reports, the filter in egress node μ Engines must have functionality to dynamically update which ports (and protocols) should be allowed to enter the cube, according to the port numbers conveyed by the RTSP SETUP exchange.

The source routing algorithm should be extended to find a new route if a link is broken during streaming. Another feature that should be implemented is the ability to reserve resources along the route, and to choose the optimal route based on recorded reservations.

Finally, there were some tests that could have been performed to support the data already collected, if time had permitted.

To get an even more accurate comparison of the routing performance compared to Linux hosts, a test could have been run monitoring the routing time for a group of directly connected Linux hosts. By eliminating the time used for switching, we could get more precise data for the time saved by processing packets on the IXP.

It is possible that, even with the kernel having priority, the packet processing time would suffer from a heavily loaded host. We could therefore have measured the time used by a Linux host to process the IP header switching when loaded with heavy work (like a tar process).

To find out which kind of load the μ Engines can take before having to drop packets, a packet generator (or several) could be used to strain the IP header switch application.

Bibliography

- [1] The Agere website. November 2005, <http://www.agere.com>.
- [2] Alex Arachnid server website. October 2005, http://www.alex.com/pages/anglais/fond_stream.html.
- [3] The AMCC website. November 2005, <http://www.amcc.com>.
- [4] The C-COR website. <http://www.c-cor.com/>.
- [5] Anton Chuvakin. Iptables linux firewall with packet string-matching support. Published at SecurityFocus: <http://www.securityfocus.com>, 2001. October 2005, <http://www.securityfocus.com/infocus/1531>.
- [6] Douglas E. Comer. *Using Network Processors, Intel IXP2xxx version*. Pearson Education, Inc., 2005.
- [7] California software laboratories: RTSP techguide. October 2005, <http://www.cswl.com/whitepr/tech/StreamingTechnology.html>.
- [8] Apple Darwin streaming server. November 2005, <http://developer.apple.com/darwin/projects/streaming/>.
- [9] The DivX website. October 2005, <http://www.divx.com>.
- [10] The Dresden real-time operating system (DROPS) website. November 2005, <http://web.inf.tu-dresden.de/SyA/labs/project/overview.html>.
- [11] DVD demystified website. November 2005, <http://www.dvddemystified.com/>.
- [12] RadiSys ENP2505 hardware reference manual. November 2005, http://www.radisys.com/files/support_downloads/007-01266-0002.ENP-2505.pdf.
- [13] RadiSys ENP2611 hardware reference manual. November 2005, http://www.radisys.com/files/support_downloads/007-01419-0003.ENP-2611HW.pdf.

- [14] Flac - free lossless audio compression website. October 2005
<http://flac.sourceforge.net>.
- [15] The Free On-line Dictionary of Computing. November 2005,
<http://foldoc.doc.ic.ac.uk/foldoc/index.html>.
- [16] The Fraunhofer website. October 2005, <http://www.iis.fraunhofer.de>.
- [17] J. Gettys, J Mogul, H. Frystyk, L. Masinter, P Leach, and T. Berners-Lee. RFC3550, 1999. October 2005, <http://www.ietf.org/rfc/rfc2616.txt>.
- [18] Moore Gordon. Cramming more components onto integrated circuits. *Electronics Magazine*, 1965. October 2005, ftp://download.intel.com/museum/Moores_Law/Articles-Press_Releases/Gordon_Moore_1965_Article.pdf.
- [19] A. Griffiths and G. Metherall. Cluster interconnection networks. Report for CSC433, Monash University, 1999. Available online: <http://www.buyya.com/csc433/ClusterNets.pdf>.
- [20] Carsten Griwodz, Michael Liepert, Michael Zink, and Ralf Steinmetz. Tune to lambda patching. *ACM Performance Evaluation Review*, 27(4):20–26, March 2000.
- [21] Pål Halvorsen. Improving I/O performance of multimedia servers. *Thesis for the Dr. Scient. (PhD) degree at University of Oslo, Published by Unipub forlag, ISSN 1501-7710, No. 161, Oslo, Norway, 2001.*
- [22] Pål Halvorsen and Carsten Griwodz. INF5060 - multimediakommunikasjon med nettverksprosessorer. IFI, UIO.
- [23] Pål Halvorsen and Carsten Griwodz. INF5070 - media storage and distribution systems. IFI, UIO.
- [24] Pål Halvorsen, Carsten Griwodz, Vera Goebel, Ketil Lund, Thomas Plagemann, and Jonathan Walpole. Storage system support for multimedia applications, part I and II. *IEEE Distributed Systems Online, Vol 5, No. 1/2, Jan/Feb 2004, 2003.*
- [25] IBM VideoCharger. October 2005, <http://www-306.ibm.com/software/data/videocharger>.
- [26] The Internet Machines website. November 2005, <http://www.internetmachines.com>.
- [27] The Intel website. November 2005, <http://www.intel.com>.
- [28] Intel IXA portability framework reference manual., 2003.

- [29] Intel IXA SDK 2.01 programmers reference manual.
- [30] Intel IXA SDK 3.51 programmers reference manual. November 2005, http://www.intel.com/design/network/products/npfamily/sdk_download.htm#351.
- [31] Intel IXP1200 hardware reference manual.
- [32] Intel IXP2400 hardware reference manual.
- [33] RadiSys ENP SDK 3.5 programmers guide.
- [34] Generic linux kernel for the IXP2xxx hardware open source initiative. November 2005, <http://ixp2xxx.sourceforge.net/>.
- [35] David B Johnson, David A. Maltz, and Yih-Chun Hu. The dynamic source routing protocol for mobile ad hoc networks (DSR), 2004. November 2005, <http://www.ietf.org/internet-drafts/draft-ietf-manet-dsr-10.txt>.
- [36] Komssys pages at sourceforge. October 2005, <http://komssys.sourceforge.net>.
- [37] The Monta Vista website. <http://www.mvista.com/>.
- [38] The Motorola website. November 2005, <http://www.motorola.com>.
- [39] The fraunhofer MP3 website. November 2005, <http://www.iis.fraunhofer.de/amm/techinf/layer3>.
- [40] Whatis.com - definition of "multimedia". "October 2005, http://searchwebservices.techtarget.com/sDefinition/0,,sid26_gci212612,00.html".
- [41] Chuck Narad. Communications and Moore's Law. Keynote at Intel 10th Academic Forum, 2005. October 2005, <http://download.intel.com/corporate/education/EMEA/academicforum/keynotes/Narad/Keynote%20Chuck%20Narad.pdf>.
- [42] The nCube website. Site decommissioned when the company was sold to c-cor. Old URL: <http://www.ncube.com/>.
- [43] The Intel network processor family website. October 2005, <http://www.intel.com/design/network/products/npfamily/index.htm>.
- [44] The Norwegian broadcasting company (NRK). November 2005, <http://www.nrk.no>.
- [45] The Ogg Vorbis website. November 2005, <http://www.vorbis.com>.
- [46] The Oracle interMedia website. November 2005, <http://www.oracle.com/technology/products/intermedia/index.html>.

- [47] T. Plagemann, V. Goebel, P. Halvorsen, and O. Anshus. Operating system support for multimedia systems, 2000.
- [48] The PMC-Sierra website. November 2005, <http://www.pmcsierra.com>.
- [49] The Qlogic website. October 2005, <http://www.qlogic.com>.
- [50] Apple Quicktime streaming server. October 2005, <http://www.apple.com/quicktime/streamingserver/>.
- [51] The RadiSys website. October 2005, <http://www.radisys.com>.
- [52] Sanjay Ranka, Youngju Won, and Sartaj Sahni. Programming a hypercube multicomputer. *IEEE Software Magazine*, 1988.
- [53] Real Networks website. November 2005, <http://www.realnetworks.com>.
- [54] Real Networks media delivery website. October 2005, http://www.realnetworks.com/products/media_delivery.html.
- [55] H Schulzrinne, S Casner, R Frederick, and V Jacobson. RFC3550, 2003.
- [56] H Schulzrinne, A Rao, and R. Lanphier. RFC2326, 1998. November 2005, <http://www.ietf.org/rfc/rfc2326.txt>.
- [57] Frank Olaf Sem-Jacobsen, Tom Anders Dalseng, and Andreas Petlund. Achieving intra-cube routing using Intel IXP cards. Assignment in INF5070: Multimedia storage and distribution systems, 2003.
- [58] The SF-anytime website. November 2005. <http://www.sf-anytime.com/>.
- [59] Dinkar Sitaram and Asit Dan. *Multimedia Servers - Applications, environments, and design*. Morgan Kaufmann Publishers, 2000.
- [60] The tcpdump/libpcap project homepage. November 2005: <http://www.tcpdump.org/>.
- [61] The Vitesse website. November 2005, <http://www.vitesse.com>.
- [62] The World of Warcraft website. November 2005, <http://www.worldofwarcraft.com/>.
- [63] George Kingsley Zipf. Human behaviour and the principle of least-effort. *Addison-Wesley, Cambridge MA*, 1949.

Appendix A

VS³ video server cube source

A.1 ICrouter microblock source

A.1.1 dl_system.excerpt.h

```
1  /*****
2  **
3  **  This file contains definitions of scratch rings      **
4  **  and shared memory areas for the VS3 application.  **
5  **  This is only an excerpt.                          **
6  **  The whole file could not be included due to      **
7  **  Intel proprietary regulations.                  **
8  **  To get the whole picture, access the file       **
9  **  dispatch_loop/dl_system.h in the build tree.    **
10 **
11 *****/
12
13 /*-----
14 * i) Base address for meta data (buffer descriptors) of packet
15 *    buffers - for free list 2 - dcube
16 *-----*/
17 #ifndef          DCUBE_NUM_BUF_HANDLES
18 #define          DCUBE_NUM_BUF_HANDLES 1024
19 #endif
20
21 #ifndef          DCUBE_SRAM_BASE
22 #define          DCUBE_SRAM_BASE      0x80000
```

```

23 #endif
24
25 /* Size of meta-data (buffer descriptor) in bytes for each buffer.
26    Should be a power of 2. */
27 #ifndef          DCUBE_META_DATA_SIZE
28 #define          DCUBE_META_DATA_SIZE          32
29 #endif
30
31 /* Total SRAM size allocated in bytes for meta data.
32    In simulation we use only 20KB. */
33 #ifndef          DCUBE_SRAM_SIZE
34 #define          DCUBE_SRAM_SIZE          (DCUBE_NUM_BUF_HANDLES * DCUBE_META_DATA_SIZE)
35 #endif
36
37 #ifndef          DCUBE_SRAM_MAX
38 #define          DCUBE_SRAM_MAX          DCUBE_SRAM_BASE + DCUBE_SRAM_SIZE
39 #endif
40
41 /*-----
42 * i) Base address for packet buffers. This is the actual packet data.
43 *    --dcube
44 *-----*/
45 #ifndef          DCUBE_SDRAM_BASE
46 #define          DCUBE_SDRAM_BASE          0x2000000
47 #endif
48
49 /*    Size of one packet buffer in bytes. */
50 #ifndef          DCUBE_BUFFER_SIZE
51 #define          DCUBE_BUFFER_SIZE          2048
52 #endif
53
54 /* Total DRAM size allocated in bytes for packet buffers. */
55 #ifndef          DCUBE_SDRAM_SIZE
56 #define          DCUBE_SDRAM_SIZE          (DCUBE_NUM_BUF_HANDLES * DCUBE_BUFFER_SIZE)
57 #endif
58
59 /*
60 * Scratch ring packet_echo --> XScale
61 */
62 #define          PE_XSCALE_COMM_RING          10
63 #define          PE_XSCALE_COMM_BASE          8192

```

```

64 #define          PE_XSCALE_COMM_SIZE          512
65
66
67 /*
68  * Scratch ring XScale --> packet_echo
69 */
70 #define          XSCALE_PE_COMM_RING          11
71 #define          XSCALE_PE_COMM_BASE          9216
72 #define          XSCALE_PE_COMM_SIZE          512

```

A.1.2 ICrouter.h

```

1  /*****
2  **                                                                 **
3  ** Includes and definitions for the ICrouter microblock **
4  **                                                                 **
5  *****/
6
7 #ifndef _ICROUTER_H
8 #define _ICROUTER_H
9
10 #include <ixp.h>
11 #include <dl_system.h>
12 #include <dl_buf.c>
13 #include <dl_meta.c>
14 #include <hardware.h>
15 #include <sig_functions.h>
16 #include <ixp_lib.h>
17 #include <ixp_crc.c>
18 #include <rtl.c>
19 #include <ix_cc_microengines_bindings.h>
20
21 #define  PKTHDR_CACHE_SIZE          5          /* Number of quadwords to read in isRTSP() */
22 #define  ETHPROT_IP                  0x0800   /* Ethernet protocol for filter */
23 #define  IPPROT_TCP                   6       /* IP protocol for filter */
24 #define  SERVER_PORT                 9070     /* Server port for filter */
25
26 extern dl_buf_handle_t dlBufHandle;          /* The current buffer handle */
27 extern dl_buf_handle_t dlEopBufHandle;      /* For large packets, this is the last buffer in the chain. */
28 extern dl_meta_t dlMeta;                   /* Metadata struct */

```

```

29 __declspec(gp_reg) int dlNextBlock;           /* Next block that should process the buffer/packet */
30
31 /* IC packet types */
32 enum {
33     ICH_FIND_FILE = 1,
34     ICH_ROUTE_FEEDBACK,
35     ICH_CONTROL_MSG,
36     ICH_RTP
37 };
38
39 #endif /* #ifndef _ICROUTER_H */

```

A.1.3 ICrouter.c

```

1  /******
2  **
3  **  The ICrouter microblock and support methods  **
4  **
5  *****/
6
7  #include "ICrouter.h"
8
9  /* Read and return the input port from metadata */
10 int getInputPort(unsigned int bufHandle){
11     /* Check if packet came in on port 0 */
12     __declspec(sram_read_reg) unsigned int metaData[5];
13     __declspec(gp_reg) unsigned int metaOffset;
14
15     SIGNAL sig_sram_rw;
16
17     metaOffset = bufHandle << 2;
18     sram_read (metaData, (volatile void __declspec(sram) *)metaOffset, 5, \
19               sig_done, &sig_sram_rw);
20     wait_for_all(&sig_sram_rw);
21     return (metaData[3] >> 16);
22 }
23
24 /* Read and return the output port from metadata */
25 int getOutputPort(unsigned int bufHandle){
26     __declspec(sram_read_reg) unsigned int metaData[5];

```

```

27  __declspec(gp_reg) unsigned int metaOffset;
28
29  SIGNAL sig_sram_rw;
30
31  metaOffset = bufHandle << 2;
32  sram_read (metaData, (volatile void __declspec(sram) *)metaOffset, 5, \
33           sig_done, &sig_sram_rw);
34  wait_for_all(&sig_sram_rw);
35  return (metaData[3] & 0xffff);
36 }
37
38 /* Pop the next port from route field of IC header */
39 int popPort(){
40     uint32_t route;
41     char port;
42     __declspec(sdram) unsigned char *p_pkt_hdr;
43
44     p_pkt_hdr = (__declspec(sdram) unsigned char *)
45         (Dl_BufGetData(dlBufHandle) + dlMeta.offset);
46
47     route = ua_get_u32(p_pkt_hdr, 4);
48
49     port = route & 0xff;
50     route >>= 8;
51     route |= 0xff000000;
52     ua_set_32(p_pkt_hdr, 4, route);
53
54     return port;
55 }
56
57 /* Decrement and return TTL in IC header */
58 int ttlDecr(){
59     __declspec(sdram) unsigned char *p_pkt_hdr;
60     int ttl;
61
62     p_pkt_hdr = (__declspec(sdram) unsigned char *)
63         (Dl_BufGetData(dlBufHandle) + dlMeta.offset);
64
65     ttl = ua_get_u8(p_pkt_hdr, 1);
66     ttl--;
67     ua_set_8(p_pkt_hdr, 1, ttl);

```

```

68
69  return ttl;
70 }
71
72 /* Get the IC packet type from IC header */
73 int getType(){
74  __declspec(s dram) unsigned char *p_pkt_hdr;
75
76  p_pkt_hdr = (__declspec(s dram) unsigned char *)
77    (DI_BufGetData(dlBufHandle) + dlMeta.offset);
78
79  return ua_get_u8(p_pkt_hdr, 0);
80 }
81
82 /* Strips the IC header (for outbound packets)*/
83 void stripICHeader(){
84  __declspec(sram) dl_meta_t *pMeta;
85
86  pMeta = (__declspec(sram) dl_meta_t *)DI_BufGetDesc(dlBufHandle);
87  pMeta->bufferSize-=8;
88  pMeta->offset +=8;
89 }
90
91 /* Remove ethernet checksum automatically generated by hardware */
92 void removeEthChecksum(){
93  __declspec(sram) dl_meta_t *pMeta;
94
95  pMeta = (__declspec(sram) dl_meta_t *)DI_BufGetDesc(dlBufHandle);
96  pMeta->bufferSize-=4;
97 }
98
99 /* Packet filter . Configured to allow TCP packets destined for port 9070 */
100 int isRTSP(unsigned int bufHandle, unsigned int dlMeta1){
101  __declspec(gp_reg) unsigned short ethtype, dstport;
102  __declspec(gp_reg) unsigned char ipprot;
103  __declspec(s dram) unsigned char *p_pkt_hdr;
104  __declspec(dram_read_reg) unsigned int pkthdr_in[10];
105  __declspec(local_mem) unsigned int temp[10];
106  SIGNAL_PAIR sig_dram_rw;
107
108  /* set the meta data accordingly */

```

```

109     dlMeta.value[1]                = dlMeta1;
110     dlBufHandle.value              = bufHandle;
111
112     p_pkt_hdr = (__declspec(s dram) unsigned char *)
113         (Dl_BufGetData(dlBufHandle) + dlMeta.offset);
114
115     dram_read(pkthdr_in,(volatile void __declspec(s dram) *)p_pkt_hdr,
116             PKTHDR_CACHE_SIZE, sig_done,&sig_dram_rw);
117     wait_for_all(&sig_dram_rw);
118
119     /* Copy the header to local memory */
120     temp[0] = pkthdr_in[0];
121     temp[1] = pkthdr_in[1];
122     temp[2] = pkthdr_in[2];
123     temp[3] = pkthdr_in[3];
124     temp[4] = pkthdr_in[4];
125     temp[5] = pkthdr_in[5];
126     temp[6] = pkthdr_in[6];
127     temp[7] = pkthdr_in[7];
128     temp[8] = pkthdr_in[8];
129     temp[9] = pkthdr_in[9];
130
131     /*Read eth header protocol type */
132     ethtype = ua_get_u16(temp, 12);
133     if(ethtype != ETHPROT_IP)
134         return 0;
135
136     /* Read IP protocol field */
137     ipprot = ua_get_u8(temp, 23);
138     if(ipprot != IPPROT_TCP)
139         return 0;
140
141     /* Check dst port */
142     dstport = ua_get_u16(temp, 36);
143     if(dstport != SERVER_PORT)
144         return 0;
145
146     return 1;
147 }
148
149 /* Main procedure containing the processing loop.

```

```

150     Upon receiving a new packet on the scratch ring,
151     the packet is processed and dropped or forwarded
152 */
153 main() {
154
155     __declspec(sram_write_reg) unsigned int txReq1, txReq2, txReq3;
156     __declspec(sram_read_reg) unsigned int rx_msg1[5], rx_msg2, rx_msg3;
157     __declspec(gp_reg) unsigned int ipsrc, ipdest, ethsrc1, ethdest1;
158     __declspec(gp_reg) unsigned short ethsrc2, ethdest2;
159
160     int sig;
161     SIGNAL sig_scr_created;
162     SIGNAL sig_new_packet;
163     SIGNAL sig_sram_read;
164     SIGNAL sig_sram_write;
165     SIGNAL sig_scr_get;
166     SIGNAL sig_scr_put;
167     SIGNAL sig_get_meta;
168     SIGNAL sig_flush_meta;
169     SIGNAL sig_never;
170     SIGNAL sig_buf_alloc;
171
172     __assign_relative_register(&sig_scr_created, 13);
173     __assign_relative_register(&sig_new_packet, 14);
174
175     /* Wait for signal from XScale before proceeding
176     This is to prevent the system from starting to
177     forward packets until the XScale functionality is ready */
178     cap_fast_write( 0, csr_thread_interrupt_a);
179     wait_for_all(&sig_scr_created);
180
181     /* Context 0: Read from RX(sr4), Filter packets.
182     Forward according to Packet type */
183     if (ctx() == 0) {
184         while(1){
185
186             /* Read from scratch ring */
187             scratch_get_ring(rx_msg1,
188                             (void*)(POS_RX_RING_OUT << 2),
189                             5,
190                             sig_done,

```



```

191         &sig_scr_get);
192     wait_for_all(&sig_scr_get);
193
194     /* If data is received, commence processing */
195     if (rx_msg1[0]){
196         /* Set metadata offset and buf handle values
197         (used by helper functions */
198         dlMeta.value[1]           = rx_msg1[2];
199         dlBufHandle.value         = rx_msg1[0];
200         txReq1 = rx_msg1[0];
201
202         /* Remove eth checksum (added by hw) by
203         decrementing bufferSize in metadata by 4 */
204         removeEthChecksum();
205
206         /* Read the input port number */
207         if(getInputPort(rx_msg1[0]) == 0){
208             /* External packet, filter */
209             if(isRTSP(rx_msg1[0], rx_msg1[2])){
210                 /* Forward to XScale */
211                 scratch_put_ring(&txReq1,
212                                 (void*)(PE_XSCALE_COMM_RING << 2),
213                                 1,
214                                 sig_done,
215                                 &sig_scr_put);
216                 wait_for_all(&sig_scr_put);
217                 /* Invoke int_a to wakeup XScale packet processing */
218                 cap_fast_write( 0, csr_thread_interrupt_a);
219             }
220         }else{
221             /* Internal packet. Route according to IC header */
222             U8 port, ttl, type;
223
224             type = getType();
225             if(type == ICH_FIND_FILE){
226                 /* File locating packet – Always forward to XScale */
227                 scratch_put_ring(&txReq1,
228                                 (void*)(PE_XSCALE_COMM_RING << 2),
229                                 1,
230                                 sig_done,
231                                 &sig_scr_put);

```

```

232     wait_for_all(&sig_scr_put);
233     cap_fast_write( 0, csr_thread_interrupt_a);
234 }else{
235     /* Other packet. Route according to route field of IC header */
236     /* Decrement ttl */
237     ttl = ttlDecr();
238     if(ttl > 0){
239         /* Pop route */
240         port = popPort();
241         /* Send to popped port */
242         if(port == 0){
243             /* If outbound, strip IC header */
244             stripICHeader();
245
246             /* Send to port 0 (external)*/
247             scratch_put_ring((void*)&txReq1,
248                             (void*)(PACKET_TX_SCR_RING_0 << 2),
249                             1,
250                             sig_done,
251                             &sig_scr_put);
252             wait_for_all(&sig_scr_put);
253         }else if(port == 1){
254             scratch_put_ring((void*)&txReq1,
255                             (void*)(PACKET_TX_SCR_RING_1 << 2),
256                             1,
257                             sig_done,
258                             &sig_scr_put);
259             wait_for_all(&sig_scr_put);
260         }else if(port == 2){
261             scratch_put_ring((void*)&txReq1,
262                             (void*)(PACKET_TX_SCR_RING_2 << 2),
263                             1,
264                             sig_done,
265                             &sig_scr_put);
266             wait_for_all(&sig_scr_put);
267         }else if(port == 0xff){
268             /* Forward to XScale */
269             scratch_put_ring(&txReq1,
270                             (void*)(PE_XSCALE_COMM_RING << 2),
271                             1,
272                             sig_done,

```

```

273                                     &sig_scr_put);
274         wait_for_all(&sig_scr_put);
275         cap_fast_write( 0, csr_thread_interrupt_a);
276     }
277     } /* If ttl == 0 -- packet is dropped */
278 }
279 } /* Packet is dropped */
280 } /*if packet*/
281 } /*while(1)*/
282 } /* else */
283 /* The second context.
284    Receives packets from the XScale, and forwards according to metadata output port
285    the route has already been popped on the XScale. */
286 if (ctx() == 2) {
287     unsigned int numPackets;
288     int i;
289
290     while(1){
291
292         scratch_get_ring(&rx_msg2,
293                         (void*)(XSCALE_PE_COMM_RING << 2),
294                         1,
295                         sig_done,
296                         &sig_scr_get);
297         wait_for_all(&sig_scr_get);
298
299         numPackets = (rx_msg2 >> 24);
300
301         for(i = 0; i < numPackets; i++){
302             txReq2 = rx_msg2 + ( i * 8 );
303             if(getOutputPort(rx_msg2) == 0){
304                 scratch_put_ring((void*)&txReq2,
305                                 (void*)(PACKET_TX_SCR_RING_0 << 2),
306                                 1,
307                                 sig_done,
308                                 &sig_scr_put);
309             }else if(getOutputPort(rx_msg2) == 1){
310                 scratch_put_ring((void*)&txReq2,
311                                 (void*)(PACKET_TX_SCR_RING_1 << 2),
312                                 1,
313                                 sig_done,

```

```

314             &sig_scr_put);
315     }else if(getOutputPort(rx_msg2) == 2){
316         scratch_put_ring((void*)&txReq2,
317             (void*)(PACKET_TX_SCR_RING_2 << 2),
318             1,
319             sig_done,
320             &sig_scr_put);
321     }
322     wait_for_all(&sig_scr_put);
323 }
324 } /*while(1)*/
325 }else{
326     /* All other contexts are sleeping */
327     wait_for_all(&sig_never);
328 }
329 } /* main*/

```

A.2 XScale source

A.2.1 dcube.h

```

1 #ifndef DCUBE_H
2 #define DCUBE_H
3
4 #include "uclo.h"
5 #include "hal_mev2.h"
6 #include "halMev2Api.h"
7 #include "hal_scratch.h"
8 #include "hal_sram.h"
9 #include "hal_dram.h"
10 #include "ix_macros.h"
11 #include "ix_ossl.h"
12 #include <stdio.h>
13 #include <errno.h>
14 #include <stdlib.h>
15 #include <string.h>
16 #include <fcntl.h>
17 #include <unistd.h>
18 #include <sys/time.h>

```

```

19 #include <sys/ioctl.h>
20 #include <enpv2_types.h>
21 #include <sys/mman.h>
22
23 /* Socket includes */
24 #include <sys/socket.h>
25 #include <sys/types.h>
26 #include <arpa/inet.h>
27
28 #ifdef DEBUG
29 #define   DBG_MSG(str,args...)           printf(str,##args)
30 #else
31 #define DBG_MSG(str, args...)
32 #endif
33
34 #define VIRT_DRAM_BASE (void *) (ix_uint32)Hal_dram_ch0_virtAddr
35 #define VIRT_SRAM_BASE (void *) (ix_uint32)Hal_sram_ch0_virtAddr
36 #define MAX_CONCURRENT_STREAMS 100
37 #define MULT_FACTOR 6 /* Used to get pointer to data from buf handle */
38 #define MAX_IC_SESSIONS MAX_CONCURRENT_STREAMS
39 #define ICH_TTL 4
40
41 /* XScale to host socket defines */
42 #define DEST_PORT 23456
43 #define DEST_IP "192.168.1.1"
44
45 enum {
46   RTSP_DESC = 1,
47   RTSP_SETUP
48 };
49
50 typedef union {
51   struct {
52     unsigned int  bufferNext;           /**< Next buffer in the chain */
53
54     unsigned short  bufferSize;         /**< amount of data currently in buffer */
55     unsigned short  offset;             /**< offset in DRAM where data begins */
56
57     unsigned int    packetSize          : 16;  /**< amount of data in the chain of buffers */
58
59     unsigned int    freeListId          : 4;    /**< Free List to which this buffer belongs to */

```

```

60  unsigned int    rxStat          : 4;    /**< Receive status */
61  unsigned int    headerType      : 8;    /**< HHeader Type: IPv4, IPv6 etc */
62
63  unsigned short  inputPort;       /**< Input port on which packet was received */
64  unsigned short  outputPort;     /**< Output port on which packet to be transmitted */
65
66  unsigned int    nextHopId       : 16;   /**< Nexthop ID */
67  unsigned int    fabricPort      : 8;    /**< Blade:Port */
68  unsigned int    reserved       : 4;    /**< reserved */
69  unsigned int    nhidType       : 4;    /**< nexthop ID type */
70
71  unsigned int    colorId         :4;
72  unsigned int    reserved1      :4;
73  unsigned int    flowId         :24;   /**< FLOW ID */
74
75  unsigned short  classId;        /**< Class ID */
76  unsigned short  reserved2;
77
78  unsigned int    packetNext;     /**< Next packet in the chain */
79  /**< (used only in Hierarchical Queuing) */
80  } __attribute__((packed));      // end of struct
81
82  unsigned int    value[8];       /**< aggregate for the above fields */
83
84  } dl_meta_t;
85
86  enum {
87    ICH_FIND_FILE    = 1,
88    ICH_ROUTE_FB,
89    ICH_CTRL_MSG,
90    ICH_CTRL_FB,
91    ICH_RTP
92  };
93
94  /* Standard IC header:
95     Common for all IC packets. */
96  typedef struct{
97    uint8_t        type; /* Type of IC packet */
98    uint8_t        ttl;  /* To avoid circulating packets */
99    uint16_t       dataLen; /* Length of IC packet data (excluding the header) */
100   uint32_t       route; /* Routing information. Src routing */

```

```

101 } ICHheader_t __attribute__((packed));
102
103 /* IC header extension:
104 Used to setup route information when assigning a new stream */
105 typedef struct{
106     uint16_t      sport; /* Client source port */
107     uint16_t      dport; /* Server port */
108
109     uint32_t      saddr; /* Client source IP addr. */
110     uint32_t      daddr; /* Egress IP addr */
111
112     char         eth_src[6]; /* Ethernet src address */
113     char         eth_dst[6]; /* Ethernet src address */
114 } ICHheader_ext_t __attribute__((packed));
115
116 enum {
117     ICS_CLOSED = 0,
118     ICS_WAIT_ROUTE,
119     ICS_ACTIVE,
120 };
121
122 typedef struct{
123     uint8_t      status;
124     uint32_t     route; /* Routing information. Src routing */
125     uint16_t     sport; /* Client source port */
126     uint16_t     dport; /* Server port */
127
128     uint32_t     saddr; /* Client source IP addr. */
129     uint32_t     daddr; /* Egress IP addr */
130
131     char         eth_src[6]; /* Ethernet src address */
132     char         eth_dst[6]; /* Ethernet src address */
133 } IC_session_t;
134 /* Structs located at the server that has the relevant file */
135 IC_session_t serv_ics[MAX_IC_SESSIONS];
136 /* Structs located at the egress (used for forwarding incoming RTSP and RTCP data) */
137 IC_session_t egr_ics[MAX_IC_SESSIONS];
138
139 extern void StartSpi3br();
140 extern void StartMacs();
141

```

```

142 /* Local includes
143     Needs structs above */
144 #include "bogus_tcp.h"
145 #include "dcube_utils.h"
146
147 #endif /* #ifndef DCUBE_H */

```

A.2.2 dcube.c

```

1 #include "dcube.h"
2
3 Hal_IntrMasks_T intMask;
4 ix_uint32 cbdata, intThread_A_handle, intThread_B_handle;
5 const char *pImageName="ICrouter.uof";
6 int hostSock, errno;
7
8
9 /*Called when the ICrouter microblock receives a
10 packet on port 0 or a packet destined for
11 this node is received*/
12 void intThread_A(){
13     unsigned int metaOffset, dataOffset, bufHandle, rtspType;
14     dl_meta_t* pMeta;
15     void* dataVAddr;
16
17     /* Get packet and vaddr for metadata and data */
18     bufHandle = SCRATCH_RING_GET(10);
19
20     DBG_MSG("Got from sr10: %0#10x\n", (int)bufHandle);
21
22     metaOffset = (bufHandle & 0xffff) << 2;
23     dataOffset = (bufHandle & 0xffff) << (MULT_FACTOR + 2);
24
25     DBG_MSG("metaOffset: %0#10x\n", (int)metaOffset);
26
27     pMeta = VIRT_SRAM_BASE + metaOffset;
28
29     dataVAddr = VIRT_DRAM_BASE + dataOffset + pMeta->offset;
30
31     /* IC packet or from port 0 0? */

```



```

32  if (pMeta->inputPort == 0){ /* External packet */
33      char data[1500];
34      int len = 0, idx = 0, status;
35
36      /* Since all other packets are filtered on the uEngines,
37       *This has to be a TCP-packet destined for port 9070 */
38
39      /* Port 0:
40          Do we have the file?
41          Yes: Create session and forward to host.
42          No: Create IC header and forward to port 1 and 2 */
43      status = tcp_recv(dataVAddr, data, &len, &idx);
44      DBG_MSG("status: %i\n", status);
45      if((status == B_TCP_DATA)){
46          DBG_MSG("tcp_recv returned the following data:\n");
47
48          /* Ensure support for creating more than one stream on each TCP connection */
49          rtspType = getRtspType(data);
50          if(rtspType == RTSP_DESC){
51              DBG_MSG("dcube.c: packet is DESC\n");
52              /* Implement later ..*/
53          }else if( rtspType == RTSP_SETUP ){
54              DBG_MSG("dcube.c: packet is SETUP\n");
55              if(localHasFile(data)){
56                  DBG_MSG("dcube.c: host has file.\n");
57                  /* Construct IC header and forward packet to host */
58                  add_IC_hdr(data, len, ICH_FIND_FILE, 0xffffffff);
59                  add_IC_ext_hdr(data, idx);
60                  ICPushPort(data, pMeta->inputPort);
61                  debug_print_ICheader(data);
62                  idx = newServICSession(data);
63                  debug_print_servICsession(idx);
64                  idx = newEgrICSession(data);
65                  egr_ics[idx].route = 0xffffffff; /* Since this is no feedback message */
66                  debug_print_egrICsession(idx);
67                  /* Send to host */
68                  errno = send(hostSock, data, len+sizeof(ICHeader_t)+sizeof(ICHeader_ext_t), 0);
69                  if(errno == -1){
70                      DBG_MSG("send() returned error\n");
71                  }else{
72                      DBG_MSG("Successfully sent %i bytes to socket\n",errno);

```

```

73     }
74 }else{
75     DBG_MSG("dcube.c: file not found.\n");
76     /* Set up IC header and forward packet
77        to port 1 and 2 */
78     /* Set route to be 0xffffffff. Ports will be pushed
79        as the packet is broadcasted through the system */
80     add_IC_hdr(data, len, ICH_FIND_FILE, 0xffffffff);
81     add_IC_ext_hdr(data, idx);
82     ICPushPort(data, pMeta->inputPort);
83     ixp_send_packet(data, len+sizeof(ICHeader_t)+sizeof(ICHeader_ext_t), 1);
84     ixp_send_packet(data, len+sizeof(ICHeader_t)+sizeof(ICHeader_ext_t), 2);
85 }
86 }else{
87     int icIdx, port;
88     /* Data for an existing stream
89        if the stream cannot be found,
90        discard the packet */
91     DBG_MSG("Received RTSP package\n");
92     debug_print_egrICsession(0);
93     icIdx = egrFindICsess(idx);
94     if(icIdx ==-1){
95         DBG_MSG("IC session not found\n");
96         /* Discard packet */
97     }else{
98         add_IC_hdr(data, len, ICH_CTRL_MSG, egr_ics[icIdx].route);
99         debug_print_buffer32(data, len+sizeof(ICHeader_t));
100
101         port = ICPopPort(data);
102         if(port == 0xff){
103             /* Send to host */
104             errno = send(hostSock, data, len+sizeof(ICHeader_t)+sizeof(ICHeader_ext_t), 0);
105             if(errno == -1){
106                 DBG_MSG("send() returned error\n");
107             }else{
108                 DBG_MSG("Successfully sent %i bytes to socket\n",errno);
109             }
110         }else{
111             ixp_send_packet(data, len+sizeof(ICHeader_t), port);
112         }
113     }

```

```

114     }
115 }else if((status == B_TCP_ERROR)){
116     DBG_MSG("tcp_recv returned B_TCP_ERROR \n");
117     return;
118 }else if((status == B_TCP_NO_DATA)){
119     DBG_MSG("tcp_recv returned B_TCP_NO_DATA \n");
120     return;
121 }
122
123 }else{ /* IC packet */
124     /* IC-packet*/
125
126     DBG_MSG("Got IC Packet, %i bytes\n", ((ICHeader_t*)dataVAddr)->dataLen);
127
128     switch(*((uint8_t*)dataVAddr)){ /* Read IC header Type */
129
130     case ICH_FIND_FILE: {
131         DBG_MSG("ICH_FIND_FILE\n");
132
133         if(localHasFile(dataVAddr+sizeof(ICHeader_t)+sizeof(ICHeader_ext_t))){
134             DBG_MSG("dcube.c - IC: host has file.\n");
135             /* If this is the first time I got this packet:
136              create stream struct - forward to host
137              If I have seen this packet before:
138              drop packet */
139             /* Find index of IC session for stream.
140              If no ssession, create a new one */
141
142             if(findICSession(dataVAddr) >= 0){
143                 DBG_MSG("ICSession found\n");
144                 /* We have already received this packet through another route:
145                  Drop packet */
146             }else{
147                 int idx, i, fbpSize =sizeof(ICHeader_t) + sizeof(ICHeader_ext_t) + 4 ;
148                 char *feedback, *rp;
149                 ICHeader_t *ich;
150
151                 feedback = (char*)malloc(fbpSize);
152                 /* Zero the packet */
153                 memset(feedback, 0, fbpSize);
154

```

```

155     DBG_MSG("Creating new IC session\n");
156     debug_print_ICheader(dataVAddr);
157     /* Push input port */
158     ICPushPort(dataVAddr, pMeta->inputPort);
159     idx = newServICSession(dataVAddr);
160     debug_print_servICsession(idx);
161     /* Send feedback packet to egress to
162        create session struct there */
163     memcpy(feedback, dataVAddr, sizeof(ICHeader_t)+sizeof(ICHeader_ext_t));
164     ich = (ICHeader_t*)feedback;
165
166     /* Make sure the packet is forwarded to XScale at egress */
167     /* Make sub of this ? */
168     rp = (uint8_t*)&(ich->route);
169     for(i=0; i<4; i++){
170         if(rp[i]!=0xff){
171             rp[i]=0xff;
172             break;
173         }
174     }
175     /* Copy the route to payload of packet */
176     memcpy(feedback+sizeof(ICHeader_t)+sizeof(ICHeader_ext_t), &(ich->route), 4);
177     /* Prepare other fields of IC header */
178     ich->type=ICH_ROUTE_FB;
179     ich->dataLen=4;
180     ich->tTL=4;
181     DBG_MSG("feedback: ich->route: %x\n", ich->route);
182
183     /* Send feedback packet */
184     ixp_send_packet(feedback, fbpSize, ICPopPort(feedback));
185     /* Send packet to server */
186     /* The program will crash (!) if the connect was not successful:
187        add handling of this problem */
188     errno = send(hostSock, dataVAddr, pMeta->bufferSize, 0);
189     if(errno == -1){
190         DBG_MSG("send() returned error\n");
191     }else{
192         DBG_MSG("Successfully sent %i bytes to socket\n",errno);
193     }
194 }
195 }else{

```

```

196     DBG_MSG("dcube.c - IC: file not found.\n");
197     /* decrement ttl, forward to port 1 and 2 */
198     if (IC_ttl_dec(dataVAddr)){
199         /* Push port on IC header */
200         ICPushPort(dataVAddr, pMeta->inputPort);
201         /* if ttl=0, drop packet */
202         /* Send to the port not received from */
203         if(pMeta->inputPort==2)
204             ixp_send_packet(dataVAddr, pMeta->bufferSize, 1);
205         else
206             ixp_send_packet(dataVAddr, pMeta->bufferSize, 2);
207     }
208     DBG_MSG("ttl is now: %i\n", *((uint8_t*)dataVAddr+1));
209 }
210
211 break;
212 }
213
214 case ICH_CTRL_FB: {
215     uint32_t saddr, daddr;
216     uint16_t sport, dport;
217     int status, size;
218     ICHheader_t *ich = (ICHheader_t*)dataVAddr;
219     ICHheader_ext_t *iche = (ICHheader_ext_t*)(dataVAddr + sizeof(ICHheader_t));
220     DBG_MSG("ICH_CTRL_FB\n");
221     debug_print_ICheader(dataVAddr);
222     /* Get port data from IC header */
223     saddr = iche->saddr;
224     daddr = iche->daddr;
225     sport = iche->sport;
226     dport = iche->dport;
227     size = ich->dataLen;
228     DBG_MSG("ICH_CTRL_FB: size = %i\n", size);
229     /* Strip IC header */
230     dataVAddr += sizeof(ICHheader_t) + sizeof(ICHheader_ext_t);
231     /* Forward payload to tcp_send */
232     status = tcp_send(dataVAddr, size, saddr, daddr, sport, dport);
233     if(status == -1){
234         DBG_MSG("tcp_send() returned an error.\n");
235     }else{
236         DBG_MSG("Successfully passed %i bytes to tcp_send.\n",status);

```

```

237     }
238
239     break;
240 }
241
242 case ICH_CTRL_MSG: {
243     DBG_MSG("ICH_CONTROL_MSG\n");
244     /* Forward to server */
245     errno = send(hostSock, dataVAddr, pMeta->bufferSize, 0);
246     if(errno == -1){
247         DBG_MSG("send() returned error\n");
248     }else{
249         DBG_MSG("Successfully sent %i bytes to socket\n",errno);
250     }
251
252     break;
253 }
254
255 case ICH_ROUTE_FB: {
256     int idx;
257     DBG_MSG("ICH_ROUTE_FEEDBACK\n");
258     debug_print_ICheader(dataVAddr);
259     debug_print_buffer32(dataVAddr, pMeta->bufferSize);
260     /* Set up egress icession with reverse route */
261     idx = newEgrICSession(dataVAddr);
262     debug_print_egrICsession(idx);
263
264     break;
265 }
266
267 default: break;
268 }
269
270 }
271
272 return;
273 }
274
275 int main(){
276     void *ucloHandle;
277     int status;

```

```

278  struct sockaddr_in dest_addr;    // will hold the destination addr
279
280  /* init ixa sdk uclo lib */
281  UcLo_InitLib();
282
283  /* init all uengs */
284  UcLo_InitLibUeng( 0xff );
285
286  /* load ueng image into memory */
287  status = UcLo_LoadObjFile( &ucloHandle, (char*) pImageName );
288  if( status != UCLO_SUCCESS ) {
289      printf("UcLo_LoadObjFile1 failed with status code :%d\n", status);
290      return 1;
291  }
292
293  /* write to microengines */
294  status = UcLo_WriteUimageAll( ucloHandle );
295  if( status != UCLO_SUCCESS ) {
296      printf("UcLo_WriteUimageAll failed\n");
297      printf("status = %i\n", status);
298      return 1;
299  }
300
301  /* verify uengine 0 for sanity check */
302  status = UcLo_VerifyUengine( ucloHandle, 0 );
303  if( status != UCLO_SUCCESS ) {
304      printf("UcLo_VerifyUengine 0 failed, status = 0x%x\n", status);
305      return 1;
306  }
307
308  printf("verify uengine passed\n");
309
310  /* delete object */
311  status = UcLo_DeleObj( ucloHandle );
312  if( status != UCLO_SUCCESS ) {
313      printf("UcLo_DeleObj failed\n");
314      return 1;
315  }
316
317  /* Enable interrupt A */
318  halMe_IntrEnable(HALME_INTR_THD_A_MASK);

```

```

319
320  /* init and start spi3br */
321  StartSpi3br();
322
323  /* init and start macs */
324  StartMacs();
325
326  halMe_Init(0xff);
327  /*Start microengines */
328  halMe_Start(0, 0xff);
329  halMe_Start(1, 0xff);
330  halMe_Start(2, 0xff);
331  printf("Microengines started\n");
332
333  /* Wait for int A Indicates that uEngines are initialized and ready */
334  halMe_IntrPoll(HALME_INTR_THD_A_MASK, &intMask);
335
336  DBG_MSG("Got int A\n");
337
338  /* Spawn callback thread for me-XScale comm */
339  status = halMe_SpawnIntrCallbackThd( HALME_INTR_THD_A_MASK, intThread_A,
340                                     &cbdata, 1, (void*)&intThread_A_handle);
341  if (status != HALME_SUCCESS){
342      printf("Error spawning intThread_A\n");
343  } else {
344      printf("Success spawning intThread_A\n");
345  }
346
347  /* Send signal 13 to me1, ctx 0, 2 and 4 */
348  status = me_signal( 1, 0, 13);
349  if (status != HALME_SUCCESS){
350      printf("Error signaling me 1, ctx 0\n");
351  }
352  status = me_signal( 1, 2, 13);
353  if (status != HALME_SUCCESS){
354      printf("Error signaling me 1, ctx 2\n");
355  }
356  status = me_signal( 1, 4, 13);
357  if (status != HALME_SUCCESS){
358      printf("Error signaling me 1, ctx 4\n");
359  }

```



```

360
361  /* Setup socket to communicate with host */
362  hostSock = socket(AF_INET, SOCK_STREAM, 0); // do some error checking!
363  if(hostSock == -1){
364      printf("socket() returned error\n");
365  }else{
366      printf("Socket successfully created\n");
367  }
368
369  dest_addr.sin_family = AF_INET;
370  dest_addr.sin_port = htons(DEST_PORT);
371  dest_addr.sin_addr.s_addr = inet_addr(DEST_IP);
372  memset(&(dest_addr.sin_zero), '\0', 8);
373
374  errno = connect(hostSock, (struct sockaddr *)&dest_addr, sizeof(struct sockaddr));
375  if(errno == -1){
376      printf("connect() returned error\n");
377  }else{
378      printf("Successfully connected to %s, port %i\n",inet_ntoa(dest_addr.sin_addr), DEST_PORT);
379  }
380  /*Show Spi3br driver information*/
381  //ShowSpi3br();
382
383  /*Show mac driver information*/
384  //ShowMac();
385
386  while(1) {
387      char data[1500], *dp;
388      int status, bytecount, dataSize;
389      /* Listen to hostSock
390         collect data, create Correct headers,
391         and forward to uEngines for transmit */
392
393      /* Assumes that IC Header is the first to arrive,
394         and that all packets arrive in the order: ICHeader->Payload->ICHeader->Payload etc. */
395
396      /* Read IC header */
397      dp = data;
398      bytecount = sizeof(ICHeader_t);
399      while(bytecount > 0){
400          status = recv(hostSock, dp, bytecount, 0);

```

```

401     if(status == -1){
402         DBG_MSG("recv() returned error\n");
403         sleep(1000);
404     }else if(status > 0){
405         DBG_MSG("Socket: Received IC header, %i bytes\n", status);
406         dp += status;
407         bytecount -= status;
408     }
409 }
410
411 /* Get size of data from IC header */
412 ICHeader_t *ich = (ICHeader_t*)data;
413 int port;
414 if((ich->type == ICH_ROUTE_FB) || (ich->type == ICH_CTRL_FB))
415     bytecount = ich->dataLen + sizeof(ICHeader_ext_t);
416 else
417     bytecount = ich->dataLen;
418
419 dataSize = bytecount;
420 DBG_MSG("dataSize: %i\n", dataSize);
421 DBG_MSG("bytecount: %i\n", bytecount);
422 while(bytecount > 0){
423     status = recv(hostSock, dp, bytecount, 0);
424     if(status == -1){
425         DBG_MSG("recv() returned error\n");
426         sleep(1000);
427     }else if(status > 0){
428         DBG_MSG("Socket: Received %i bytes of data\n", status);
429         dp += status;
430         bytecount -= status;
431     }
432 }
433 dataSize+=sizeof(ICHeader_t);
434
435 debug_print_buffer32(data, dataSize);
436 port = ICPopPort(data);
437 if(port == 0xff){ /* This is an RTSP message on egress machine */
438     ICHeader_ext_t *iche = (ICHeader_ext_t*)(data+sizeof(ICHeader_t));
439
440     status = tcp_send(data+sizeof(ICHeader_t)+sizeof(ICHeader_ext_t), ich->dataLen,
441                     iche->saddr, iche->daddr, iche->sport, iche->dport);

```

```

442     if(status == -1){
443         DBG_MSG("tcp_send() returned an error.\n");
444     }else{
445         DBG_MSG("Successfully passed %i bytes to tcp_send.\n",status);
446     }
447 }else{
448     dp = data;
449     if(port == 0){ /* If data is at egress, strip IC header */
450         dp += sizeof(ICHeader_t);
451         dataSize -= sizeof(ICHeader_t);
452     }
453     ixp_send_packet(dp, dataSize, port);
454 }
455 }
456
457 return 0;
458 }
459
460
461
462
463
464
465
466
467
468
469

```

A.2.3 dcube_utils.h

```

1 #ifndef _DCUBE_UTILS_H
2 #define _DCUBE_UTILS_H
3
4 #include "dcube.h"
5
6 void tcp_udp_checksum(void*, int);
7 void ip_checksum(void*);
8 ix_error me_signal( ix_uint32, ix_uint32, ix_uint32);
9 void ixp_send_packet(uint8_t*, int, uint16_t);

```

```

10 void debug_print_metadata(dl_meta_t*);
11 void debug_print_buffer32(void*, int);
12 int getRtspType(char *data);
13 int localHasFile (char *descString);
14 void add_IC_hdr(char *data, int len, int type, uint32_t route);
15 void add_IC_ext_hdr(char *data, int idx);
16 int IC_ttl_dec(void* dataVAddr);
17 int findICSession(void* ICpacket);
18 int newServICSession(void* ICpacket);
19 int newEgrICSession(void* ICpacket);
20 void debug_print_egrICsession(int idx);
21 void debug_print_servICsession(int idx);
22 void debug_print_ICheader(void* ICpacket);
23 void ICPushPort(void* ICpacket, uint8_t port);
24 int ICPopPort(void* ICpacket);
25 int egrFindICsess(int idx);
26 int timeval_subtract (struct timeval *result, struct timeval *x, struct timeval *y);
27
28 #endif /* #ifndef _DCUBE_UTILS_H */

```

A.2.4 dcube_utils.c

```

1 #include "dcube_utils.h"
2
3 #define IX_RM_SAME_ME_SIGNAL_OFFSET          0x108
4 #define VIRT_DRAM_BASE (void *) (ix_uint32)Hal_dram_ch0_virtAddr
5 #define VIRT_SRAM_BASE (void *) (ix_uint32)Hal_sram_ch0_virtAddr
6
7 int bufCount = 0; /* what buffer handle number is currently "active */
8 uint8_t packetcount = 0; /* How many packets currently scheduled for batch send */
9 int bufBase; /* The base buffer handle for the next batch */
10 struct timeval curTime, lastBatch;
11
12 ix_error me_signal( ix_uint32 arg_MENumber,
13                   ix_uint32 arg_ContextNumber,
14                   ix_uint32 arg_SignalNumber ) {
15     ix_error err = IX_SUCCESS;
16
17     WRITE_LWORD(((ix_uint32)Hal_cap_me_local_csr_virtAddr
18                + (arg_MENumber << 10) + IX_RM_SAME_ME_SIGNAL_OFFSET),

```

```

19             ((arg_SignalNumber << 3) + arg_ContextNumber));
20
21     return err;
22 }
23
24 /******
25 * To calculate tcp checksum :
26 * - Construct pseudo header:
27 *   +-----+-----+-----+-----+
28 *   |           Source Address           |
29 *   +-----+-----+-----+-----+
30 *   |           Destination Address       |
31 *   +-----+-----+-----+-----+
32 *   | zero | PTCL |   TCP Length   |
33 *   +-----+-----+-----+-----+
34 * - Zero out checksum field
35 * - Pad one zeroed byte if length is an odd number
36 * - Do a ones complement sum of the whole thing
37 *   -If a carry occur, add one, and return
38 * - Return the ones complement of the sum
39 *****/
40
41 /* Calculates the 16 bit ones-complement sum of a given
42 buffer. Pads the last byte with 0 if odd size
43 Input: data: Pointer to buffer start
44 len: length of data */
45 uint32_t partial_csum(uint8_t* data, int len) {
46     uint32_t sum = 0;
47     uint16_t last = 0;
48     int odd = 0;
49
50     if(len & 1)
51         odd = 1;
52
53     len >>= 1;
54     while (len > 0) {
55         sum += *((uint16_t*)data);
56         data += sizeof(uint16_t);
57         len--;
58     }
59     /* odd len */

```

```

60  if (odd){
61      last = *data << 8;
62      sum += last;
63  }
64
65  return sum ;
66 }
67
68 /* Calculate checksum for TCP or UDP header, and
69 write it to the packet
70 Input: ipHdrStart – Pointer to beginning of IP header
71 size – total size of packet */
72 void tcp_udp_checksum(void* ipHdrStart, int size){
73     uint32_t sum = 0;
74     uint16_t finalsum;
75     uint8_t protocol;
76     uint32_t zero_ptcl_tcpsize;
77
78     /* Extract protocol type from IP header */
79     memcpy(&protocol, ipHdrStart+9, 1);
80     zero_ptcl_tcpsize= (protocol << 16) | ((size-20) & 0xffff);
81
82     /* zero out checksum field */
83     memset(ipHdrStart+36, 0, 2);
84     /* Pseudo header calculations */
85     /* Calculate checksum of src & dest */
86     sum = partial_csum(ipHdrStart+12, 8);
87     /* Calculate checksum of protocol nr and TCP size*/
88     sum += partial_csum((uint8_t*)&zero_ptcl_tcpsize, 4);
89
90     /* Calculate checksum of tcp header and data*/
91     sum += partial_csum(ipHdrStart+20, size-20);
92
93     /* Add carries */
94     while (sum & 0xffff0000)
95         sum = (sum >> 16) + (sum & 0xffff);
96
97     /* Write the one's complement of the sum to
98 the correct spot in the TCP header */
99     finalsum = ~sum;
100    memcpy(ipHdrStart+36, &finalsum, 2);

```

```

101 }
102
103 /* Calculate checksum for IP header, and
104 write it to the packet
105 Input: Pointer to beginning of IP header */
106 void ip_checksum(void* ipHdrStart){
107     uint32_t sum;
108     uint16_t finalsum;
109
110     memset(ipHdrStart+10, 0, 2);
111     sum = partial_csum(ipHdrStart, 20);
112
113     /* Add carries */
114     while (sum & 0xffff0000)
115         sum = (sum >> 16) + (sum & 0xffff);
116
117     finalsum = ~sum;
118     memcpy(ipHdrStart+10, &finalsum, 2);
119 }
120
121 /* Send packet buffer to port specified
122 todo: implement batch sending and sending to
123 different ports */
124 void ixp_send_packet(uint8_t* buffer, int size, uint16_t port){
125     int curBufHandle, bufHandle, metaOffset, dataOffset, status;
126     struct timeval timediff;
127     void* dataVAddr;
128     dl_meta_t* pMeta;
129     DBG_MSG("Send packet of %i bytes to port %i\n", size, port);
130
131     packetcount++;
132     curBufHandle = 0x20000 + (bufCount * 8);
133     /* If we're at the beginning of a new batch */
134     if( packetcount == 1){
135         bufBase = curBufHandle;
136     }
137
138     /* Prepare the current packet */
139     bufHandle = curBufHandle;
140     metaOffset = (bufHandle & 0xfffff) << 2;
141     dataOffset = (bufHandle & 0xfffff) << (MULT_FACTOR + 2);

```

```

142 pMeta = VIRT_SRAM_BASE + metaOffset;
143 /* Zero metadata */
144 memset(pMeta, 0, sizeof(dl_meta_t));
145
146 pMeta->bufferNext = 0xff;
147 pMeta->value[2] = 0x00001001;
148 pMeta->bufferSize = size;
149 pMeta->packetSize = size;
150 pMeta->nextHopId = 0xff;
151 pMeta->offset = 0x100;
152 pMeta->outputPort = port;
153
154 dataVAddr = VIRT_DRAM_BASE + dataOffset + pMeta->offset;
155
156 memcpy(dataVAddr, buffer, size);
157
158 /* Increment or wrap buffer count */
159 if(bufCount == 1024){
160     bufCount = 0;
161 } else
162     bufCount++;
163
164 status = gettimeofday(&curTime, NULL);
165 if(status == -1){
166     DBG_MSG("gettimeofday returned an error\n");
167 }
168
169 timeval_subtract (&timediff, &curTime, &lastBatch);
170 DBG_MSG("timediff.tv_sec: %u - timediff.tv_usec: %u\n", timediff.tv_sec, timediff.tv_usec);
171
172 //if( timediff.tv_sec > 0 || timediff.tv_usec > 500 || packetcount >= 50){
173     /*
174         |   31 - 24   | 23   -   0 |
175         num packets   bufHandle
176     */
177     bufHandle = (packetcount << 24) | curBufHandle;
178     DBG_MSG("bufHandle : %0#10x\n", bufHandle);
179     SCRATCH_RING_PUT(11, bufHandle);
180     /* Reset batch packet counter */
181     packetcount = 0;
182     status = gettimeofday(&lastBatch, NULL);

```



```

183     if(status == -1){
184         DBG_MSG("gettimeofday returned an error\n");
185     }
186     //}
187
188     //status = me_signal( 1, 2, 14);
189     //if (status != HALME_SUCCESS){
190     //    DBG_MSG("Error signaling me 1, ctx 2\n");
191     // }
192 }
193
194 /* Print relevant metadata */
195 void debug_print_metadata(dl_meta_t* pMeta){
196     DBG_MSG("bufferNext: %0#10x\n", pMeta->bufferNext);
197     DBG_MSG("bufferSize: %0#10x\n", pMeta->bufferSize);
198     DBG_MSG("offset: %0#10x\n", pMeta->offset);
199     DBG_MSG("packetSize: %0#10x\n", pMeta->packetSize);
200     DBG_MSG("freeListId: %0#10x\n", pMeta->freeListId);
201     DBG_MSG("rxStat: %0#10x\n", pMeta->rxStat);
202     DBG_MSG("headerType:: %0#10x\n", pMeta->headerType);
203     DBG_MSG("inputPort: %0#10x\n", pMeta->inputPort);
204     DBG_MSG("outputPort: %0#10x\n", pMeta->outputPort);
205     DBG_MSG("nextHopId: %0#10x\n", pMeta->nextHopId);
206     DBG_MSG("fabricPort: %0#10x\n", pMeta->fabricPort);
207     DBG_MSG("flowId: %0#10x\n", pMeta->flowId);
208     DBG_MSG("classId: %0#10x\n", pMeta->classId);
209     DBG_MSG("classId: %0#10x\n\n", pMeta->packetNext);
210 }
211
212 /* Print buffer in 32 bit words hex in byte order*/
213 void debug_print_buffer32(void* buf, int size){
214     int i;
215
216     for(i=0; i < size; i++){
217         if(i % 4==0)
218             DBG_MSG("\nbuf[%i]: 0x", i);
219
220         if(*((uint8_t*)(buf+i)) == 0)
221             DBG_MSG("00");
222         else if(*((uint8_t*)(buf+i)) < 15)
223             DBG_MSG("0%x", *((uint8_t*)(buf+i)));

```

```

224     else
225         DBG_MSG("%x", *((uint8_t*)(buf+i)));
226     }
227     DBG_MSG("\n");
228 }
229
230 /* Check if the payload is a SETUP packet
231 Implement check for other types when needed */
232 int getRtspType(char *data){
233     char *setupStr = "SETUP";
234     if(strncmp(data, setupStr, strlen(setupStr)) == 0)
235         return RTSP_SETUP;
236
237     return 0;
238 }
239
240 /* Check if the requested file exists on the local server */
241 int localHasFile (char *descString){
242     char fname[50] = "/opt/storage/\0";
243     char *start;
244     FILE *status;
245
246     start = strstr(descString, "9070");
247     if(start == NULL)
248         return 0;
249     start += 5;
250
251     strncat(fname, start, strchr(start, ' ')-start);
252     DBG_MSG("fname = %s\n", fname);
253
254     status = fopen (fname, "rb");
255     if(status == NULL)
256         return 0;
257     else
258         fclose(status);
259     return 1;
260 }
261
262 /* Input: data: pointer to buffer containing payload
263     len: pointer to length of payload
264     idx: index of tcp stream in the ts-array

```

```

265         Prepends the payload with an IC header */
266 void add_IC_hdr(char *data, int len, int type, uint32_t route){
267     ICHheader_t ich;
268     /* Zero IC header */
269     memset(&ich, 0, sizeof(ICHheader_t));
270     ich.type = type;
271     ich.ttl = ICH_TTL;
272     ich.dataLen = len;
273     ich.route = route;
274     memmove(data+sizeof(ICHheader_t), data, len);
275     memcpy(data, &ich, sizeof(ICHheader_t));
276 }
277
278 /* Assumes that an IC header already has been added */
279 void add_IC_ext_hdr(char *data, int idx){
280     ICHheader_t *ich = (ICHheader_t*)data;
281     ICHheader_ext_t iche;
282     /* Zero IC header ext */
283     memset(&iche, 0, sizeof(ICHheader_ext_t));
284     iche.sport = ts[idx].sport;
285     iche.dport = ts[idx].dport;
286     iche.saddr = ts[idx].saddr;
287     iche.daddr = ts[idx].daddr;
288     memcpy(&iche.eth_src, &ts[idx].eth_src, 6);
289     memcpy(&iche.eth_dst, &ts[idx].eth_dst, 6);
290     memmove(data+sizeof(ICHheader_t)+sizeof(ICHheader_ext_t), data+sizeof(ICHheader_t), ich->dataLen);
291     memcpy(data+sizeof(ICHheader_t), &iche, sizeof(ICHheader_ext_t));
292 }
293
294 /*Input: pointer to beginning of IC packet
295 Decrements IC ttl by one, and returns the value */
296 int IC_ttl_dec(void* dataVAddr){
297     return --(*((uint8_t*)dataVAddr+1));
298 }
299
300 /* Check if IC session exists.
301 Input: pointer to start of IC packet
302 Output: position of IC session (or -1 if it does not exist) */
303 int findICSession(void* ICpacket){
304     ICHheader_ext_t *iche = (ICHheader_ext_t*)(ICpacket+sizeof(ICHheader_t));
305     int i, pos = -1;

```

```

306
307   DBG_MSG("findICSession: sport= %i\n", iche->sport);
308
309   for (i=0; i<MAX_IC_SESSIONS; i++) {
310       if(serv_ics[i].status != ICS_CLOSED && serv_ics[i].saddr==iche->saddr &&
311           serv_ics[i].daddr==iche->daddr && serv_ics[i].sport == iche->sport && serv_ics[i].dport == iche->dport) {
312           /* We have found an active session that corresponds to the ic header */
313           pos = i;
314           DBG_MSG("ICSession: found session\n");
315           break;
316       }
317   }
318   return pos;
319 }
320
321 uint32_t routeInverse(uint32_t route){
322     uint32_t inv_route;
323     uint8_t *rp= (uint8_t*)&route, *irp=(uint8_t*)&inv_route;
324     int i, j=3;
325
326     for(i=0; i < 4; i++){
327         if(*(rp+i)==0xff)
328             *(irp+i) = 0xff;
329         else
330             *(irp+(j--)) = *(rp+i);
331     }
332     return inv_route;
333 }
334
335 /* Initialize IC session from data
336    Input: Pointer to IC packet
337    Output: index of new IC session (-1 if no free slot)
338 */
339 int newEgrICSession(void* ICpacket){
340     ICHeader_ext_t *iche = (ICHeader_ext_t*)(ICpacket + sizeof(ICHeader_t));
341     /* The payload of IC_ROUTE_FB is rout to server with the file */
342     uint32_t route_data = *(uint32_t*)(ICpacket+sizeof(ICHeader_t)+sizeof(ICHeader_ext_t));
343     DBG_MSG("newEgrICSession: route_data: %x\n", route_data);
344     int i, pos = -1;
345     /* locate first free slot */
346     for (i=0; i<MAX_IC_SESSIONS; i++) {

```

```

347     if(egr_ics[i].status == ICS_CLOSED){
348         pos = i;
349         DBG_MSG("newEgrICSession: first free: %i\n", pos);
350         break;
351     }
352 }
353
354 if(pos != -1){
355     egr_ics[pos].status = ICS_ACTIVE;
356     egr_ics[pos].route = routeInverse(route_data);
357     egr_ics[pos].sport = iche->sport;
358     egr_ics[pos].dport = iche->dport;
359     egr_ics[pos].saddr = iche->saddr;
360     egr_ics[pos].daddr = iche->daddr;
361     memcpy(&egr_ics[pos].eth_src, &iche->eth_src, 6);
362     memcpy(&egr_ics[pos].eth_dst, &iche->eth_dst, 6);
363 }
364 return pos;
365 }
366
367 /* Initialize IC session from data
368    Input: Pointer to IC packet
369    Output: index of new IC session (-1 if no free slot)
370 */
371 int newServICSession(void* ICpacket){
372     ICHeader_t *ich = (ICHeader_t*)ICpacket;
373     ICHeader_ext_t *iche = (ICHeader_ext_t*)(ICpacket+sizeof(ICHeader_t));
374     int i, pos = -1;
375     /* locate first free slot */
376     for (i=0; i<MAX_IC_SESSIONS; i++) {
377         if(serv_ics[i].status == ICS_CLOSED){
378             pos = i;
379             DBG_MSG("newServICSession: first free: %i\n", pos);
380             break;
381         }
382     }
383
384     if(pos != -1){
385         serv_ics[pos].status = ICS_ACTIVE;
386         serv_ics[pos].route = ich->route;
387         serv_ics[pos].sport = iche->sport;

```

```

388     serv_ics[pos].dport = iche->dport;
389     serv_ics[pos].saddr = iche->saddr;
390     serv_ics[pos].daddr = iche->daddr;
391     memcpy(&serv_ics[pos].eth_src, &iche->eth_src, 6);
392     memcpy(&serv_ics[pos].eth_dst, &iche->eth_dst, 6);
393 }
394 return pos;
395 }
396
397 void debug_print_servICsession(int idx){
398     if(idx < MAX_IC_SESSIONS){
399         DBG_MSG("IC-session:\n");
400         DBG_MSG("index: %i\n", idx);
401         DBG_MSG("status: %u\n", serv_ics[idx].status);
402         DBG_MSG("route: %x\n", serv_ics[idx].route);
403         DBG_MSG("sport: %u\n", serv_ics[idx].sport);
404         DBG_MSG("dport: %u\n", serv_ics[idx].dport);
405         DBG_MSG("saddr: %x\n", serv_ics[idx].saddr);
406         DBG_MSG("daddr: %x\n", serv_ics[idx].daddr);
407         DBG_MSG("eth_src: ");
408         debug_print_buffer32(serv_ics[idx].eth_src, 6);
409         DBG_MSG("eth_dst: ");
410         debug_print_buffer32(serv_ics[idx].eth_dst, 6);
411     }
412 }
413
414 void debug_print_egrICsession(int idx){
415     if(idx < MAX_IC_SESSIONS){
416         DBG_MSG("IC-session:\n");
417         DBG_MSG("index: %i\n", idx);
418         DBG_MSG("status: %u\n", egr_ics[idx].status);
419         DBG_MSG("route: %x\n", egr_ics[idx].route);
420         DBG_MSG("sport: %u\n", egr_ics[idx].sport);
421         DBG_MSG("dport: %u\n", egr_ics[idx].dport);
422         DBG_MSG("saddr: %x\n", egr_ics[idx].saddr);
423         DBG_MSG("daddr: %x\n", egr_ics[idx].daddr);
424         DBG_MSG("eth_src: ");
425         debug_print_buffer32(egr_ics[idx].eth_src, 6);
426         DBG_MSG("eth_dst: ");
427         debug_print_buffer32(egr_ics[idx].eth_dst, 6);
428     }

```

```

429 }
430
431 void debug_print_ICheader(void* ICpacket){
432     ICHheader_t *ich = (ICHheader_t*)ICpacket;
433
434     DBG_MSG("ICheader:\n");
435     DBG_MSG("type: %u\n", ich->type);
436     DBG_MSG("ttl: %u\n", ich->ttl);
437     DBG_MSG("dataLen: %u\n", ich->dataLen);
438     DBG_MSG("route: %x\n", ich->route);
439
440     if((ich->type == ICH_FIND_FILE) || (ich->type == ICH_ROUTE_FB) || (ich->type == ICH_CTRL_FB)){
441         ICHheader_ext_t *iche = (ICHheader_ext_t*)(ICpacket+sizeof(ICHheader_t));
442         DBG_MSG("sport: %u\n", iche->sport);
443         DBG_MSG("dport: %u\n", iche->dport);
444         DBG_MSG("saddr: %x\n", iche->saddr);
445         DBG_MSG("daddr: %x\n", iche->daddr);
446         DBG_MSG("eth_src: ");
447         debug_print_buffer32(iche->eth_src, 6);
448         DBG_MSG("eth_dst: ");
449         debug_print_buffer32(iche->eth_dst, 6);
450     }
451 }
452
453 /* Input: ICpacket: Pointer to IC packet
454          port: input port number
455          Pushes the given port number on the right side of the 4 byte route */
456 void ICPushPort(void* ICpacket, uint8_t port){
457     ICHheader_t *ich = (ICHheader_t*)ICpacket;
458
459     ich->route <<= 8;
460     ich->route |= port;
461     DBG_MSG("ICPushPort: ich->route: %x\n", ich->route);
462 }
463
464 /* Input: ICpacket: Pointer to IC packet
465          Pops a port from the IC header route, and returns it. */
466 int ICPopPort(void* ICpacket){
467     ICHheader_t *ich = (ICHheader_t*)ICpacket;
468     int port;
469

```

```

470 port = ich->route & 0xff;
471 ich->route >>= 8;
472 ich->route |= 0xff000000;
473 DBG_MSG("ICPopPort: ich->route: %x\n", ich->route);
474 DBG_MSG("ICPopPort: popped port: %u\n", port);
475 return port;
476 }
477
478 /* Input: idx: Index of TCP session
479      Output: index of IC session (or -1 if not found) */
480 int egrFindICsess(int idx){
481     int i;
482
483     for (i=0; i<MAX_IC_SESSIONS; i++) {
484         if(egr_ics[i].status != ICS_CLOSED && egr_ics[i].saddr==ts[idx].saddr &&
485            egr_ics[i].daddr==ts[idx].daddr && egr_ics[i].sport == ts[idx].sport &&
486            egr_ics[i].dport == ts[idx].dport){
487             /* We have found an active IC session that corresponds to the received packet */
488             return i;
489             DBG_MSG("egrFindICsess: found session\n");
490         }
491     }
492     return -1;
493 }
494
495 /* Subtract the 'struct timeval' values X and Y,
496      storing the result in RESULT.
497      Return 1 if the difference is negative, otherwise 0. */
498 int timeval_subtract (struct timeval *result,
499                      struct timeval *x, struct timeval *y){
500     /* Perform the carry for the later subtraction by updating y. */
501     if (x->tv_usec < y->tv_usec) {
502         int nsec = (y->tv_usec - x->tv_usec) / 1000000 + 1;
503         y->tv_usec -= 1000000 * nsec;
504         y->tv_sec += nsec;
505     }
506     if (x->tv_usec - y->tv_usec > 1000000) {
507         int nsec = (x->tv_usec - y->tv_usec) / 1000000;
508         y->tv_usec += 1000000 * nsec;
509         y->tv_sec -= nsec;
510     }

```



```

511
512  /* Compute the time remaining to wait.
513     tv_usec is certainly positive. */
514  result->tv_sec = x->tv_sec - y->tv_sec;
515  result->tv_usec = x->tv_usec - y->tv_usec;
516
517  /* Return 1 if result is negative. */
518  return x->tv_sec < y->tv_sec;
519 }
520
521 /* To do:
522     Make linked lists of TCP structs and IC session structs */

```

A.2.5 bogus_tcp.h

```

1  #ifndef _BOGUS_TCP_H
2  #define _BOGUS_TCP_H
3
4  #include "dcube.h"
5
6  #define MAX_TCP_SESSIONS          MAX_CONCURRENT_STREAMS
7  #define ETH_HDR_SIZE              14
8  #define ETH_TYPE_IP               0x0800
9  #define IP_HDR_SIZE              20
10 #define TCP_HDR_SIZE              20
11
12 /* Return states */
13 #define B_TCP_ERROR                0
14 #define B_TCP_NO_DATA             1
15 #define B_TCP_DATA                2
16
17 /* TCP FLAGS */
18 #define TCP_FLAGS_FIN             1
19 #define TCP_FLAGS_SYN             1<<1
20 #define TCP_FLAGS_RST             1<<2
21 #define TCP_FLAGS_PSH             1<<3
22 #define TCP_FLAGS_ACK             1<<4
23 #define TCP_FLAGS_URG             1<<5
24
25 /* TCP states */

```

```

26 enum {
27     TCP_CLOSED = 0,
28     TCP_SYN_RCVD,
29     TCP_ESTABLISHED,
30     TCP_CLOSE_WAIT,
31     TCP_FIN_WAIT,
32     TCP_TIME_WAIT
33 };
34
35 /* TCP Change options */
36 enum {
37     TCP_NO_CHANGE = 0,
38     TCP_RST,
39     TCP_FIN
40 };
41
42 typedef struct{
43     uint8_t          status;
44     uint32_t         saddr;
45     uint32_t         daddr;
46     uint16_t         sport;
47     uint16_t         dport;
48     uint32_t         local_num;
49     uint32_t         remote_num;
50     time_t          last_used;
51     char            eth_dst[6];
52     char            eth_src[6];
53 } tcp_session_s;
54 tcp_session_s ts[MAX_TCP_SESSIONS];
55
56 /* 14 bytes Ethernet header */
57 typedef struct{
58     char            dst[6];
59     char            src[6];
60     uint16_t type;
61 } eth_hdr __attribute__((packed));
62
63 /* 20 bytes IP Header */
64 typedef struct{
65     uint8_t version    : 4 ;          /* Version */
66     uint8_t hlen       : 4 ;          /* Header length */

```

```

67  uint8_t  tos;                /* Type of service */
68  uint16_t length;            /* Total length */
69  uint16_t ident;             /* Identification */
70  uint16_t flags    : 3 ;     /* Flags */
71  uint16_t offset   : 13;    /* Fragment offset */
72  uint8_t  ttl;           /* Time to live */
73  uint8_t  protocol;     /* Protocol */
74  uint16_t checksum;     /* Header checksum */
75  uint32_t src;          /* Source address */
76  uint32_t dst;          /* Destination address */
77  }ip_hdr __attribute__((packed));
78
79 /* 20 bytes TCP Header */
80 typedef struct{
81  uint16_t sport;        /* Source port */
82  uint16_t dport;        /* Destination port */
83  uint32_t seq;          /* Sequence Number */
84  uint32_t ack;          /* Acknowledgement number */
85  uint16_t hdrlen    : 4; /* TCP Header length */
86  uint16_t reserved : 6; /* Reserverd - Zero */
87  uint16_t flags     : 6; /* Flags */
88  uint16_t win;      /* Window size */
89  uint16_t checksum; /* Header Checksum */
90  uint16_t urgptr;   /* Urgent pointer */
91  }tcp_hdr __attribute__((packed));
92
93 /* Exported methods */
94 int tcp_recv(void *dataVAddr, char *data, int *len, int *idx);
95 int tcp_send(void *data, int size, uint32_t saddr, uint32_t daddr, uint16_t sport, uint16_t dport);
96 void set_tcp_data(int ts_id, uint32_t sa, uint32_t da, uint16_t sp, uint16_t dp, uint32_t seq, uint8_t flags);
97 int find_create_tcp(uint32_t sa, uint32_t da, uint32_t sp, uint16_t dp, uint32_t seq, uint8_t flags);
98 int find_tcp(uint32_t sa, uint32_t da, uint32_t sp, uint16_t dp);
99 void get_eth_src(void *dataVAddr, char *eth_src);
100 void get_eth_dst(void *dataVAddr, char *eth_dst);
101 void get_ip_hlen(void *dataVAddr, uint8_t *ip_hlen);
102 void get_ip_data_len(void *dataVAddr, uint16_t *datagram_len);
103 void get_ip_src(void* dataVAddr, uint32_t *ip_sa);
104 void get_ip_dst(void* dataVAddr, uint32_t *ip_da);
105 void get_tcp_sport(void* dataVAddr, uint8_t ip_hlen, uint16_t *tcp_sp);
106 void get_tcp_dport(void* dataVAddr, uint8_t ip_hlen, uint16_t *tcp_dp);
107 void get_tcp_seq(void* dataVAddr, uint8_t ip_hlen, uint32_t *tcp_seq);

```

```

108 void get_tcp_ack(void* dataVAddr, uint8_t ip_hlen, uint32_t *tcp_ack);
109 void get_tcp_flags(void* dataVAddr, uint8_t ip_hlen, uint8_t *tcp_flags);
110 void get_tcp_hlen(void* dataVAddr, uint8_t ip_hlen, uint8_t *tcp_hlen);
111 void prepare_iph(ip_hdr *iph);
112 void prepare_tcph(tcp_hdr *tcph);
113 void debug_print_TCPsess(tcp_session_s tcps);
114
115 #endif /* #ifndef _BOGUS_TCP_H */

```

A.2.6 bogus_tcp.c

```

1 #include "bogus_tcp.h"
2
3 /*
4  * (C)2003, Matija Puzar <matija@ifi.uio.no>
5  * Adapted by Andreas Petlund 2005
6  */
7
8 /* Input: ipData: Pointer to beginning of IP packet
9  Returns: B_TCP_ERROR if error, B_TCP_NO_DATA if no data, B_TCP_DATA if data is received.
10 data: pointer to packet payload (should be at least 1500 bytes long).
11 len: length of packet payload.
12 idx: index into tcp_struct_array (to get port and IP data for IC sessions)*/
13 int tcp_recv(void *ipData, char *data, int *len, int *idx){
14 /* We already know this is a TCP packet destined for port 9070 */
15
16 char eth_src[6], eth_dst[6], resp_data[54];
17 uint8_t tcp_flags, new_tcp_flags, ip_hlen, tcp_hlen;
18 uint16_t tcp_sp, tcp_dp, datagram_len, tcp_len, data_len;
19 uint32_t ip_sa, ip_da, tcp_seq, tcp_ack;
20 int ts_id, tcp_ack_incr = 0, fin_incr = 0;
21
22 *idx = -1;
23 /* Get necessary data */
24 /* Get eth src and dst addr */
25 get_eth_src(ipData, eth_src);
26 debug_print_buffer32(eth_src, 6);
27 get_eth_dst(ipData, eth_dst);
28 debug_print_buffer32(eth_dst, 6);
29

```

```

30  /* Get necessary data from IP header */
31  get_ip_hlen(ipData, &ip_hlen);
32  get_ip_data_len(ipData, &datagram_len);
33  get_ip_src(ipData, &ip_sa);
34  get_ip_dst(ipData, &ip_da);
35
36  /* Let's get all the necessary data from the TCP header */
37  get_tcp_sport(ipData, ip_hlen, &tcp_sp);
38  get_tcp_dport(ipData, ip_hlen, &tcp_dp);
39  get_tcp_seq(ipData, ip_hlen, &tcp_seq);
40  get_tcp_ack(ipData, ip_hlen, &tcp_ack);
41  get_tcp_flags(ipData, ip_hlen, &tcp_flags);
42  get_tcp_hlen(ipData, ip_hlen, &tcp_hlen);
43  tcp_len = datagram_len - (ip_hlen * 4); /* hlen is given in 32 bit words */
44  data_len = tcp_len - (tcp_hlen * 4); /* hlen is given in 32 bit words */
45
46  /* Let's find the TCP session which corresponds to the given data (or create a new one) */
47  ts_id = find_create_tcp(ip_sa, ip_da, tcp_sp, tcp_dp, tcp_seq, tcp_flags);
48  /* Copy MAC addresses into struct
49      --- Do this only once (to be implemented :) --- */
50  memcpy(ts[ts_id].eth_src, eth_src, 6);
51  memcpy(ts[ts_id].eth_dst, eth_dst, 6);
52
53  /* Now, let's handle the received flags and set the sending ones.
54   * By default, we set the ACK flag
55   * tcp_ack_incr is set to 1 if we got a SYN or FIN flag (each of them "consumes" 1 sequence number)
56   */
57  new_tcp_flags = TCP_FLAGS_ACK;
58  tcp_ack_incr = 0;
59
60  /* If we got a SYN packet, we will respond also with a SYN packet */
61  if (tcp_flags & TCP_FLAGS_SYN) {
62      new_tcp_flags |= TCP_FLAGS_SYN;
63      tcp_ack_incr = 1;
64  }
65
66  /* By default we push all received data immediately, so here we don't do anything */
67  if (tcp_flags & TCP_FLAGS_PSH) {
68      DBG_MSG("TCP_FLAGS_PSH\n");
69  }
70

```

```

71  /* If we get an ACK, let's check if we have to switch to a different state */
72  if (tcp_flags & TCP_FLAGS_ACK) {
73      switch (ts[ts_id].status) {
74
75          /* We received the third TCP packet (last in the 3-way handshake operation) */
76          case TCP_SYN_RCVD: {
77              //fprintf(stderr, "Received 3hs ACK, dropping packet\n");
78              ts[ts_id].status = TCP_ESTABLISHED;
79              return B_TCP_NO_DATA;
80          }
81
82          /* We received the last ACK, after already closing the connection */
83          case TCP_CLOSE_WAIT: {
84              //fprintf(stderr, "Received last ACK, dropping packet\n");
85              ts[ts_id].status = TCP_CLOSED;
86              return B_TCP_NO_DATA;
87          }
88
89          /* We received the last ACK but still need to send ours */
90          case TCP_FIN_WAIT: {
91              ts[ts_id].status = TCP_TIME_WAIT;
92              break;
93          }
94
95          default: break;
96      }
97  }
98
99  /* If we get a FIN, we should respond with an ACK and, if we didn't initiate
100  * a FIN before, we will do it now and pass to the CLOSE_WAIT state
101  */
102 if (tcp_flags & TCP_FLAGS_FIN) {
103     if (ts[ts_id].status == TCP_ESTABLISHED) {
104         new_tcp_flags |= TCP_FLAGS_FIN;
105         ts[ts_id].status = TCP_CLOSE_WAIT;
106     }
107     tcp_ack_incr = 1;
108 }
109
110 /* If we get a seq. number that we already ACK-ed, we don't do it again */
111 if (ts[ts_id].status == TCP_ESTABLISHED && data_len == 0 && tcp_seq == ts[ts_id].remote_num) {

```

```

112     //fprintf(stderr, "Already ACK-ed, dropping packet\n");
113     return B_TCP_NO_DATA;
114 }
115
116 /* If any payload, copy it into buffer pointed to by data */
117 if ((data_len > 0) && (ts[ts_id].status != TCP_FIN_WAIT)) {
118     memcpy(data, ipData+ETH_HDR_SIZE+(ip_hlen*4)+(tcp_hlen*4), data_len);
119     *len = data_len;
120     *idx = ts_id;
121 }
122
123 /* Now, let's send a response */
124 memset(resp_data, 0, sizeof(resp_data));
125 if (tcp_seq == ts[ts_id].remote_num && ts[ts_id].status != TCP_CLOSED) {
126
127     /* Prepare ethernet header */
128     eth_hdr ethh;
129     memset(&ethh, 0, sizeof(ethh)); /* make sure struct is zeroed */
130     memcpy(&ethh.dst, eth_src, sizeof(eth_src));
131     memcpy(&ethh.src, eth_dst, sizeof(eth_dst));
132     ethh.type = ETH_TYPE_IP;
133
134     /* Let's prepare the IP header as well */
135     ip_hdr iph;
136     prepare_iph(&iph); /* Fill in predefined fields */
137     iph.src = ip_da;
138     iph.dst = ip_sa;
139     iph.length = sizeof(resp_data)-ETH_HDR_SIZE;
140
141     /* remote_num is the remote sequence number we expect next */
142     ts[ts_id].remote_num = tcp_seq + data_len + tcp_ack_incr;
143
144     /* Now, we prepare the TCP header as well */
145     tcp_hdr tcp_h;
146     prepare_tcp_h(&tcp_h); /* Fill in predefined fields */
147     iph.src = ip_da;
148     tcp_h.sport = tcp_dp;
149     tcp_h.dport = tcp_sp;
150     tcp_h.seq = ts[ts_id].local_num;
151     tcp_h.ack = ts[ts_id].remote_num;
152     tcp_h.flags = new_tcp_flags;

```

```

153
154     /* Copy headers to buffer */
155     memcpy(resp_data, &ethh, sizeof(ethh));
156     memcpy(resp_data+ETH_HDR_SIZE, &iph, IP_HDR_SIZE);
157     memcpy(resp_data+ETH_HDR_SIZE+IP_HDR_SIZE, &tcph, TCP_HDR_SIZE);
158
159     /* Calculate checksums */
160     ip_checksum(resp_data + ETH_HDR_SIZE);
161     tcp_udp_checksum(resp_data + ETH_HDR_SIZE, IP_HDR_SIZE + TCP_HDR_SIZE);
162
163     /* local_num has the next sequence number on our side */
164     ts[ts_id].local_num += tcp_ack_incr + fin_incr;
165
166     } else {
167         /* Something went wrong */
168         DBG_MSG("B_TCP_ERROR\n");
169         return B_TCP_ERROR;
170     }
171
172     /* In this app, only port 0 (external) will need TCP */
173     ixp_send_packet(resp_data, sizeof(resp_data), 0);
174
175     /* If we were in the TIME_WAIT state and sent our last ACK,
176     * we close the connection totally */
177     if (ts[ts_id].status == TCP_TIME_WAIT)
178         ts[ts_id].status = TCP_CLOSED;
179
180     /* If we get a RST flag, we reset the connection and discard the packet */
181     if (tcp_flags & TCP_FLAGS_RST) {
182         ts[ts_id].status = TCP_CLOSED;
183     }
184
185     /* Give proper return value */
186     if(data_len > 0)
187         return B_TCP_DATA;
188     else
189         return B_TCP_NO_DATA;
190 }
191
192 /* Input: data:    Pointer to payload.
193     size:         Size of payload in bytes.

```



```

194         saddr: Source IP address.
195         daddr: Destination IP address.
196         sport: Source port.
197         dport: Destination port.
198         change: Status change (i.e. FIN or RST) (0 if no change)
199         Output: Nuber of bytes sent, -1 if error. */
200 int tcp_send(void *data, int size, uint32_t saddr, uint32_t daddr, uint16_t sport, uint16_t dport){
201     int tcpIdx, psize = size + TCP_HDR_SIZE + IP_HDR_SIZE + ETH_HDR_SIZE;
202     char *packet = malloc(psize);
203
204     tcpIdx = find_tcp(saddr, daddr, sport, dport);
205     if(tcpIdx == -1)
206         return tcpIdx;
207
208     DBG_MSG("tcp_send(): Found TCP stream. Index: %i\n", tcpIdx);
209
210     /* Prepare ethernet header */
211     eth_hdr eh;
212     memcpy(&eh.dst, &ts[tcpIdx].eth_src, 6);
213     memcpy(&eh.src, &ts[tcpIdx].eth_dst, 6);
214     eh.type = ETH_TYPE_IP;
215
216     DBG_MSG("After ETH header prepare");
217
218     /* Let's prepare the IP header as well */
219     ip_hdr iph;
220     prepare_iph(&iph); /* Fill in predefined fields */
221     iph.src = daddr;
222     iph.dst = saddr;
223     iph.length = (uint16_t)(psize-ETH_HDR_SIZE);
224
225     DBG_MSG("After IP header prepare");
226
227     /* Now, we prepare the TCP header as well */
228     tcp_hdr tcpH;
229     prepare_tcpH(&tcpH); /* Fill in predefined fields */
230     tcpH.sport = ts[tcpIdx].dport;
231     tcpH.dport = ts[tcpIdx].sport;
232     tcpH.seq = ts[tcpIdx].local_num;
233     tcpH.ack = ts[tcpIdx].remote_num;
234     tcpH.flags = TCP_FLAGS_PSH | TCP_FLAGS_ACK;

```

```

235
236   DBG_MSG("After TCP header prepare");
237
238   /* Increment local_num with number of bytes sent */
239   ts[tcpIdx].local_num = ts[tcpIdx].local_num + size;
240
241   /* Copy headers to packet */
242   memcpy(packet, &eh, ETH_HDR_SIZE);
243   memcpy(packet+ETH_HDR_SIZE, &iph, ETH_HDR_SIZE+IP_HDR_SIZE);
244   memcpy(packet+ETH_HDR_SIZE+IP_HDR_SIZE, &tcph, ETH_HDR_SIZE+IP_HDR_SIZE+TCP_HDR_SIZE);
245   /* Copy payload to packet */
246   memcpy(packet+ETH_HDR_SIZE+IP_HDR_SIZE+TCP_HDR_SIZE, data, size);
247
248   debug_print_buffer32(packet, psize);
249   /* Calculate checksums */
250   ip_checksum(packet + ETH_HDR_SIZE);
251   tcp_udp_checksum(packet + ETH_HDR_SIZE, size + IP_HDR_SIZE + TCP_HDR_SIZE);
252
253   ixp_send_packet(packet, psize, 0);
254   free(packet);
255
256   return size;
257 }

```

A.2.7 bogus_tcp_utils.c

```

1 #include "bogus_tcp.h"
2
3 /* Look up our table of established TCP connections – if found, return the index in the table
4  * If not, create a new entry */
5 int find_create_tcp(uint32_t sa, uint32_t da, uint32_t sp, uint16_t dp, uint32_t seq, uint8_t flags){
6   int i, pos = -1, min_time_id = 0, first_free = -1;
7   time_t min_time;
8
9   min_time = time(NULL);
10
11  for (i=0; i<MAX_TCP_SESSIONS; i++) {
12    if(ts[i].status != TCP_CLOSED && ts[i].saddr==sa && ts[i].daddr==da && ts[i].sport == sp && ts[i].dport ==
13      /* We have found an active stream that corresponds to the packet */
14      pos = i;

```

```

15     DBG_MSG("b_tcp: found stream\n");
16     break;
17 } else if(ts[i].status != TCP_CLOSED){
18     /* This is an active stream that does not correspond to the packet.
19        Check when it was last used */
20     if (ts[i].last_used < min_time) {
21         min_time = ts[i].last_used;
22         min_time_id = i;
23     }
24 } else if(first_free == -1){
25     /* First free slot */
26     first_free = i;
27 }
28 }
29
30 /* Use the first free slot, if none, reuse the least recently used */
31 if (pos == -1){
32     if(first_free != -1){
33         pos = first_free;
34         DBG_MSG("b_tcp: using first free\n");
35     }else{
36         pos = min_time_id;
37         DBG_MSG("b_tcp: using lru\n");
38     }
39 }
40
41 set_tcp_data(pos, sa, da, sp, dp, seq, flags);
42
43 return pos;
44 }
45
46 /* Look up our table of established TCP connections.
47    If found, return the index in the table.
48    If not found, return -1 */
49 int find_tcp(uint32_t sa, uint32_t da, uint32_t sp, uint16_t dp){
50     int i;
51
52     for (i=0; i < MAX_TCP_SESSIONS; i++) {
53         if(ts[i].status != TCP_CLOSED && ts[i].saddr == sa
54            && ts[i].daddr == da && ts[i].sport == sp && ts[i].dport == dp) {
55             /* We have found an active stream that corresponds to the packet */

```

```

56     DBG_MSG("b_tcp: found stream\n");
57     return i;
58 }
59 }
60 DBG_MSG("find_tcp: No matching stream found\n");
61 return -1;
62 }
63
64 void set_tcp_data(int ts_id, uint32_t sa, uint32_t da, uint16_t sp, uint16_t dp, uint32_t seq, uint8_t flags){
65     if (flags & TCP_FLAGS_SYN){
66         ts[ts_id].remote_num = seq;
67         ts[ts_id].local_num = 1;
68         ts[ts_id].status = TCP_SYN_RCVD;
69     }
70     if (ts[ts_id].status == TCP_SYN_RCVD){
71         ts[ts_id].saddr = sa;
72         ts[ts_id].daddr = da;
73         ts[ts_id].sport = sp;
74         ts[ts_id].dport = dp;
75     }
76     ts[ts_id].last_used = time(NULL);
77 }
78
79 /* Input: Pointer to start of eth packet
80    Output: Eth src - 6 bytes */
81 void get_eth_src(void *dataVAddr, char *eth_src){
82     memcpy(eth_src, dataVAddr+6, 6);
83 }
84
85 /* Input: Pointer to start of eth packet
86    Output: Eth dst - 6 bytes */
87 void get_eth_dst(void *dataVAddr, char *eth_dst){
88     memcpy(eth_dst, dataVAddr, 6);
89 }
90
91 /* Input: Pointer to start of eth packet
92    Output: ip src - 4 bytes */
93 void get_ip_src(void* dataVAddr, uint32_t *ip_sa){
94     memcpy(ip_sa, dataVAddr+ETH_HDR_SIZE+12, 4);
95 }
96

```

```

97 /* Input: Pointer to start of eth packet
98   Output: ip dst - 4 bytes */
99 void get_ip_dst(void* dataVAddr, uint32_t *ip_da){
100   memcpy(ip_da, dataVAddr+ETH_HDR_SIZE+16, 4);
101 }
102
103 /* Input: Pointer to start of eth packet
104   Output: IP Hlen - 4 bit */
105 void get_ip_hlen(void *dataVAddr, uint8_t *ip_hlen){
106   memcpy(ip_hlen, dataVAddr+ETH_HDR_SIZE, 1);
107   *ip_hlen &= 0xf; /* Hlen is the 4 lsb */
108 }
109
110 /* Input: Pointer to start of eth packet
111   Output: IP length - 2 bytes */
112 void get_ip_data_len(void *dataVAddr, uint16_t *datagram_len){
113   memcpy(datagram_len, dataVAddr+ETH_HDR_SIZE+2, 2);
114 }
115
116 /* Input: dataVAddr: Pointer to start of eth packet
117   ip_hlen: length of IP header in 32 bit words
118   Output: TCP source port - 2 bytes */
119 void get_tcp_sport(void* dataVAddr, uint8_t ip_hlen, uint16_t *tcp_sp){
120   int tcp_hdr_off = ip_hlen * 4; /* ip_hlen is given in 32 bit words */
121   memcpy(tcp_sp, dataVAddr+ETH_HDR_SIZE+tcp_hdr_off, 2);
122 }
123
124 /* Input: dataVAddr: Pointer to start of eth packet
125   ip_hlen: length of IP header in 32 bit words
126   Output: TCP dest port - 2 bytes */
127 void get_tcp_dport(void* dataVAddr, uint8_t ip_hlen, uint16_t *tcp_dp){
128   int tcp_hdr_off = ip_hlen * 4; /* ip_hlen is given in 32 bit words */
129   memcpy(tcp_dp, dataVAddr+ETH_HDR_SIZE+tcp_hdr_off+2, 2);
130 }
131
132 /* Input: dataVAddr: Pointer to start of eth packet
133   ip_hlen: length of IP header in 32 bit words
134   Output: TCP seq - 4 bytes */
135 void get_tcp_seq(void* dataVAddr, uint8_t ip_hlen, uint32_t *tcp_seq){
136   int tcp_hdr_off = ip_hlen * 4; /* ip_hlen is given in 32 bit words */
137   memcpy( tcp_seq, dataVAddr+ETH_HDR_SIZE+tcp_hdr_off+4, 4);

```

```

138 }
139
140 /* Input: dataVAddr: Pointer to start of eth packet
141 ip_hlen: length of IP header in 32 bit words
142 Output: TCP ack - 4 bytes */
143 void get_tcp_ack(void* dataVAddr, uint8_t ip_hlen, uint32_t *tcp_ack){
144     int tcp_hdr_off = ip_hlen * 4; /* ip_hlen is given in 32 bit words */
145     memcpy( tcp_ack, dataVAddr+ETH_HDR_SIZE+tcp_hdr_off+8, 4);
146 }
147
148 /* Input: dataVAddr: Pointer to start of eth packet
149 ip_hlen: length of IP header in 32 bit words
150 Output: TCP flags - 6 bits */
151 void get_tcp_flags(void* dataVAddr, uint8_t ip_hlen, uint8_t *tcp_flags){
152     int tcp_hdr_off = ip_hlen * 4; /* ip_hlen is given in 32 bit words */
153     memcpy( tcp_flags, dataVAddr+ETH_HDR_SIZE+tcp_hdr_off+13, 1);
154     *tcp_flags &= 0x3F; /* all but the 6 lsb's */
155 }
156
157 /* Input: dataVAddr: Pointer to start of eth packet
158 ip_hlen: length of IP header in 32 bit words
159 Output: TCP header length - 4 bits */
160 void get_tcp_hlen(void* dataVAddr, uint8_t ip_hlen, uint8_t *tcp_hlen){
161     int tcp_hdr_off = ip_hlen * 4; /* ip_hlen is given in 32 bit words */
162     memcpy( tcp_hlen, dataVAddr+ETH_HDR_SIZE+tcp_hdr_off+12, 1);
163     *tcp_hlen >>= 4; /* Hlen is in the 4 msb's */
164 }
165
166 void prepare_iph(ip_hdr *iph){
167     memset(iph, 0, sizeof(ip_hdr)); /* make sure struct is zeroed */
168     iph->ttl = 64;
169     iph->version = 4;
170     iph->hlen = (sizeof(ip_hdr) / 4 ); /* hlen is given in 32 bit words */
171     iph->flags = 2;
172     iph->protocol = 6;
173 }
174
175 void prepare_tcp(tcp_hdr *tcp){
176     memset(tcp, 0, sizeof(tcp_hdr)); /* make sure struct is zeroed */
177     tcp->hdrlen = (sizeof(tcp_hdr) / 4); /* hlen is given in 32 bit words */
178     tcp->win = 0x16d0;

```

```

179 }
180
181 void debug_print_TCPsess(tcp_session_s tcps){
182     DBG_MSG("TCP session:\n");
183     DBG_MSG("tcps.status: %i\n", tcps.status);
184     DBG_MSG("tcps.saddr: %x\n", tcps.saddr);
185     DBG_MSG("tcps.daddr: %x\n", tcps.daddr);
186     DBG_MSG("tcps.sport: %i\n", tcps.sport);
187     DBG_MSG("tcps.dport: %i\n", tcps.dport);
188     DBG_MSG("tcps.local_num: %x\n", tcps.local_num);
189     DBG_MSG("tcps.remote_num: %x\n", tcps.remote_num);
190     DBG_MSG("tcps.last_used: %x\n", (uint32_t)tcps.last_used);
191     DBG_MSG("eth_dst: 0x%x\n", *((uint32_t*)&(tcps.eth_dst)));
192     DBG_MSG("eth_src: 0x%x\n", *((uint32_t*)&(tcps.eth_src)));
193 }

```

A.3 Host source

A.3.1 srtsp.h

```

1 #ifndef _SRTSP_H
2 #define _SRTSP_H
3
4 /* felles */
5 #include <pthread.h>
6 #include <stdio.h>
7 #include <linux/types.h>
8 #include <time.h>
9 #include <sys/types.h>
10 #include <sys/socket.h>
11 #include <sys/time.h>
12 #include <sys/sendfile.h>
13 #include <netinet/in.h>
14 #include <netinet/tcp.h>
15 #include <errno.h>
16 #include <ctype.h>
17 #include <netdb.h>
18 #include <string.h>
19 #include <arpa/inet.h>

```

```

20 #include <fcntl.h>
21 #include <stdlib.h>
22 #include <sys/stat.h>
23 #include <ctype.h>
24 #include <sys/resource.h>
25 #include <sys/mman.h>
26 #include <unistd.h>
27 #include <linux/unistd.h>
28 #include <syscall.h>
29
30 /* "Local" includes */
31 #include "dbg_msg.h"
32 #include "rtp.h"
33 #include "ic.h"
34
35
36 #define UDP_CORK 1
37 #define NSEC_PER_USEC (1000L)
38
39 #define MYPORT 23456 // the port users will be connecting to
40 #define BACKLOG 5 // how many pending connections queue will hold
41 #define MAX_RTP_STREAMS 100
42 #define ETH_TYPE_IP 0x0800
43 #define BIT_WRAP 0xffff
44
45 #define MEDIASTREAM_PLAY 1
46 #define MEDIASTREAM_PAUSE 2
47 #define MEDIASTREAM_FF 3
48 #define MEDIASTREAM_REWIND 4
49 #define MEDIASTREAM_TEARDOWN 5
50 #define MEDIASTREAM_STATISTICS_INIT 6
51 #define MEDIASTREAM_STATISTICS 7
52 #define MEDIASTREAM_STATUS 8
53
54 #define SESSION_STATUS_INIT 201
55 #define SESSION_STATUS_RUNNING 202
56 #define SESSION_STATUS_PAUSED 203
57 #define SESSION_STATUS_FF_REWIND 204
58 #define SESSION_STATUS_STOPPED 205
59
60

```



```

61 typedef struct {
62     uint32_t ipsrc;
63     uint32_t ipdst;
64     uint8_t notused; /* always zero */
65     uint8_t proto; /* protocol used */
66     uint16_t len; /* UDP len */
67 } pseudo_udp __attribute__((packed));
68
69 /* 14 bytes Ethernet header */
70 typedef struct{
71     char        dst[6];
72     char        src[6];
73     uint16_t type;
74 } eth_hdr __attribute__((packed));
75
76 /* 20 bytes IP Header */
77 typedef struct{
78     uint8_t version    : 4 ;           /* Version */
79     uint8_t hlen       : 4 ;           /* Header length */
80     uint8_t tos;        /* Type of service */
81     uint16_t length;    /* Total length */
82     uint16_t ident;    /* Identification */
83     uint16_t flags     : 3 ;           /* Flags */
84     uint16_t offset    : 13;          /* Fragment offset */
85     uint8_t ttl;       /* Time to live */
86     uint8_t protocol;  /* Protocol */
87     uint16_t checksum; /* Header checksum */
88     uint32_t src;      /* Source address */
89     uint32_t dst;      /* Destination address */
90 }ip_hdr __attribute__((packed));
91
92 /* 8 bytes UDP Header */
93 typedef struct{
94     uint16_t sport;    /* Source port */
95     uint16_t dport;    /* Destination port */
96     uint16_t len;      /* Length */
97     uint16_t csum;     /* Checksum */
98 }udp_hdr __attribute__((packed));
99
100 typedef struct {
101     uint32_t secs;

```

```

102  uint32_t frac;
103 } ntp64_t;
104
105 /* supported RTSP message types: */
106 enum {
107     RTSP_SETUP = 1,
108     RTSP_PLAY,
109     RTSP_UNKNOWN
110 };
111
112 struct playinfo {
113     pid_t pid;           /* initielt 0, settes etter retur fra første play kall. */
114     int in_fd;          /* fil som skal streames. */
115     int out_fd;         /* socket. */
116     loff_t start_offset; /* start offset i in_fd. */
117     loff_t stop_offset; /* stop offset i in_fd. */
118     long trans_interval; /* transmission interval:
119                             - mpeg ts: antall usec mellom hver pakke.
120                             - mpeg es: antall usec mellom hver frame.
121                             */
122     uint16_t max_pkt_len; /* max pakke lengde:
123                             - mpeg ts payload:
124                             - max_pkt_len-rtp_hdr_size.
125                             - max_pkt_len må være (188*n)+rtp_hdr_size.
126                             - mpeg es payload:
127                             - payload > 1 && payload < max_pkt_len-rtp_hdr_size-mpeg_hdr_size.
128                             */
129     uint32_t tsinc;     /* timestamp increment:
130                             - mpeg ts: increment per pakke
131                             - mpeg es: increment per frame
132                             */
133     uint8_t type;       /* type media */
134     /*RTP*/
135     uint32_t ssrc;     /* synchronization source */
136
137     unsigned long head;
138     unsigned long tail;
139     unsigned long datasize;
140 };
141
142

```

```

143 struct ff_rewind_info {
144     loff_t start_offset;
145     loff_t stop_offset;
146 };
147
148 struct ses_stat {
149     ntp64_t ntp_timestamp;
150     uint32_t rtp_timestamp;
151     uint32_t packet_count;
152     uint32_t octet_count;
153     uint16_t rtp_seq;
154 };
155
156 struct ses_stat_first {
157     uint32_t rtp_timestamp;
158     ntp64_t ntp_timestamp;
159 };
160
161 typedef struct {
162     //struct list_head list;
163     uint32_t sid;
164     pid_t pid;
165     uint8_t status;
166
167     int in_fd;
168     int out_fd;
169     loff_t start_offset;
170     loff_t stop_offset;
171     unsigned long buffaddr;
172
173     long trans_interval; //i usec
174     uint32_t tsinc;
175     uint16_t max_pkt_len;
176
177     /*RTP*/
178     uint16_t rtp_seq;
179     uint32_t rtp_timestamp;
180     uint32_t ntp_timestamp;
181     uint32_t ssrc;
182
183     /* The value of the RTP timestamp when the session opened (chosen randomly)

```

```

184  */
185  uint16_t  init_RTP_seq;
186  uint32_t  init_RTP_timestamp;
187  ntp64_t   init_NTP_timestamp;
188
189  struct timeval time_start; // settes til timeofday etter pause
190  /* current rtp timestamp offset*/
191  struct timeval time_elapsed; // settes til 0 ved pause
192
193  //DISSE NULLES UT FOR HVER GANG STAT KALLES
194  uint32_t  packet_count;
195  uint32_t  octet_count;
196
197  unsigned long head;
198  unsigned long tail;
199  unsigned long datasize;
200
201 } rtp_session_info_t;
202
203 enum {
204     RTPS_CLOSED = 0,
205     RTPS_ACTIVE
206 };
207
208 enum {
209     R_EGRESS = 0,
210     R_OUTBOUND
211 };
212
213 typedef struct{
214     uint8_t status;
215     int thdId; /* Id of streaming thread */
216     int client_rtp_port;
217     int client_rtcp_port;
218     int server_rtp_port;
219     int server_rtcp_port;
220     int icIdx; /* Index of corresponding IC session */
221     char filename[100];
222 } RTP_stream_t;
223 RTP_stream_t rtps[MAX_RTP_STREAMS];
224

```

```

225 /* Socket descriptor for comm with XScale */
226 int xscale_sock;
227
228
229 void debug_print_ICheader(void* ICpacket);
230 void debug_print_buffer32(void* buf, int size);
231 void debug_print_ICsession(int idx);
232 struct timeval usec2tv(long usec);
233 struct timeval AddTimes(struct timeval *time1, struct timeval *time2);
234 struct timeval SubTimes(struct timeval *time1, struct timeval *time2);
235 int getRtspType(char* data);
236 int getRtspSession(char *data);
237 void getFileName(char *data, char* fn);
238 int newICSession(void* ICheader, void* icHdrExt);
239 int newRtpStream(int icIdx, ICHeader_t *ich, char *data);
240 void getClientPorts(char *data, int *c_rtp_p, int *c_rtcp_p);
241 void createRtspReply(int ses, char *data, int type, char *reply);
242 int addICheader(IC_session_t ics, char *data, int size, int ICtype, int routeType);
243 void addICEExtHeader(IC_session_t ics, char *data, int size);
244 void sendToXScale(char *reply, int size);
245 uint32_t egressRoute(uint32_t route);
246
247 void partial_csum(uint32_t *sum, void* data, int len);
248 void udp_checksum(void* ipHdrStart, int size);
249 void ip_checksum(void* ipHdrStart);
250
251 #endif /* #ifndef _SRTSP_H */

```

A.3.2 srtsp.c

```

1 #include "srtsp.h"
2
3 /* Will send the whole file with packets in a predefined interval */
4 void *mpeg_ps(RTP_stream_t *streamInfo) {
5     char *path = "/opt/hardhat/previewkit/arm/xscale_be/target/opt/storage/", fullpath[100];
6     char *packet = malloc(1500); //char packet[1500];
7     int hdrs_size= sizeof(ICHeader_t) + sizeof(eth_hdr) + sizeof(ip_hdr) + sizeof(udp_hdr) + sizeof(rtp_hdr_t);
8     ICHeader_t *ich=(ICHeader_t*)packet;
9     eth_hdr *eth=(eth_hdr*)(packet+sizeof(ICHeader_t));
10    ip_hdr *iph=(ip_hdr*)(packet+sizeof(ICHeader_t) + sizeof(eth_hdr));

```

```

11  udp_hdr *udph=(udp_hdr*)(packet+sizeof(ICHeader_t) + sizeof(eth_hdr) + sizeof(ip_hdr));
12  rtp_hdr_t *rtph = (rtp_hdr_t*)(packet+sizeof(ICHeader_t) + sizeof(eth_hdr) + sizeof(ip_hdr)+sizeof(udp_hdr));
13  uint32_t new_rtp_seq, new_rtp_ts;
14  int fd, bytes_read;
15
16  DBG_MSG("*****NEW THREAD - ID: %i *****\n", (int)pthread_self());
17
18  /* Generate full path to file */
19  sprintf(fullpath, "%s%s", path, streamInfo->filename);
20
21  /* Open file */
22  fd = open(fullpath, O_RDONLY);
23  if(fd == -1){
24      DBG_MSG("open() failed\n");
25      DBG_MSG("Couldn't open file. error: %s. \n", strerror (errno));
26      goto out;
27  }else{
28      DBG_MSG("Successfully opened %s, fd: %i\n", fullpath, (int)fd);
29  }
30  /* Init IC header */
31  memset(ich, 0, sizeof(ICHeader_t));
32  ich->type = ICH_RTP;
33  ich->ttl = 4;
34  ich->route = htonl(ics[streamInfo->icIdx].route);
35
36  /* Init ETH header */
37  memset(eth, 0, sizeof(eth_hdr));
38  memcpy(eth->dst, ics[streamInfo->icIdx].eth_src, 6);
39  memcpy(eth->src, ics[streamInfo->icIdx].eth_dst, 6);
40  eth->type = htons(ETH_TYPE_IP);
41
42  /* Init IP header */
43  memset(iph, 0, sizeof(ip_hdr));
44  //iph->version = 4;
45  //iph->hlen = 5;
46  *((uint8_t*)iph) = 0x45;
47  //iph->flags = 2;
48  *((uint8_t*)iph+6) = 0x40;
49
50  iph->ttl = 64;
51  iph->protocol = 0x11;

```

```

52 iph->src = htonl(ics[streamInfo->icIdx].daddr);
53 iph->dst = htonl(ics[streamInfo->icIdx].saddr);
54
55 /* Init UDP header */
56 memset(udph, 0, sizeof(udp_hdr));
57 udph->sport = htons(streamInfo->server_rtp_port);
58 udph->dport = htons(streamInfo->client_rtp_port);
59 udph->csum = 0xffff;
60
61 /* Init RTP header */
62 memset(rtph, 0, sizeof(rtp_hdr_t));
63 rtph->version = RTP_VERSION;
64 //rtp_hdr.p = 0;
65 //rtp_hdr.x = 0;
66 //rtp_hdr.cc = 0;
67 //rtp_hdr.m = 0;
68 rtph->pt = PT_MP2P;
69 rtph->seq = htons((uint16_t)rand());
70 rtph->ts = htonl(rand());
71 rtph->ssrc = htonl(rand());
72
73 /******
74 /*          STREAM FILE          */
75 /******
76 while(streamInfo->status==RTPS_ACTIVE) {
77
78 /* Read file data into buffer+headers */
79 bytes_read = read( fd, packet + hdrs_size, 1400);
80 if(bytes_read > 0){
81
82 /* Update IC header */
83 ich->dataLen = htons(bytes_read + sizeof(eth_hdr) + sizeof(ip_hdr) + sizeof(udp_hdr) + sizeof(rtp_hdr_t));
84 /* Update IP header */
85 iph->length = htons(bytes_read + sizeof(ip_hdr) + sizeof(udp_hdr) + sizeof(rtp_hdr_t));
86 /* Update UDP header */
87 udph->len = htons(bytes_read + sizeof(udp_hdr) + sizeof(rtp_hdr_t));
88
89 /* Update RTP header */
90 new_rtp_seq = ntohs(rtph->seq);
91 rtph->seq = htons(new_rtp_seq++);
92 new_rtp_ts = ntohl(rtph->ts);

```

```

93     rtp->ts = htonl(new_rtp_ts++); /* Update this to correct values */
94
95     /* Calculate checksums */
96     ip_checksum(iph);
97     udp_checksum(iph, bytes_read + sizeof(rtp_hdr_t)+sizeof(udp_hdr));
98
99     /* Send packet to socket */
100    sendToXScale(packet, bytes_read + hdrs_size);
101
102    /* Wait appropriate amount of time */
103    usleep(2000);
104    }else{
105        close(fd);
106        DBG_MSG("*****BREAK*****\n");
107        break;
108    }
109 }
110 out:
111     DBG_MSG("***** EXIT THREAD *****\n");
112     free(packet);
113     pthread_exit(0);
114 }
115
116 int main(){
117
118     int sockfd; // listen on sock_fd, new connection on xscale_sock
119     struct sockaddr_in my_addr; // my address information
120     struct sockaddr_in their_addr; // connector's address information
121     int sin_size, dataLen, rtspType, icType;
122     char data[1500], icheader[20];
123     int status, recv_size, threadcount=0;
124     pthread_t thread[100];
125
126     sockfd = socket(AF_INET, SOCK_STREAM, 0); // do some error checking!
127     if(sockfd == -1){
128         printf("main(): socket returned error\n");
129     }else{
130         DBG_MSG("main(): Socket successfully created\n");
131     }
132
133     my_addr.sin_family = AF_INET; // host byte order

```



```

134 my_addr.sin_port = htons(MYPORT); // short, network byte order
135 my_addr.sin_addr.s_addr = INADDR_ANY; // auto-fill with my IP --endre til 192.168.1.1
136 memset(&(my_addr.sin_zero), '\0', 8); // zero the rest of the struct
137
138 status = bind(sockfd, (struct sockaddr *)&my_addr, sizeof(struct sockaddr));
139 if(status == -1){
140     printf("main(): bind returned error\n");
141     exit(1);
142 }else{
143     printf("main(): Port %i successfully bound\n", MYPORT);
144 }
145
146 status = listen(sockfd, BACKLOG);
147 if(status == -1){
148     printf("main(): listen() returned error\n");
149 }else{
150     printf("main(): Listening on port %i\n", MYPORT);
151 }
152
153 sin_size = sizeof(struct sockaddr_in);
154
155 status = accept(sockfd, (struct sockaddr *)&their_addr, &sin_size);
156 if(status == -1){
157     printf("main(): accept() returned error\n");
158 }else{
159     printf("main(): accepting connection on port %i\n", MYPORT);
160     xscale_sock = status;
161 }
162
163 /* Close listen-socket (will only need one connection) */
164 close(sockfd);
165
166 while(1){
167     /* Assumes that IC Header is the first to arrive,
168     and that all packets arrive in the order: ICHeader->Payload->ICHeader->Payload etc. */
169
170     /* Read IC header */
171     status = recv(xscale_sock, icheader, sizeof(ICHeader_t), 0);
172     if(status == -1){
173         DBG_MSG("main(): recv() returned error\n");
174     }else if(status > 0){

```

```

175     DBG_MSG("main(): Received IC header, %i bytes\n", status);
176 }
177
178 /* Get size of data from IC header
179     If ICH type is ICH_FIND_FILE, we also have to
180     receive an IC header extension with address-info */
181 ICHeader_t *ich = (ICHeader_t*)icheader;
182 icType = ich->type;
183 dataLen = ntohs(ich->dataLen);
184 if((icType == ICH_FIND_FILE))
185     recv_size = dataLen + sizeof(ICHeader_ext_t);
186 else
187     recv_size = dataLen;
188
189 status = recv(xscale_sock, data, recv_size, 0);
190 if(status == -1){
191     DBG_MSG("main(): recv() returned error\n");
192 }else if(status > 0){
193     DBG_MSG("main(): Received %i bytes of data.\n", status);
194 }
195
196 if(icType == ICH_FIND_FILE)
197     DBG_MSG("main(): Data: %s\n", data+sizeof(ICHeader_ext_t));
198 else
199     DBG_MSG("main(): Data: %s\n", data);
200
201
202 /* Get RTSP message type
203     take appropriate action
204     IF type==ICH_FIND_FILE data resides after IC header ext*/
205 if(icType == ICH_FIND_FILE){
206     DBG_MSG("ICH_FIND_FILE\n");
207     rtspType = getRtspType(data + sizeof(ICHeader_ext_t));
208 }else{
209     rtspType = getRtspType(data);
210 }
211 if(rtspType == RTSP_SETUP){
212     int rtpIdx, icIdx, packetSize, replysize; /* RTP-session index and IC-session index*/
213     char reply[200];
214     DBG_MSG("main(): Received RTSP_SETUP\n");
215

```

```

216     /* Setup icsession struct */
217     icIdx = newICSession(ich, data);
218     /* Create RTP stream struct */
219     rtpIdx = newRtpStream(icIdx, ich, data+sizeof(ICHeader_ext_t));
220
221     /* Give RTSP feedback */
222     createRtspReply(rtpIdx, (char*)(data+sizeof(ICHeader_ext_t)), RTSP_SETUP, reply);
223     replysize = strlen(reply);
224     packetSize = addICheader(ics[icIdx], reply, replysize, ICH_CTRL_FB, R_EGRESS);
225     addICExtHeader(ics[icIdx], reply, replysize);
226     packetSize+=sizeof(ICHeader_ext_t);
227     sendToXScale(reply, packetSize);
228
229 }else if(rtspType == RTSP_PLAY){
230     int rtpIdx, packetSize, replysize;
231     char reply[200];
232     DBG_MSG("main(): Received RTSP PLAY\n");
233
234     /* Give rtsp feedback */
235     rtpIdx = getRtspSession(data);
236
237     createRtspReply(rtpIdx, (char*)data, RTSP_PLAY, reply);
238     replysize = strlen(reply);
239     packetSize = addICheader(ics[rtps[rtpIdx].icIdx], reply, replysize, ICH_CTRL_FB, R_EGRESS);
240     addICExtHeader(ics[rtps[rtpIdx].icIdx], reply, replysize);
241     packetSize+=sizeof(ICHeader_ext_t);
242
243     sendToXScale(reply, packetSize);
244
245     /* Create rtp stream thread */
246     /* Implement handling of thread_exit */
247     pthread_create( &(thread[threadcount++]), NULL, (void*)mpeg_ps, &rtps[rtpIdx]);
248     //threadcount++;
249
250 }else if(rtspType == RTSP_UNKNOWN){
251     DBG_MSG("main(): Unknown RTSP type received, packet discarded.\n");
252 }
253 }
254 }

```

A.3.3 srtsp_utils.h

```
1 #ifndef _SRTPS_UTILS_H
2 #define _SRTPS_UTILS_H
3
4 #include <netinet/in.h>
5 #include <string.h>
6 #include <sys/time.h>
7 #include <time.h>
8 #include "dbg_msg.h"
9 #include "ic.h"
10
11 void debug_print_ICheader(void* ICpacket);
12 void debug_print_buffer32(void* buf, int size);
13
14 #ifndef _RTSP_TYPE
15 #define _RTSP_TYPE
16 enum {
17     RTSP_SETUP = 1,
18     RTSP_PLAY,
19     RTSP_UNKNOWN
20 };
21 #endif
22
23 struct timeval usec2tv(long usec);
24 struct timeval AddTimes(struct timeval *time1, struct timeval *time2);
25 struct timeval SubTimes(struct timeval *time1, struct timeval *time2);
26 int getRTSPType(char* data);
27
28 #endif /* #ifndef _SRTPS_UTILS_H */
```

A.3.4 srtsp_utils.c

```
1 #include "srtsp.h"
2
3 pthread_mutex_t socklock = PTHREAD_MUTEX_INITIALIZER;
4
5 void debug_print_ICheader(void* ICpacket){
6     ICHeader_t *ich = (ICHeader_t*)ICpacket;
7
```

```

8   DBG_MSG("IHeader:\n");
9   DBG_MSG("type: %u\n", ich->type);
10  DBG_MSG("ttl: %u\n", ich->ttl);
11  DBG_MSG("dataLen: %u\n", ich->dataLen);
12  DBG_MSG("route: %x\n", ich->route);
13
14  if((ich->type == ICH_FIND_FILE) || (ich->type == ICH_ROUTE_FB) || (ich->type == ICH_CTRL_FB)){
15      IHeader_ext_t *iche = (IHeader_ext_t*)(ICpacket+sizeof(IHeader_t));
16      DBG_MSG("sport: %u\n", iche->sport);
17      DBG_MSG("dport: %u\n", iche->dport);
18      DBG_MSG("saddr: %u\n", iche->saddr);
19      DBG_MSG("daddr: %u\n", iche->daddr);
20      DBG_MSG("eth_src: ");
21      debug_print_buffer32(iche->eth_src, 6);
22      DBG_MSG("eth_dst: ");
23      debug_print_buffer32(iche->eth_dst, 6);
24  }
25 }
26
27 void debug_print_ICsession(int idx){
28     if(idx < MAX_IC_SESSIONS){
29         DBG_MSG("IC-session:\n");
30         DBG_MSG("index: %i\n", idx);
31         DBG_MSG("status: %u\n", ics[idx].status);
32         DBG_MSG("route: %x\n", ics[idx].route);
33         DBG_MSG("sport: %u\n", ics[idx].sport);
34         DBG_MSG("dport: %u\n", ics[idx].dport);
35         DBG_MSG("saddr: %x\n", ics[idx].saddr);
36         DBG_MSG("daddr: %x\n", ics[idx].daddr);
37         DBG_MSG("eth_src: ");
38         debug_print_buffer32(ics[idx].eth_src, 6);
39         DBG_MSG("eth_dst: ");
40         debug_print_buffer32(ics[idx].eth_dst, 6);
41
42     }
43 }
44
45 /* Print buffer in 32 bit words hex in byte order*/
46 void debug_print_buffer32(void* buf, int size){
47     int i;
48

```

```

49  for(i=0; i < size; i++){
50      if(i % 4 == 0)
51          DBG_MSGG("\nbuf[%i]: 0x", i);
52
53      if((uint8_t*)(buf+i) == 0)
54          DBG_MSGG("00");
55      else if((uint8_t*)(buf+i) < 15)
56          DBG_MSGG("0%x", *((uint8_t*)(buf+i)));
57      else
58          DBG_MSGG("%x", *((uint8_t*)(buf+i)));
59  }
60  DBG_MSGG("\n");
61  }
62
63  /* Return RTSP message type, based on the first characters in the packet */
64  int getRtspType(char* data){
65      char *setup = "SETUP";
66      char *play = "PLAY";
67
68      if(strncmp(data, setup, strlen(setup)) == 0)
69          return RTSP_SETUP;
70      else if(strncmp(data, play, strlen(play)) == 0)
71          return RTSP_PLAY;
72      else
73          return RTSP_UNKNOWN;
74
75      return 0;
76  }
77
78  /* Initialize IC session from data
79     Input: Pointer to IC header
80     Output: index of new IC session (-1 if no free slot)
81  */
82  int newICSession(void* IChheader, void* icHdrExt){
83      ICHeader_t *ich = (ICHeader_t*)IChheader;
84      ICHeader_ext_t *iche = (ICHeader_ext_t*)icHdrExt;
85      int i, pos = -1;
86      /* locate first free slot */
87      for (i=0; i<MAX_IC_SESSIONS; i++) {
88          if(ics[i].status == ICS_CLOSED){
89              pos = i;

```

```

90     DBG_MSG("newICSession: first free: %i\n", pos);
91     break;
92 }
93 }
94
95 if(pos != -1){
96     ics[pos].status = ICS_ACTIVE;
97     ics[pos].route = ntohl(iche->route);
98     ics[pos].sport = ntohs(iche->sport);
99     ics[pos].dport = ntohs(iche->dport);
100    ics[pos].saddr = ntohl(iche->saddr);
101    ics[pos].daddr = ntohl(iche->daddr);
102    memcpy(&ics[pos].eth_src, &iche->eth_src, 6);
103    memcpy(&ics[pos].eth_dst, &iche->eth_dst, 6);
104 }
105 return pos;
106 }
107
108
109 /*Input: icIdx: index int array of ic-session structs
110         ich: pointer to start of IC header
111         data: pointer to packet data
112
113         Allocate RTP-stream struct
114         return index (session nr */
115 int newRtpStream(int icIdx, ICHeader_t *ich, char *data){
116     int i, pos=-1, c_rtp_p, c_rtcp_p;
117     char fn[100];
118
119     /* locate first free slot */
120     for (i=0; i<MAX_RTP_STREAMS; i++) {
121         if(rtps[i].status == RTPS_CLOSED){
122             pos = i;
123             DBG_MSG("newRTPStream: first free: %i\n", pos);
124             break;
125         }
126     }
127
128     getClientPorts(data, &c_rtp_p, &c_rtcp_p);
129     getFileName(data, fn);
130

```

```

131  if(pos != -1){
132      rtps[pos].status = RTPS_ACTIVE;
133      rtps[pos].client_rtp_port = c_rtp_p;
134      rtps[pos].client_rtcp_port = c_rtcp_p;
135      rtps[pos].server_rtp_port = 7654; /* bogus value */
136      rtps[pos].server_rtcp_port = 7655; /* bogus value */
137      rtps[pos].icIdx = icIdx;
138      strcpy(rtps[pos].filename, fn);
139  }
140  return pos;
141 }
142
143 /* Input: data: Pointer to RTSP SETUP message
144    Output: c_rtp_p: Client RTP port
145           c_rtcp_p: Client RTCP port
146 */
147 void getClientPorts(char *data, int *c_rtp_p, int *c_rtcp_p){
148     char *sp, *ep;
149
150     sp = strstr(data, "client_port=");
151     sp+=12; /* Move pointer to beginning of port number */
152     *c_rtp_p = strtol(sp, &ep, 10);
153     *c_rtcp_p = strtol(ep+1, NULL, 10);
154 }
155
156 /* Input: data: Pointer to RTSP SETUP message
157    Output: fn pointer to filename (\0 terminated)
158 */
159 void getFileName(char *data, char* fn){
160     char *sp, *ep;
161     int len;
162
163     sp = strstr(data, "9070/");
164     sp+=5;
165     ep = strstr(sp, "RTSP/1.0");
166     ep-=1;
167     len = ep-sp;
168     strncpy(fn, sp, len);
169     memset(fn+len, 0, 1); /* Zero-terminate string */
170 }
171

```



```

172 void createRtspReply(int ses, char *data, int type, char *reply){
173     char *cseq;
174
175     switch(type){
176
177     case RTSP_SETUP: {
178         cseq= strstr(data, "Cseq: ");
179         cseq+= 6;
180         sprintf(reply, "RTSP/1.0 200 OK\nCseq: %i\nSession: %i\nTransport: RTP/AVP;unicast;client_port=%i-%i\n",
181                 (int)strtol(cseq, NULL, 10), ses, rtps[ses].client_rtp_port, rtps[ses].client_rtcp_port,
182                 rtps[ses].server_rtp_port, rtps[ses].server_rtcp_port);
183
184         break;
185     }
186
187     /* May supply more info here */
188     case RTSP_PLAY: {
189         cseq= strstr(data, "Cseq: ");
190         cseq+= 6;
191         sprintf(reply, "RTSP/1.0 200 OK\nCseq: %i", (int)strtol(cseq, NULL, 10));
192
193         break;
194     }
195
196     default: break;
197
198     }
199 }
200
201 /* Assumes there is room in data buffer for header
202    Adds values in network byte order */
203 int addICheader(IC_session_t ics, char *data, int size, int ICtype, int routeType){
204     ICHeader_t *ich;
205
206     /* Make space for header */
207     memmove(data + sizeof(ICHeader_t), data, size);
208     ich = (ICHeader_t*)data;
209     /* Zero IC header */
210     memset(ich, 0, sizeof(ICHeader_t));
211
212     ich->type = ICtype;

```

```

213     ich->ttl = 4;
214     ich->dataLen = htons(size);
215     if (routeType == R_EGRESS)
216         ich->route = htonl(egressRoute(ics.route));
217     else
218         ich->route = htonl(ics.route);
219
220
221     return size + sizeof(ICHeader_t);
222 }
223
224 void addICExtHeader(IC_session_t ics, char *data, int size){
225     ICHeader_ext_t *iche;
226
227     /* Make space for header */
228     memmove(data + sizeof(ICHeader_t) + sizeof(ICHeader_ext_t), data+sizeof(ICHeader_t), size);
229     iche = (ICHeader_ext_t*)(data + sizeof(ICHeader_t));
230     /* Zero ICExt header */
231     memset(iche, 0, sizeof(ICHeader_ext_t));
232     iche->sport = htons(ics.sport);
233     iche->dport = htons(ics.dport);
234     iche->saddr = htonl(ics.saddr);
235     iche->daddr = htonl(ics.daddr);
236     memcpy(&iche->eth_src, &ics.eth_src, 6);
237     memcpy(&iche->eth_dst, &ics.eth_dst, 6);
238 }
239
240 void sendToXScale(char *data, int size){
241
242     pthread_mutex_lock( &socklock );
243     errno = send(xscale_sock, data, size, 0);
244     pthread_mutex_unlock( &socklock );
245     if(errno == -1){
246         DBG_MSG("send() returned error\n");
247     }else{
248         //DBG_MSG("Successfully sent %i bytes to socket\n", errno);
249     }
250 }
251
252 /* Removes the outgoing port from a route to make
253    the packet go to the XScale on the egress node

```

```

254     to be transmitted with in the correct TCP stream */
255 uint32_t egressRoute(uint32_t route){
256     uint8_t *rp;
257     int i;
258
259     rp = (uint8_t*)&route;
260     for(i=3; i>=0; i--){
261         if(rp[i]!=0xff && rp[i]==0){
262             rp[i]=0xff;
263             break;
264         }
265     }
266     return route;
267 }
268
269 /* Input: pointer to rtsp packet
270    Output: Session number for This RTP stream. */
271 int getRtspSession(char *data){
272     int ses;
273     char *sp;
274     sp = strstr(data, "Session: ");
275     sp+=9; /* Move pointer to beginning of session number */
276     ses = strtol(sp, NULL, 10);
277
278     return ses;
279 }
280
281 /* Calculates the 16 bit ones-complement sum of a given
282    buffer. Pads the last byte with 0 if odd size
283    Input: data: Pointer to buffer start
284           len: length of data */
285 void partial_csum(uint32_t *sum, void* data, int len) {
286     int i;
287
288     for(i = len>>1; i>0; i--){
289         *sum += *((uint16_t*)data);
290         data += sizeof(uint16_t);
291         /* Add carries */
292         *sum = (*sum >> 16) + (*sum & 0xffff);
293         *sum += *sum >> 16;
294     }

```

```

295
296 }
297
298 void udp_checksum(void* ipHdrStart, int size){
299     uint32_t sum = 0;
300     uint16_t finalsum;
301     pseudo_udp ph;
302
303     /* Construct pseudo header */
304     memcpy(&ph.ipsrc, ipHdrStart+12, 4);
305     //ph.ipsrc = ntohl(ph.ipsrc);
306     memcpy(&ph.ipdst, ipHdrStart+16, 4);
307     //ph.ipdst = ntohl(ph.ipdst);
308     ph.notused = 0;
309     ph.proto = 0x11;
310     ph.len = htons(size);
311
312     /* Calculate checksum of pseudo header */
313     partial_csum(&sum, &ph, 12);
314     /* Zero checksum field */
315     memset(ipHdrStart+26, 0, 2);
316     /* Calculate checksum of udp header and data*/
317     partial_csum(&sum, ipHdrStart+20, size);
318
319     /* Write the one's complement of the sum to
320        the correct spot in the TCP header */
321     finalsum = ~sum;
322     memcpy(ipHdrStart+26, &finalsum, 2);
323 }
324
325 void ip_checksum(void* ipHdrStart){
326     uint32_t sum=0;
327     uint16_t finalsum;
328
329     memset(ipHdrStart+10, 0, 2);
330     partial_csum(&sum, ipHdrStart, 20);
331
332     finalsum = ~sum;
333     memcpy(ipHdrStart+10, &finalsum, 2);
334 }

```

A.3.5 ic.h

```
1 #ifndef _IC_H
2 #define _IC_H
3
4 #include <inttypes.h>
5
6 #define MAX_IC_SESSIONS 100
7
8 enum {
9     ICH_FIND_FILE = 1, /* Packets used to locate file */
10    ICH_ROUTE_FB, /* Used to setup egress IC sessions */
11    ICH_CTRL_MSG, /* RTSP messages (except DESC and SETUP */
12    ICH_CTRL_FB, /* RTSP replies */
13    ICH_RTP /* RTP data Packets */
14 };
15
16 /* Standard IC header:
17    Common for all IC packets. */
18 typedef struct{
19     uint8_t      type; /* Type of IC packet */
20     uint8_t      ttl; /* To avoid circulating packets */
21     uint16_t     dataLen; /* Length of IC packet data (excluding the header) */
22     uint32_t     route; /* Routing information. Src routing */
23 } ICHHeader_t __attribute__((packed));
24
25 /* IC header extension:
26    Used to setup route information when assigning a new stream */
27 typedef struct{
28     uint16_t     sport; /* Client source port */
29     uint16_t     dport; /* Server port */
30
31     uint32_t     saddr; /* Client source IP addr. */
32     uint32_t     daddr; /* Egress IP addr */
33
34     char         eth_src[6]; /* Ethernet src address */
35     char         eth_dst[6]; /* Ethernet src address */
36 } ICHHeader_ext_t __attribute__((packed));
37
38 enum {
39     ICS_CLOSED = 0,
```

```

40  ICS_WAIT_ROUTE,
41  ICS_ACTIVE,
42 };
43
44 typedef struct{
45  uint8_t      status;
46  uint32_t     route; /* Routing information. Src routing */
47  uint16_t     sport; /* Client source port */
48  uint16_t     dport; /* Server port */
49
50  uint32_t     saddr; /* Client source IP addr. */
51  uint32_t     daddr; /* Egress IP addr */
52
53  char        eth_src[6]; /* Ethernet src address */
54  char        eth_dst[6]; /* Ethernet src address */
55 } IC_session_t;
56
57 /* Information about the IC sessions */
58 IC_session_t  ics[MAX_IC_SESSIONS];
59
60 #endif /* #ifndef _IC_H */

```

A.3.6 rtp.h

```

1  #ifndef _RTP_H
2  #define _RTP_H
3
4  #include <inttypes.h>
5
6  #define RTP_VERSION 2
7  #define PT_MP2T          33 /* MPEG2 TS A/V */
8  #define PT_MP2P          567 /* MPEG2 PS A/V */
9
10 /* RTP data header */
11 typedef struct {
12  uint16_t version:2;      /* protocol version */
13  uint16_t p:1;           /* padding flag */
14  uint16_t x:1;           /* header extension flag */
15  uint16_t cc:4;          /* CSRC count */
16  uint16_t m:1;           /* marker bit */

```

```

17  uint16_t pt:7;           /* payload type */
18  uint16_t seq;           /* sequence number */
19  uint32_t ts;            /* timestamp */
20  uint32_t ssrc;         /* synchronization source */
21  //u_int32 csrc[1];     /* optional CSRC list */
22 } rtp_hdr_t __attribute__((packed));
23
24 #endif /* #ifndef _RTP_H */

```

A.3.7 dbg_msg.h

```

1  #ifndef _DBG_MSG_H
2  #define _DBG_MSG_H
3
4  #define DEBUG
5
6  #include <stdio.h>
7
8  #ifdef DEBUG
9  #define   DBG_MSG(str,args...)           printf(str,##args)
10 #else
11 #define DBG_MSG(str, args...)
12 #endif
13
14 #endif /* #ifndef _DBG_MSG_H */

```