A Scalable, High Performance Active Network Node

Dan Decasper¹, Guru Parulkar², Sumi Choi², John DeHart², Tilman Wolf², Bernhard Plattner¹ [dan|plattner]@tik.ee.ethz.ch [guru|syc1|jdd|wolf]@arl.wustl.edu ¹Computer Engineering and Networks Laboratory, ETH Zurich, Switzerland Phone: +41-1-632 7019 Fax: +41-1-632 1035 ²Applied Research Laboratory, Washington University, St. Louis, USA Phone: +1-314-935 4586 Fax: +1-314-935 7302

October 1998

Abstract

Active networking in environments built to support link rates up to several gigabits per second poses many challenges. One such challenge is that the memory bandwidth and individual processing power of the router's microprocessors limit the total available processing power of a router. In this paper, we identify and describe three key components, which promise a high performance active network solution. This solution implements the key features typical to active networking, such as automatic protocol deployment and application specific processing, and it is suitable for a gigabit environment. First, we describe the hardware of the Active Network Node (ANN), a scalable, high performance platform based on off-the-shelf CPUs connected to a gigabit ATM switch backplane. Second, we introduce the ANN's modular, extensible and highly efficient operating system (NodeOS). Third, we describe an Execution Environment running on top of the NodeOS, which implements a novel large-scale active networking architecture called "Distributed Code Caching".

Key words: active networks; distributed code caching; gigabit active networking; scalable active network node

1 Introduction

Active networks [27] are packet-switched networks in which packets can contain code fragments that are executed on the intermediary nodes. The code carried by a packet may extend and modify the network infrastructure. The goal of active network research is to develop mechanisms to increase the flexibility and customizability of the network and to accelerate the pace at which network software is deployed. Applications running on end systems are allowed to inject code into the network to change the network's behavior to their favor.

Until recently, active networking research concentrated on two distinct approaches: "programmable switches" [4, 6] and "capsules" [23, 29]. These two approaches can be viewed as the two extremes in terms of how program code is injected into network nodes. Programmable switches typically upgrade by implicit injection of code by a network administrator. Research in the area of programmable switches focuses on how to upgrade network devices at run time, on upgrades introduced by administrators which support end system applications (e.g. congestion control for real-time data streams), or on a combination of both. Example applications include self-learning web caches, congestion control algorithms, on-line auctions, and sensor data mixing. Since the code is injected out-of-band, programmable switches provide no automated, on-the-fly upgrading functionality. Capsules, on the other hand, are packets carrying small amounts of program code, which is transported, in-band and executed on every node along a packet's path. This approach introduces a totally new paradigm to packet switched networks. Instead of "passively" forwarding data packets, routers execute the packet's code. The result of that computation determines what happens next to the packet. Applications include simple proof-of-concept ping applications, network diagnostic tools, active multicasting and more. This approach has the potential for an enormous impact on the future of networking. However, in the near future, security constraints will cause severe performance problems for capsule-based solutions. Capsules commonly make use of a virtual machine that interprets the capsule's code to safely execute it on a node. In order to ensure security, the virtual machines must restrict the address space a particular capsule might access, thus restricting the application of capsules. We expect network links to be 10 Gb/s or faster in the near future. With an optimistic average packet size of 512 bytes for IP traffic, a router has to process 2.6 million packets per second on every port, which is less than 380 nanoseconds per packet. A 300 MHz PentiumTM processor can therefore not spend on average more than 114 cycles to receive, process, and forward a packet just to keep up with the link speed. Even if we assume that a significant fraction of the packets forwarded do not require active processing and can be handled in hardware, it seems obvious that active network architectures based on virtual machines are not well suited to a multi-gigabit scenario. They may, however, be relevant to network management.

Recently, convergence between the pure "programmable switch" and the pure "capsule" approach became visible. Most of the research groups involved agree that some sort of code caching makes a lot of sense. The main motivation for this convergence is the realization that potential capsule code is more application specific than user specific. In the same way, users usually do not write their own applications but use off-the-shelf software. They are *not* expected to inject their *own* programs into the network, but use code from a set of code modules written by specialists. This allows for various optimizations in the form of caching, as we will see in the related work section. We will also show how our approach aggressively builds upon this same realization.

Another very important observation is that the deployment of multimedia data sources and applications (e.g. real-time audio/video, IP telephony) will produce longer lived packet streams (flows[†]) with more packets per session than is common in today's Internet. Especially for these kinds of applications, active networking offers very promising possibilities: media gateways; data fusion and merging; and sophisti-

[†] Flows are sequences of packets with a common five-tuple of IP header fields consisting of source address, destination address, source port number, destination port number, and protocol.

cated application specific congestion control. Both our hardware and software architectures support the notion of flows. In particular, the locality properties of flows are effectively exploited to provide for a highly efficient data path.

This paper describes the design of a high performance Active Network Node (ANN) that supports network traffic at gigabit rates and provides the required flexibility of Active Network technology for automatic, rapid protocol deployment and application specific data processing and forwarding. In section 2, we look at ongoing research conducted in other labs. Section 3 describes the hardware of our ANN, and section 4 describes the software platform running on top of that hardware. Section 5 shows how we integrate our Distributed Code Caching (DAN, [10]) architecture on top of our platform. Section 6 summarizes our ideas and takes a look at what we plan to do in the future.

2 Related Work

Active networking research has been ongoing for several years. Various research labs have described and implemented interesting approaches. In this section, we give an overview of some of these efforts.

2.1 MIT

Tennenhouse *et al.* [27] proposed "capsules", that is datagrams carrying small fragments of code, and an implementation in the form of an IP option [29]. The TCL language and a stripped-down TCL interpreter are used to provide safe execution of the code. Some simple, well-known network utilities (e.g. *traceroute*) have been implemented using capsules. So far, this work is mainly focused on a proof-of-concept for the capsule idea. To overcome security issues, the capsules in this approach are interpreted by a virtual machine. This method of execution introduces performance problems.

Further, this group proposed the ANTS [30] toolkit (downloadable code available). The main goal of ANTS is to provide an architecture for dynamic network protocol deployment. It introduces an optimization to the traditional capsule model. Instead of carrying code in every packet, packets carry pointers to code. This code is then loaded the first time it is needed from the previous hop along a packet's path. Java is used as a programming language for active code. This approach optimizes the bandwidth usage with the drawback of a considerable initial delay. Further, the usefulness of active reliable multicast is shown [18] using the ANTS platform. Results are measured in [17] and show significant performance gains. [17] also proposes various applications of active network technologies like sensor data mixing and inspired us to follow a similar approach. Recently, this group presented PAN [20], an approach similar to ANTS with the difference that instead of Java, machine code is transported in packets. This gives better forwarding performance, but security and interoperability issues are not addressed. Therefore, this approach has yet to prove its practical usefulness.

2.2 BBN

Smart Packets [23] is another capsule approach. The main focus is on implementation of extended diagnostic functionality in the network. A new compact programming language called "Sprocket" is specified and implemented. Its goal is to produce code that is compact enough to fit into an Ethernet packet. Sprocket programs are compiled into "Spanner" code that represents the assembly language for Smart Packets. Spanner code is interpreted in a virtual machine on the node receiving the packet. The programs are authenticated before interpretation and run-time limited during execution. With its clear focus on network management, SmartPackets provide a very powerful improvement over SNMP, which is used for management of conventional networks. The group is in the process of deploying SmartPackets on the CAIRN network.

2.3 Georgia Tech

Zegura *et al.* [6, 7, 32] introduced a generic view of network code as a set of functions which are called depending on identifiers found in data packets. Application specific data processing is implemented, as an example, for congestion control for MPEG video streams. The functions referred to in the data packets are loaded out-of-band into the network nodes.

This group also showed how self-organizing network caches could be built using active network technologies. They used simulations and analytical models to evaluate the performance gains offered by caching. This work shows that network caching is an application of active networks worth pursuing. Our system builds in part on the theoretical background and terminology introduced in [6], but significantly extends the system's capabilities.

2.4 University of Pennsylvania

The SwitchWare [3] project uses three important components: active packets, switchlets, and a secure active router infrastructure. Active packets are similar to MIT's capsules. Switchlets are dynamically loadable programs that provide specific services on the network nodes. Active packets are programmed in a simple language called PLAN (Programming Language for Active Networks). PLAN programs are strongly typed and statically type-checked to provide safety before being injected into the network. Further, PLAN programs are made secure by restricting their actions (e.g. a PLAN program cannot manipulate node-resident state). To compensate for these limitations, PLAN programs can call switchlets. Switchlet modules are written in a language (CAML) which supports formal methodologies to prove security properties of the modules at compile time and no interpretation is needed. The code fragments are authenticated by the developer and explicitly (and not on-demand) loaded into the switch. At the lowest layer, the Secure Active Network Environment (SANE) ensures the integrity of the entire environment. SANE identifies a minimal set of system elements (e.g. a small area of the BIOS) upon which system integrity is dependent and builds an integrity chain with cryptographic hashes on the image of the succeeding layer in the system, before passing control to that image. If an image is corrupted, it is automatically recovered from an authenticated copy over the network. Although the project shows very interesting properties, the main problem with this architecture seems to be that PLAN programs are not powerful enough for many applications. Therefore, switchlets have to be installed out-of-band to provide "handles" for the PLAN programs. This makes the system less flexible. Active Bridging [4] is an application of SwitchWare which shows reprogramming of a bridge with switchlets.

Another group at the University of Pennsylvania works on the Programmable Protocol Processing Pipeline (P4, [14]) project. They use reconfigurable Field Programmable Gate Arrays (FPGAs) to implement datagram processing functionality (forward error correction in their case) in hardware. Using this kind of hardware support looks very promising for the future of active networking. As far as we know, this is the only other group besides us to consider hardware support for active networking at this point in time.

2.5 University of Arizona

Scout [19] is a communication-oriented operating system. The kernel is a customized composition of low-level communication primitives that are implemented as modules. Modules implement independent functionality, like IP, UDP, or TCP protocols. Modules can be combined to form *paths* that build a logical channel over which I/O data flows. *Joust* [15] runs on top of Scout and consists of an implementation of the Java virtual machine (VM) including both the runtime system and a just-in-time compiler. The VM's API has been extended to interact closely with Scout and to allow applications to access low level resources. All fixed components are written in C or Java and compiled to machine code ahead of time. The Joust/Scout implementation of ANTS performs two to three times faster than an implementation using Sun's JDK and an of-the-shelf operating system like Linux. The Scout/Joust combination provides the fastest Java environment for active networks documented thus far. However, it still looks like it is not suited for high-volume, high-bandwidth traffic.

2.6 Columbia University

Netscript [31] is middleware for programming functions of intermediate network nodes. The Netscript programming language allows script processing of packet streams with a focus on routing, packet analyzers and signaling functions. Netscript programs are organized as mobile agents that are dispatched to remote systems and executed under local or remote control. The goal of Netscript is to simplify the development of networked systems and to enable their remote programming. The Netscript project envisions networks that support flexible programmability and dynamic deployment of software at all nodes. The Netscript language includes constructs and abstractions that greatly simplify the design of traffic-handling software. These abstractions hide the heterogeneous details of networked systems. Protocol messages are defined and encoded as high-level Netscript objects. Netscript programs are message interpreters that operate on streams of messages. Messages can be encoded either as high-level Netscript objects or in a format compatible with existing standards.

2.7 Implications on our work

The use of interpreted capsules has a lot of potential in areas like network management, where performance is not a primary concern. As demonstrated by others, using an interpretation-based approach delivers far more flexible mechanisms for network management than traditional approaches (e.g. SNMP). Java as the language of choice for capsules provides the advantages of having a lot of market and research momentum. This leads to a variety of available execution environments for Java capsules and increasingly higher quality of these environments. We therefore decided to support the ANTS execution environment on our ANN node architecture. However, for applications requiring a maximal amount of computation performance, minimal latency and high bandwidth, our DAN architecture promises to be better suited. We describe the ANN's hardware and software architecture in the next two sections, before we elaborate on the implementation of DAN on our platform.

3 The Active Networking Node Hardware

The current trend in Internet router technology is to implement an increasingly higher amount of datagram processing in hardware ASICs. Most modern high-performance IP routers are capable of forwarding standard datagrams (without special features like IP options) entirely in hardware. This is required for large backbone routers to keep up with link speeds. These routers typically use ASICs on every port, which have high-bandwidth access to a local table of routes. The central CPU is only involved in processing of non-standard datagrams and to implement control-path functionality like routing protocols.

By definition, active networking extends the amount of processing spent on a single packet and since the processing is application-specific for a potentially significant variety of applications, it cannot be implemented in ASICs. We believe, however, that with the following set of hardware design measures we can optimally address the problem. Active networking router hardware designed for high performance requires:

- A high number of processing elements compared to the number of router ports. A single central CPU attached to a backplane serving all ports is unable to keep up with link speed even for a small number of ports and relatively low bandwidth links (e.g. 10 Mbits/s). We use a general purpose CPU and a Field Programmable Gate Array (FPGA) *on every port* of a switch backplane. We call the combination of CPU and FPGA the *processing engine*. The CPU takes care of the majority of active functions applied to a packet, while the FPGA implements functions which are particularly performance critical in hardware. Both can be programmed on-the-fly.
- Tight coupling between a processing engine and the network, as well as between the processing engine and a switch backplane. Since the main limiting factors are processing power and memory bandwidth, one has to make sure that these valuable resources are used in the most effective fash-

ion. Two selective measures are applied. First, we benefit from the fact that most network traffic is flow-oriented. Bursts of packets share important forwarding properties that are, once determined, common to all packets of a particular flow. Thus, the majority of non-active packets allow cut-through routing directly through the switch backplane without CPU intervention. Second, by tightly coupling the processing engine to the link, packets arrive at an ANN with minimal overhead through zero-copy DMA.

• Scalable processing power to meet the demands of active processing of packets. Computation on active flows must be evenly distributed over the processing engines available.



Figure 1: Active Network Node (ANN)

The top level hardware architecture for the ANN is shown in Figure 1. It is derived from our high performance IP routing architecture [22] and has been refined and optimized for the purpose of active networks. The node consists of a set of Active Network Processing Elements (ANPE, four in Figure 1) connected to an ATM switch fabric [8]. ANNs are interconnected only through ANPES. The scalable switch fabric currently supports eight ports with data rates as high as 2.4 Gb/s on each port. The ANPE comprises a general-purpose processor, a

large FPGA (100,000 gates), and memory. The ANPEs are connected to the backplane via the ATM Port Interconnect Controller (APIC, [12]) chip. Other devices, like workstations and servers, are connected through a line card directly to the switch fabric (not shown in Figure 1).

Scalability is guaranteed through (1) the ability to configure any number of ANPEs which can be added to the ANN; (2) a scalable switch backplane; (3) a load sharing algorithm which dynamically distributes active flows over the ANPEs by configuring the corresponding APICs (setting/resetting cut-through switching of selected VCs) in order to move active flows from heavily loaded ANPEs to less loaded ones. Figure 1 shows an example data flow coming into the ANN at ANPE A and going out at ANPE D. The active processing is done in ANPE C since ANPE A is heavily loaded and the load-sharing algorithm directed the flow to ANPE C which finally directs the flow to the ANN connected to ANPE D (ANPE A and D switch the flow in hardware without CPU intervention through the APIC). We are developing an intra-ANN protocol to communicate the status of processing engine load between ANPEs on a reserved VC.

3.1 ANPE Architecture

The ANPE consists of an APIC ATM host adapter chip, a PentiumTM CPU, a large FPGA, and up to four gigabytes of DRAM.

The APIC (ATM Port Interconnect Controller) is an ATM host-network interface device with two ATM ports and a built-in PCI (Peripheral Component Interconnect) bus interface. Each of the ATM ports can be independently operated at full duplex rates ranging from 155 Mb/s to 1.2 Gb/s. The ATM cell handling is done entirely in hardware and structured so as not to affect the active networking software subsystem. The

APIC further implements VC switching in hardware and is capable of forwarding cells directly without passing them to the processing engine. This is used to implement the load-sharing algorithm as shown in the previous section. It allows the ANPE to forward plain (non-active) IP traffic without touching the processing engine (also called cut-through forwarding), which leaves valuable cycles available for active processing. The APIC provides a powerful host system interface featuring different modes of scatter-gather Direct Memory Access (DMA) to provide true zero-copy protocol processing. This allows the implementation of a very high-performance I/O subsystem that supports high-bandwidth and low latency.

The processing engine consists of an Intel PentiumTM CPU and a 240-pin FPGA, which can have a very large number of gates (up to 100,000). The CPU runs our NodeOS, which is an optimized version of NetBSD derived from our Router Plugins [9] architecture. We will describe this architecture in the next section. The FPGA can be programmed by the CPU on-the-fly to implement the most performance-critical algorithms in hardware. The APIC can distribute individual flows to the CPU or the FPGA on a per VC basis. A packet can first go to the FPGA and then either be passed to the CPU or forwarded straight through the APIC to the link, depending on whether there is additional processing required. We expect this combination of FPGA/CPU/cut-through processing to provide excellent performance for software based packet forwarding.

One or multiple ANPEs can be physically placed on one ANPE card (Figure 1 shows only one-CPU/FPGA/APIC ANPE cards). The main advantage of having multiple CPU/FPGA/APIC combinations on one card is that the loadsharing algorithm can save switch backplane bandwidth by distributing the active flows to processing engines on the same ANPE card. Figure 2 shows an ANPE card with two ANPEs and three flows of data packets traversing the card. The three flows are drawn as a dotted, a dashed, and a solid line. The lines show the three different modes of operation of an ANPE card. The flow shown as dashed line is routed through the first APIC into the first processing engine. The processing engine processes the packets and forwards them to one of the ports of the attached switch backplane. Note that the second APIC on the ANPE card routes this flow in cut-through mode. The flow shown as solid line is an example of a flow that is diverted to the



Figure 2: ANPE card

second processing engine on the same ANPE card. This happens in case the first processing engine is heavily loaded. In this case, the load-sharing algorithm picks the second processing engine to process the packets and configures both APICs appropriately. Finally, the dotted flow is cut-through routed through both the APICs without any CPU intervention. This could either be a regular IP flow not requiring active processing or an active flow that is diverted to another anpe card attached to another switch port. This would happen if both processing engines on the card are heavily used by active flows and there are other CPUs on other ANPE cards in the same ANN which are lightly loaded. The usage of these different modes of operation of the individual ANPE cards lead to optimal performance and scalability of the ANN as a whole.

4 ANN Software Infrastructure

To utilize the hardware architecture described in the previous section in an efficient way, our software architecture must be optimized to a similar extent for high performance. The main design goal is to provide a highly efficient data path and a flexible control path. All high-bandwidth data path components are implemented in the system's kernel, whereas all management components are implemented in user space. We embed the architecture described here into our Router Plugins research platform. We showed in [9] that a highly modular router software architecture could be implemented without any significant performance penalties. The architecture described here leverages the results from the Router Plugins work.

Within the active networking community, it is common to distinguish between the "NodeOS" [1] and "Execution Environments" (EEs) for active networking software architectures. The NodeOS represents the operating system components implementing services like packet scheduling, resource management, and packet classification, which are independent of a specific active networking implementation. The NodeOS offers these services to the EEs running on top of it. An EE implements active networking protocol specific processing. For example, there can be an independent EE for ANTS, for SmartPackets or for SwitchWare.

Our software architecture is shown in Figure 3. It supports two EEs, namely the ANTS EE and the Distributed Code Caching for Active Networks (DAN, [10]) EE. IP can be viewed as another EE with the distinguishing property that the other EEs typically can not work without IP since they use it for routing and forwarding.

We target the use of ANTS to network management tasks and experimental prototyping of network protocols. ANTS is described in [30]. DAN is described in more detail in the next section. It represents our own active network archi-



Figure 3: Software architecture

tecture specially targeted at high-bandwidth, low-latency applications.

Before we give an overview of the individual NodeOS components as shown in Figure 3, we introduce some of the general concepts of this architecture.

In the context of our architecture, we call code blocks implementing application-specific network functions *active plugins*. Active plugins contain code that is downloaded and installed on the node. The downloading is triggered by the occurrence of a reference in a datagram as shown in the section on DAN, by a special configuration packet, or by an administrator. Active plugins can create instances. The terminology is intentionally derived from object-oriented programming since the semantics are similar. Instances are flow-specific configurations of active plugins. The individual properties of instances are EE and plugin specific. For example, an IP instance consists of the code that forwards the packet and the required information about the interface on which the packet has to be forwarded. However, all instances use the same well-

defined API that embeds them into the system. The API consists of an *entry()* and an *exit()* function among others. The *entry()* function is called to pass a packet to the instance. The *exit()* function is called by the instance when it is done with packet processing.

We make an important distinction between the first few packets of a flow and subsequent packets. The reception of the first packets of a flow usually causes the plugins to create an instance for the new flow. If the packet is passed to multiple instances, these instances are chained together by making the first plugin's *exit()* function call the second plugins *entry()* function and so on. Every



Figure 4: Instance chaining and labeling

EE is allowed to request a *selector* from the NodeOS to label the chain of instances. The selector is sent to the upstream node by the NodeOS. The upstream node puts the selector in subsequent packets of the same flow, which allows the downstream node to efficiently lookup the state information using hashing, and directly assign the flows state information to the packet. Note that this is similar to tag switching. We introduced the Simple Active Packet Format (SAPF), which describes the format of the selector, in [28]. Figure 4 schematically depicts instance chaining and labeling using a selector. As described below, the Selector Dispatcher implements flow lookups based on selectors. The Packet Scheduler is called last to send the packet off to the network.

While per flow instance creation and management introduces a certain amount of overhead, the payoff in the context of active networks is dramatic for subsequent packets. Other than the selector, no demultiplexing has to be performed and the operations of instances are reduced to only those which vary from packet to packet of the same flow (e.g. there is no routing lookup). Note that all flow specific information has soft-state characteristics: it is automatically removed when no packets of a given flow are received for a configurable amount of time.

For the rest of this section, we give a high level overview of the NodeOS building blocks. The DAN EE is discussed in section 5, the ANTS EE in [30].

4.1 NodeOS components

The kernel consists of the following components (from bottom to top):

- **Device Drivers / Layer 2 processing (DD)**: the DD are standard NetBSD device drivers implementing network hardware specific send and receive functions. They are modified only in two ways: First, to allow packet scheduling, they do not implement packet queues on their own. Second, they pass incoming packets to the Selector Dispatcher instead of the IP stack if the packet contains a selector. If no selector is present, the packet is first passed to the packet classifier and then to the IP stack.
- **Packet Classifier (PC)**: all packets not carrying a selector are passed to the packet classifier. We implemented a highly efficient packet classifier based on a Directed Acyclic Graph (DAG) in the context of the Router Plugins platform. The packet is classified on a five-tuple of IP header fields and the interface on which it is received. The five-tuple consists of a pair of IP addresses, port numbers and the protocol used. The PC allocates a flow record for every packet of a new flow. It tags all incoming packets with a flow index (FIX) which is carried in the packet's mbuf¹ and points to the

¹ The mbuf is a data structure that is used to store packets and packet related information efficiently in BSD derived operating system kernels.

packet's flow record. Plugins can access the flow record through the FIX. It stores all flow related information.

- Selector Dispatcher (SD): the SD scans a data packet for its SAPF selector. Using the selector, it hashes into a table to find the flow's FIX. The same hash table stores the outgoing selector, which will replace the incoming selector in the packet before it, is sent to the downstream node. Further, it stores a pointer to the first instance in the flow's instance chain as previously described and shown in Figure 4. To negotiate new selector values, the selector dispatcher offers an API to send and receive messages from neighbor nodes.
- **Packet Scheduler** (PS): the PS in use is a modified version of a Deficit Round Robin (DRR [25]) scheduler which allows bandwidth reservations using filters in addition to fair queueing. Besides DRR, we plan to use a port of CMU's Hierarchical Fair Service Curves (HFSC, [26]) scheduler which represents the state-of-the-art in flow based packet scheduling providing hierarchical scheduling and decoupling of bandwidth and delay. Both schedulers can be programmed through the PS API thus allowing plugins as well as the administrator to reserve resources.
- **Resource Controller (RC):** the RC keeps track of the CPU cycles and memory consumed by active plugin instances. The RC is responsible for fair CPU time sharing between different instances. Since NetBSD is not a real time operating system and does not have a preemptive or multithreaded kernel, we implement this by keeping track of the consumed CPU cycles on a per flow basis. Incoming packets are enqueued in an input queue associated with the flow. Input queues are served in a round robin fashion. On reception of a packet, the system goes through the following steps:
 - get packet from network card
 - find corresponding flow (call SD or PC)
 - enqueue packet in flow's input queue
 - pick packet to process from the set of all input queues
 - timestamp selected packet with CPU TCS^2
 - continue processing with that packet
 - before enqueueing the packet on the output queue, read TCS again and add difference to flow total
 - enqueue packet in output queue

The RC implements the selection of the queue. We will investigate different schemes to pick the right packet. Besides per-flow CPU distribution, we will measure what fraction of the total time the system spends processing packets. This will give us an idea of how heavily our system is loaded at any given time. All RCs in an ANN periodically exchange this quantity with each other on a reserved VC to serve as input parameters for the load-sharing algorithm.

The second quantity worth observing is the active plugin's memory consumption. The plugin must be restricted to an upper limit of memory usage by policy or depending on the current average utilization of the networking subsystem. We modify the kernel's memory management to keep track of memory usage on a per instance basis and possibly deny additional memory to greedy instances. An important resource management issue has to do with the ANN running out of resources such as CPU or memory capacity. Besides implementing the load-sharing algorithm, we explore both policies and mechanisms which can do: (1) effective "admission control" to ensure sufficient resources for admitted active connections; (2) cache management for active plugins to decide which active plugins to replace to create room for the active plugins fetched on demand to be used immediately.

• **Plugin Control Unit (PCU):** the PCU manages plugins, and is responsible for forwarding control path messages (e.g. for instance creation and registration messages) to individual plugins from

 $^{^{2}}$ The Pentiums TCS register is a 64 bit register which is incremented by one on every clock cycle.

other kernel components, as well as from user space programs using ANN library calls (we provide the ANN library with our system). Plugins register themselves with a callback function and an identifier (plugin code) when they get loaded into the kernel. All plugins must reply to a defined set of control messages, (e.g. messages to create and free instances).

• **Plugin Manager (PM):** the PM is a user-space utility for configuring the system. It is a simple application that takes arguments from the command line and translates them into calls to the user-space library. This library implements the function calls needed to configure all kernel level components. In most cases, the plugin manager is invoked from a configuration script during system initialization, but it can also be used to manually issue commands to various plugins. This is especially useful to test new active plugins.

The design of this software architecture is clearly driven by the goal of allowing efficient implementations of EEs and active plugins. While we favor our DAN architecture with this design, we make all components accessible to other EEs through a well-documented API.

5 The Distributed Code Caching Approach to Active Networking

Distributed Code Caching for Active Networks (DAN) has been described in [10]. Since [10] has been published, the architecture has matured significantly. We review the important ideas here, give an update reflecting the latest development, and elaborate on how DAN is embedded into the NodeOS described in the previous section.

Before we present our approach, we review the basic requirement for active networking. It is to allow users and applications to control networking nodes and how their packets are processed and forwarded. This necessitates computing and programmability at each network node. However, this requirement should not considerably degrade the performance of an EE through excessively complex and inefficient security mechanisms. In other words, per-packet processing should not require a long and inefficient software path. Thus, the fundamental challenge that high performance active networking poses can be summarized as follows:

Allow relocating part of the processing from the endsystems into the network, however minimize the amount of processing on a single node and make the processing as efficient as possible while keeping the flexibility and customizability that the active networking paradigm introduces.

We believe that our architecture, which we call "Distributed Code Caching", does just that.

5.1 Distributed Code Caching

To overcome the performance-related problems that will exist, at least in the near term, for capsules, we think that a combination of the programmable switch and the capsule approaches is very appealing. We replace the capsules' program code by a reference to an active plugin stored on a code server. On a reference to an unknown code segment in a router or an end system, the code is automatically downloaded from a code server. It is important to note that the code fragment or plugin is dynamically linked and executes like native code on the router/node, and thus, it runs as fast as any other code. The security issues are addressed by usage of well-known cryptography techniques (as explained later in this section), and thus, our scheme does not require slow virtual machines. This introduces some restrictions regarding the authorship and the source of active network code for the benefit of security and performance but we believe this to be an appropriate compromise. To explain our idea, we first analyze the layout of data packets and network nodes common in today's networks.

Each network node typically supports a particular set of functions that may be applied to data packets. One or more unique identifiers in the packet's header identify these functions. When a packet is processed, the referenced functions are applied to the data of the packet. The packet's data can be viewed as the input parameter to a function. An Ethernet packet, for example, contains a unique identifier for the upper layer protocol (0x0800 for IPv4, 0x08dd for IPv6). By demultiplexing an incoming packet on this value, the kernel decides to which function or set of functions the packet gets passed next. Packets consist of a finite sequence of such identifiers for functions and input parameters. The functions are normally daisy-chained in a sense that one function calls the next according to the order of the identifiers in the data packet. The first function is determined by the hardware (the interface on which the packet is received) and the last function or set of functions is implemented in the application consuming the packet. Each of the functions may also decide not to call the next function for several reasons (e.g. forwarding the packet to the next hop and thereby skipping over the higher layer data or detection of errors). Depending on the type of node and the packet's content, only a subset of these functions may be called. It is possible to think of these function identifiers in data packets as "pointers" to code fragments. In today's systems, the code that implements these functions must be available on the node processing the packet. In our system, the node contacts a "code server" for the necessary code in case the node does not already have the required code locally. In contrast to data servers, which provide a client with "passive" data, code servers provide active plugins stored in a database of code fragments. A code server is a well-known node in the network that provides a library of possibly unrelated functions for different types of operating systems from various developers.

Figure 5 shows an example of a client downloading real-time video through an Active Network Node (ANN) which involves several steps: (1) The ANN receives the connection setup request and forwards it to the video sever; (2) the video server replies with a packet referencing a function for congestion control of the video stream; (3) the ANN does not have the code referenced in its local cache and therefore contacts a code server for the plugin; (4) the ANN receives the active plugin, dynamically links it in its network-



Figure 5: ANN downloading an active plugin

ing subsystem, possibly applies the data to the congestion control function, and forwards the packet to the client. Once the plugin is downloaded, it is stored locally on the ANN removing the need to download the same active plugin in the future. Distributed code caching features the following important properties:

- Active plugins in object code: It is important to note that the active plugins offered by the code server are programmed in a higher level language such as 'C' and compiled into object code for the ANN platform. Once the node loads the functions, they are in no way different than the ones compiled into the network subsystem at build-time. For example, the functions have as much control over the network subsystem's data structures as any other function in the same context, and they are executed as fast as any other code.
- Security addressed by usage of well known cryptography techniques: All active plugins stored on code servers are digitally signed by their developers. Code servers are well known network nodes that authenticate the active plugin when sending them to ANNs. ANNs load only authenticated, digitally signed active plugins and have the capability to check the plugin's sources and developer before installing and running the plugin locally. The security problem is reduced to the implementation of a simple policy rule on the node which lets it choose the right code server and a database of public keys to check the developer's signature and the code server's authentication.
- **Minimization of code download time**: Downloading of active plugins from code servers to ANNs happens infrequently, since this is necessary only on the first occurrence of a new function identifier. Still, some attention has to be paid to the minimization of the download delay. Download time

can be minimized by the following three architectural considerations:

(1) "Probe" packet: By using one or multiple "probe" packet sent from the server to the client before sending data packets. The probe packet triggers the downloading of active plugins in parallel on all routers along the packet's path. The total end-to-end code download delay can be reduced approximately to the time a single ANN requires for the download.

(2) Optimal code server arrangement. Code servers should be as "close" as possible to the ANN. They can be put into a hierarchy similar to DNS servers where the root code servers get their active plugins from the programmers of the plugins.

(3) Minimizing the distance between the ANN and the code server: The ANN can reach the code server through different ways. One or multiple unicast addresses of code servers can be configured on the ANN (again similar to DNS). The responsibility for finding a suitable code server is up to the administrator selecting the unicast address. Another possibility is the usage of anycast or multicast addresses, which would delegate the responsibility for finding the best code server to the anycast or multicast routing. Last, the data server itself could maintain a database of active plugins and the ANN could query the data server for the plugin. This solution has the advantage that no particular configuration information for code servers must be present on the node and there is no need for a particular active plugin distribution infrastructure. It seems very natural that the organization providing a data server makes sure that not only end systems (e.g. by offering a plugin for a web browser) but also all nodes along the data packet's path are able to process the data offered in the best possible way. One disadvantage of this solution is that it allows only one level of authentication (the developer's digital signature). Also, code plugins may come from "non-optimal" sources in respect to bandwidth and delay since all routers along the packet's path might access the same data server instead of possibly utilizing parallel active plugin downloading through a hierarchical infrastructure.

• **Policies**: We support policies for at least two important system properties: Acceptance of specified active plugins and plugin caching behavior.

(1) Acceptance policies: Policies regarding acceptance of active plugins on nodes are desirable. Even if plugin sources and the plugins themselves are authenticated, network administrators may wish to restrict the set of developers they accept active plugins from or exclude certain specific active plugins because of undesired behavior.

(2) Caching policies: Developers are able to set time-outs for active plugins. When a time-out is reached for a plugin on an ANN, the ANN would delete it and refetch it on the next reference in a data packet. This mechanism can be used to deploy prototype versions of new network protocol implementations. Time-outs can be set to infinity for non-expiring plugins. In addition to these developer specified time-outs, the administrator of an ANN can set time-outs for an individual plugin, for sets of plugins or for all plugins in the ANN. These time-outs force a periodic refetch of specified plugins independent of a developer's settings. Such a refetch can be set to happen out-of-band to provide the node with the most recent plugin version independent of references in data packets.

By installing a set of rules on the ANN we will enable both mechanisms that implement these policies to be configured.

• **Integration with existing network protocols**: Our active networking support can be provided in existing protocols by introducing new function identifiers at different layers. We will briefly look into the three possibilities:

(1) Data link layer: Using ATM as a link layer, an application of function identifiers could be to use them in the LLC SNAP field.

(2) Network layer: IP options, which are defined for both IPv4 and IPv6, might be the most common way to introduce new function identifiers ([6] and [29] also describe a way to use IP option fields for AN). Options are commonly used to specify unusual datagram processing, e.g. source

routing. Whereas option usage in IPv4 is very limited because the total option length is 40 bytes, IPv6 introduces a very flexible option concept by allowing very long and unlimited numbers of options wrapped into either Hop-by-Hop or Destination option extension headers. Using IPv6 options has the further advantage that the system can benefit from the option type semantics which specifies the node's behavior in case it does not recognize the option type (skipping over

option/discarding packet/sending ICMP message to source). For function identifiers in IP options in the context of connection-oriented protocols like TCP, active plugin download can take place on connection setup. When the data server replies with a SYN back to the client requesting the connection, it may include a packet containing the "probe" function identifier and optional configuration information and force the nodes along the path to fetch the active plugins. In BSD 4.4, the retransmission delay for the client initiating the SYN is approximately 6 seconds before the next SYN is sent out and the client waits 76 seconds before considering the request as failed. Thus on-demand loading of the corresponding code should be possible without the need for changing the end node's TCP.

(3) Transport layer: For functions to be executed on end systems only, function identifiers can occur in addition to or instead of the usual transport layer function identifier for TCP/UDP.

We described the most important properties of distributed code caching in this section. The usage of caching techniques and active plugins in machine code promises to deliver a significant performance improvement over the traditional, interpretation-based capsules. Next, we describe the implementation of our DAN EE on top of the ANN NodeOS described in section 4.

5.2 The DAN Execution Environment

The DAN EE is shown in Figure 6. It consists of the Active Function Dispatcher (AFD) in the kernel and the DAN Plugin Management (DPMgmt) in user space. The DPMgmt consists of the Active Plugin Loader, the Policy Controller, the Security Gateway, the Plugin Database Controller, and the Plugin Requester. Next follows a description of the individual components.

> Function Active Dispatcher (AFD): the AFD scans a data packet for function identifiers and passes the packet to the corresponding active plugins. Although, as shown in the previous section, it is possible to integrate DAN function identifiers at various places in datagrams, in its current implementation the AFD looks for DAN function identifiers only in IPv6 hop-by-hop options. We embed DAN function identifiers in ANEP [2] packets. ANEP is a packet header defined by the active net-



works working group to precede EE specific packet headers. The IPv6 option processing code passes the packet to the AFD when it finds DAN function identifiers following an ANEP header. While scanning the packet, the AFD calls the functions referenced until they have all been called or

one of them has dropped the packet. The AFD keeps track of all known function identifiers and a pointer to their corresponding instances on a per flow basis. As the instances are called, the AFD chains them together as shown in section 4. Note however, that the AFD calls the instances *individually*. The chaining has no effect on the AFD. It is only considered if the packet contains a selector. In that case, it is passed directly to the first instance of the chain without going through the AFD. In case of a previously unknown function identifier, the AFD contacts the Active Plugin Loader (APL, described next) in order to request the corresponding active plugin. It temporarily suspends packet processing for the packet received. The AFD maintains its own queue of active packets with previously unknown function identifiers. On a call from the APL, the AFD resumes the processing of the enqueued packet by calling the newly installed active plugin.

- Active Plugin Loader (APL): this component interfaces with the networking subsystem in the kernel through a dedicated socket interface similar to the way *routed* does in BSD Unix. On the occurrence of an identifier for a previously unknown active plugin, the AFD requests the corresponding active plugin from the APL. The APL talks first to the Policy Controller to find out whether the request for the plugin is permitted. If the Policy Controller positively acknowledges the request, the APL requests the plugin from the Plugin Database Controller which maintains the database of local active plugins. If the plugin is locally available, it is immediately loaded into the networking subsystem through the Plugin Control Unit. If not, the APL contacts the Plugin Requestor to send out a request to a code server. On reception of the plugin from a code server, the APL passes it to the Security Gateway for origin and signature control. If the active plugin's signature is valid and its origin proven, it gets passed down to the Plugin Control Unit for integration into the networking subsystem. Previously suspended packet processing then resumes. Finally, the plugin is passed to the Plugin Database Controller, which includes it in its local database of active plugins.
- **Policy Controller (PC)**: the PC maintains policy rules set up by the node's administrator. As previously described, we implement both acceptance policies as well as caching policies for active plugins.
- Security Gateway (SG): the SG is responsible for checking the integrity and origin of active plugins. Which security checks are required is determined by the configuration of the ANN in question. The SG maintains a database of public keys. We implement full RSA public-key encryption using the RSAREF [24] library as a basis. This library provides both MD5 one-way hashing as well as RSA public key encryption. MD5 one-way hashing will be used to generate a plugin specific hash key that is then digitally signed with RSA using the developer's private key. The code server transmits the active plugin together with the signed hash to the ANN. On reception of the plugin, the ANN calculates the plugin's MD5 hash, decrypts the received hash with the developers public key and compares both hashes. If they match, the plugin is assumed to be valid. Security extensions [13] to the Domain Name System (DNS) provide support for a general public key distribution service which we use to distribute the developer's and code server's public keys. The stored keys enable ANNs to learn the authenticating keys of code servers in addition to those for which they are initially configured. Keys associated with the code server's and developer site's DNS names can be retrieved to support our system. An ANN can learn the public key of a code server or a developer either by reading it from the DNS or by having it statically configured. To reliably learn the public key by reading it from the DNS, the key itself must be signed with a key the ANN trusts. The ANN must be initially configured with at least the public key of one code server and developer. This is typically the network subsystem developer's key. From there, the node can securely read the public keys of other code servers and developers. As an alternative to the DNS security scheme, we will use IP security [5], which is mandatory for IPv6 implementations, for code server authentication. This allows simple and streamlined security for ANNs which do not check developer signatures but depend on code

server authentication only.

- **Plugin Database Controller (PDC)**: the PDC efficiently administers the local database of active plugins. Plugins are indexed by developer codes and function identifiers for fast access. If the ANN offers code server service to other ANNs, the database may contain active plugins for foreign hardware and software architectures. The active plugins are stored together with the developer's digital signature and the originating code server authentication. Typically, the plugins are stored on a non-volatile storage, like disks or flash RAM. This is not required for regular ANNs since they can refetch active plugins from code servers on system startup, thus saving the download time during packet processing. Plugins come with an expiration date that can be set to "infinite". Administrators are able to set global expiration time for unused plugins independent of plugin specific settings. On expiration of a code plugin, the PDC deletes it from its non-volatile storage and reloads it on request.
- **Plugin Requester (PR)**: the PR is responsible for requesting active plugins from code servers and replying to such requests. The request is either unicasted, multicasted or anycasted depending on the local configuration as described using a lightweight non-connection oriented protocol (e.g. UDP/IP) to ensure stable operation under heavy load. Since we use a datagram oriented protocol instead of a connection-oriented protocol, both loss of the active plugin request or loss of the reply (the active plugin itself) may occur. In this case, the packet causing the request is dropped by the networking subsystem with a possible error message sent to the source. The download of the active plugin is reinitiated the next time the same function reference occurs in a packet.

5.3 Code Server

Code servers feature a database of active plugins for possibly different operating systems and hardware architectures. They get their plugins either manually by configuration through a system administrator or automatically from an upper level code server in the code server hierarchy.

Code servers are network nodes running a version of the DPMgmt. We are carefully separating ANN NodeOS-dependent layers (e.g. the part that communicates with the AFD) from the rest of the DPMgmt to be able to port the DPMgmt to a wide range of platforms. Since most of today's router hardware lacks large mass storage, end systems similar to database servers are better suited to be configured as code servers. We are exploring the usability of publicly available relational and object-oriented database technologies to efficiently store active plugins.

5.4 Plugin Packages

The code for active plugins is stored on code servers and in the local active plugin database on the individual ANNs together with additional data. The code for multiple active functions can be wrapped together into one active plugin package. A download of such a package would install multiple active functions on the ANN. This is useful for strongly correlated active functions such as the four options for IPv6 mobility support [16]. On occurrence of a function identifier in a data packet, not only the referenced active function implementation is downloaded and installed but also one or more others, since they will very likely be referenced in the future. This mechanism requires only one download cycle and one application of related functions such as security checks for the whole package. A package contains at least:

- the code for one or more active functions
- the developer's digital signature
- the code server's authentication information
- configuration information for the ANN which will be passed to the plugin after installation and a set of rules regarding storage of the plugin package, like its expiration date

6 Conclusions and Future Work

We elaborated on three key factors critical to pave the way for active networking in a gigabit environment. First, we elaborated on a gigabit hardware platform that allows high-performance active networking in a scalable fashion combining off-the-shelf and customized hardware components. Second, we described our NodeOS supporting both the active network paradigm as well as the hardware in a highly optimized way. Third, we reviewed a new active networking execution environment called "Distributed Code Caching" which we believe to be especially well suited for our high-performance hardware and software environment.

We are currently in the process of implementing the system described here and expect to have a prototype of the system up and running by the time this paper is published. As a next step, we will start working on a variety of applications.

One of the most promising applications of our environment is automatic protocol deployment. We plan to show automatic upgrading of IPv4 nodes to IPv6 nodes as well as on-the-fly revision of IPv6 implementations. Without active networking, it is extremely hard (if not impossible) to change a protocol once it is deployed. What is needed, is a fully automated way to deploy and revise new protocols. This would allow for incremental refinement of specifications and implementations based on real-world experience, which has not been possible so far. Consider for example IPv6 options: in IPv6, only a very small set of IP options is specified in the base specification [11]. These options are mainly used to pad data packets to certain sizes in order to align them at word boundaries. However, the protocol supports new options in a modular way. An arbitrary number of IPv6 options can follow the IPv6 header in the form of Hop-by-Hop or Destination options. It is expected that these new options are 'hardwired' into an IPv6 implementation. To support new options, such an implementation would require recompiling, which is difficult and time consuming to do in an operational network. With the system described here, a new active plugin for an IPv6 option is downloaded on demand from a code server the first time the new option is referenced and the active plugin is stored for later use in the local cache. We will show this feature for options required for IPv6 mobility support and possibly others. Other applications planned for implementation on the platform include application-specific reliable multicast, congestion-control for real-time audio/video, media gateways and sensor data mixing.

References

- [1] Active Networks Working Group, "Architectural Framework for Active Networks", at *http://www.cc.gatech.edu/projects/canes/arch/arch-0-9.ps*, August 1998
- [2] Alexander, D., et al., "Active Network Encapsulation Protocol (ANEP)", RFC DRAFT, July 1997
- [3] Alexander, D., et al., "The SwitchWare Active Network Architecture", In *IEEE Network Special Issue on Active and Programmable Networks*, May/June 1998
- [4] Alexander, D., Shaw, M., Nettles, S., Smith, J., "Active Bridging", In *Proceedings of SIGCOMM '97*, September 1997
- [5] Atkinson, R., "Security Architecture for the Internet Protocol", RFC 1825, August 1995.
- [6] Bhattacharjee, S., et al., "An Architecture for Active Networking", In *Proceedings of INFOCOM* '97, April 1997
- Bhattacharjee, S., Calvert K., Zegura, E., "Self-Organizing Wide-Area Network Caches", In Proceedings of INFOCOM '98, April 1998
- [8] Chaney, T., et al., "Design of a Gigabit ATM Switch", In Proceedings of INFOCOM'97, April 1997.
- [9] Decasper, D., Dittia, Z., Parulkar, G., Plattner, B., "Router Plugins A Software Architecture for Next Generation Routers", In *Proceedings of SIGCOMM'98*, September 1998
- [10] Decasper, D., Plattner, B., "DAN: Distributed Code Caching for Active Networks", In Proceedings of INFOCOM'98, April 1998
- [11] Deering, S., Hinden, R., "Internet Protocol, Version 6 (IPv6), Specification", RFC 1883, December 1995
- [12] Dittia, Zubin, Jerome R. Cox, Jr., and Guru Parulkar. "Design of the APIC: A High Performance ATM Host-Network Interface Chip", In *Proceedings of INFOCOM*'95, April 95.
- [13] Eastlake, D.E., "Domain Name System Security Extensions", *draft-ietf-dnssec-secext2-02.txt*, November 1997
- [14] Hadzic, I., Smith, J., "On-the-fly Programmable Hardware for Networks", In Proceedings of Globecom 1998
- [15] Hartman, J., et al., "Joust: A Platform for Liquid Software", In IEEE Networks, July 1998
- [16] Johnson, D., Perkins, C., "Mobility Support in IPv6", draft-ietf-mobileip-ipv6-03.txt, July 1997
- [17] Legedza, U., Wetherall D., Guttag, J., "Improving the Performance of Distributed Applications Using Active Networks", *In Proceedings of INFOCOM'98*, April 1998
- [18] Lehman, L., Garland, S.J., and Tennenhouse, D., "Active Reliable Multicast", In *Proceedings of INFOCOM'98*, April 1998
- [19] Mosberger, D., Peterson, D., "Making paths explicit in the Scout operating system", In *Proceedings of the Second Symposium on Operating Systems Design and Implementation*, October 1996
- [20] Nygren, E., "The Design and Implementation of a High Performance Active Network Node", MIT Master's thesis, February 1998
- [21] Papadopoulos, C., Parulkar, G., and Varghese, G., "An Error Control Scheme for Large-Scale Multicast Applications", In *Proceedings of INFOCOM*'98, April 1998
- [22] Parulkar, G.M., Schmidt, D.C., Turner, J.S., "a^It^Pm : A Strategy for Integrating IP with ATM," In *Proceedings of SIGCOMM*'95, August 1995.
- [23] Schwartz, B., et al., "Smart Packets for Active Networks", available from *http://www.net-tech.bbn.com/smtpkts/smart.ps.gz*, January 1998
- [24] RSA Library, *ftp://ftp.rsa.com/rsaref*
- [25] Shreedar, M., Varghese, G., "Efficient Fair Queueing using Deficit Round Robin", In *Proceedings of SIGCOMM '95*, August 1995
- [26] Stoica, I., Zhang, H., and Ng, T.S.E., "A Hiearchical Fair Service Curve Algorithm for Link-Sharing, Reail-Time and Priority Services", In *Proceedings of SIGCOMM*'97, September 1997.
- [27] Tennenhouse, D., et al. "A Survey of Active Network Research", *IEEE Communications*, January 1997

- [28] Tschudin, C., Decasper, D., "Simple Active Packet Format (SAPF)", Experimental RFC, August 1998
- [29] Wetherall, D., Tennenhouse, D., "The ACTIVE IP Options", In *Proceedings of the 7th ACM SIGOPS European Workshop*, September 1996
- [30] Wetherall, D., et al, "ANTS: A Toolkit for Building and Dynamically Deploying Network Protocols", In *Proceedings of IEEE OPENARCH'98, April 1998*
- [31] Yemini, Y., daSilva, S., "Towards programmable networks", In *IFIP/IEEE International Workshop* on *Distributed Systems: Operation and Management*, October 1996
- [32] Zegura, E., "CANEs: Composable Active Network Elements", Georgia Institute of Technology, http://www.cc.gatech.edu/projects/canes/