

---

# Public Review for Flow Processing and the Rise of Commodity Network Hardware

Adam Greenhalgh, Mark Handley, Mickael Hoerd, Felipe Huici,  
Laurent Mathy, and Panagiotis Papadimitriou

Network functionalities such as intrusion detection and load balancing are often implemented in specialized expensive middleboxes plugged inside the network. But, with the advent of commodity hardware and network switches, it is time to think about leveraging these new and cheap resources to support the same functionalities with lower cost without compromising efficiency. This is in the same spirit that software radio, virtual machines and virtual routers, have been introduced. The implementation of network functionalities in a kind of software environment has the further advantage of making them easily manageable and extendable to other applications (on software timescales).

The architecture introduced in this paper is called Flowstream. It proposes the implementation of network functionalities in virtualized machines/servers/routers run on top of commodity PCs. The flow of traffic among these virtual network entities is controlled by a programmable network switch implementing Openflow. The paper motivates the problem and discusses the architecture and its main components, plus a description of some potential applications. Even though there are no validation results, all reviewers appreciate the idea and agree on the fact that it will trigger discussions among CCR readers and the members of the networking community. This is a new research area that involves several tradeoffs (technical vs. economical, reliability vs. programmability) to be clearly understood and evaluated.

Programmable flow forwarding using Openflow has been already proposed in an operating system context as for example in the NOX architecture that has appeared as an editorial note in the CCR July 2008 issue. The novelty of this new paper is in combining flow forwarding and virtualization to replace network middlebox functionalities.

*Public review written by*  
**Chadi Barakat**  
*Planète Research Group*  
*INRIA Sophia Antipolis*



# Flow Processing and the Rise of Commodity Network Hardware

Adam Greenhalgh  
University College London, UK  
a.greenhalgh@cs.ucl.ac.uk

Felipe Huici  
NEC Europe Ltd, Germany  
felipe.huici@nw.neclab.eu

Mickael Hoerd  
Lancaster University, UK  
m.hoerd@lancaster.ac.uk

Panagiotis Papadimitriou  
Lancaster University, UK  
p.papadimitriou@lancaster.ac.uk

Mark Handley  
University College London, UK  
m.handley@cs.ucl.ac.uk

Laurent Mathy  
Lancaster University, UK  
l.mathy@lancaster.ac.uk

## ABSTRACT

The Internet has seen a proliferation of specialized middlebox devices that carry out crucial network functionality such as load balancing, packet inspection and intrusion detection. Recent advances in CPU power, memory, buses and network connectivity have turned commodity PC hardware into a powerful network platform. Furthermore, commodity switch technologies have recently emerged offering the possibility to control the switching of flows in a fine-grained manner. Exploiting these new technologies, we present a new class of network architectures which enables flow processing and forwarding at unprecedented flexibility and low cost.

## Categories and Subject Descriptors

C.2.1 [Computer Communication Networks]: [Network Architecture and Design]

## General Terms

Design

## Keywords

Architecture, Flow processing, Virtualization, Internet

## 1. INTRODUCTION

In the last few years two trends have started to reshape the Internet. The first of these is steady encroachment of economic reality on the architecture of the network itself; primarily this takes the form of embedding higher level knowledge *inside* the network to enhance the ability to manage services, make better use of limited resources and control costs. Usually this means using middleboxes such as firewalls[11], traffic shapers[15], load balancers[14], intrusion detection (IDS) and prevention systems (IPS)[6], and application enhancement boxes[17, 18, 19]. It has become rare to find an end-to-end path that does not encounter at least one such device.

Middleboxes have become a multi-billion dollar market, but network operators do not spend all this money without good reason: such technologies have become essential to providing high levels of service for key applications. The great merit of the original Internet architecture was its ability to

support as-yet-unforeseen applications, but the downside is that the network does not know when the applications it supports are actually working. To prosper, enterprises need additional control, and middleboxes provide this.

The second trend lies in the commoditization of hardware. Over many years components and systems designed primarily for the mass market have achieved such large volumes of sales that their capabilities increased to the point of displacing high-end products. The canonical example is the rise of the Intel x86 CPU architecture, first displacing high-end Unix workstations running on RISC processors, and now breaking into the very top of the supercomputer league tables. In the data center, the commoditization of the 1U form-factor x86 server combined with drastically reduced CPU costs has greatly narrowed the price gap between server and desktop systems: a rack-mount case still costs more than a desktop case, but very capable servers can be bought for \$1500 complete with multiple on-board Gigabit Ethernet network interface cards.

It is, however, not only computers that have become commodity items. A few years ago, many people were forecasting that network processors would displace custom silicon in high-end router platforms. However, something different happened. The combination of a rise in very capable and cheap chipsets for Gigabit Ethernet from the likes of Broadcom and Marvell, huge volume shipments from companies such as Dell and Netgear, and bulk manufacturing from manufacturers such as Quanta (who also make many of the world's laptops) has caused switching to become commoditized. As with x86 processors, the low end has started to increase in capability and displace high-end specialist products. Today's 48-port Gigabit switches support both layer 2 and layer 3 forwarding, access control lists (ACLs) and other features at a price of around \$20 per port; commodity 10 Gigabit switches are now starting to emerge.

Researchers and middlebox manufacturers are both well-aware of the capabilities of x86 commodity hardware, but the commoditization of switch hardware and the potential to rewrite their control software has not received quite the same attention. While switches have become more powerful, they are still relatively inflexible devices: as a platform, they are rather limited in capability. Things only become interesting when you combine switches with servers.

Consider now the confluence of these two trends. There

is a huge proliferation of middleboxes, each servicing a single role performing L4-L7 functionality on data flows. At the same time, we now have cheap and extremely capable switching and processing components. However, the switches are too dumb and the servers have their limitations (despite their pretty good performance, there is only so much you can do with one box before memory bottlenecks start to kick in[5]). The clear solution is to build a generic network control, forwarding and flow processing platform from commodity switch hardware unified with a small cluster of servers, all managed as a single platform. Such a platform is inexpensive, flexible, scalable, and tolerant of failures<sup>1</sup>.

Perhaps more importantly, the rise of such platforms would open up the possibility of a commodity market for high-performance middlebox software, where a network operator might be able to mix and match control and management software in a way which is currently difficult at best.

The biggest downside of middleboxes is that they embed into the network knowledge of today's applications at the expense of tomorrow's innovations. It might seem like we are attempting to encourage this process, but the reality is that it has already happened. Once a middlebox is deployed, the cost of changing is substantial. We hope that by encouraging a common platform for such capabilities, and by making this market one for software rather than for appliances, the additional flexibility and reduced time to deployment might remove some of the barriers faced by innovative applications of the future.

The design of such a platform is the object of the rest of this paper, which is organized as follows. Section 2 describes the building blocks or technologies we base the platform on in greater detail. Section 3 provides an overview of *Flowstream*, our proposed flow processing platform, including usage scenarios and applications. In section 4, we discuss the consolidation of the platform and section 5 covers related work. Finally, section 6 concludes.

## 2. BUILDING BLOCKS

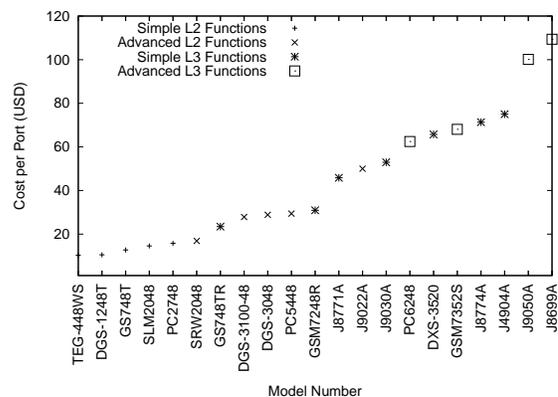
So far we have identified two trends, the middleboxes in the network and the commoditization of servers and switches. We also advocate that the solution is to build a generic network control, forwarding and flow processing platform from these commoditized elements. In this section, we describe the building blocks of such a platform in greater detail before discussing how these might be put together in the next section.

### 2.1 Commodity Switches

How cheap have network switches become? In order to answer this question, we conducted a survey (25/11/2008) of a range of lower-end gigabit switches, calculating for each the retail price per port. While this survey is by no means thorough, it gives a good idea of what the costs are when purchasing a lower-end network switch.

Figure 1 shows the results of the survey. We classified the switches into four groups: those with the simplest layer-2 forwarding (simple L2), those with layer-2 forwarding and advanced features such as ACLs (advanced L2), those with

<sup>1</sup>It is worth noting that while cheap embedded devices and processors exist, they have limited resources and flexibility and therefore they are not a good building block for a general flow processing platform.



vanced L3 forwarding should be powerful enough to comply with the Openflow specification, giving an indication that the cost of these switches should also be reasonable.

## 2.2 Commodity PCs

Even though commodity PCs have been used for a while to process network traffic, it is only in recent years that improvements in various technologies have allowed them to become a powerful network platform. The introduction of PCI Express, for example, removed the bottleneck presented by its predecessor, PCI-X[10]. Further, the availability of an increasing number of CPU cores allows a PC to run several network processes concurrently while providing high performance to each of them (as long as memory hierarchy issues are carefully considered). Ethernet port density has also increased: quad-port cards are now commonplace, which combined with motherboard interfaces allow a server to have as many as 15 or more ports.

The combination of these technologies, along with the drop in prices, has rendered the PC a viable platform for network processing. But exactly how powerful can a commodity PC be? In previous work [5] we used a relatively inexpensive Dell 2950 server with 8 processor cores, 12 Ethernet ports, and standards-compliant IP forwarding paths implemented with the Click Modular Router [9] software. With this setup we were able to forward IP packets of most sizes at line rate, and minimum-sized packets at a very reasonable 4.9 million packets per second.

## 2.3 Virtualization and Virtual Routers

Virtualization techniques enable a PC to run multiple OSes concurrently, giving them access to the underlying hardware while isolating them from each other. In addition, virtualization makes it relatively easy to migrate these OSes to another PC, a mechanism that we will exploit later. Related work in [21] shows how to prevent network traffic disruption during the migration of virtual routers.

In [4] we tackled basic fairness issues and limitations of a modern PC for software packet forwarding, exploring alternative virtualization technologies and different forwarding scenarios. From these findings we designed a *virtual router* that has highly configurable forwarding planes for advanced programmability, optimized core scheduling for high performance, and hardware multi-queueing for sharing interfaces among virtual routers. In [5] we analyzed the virtual router's performance, showing that virtual router solutions based on current commodity hardware represent a powerful, flexible, practical and inexpensive proposition.

## 3. FLOWSTREAM ARCHITECTURES

Given the building blocks described above, we can now discuss in more detail a class of system architectures for building in-network processing platforms which represent a "sweet spot" between performance, scalability and flexibility. We call such platforms "Flowstream Architectures", for reasons that should be clear shortly. Platforms built according to the Flowstream architecture can be characterized by the following properties:

- The core of the platform consists of an Ethernet switch configured to route flows. A flow is defined in the Openflow sense, as packets that match a tuple of source and destination addresses and ports.

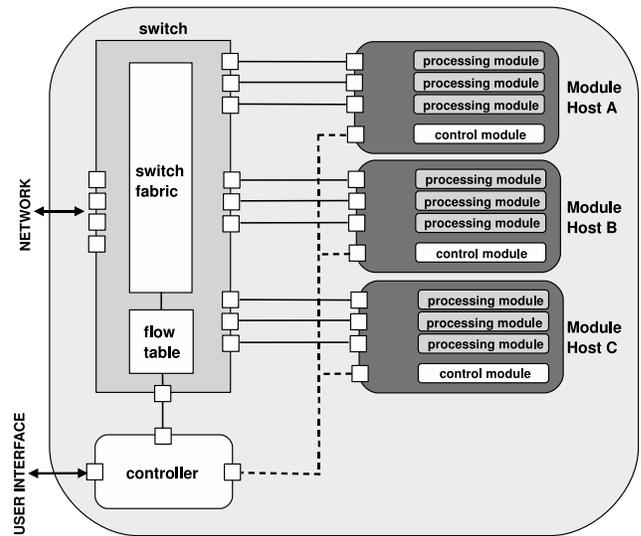


Figure 2: Overview of a Flowstream platform.

- Streams of data from these flows are then routed to one of a number of attached commodity server boxes for additional processing, before being forwarded on to the final destination.
- Software running on the server boxes can be composed to provide processing pipelines of modules.
- These modules are virtualized, in the sense that they can be moved between the servers to balance load and provide robust service in the presence of failures.
- The switch and servers are managed as a single platform from the operators' point of view by a controller.

### 3.1 Description of a Platform

Figure 2 illustrates how the server boxes (we call them *module hosts* to distinguish them from traditional servers) and flow-based switch are connected together with a controller host to comprise a Flowstream platform. Each host runs a number of *processing modules* where all of the actual flow processing takes place (except for basic forwarding, which can be done by the switch). Further, hosts contain a special module called a *control module*, which receives commands from the platform's controller to remove, install or migrate modules, as well as to provide monitoring information about the host's current load and performance.

There are three main technologies available to us for implementing a module:

- A virtual machine running its own OS and module application.
- A process running on a virtual machine shared with other modules.
- A set of kernel forwarding elements instantiated in the kernel of the device driver domain on one of the module hosts.

The first of these options is the most general and provides the best inter-module isolation, whereas the third will provide the highest performance for traffic that needs to traverse several modules in the same module host. We envis-

age different applications will use different implementation options, often on the same Flowstream platform.

For composing kernel forwarding elements, the Click modular router [9] provides a suitable set of building blocks. For example, a module can be composed of a predefined set of Click elements, and under the control of the operator, cascades of such modules can be plumbed together at run-time.

A Flowstream platform's second main component is the Openflow switch, providing the basic connectivity between module hosts and the network. In addition to this, the switch contains a flow table that is configured by the controller at runtime, allowing different flows to be directed to any of the ports on the switch. It is worth pointing out that while figure 2 shows a single switch, it would be certainly possible to scale the platform's port density by including additional switches.

The final component is the controller. Essentially, this is the brains of the platform and also its user interface to the outside world. When the operator makes a request (for instance, running an IDS on flows to a particular web server), the controller begins by choosing the module host or hosts to install the processing module(s) on. Such a decision could be based on the hosts' current load, information that the controller retrieves periodically from the control modules. Having selected a host, the controller then instructs the control module to install the requested processing module. Once this is done, the controller configures the switch's flow table so that the corresponding flows are directed to the right processing module.

With all of these components in place, a Flowstream architecture provides a powerful platform for flow processing. The fact that it is built upon commodity yet, as shown in previous work, high performance hardware should result in significant cost savings. In addition, a Flowstream setup can be easily expanded and contracted dynamically by adding or removing module hosts, something that cannot be easily accomplished on conventional routers or middleboxes. Further, when required, the isolation provided by virtualized module hosts means that several different flow processing operations can be performed simultaneously while minimizing negative interactions. The controller can migrate modules as required to ensure that a processing task does not significantly degrade the performance of others. Last but not least, using general-purpose processors and allowing operators to install their own flow processing modules yields great flexibility. So long as modules have access to well-defined flow APIs, a Flowstream platform can accommodate a wide range of existing and even future network applications. It is precisely the usage of the platform and its potential applications that we discuss next.

### 3.2 Usage Scenarios

In the most basic case, the operator submits a request to a Flowstream platform's controller asking it to apply a certain processing module to a subset of the traffic being forwarded. The controller then chooses a module host with appropriate load levels and installs the module on it, then configures the switch's flow table. The flow then travels from the switch to the module for processing before being sent back to the switch and subsequently out onto the network<sup>2</sup>.

<sup>2</sup>Note that while so far we have described modules as receiving flows, processing them and then forwarding them, it is certainly possible for a module to act as a traffic sink

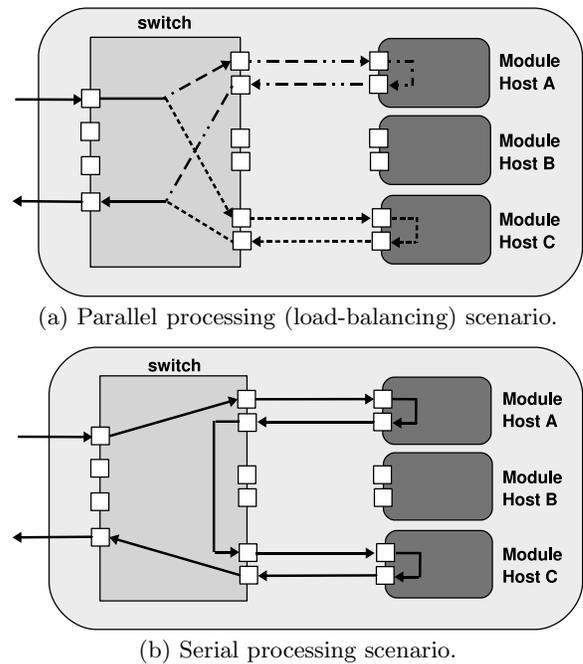
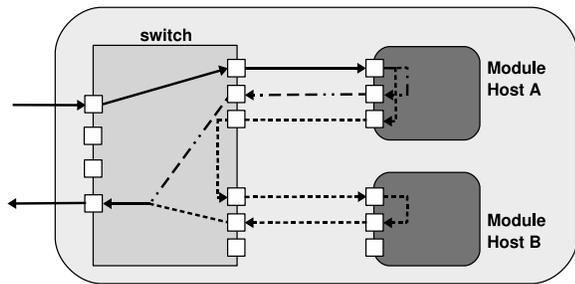


Figure 3: Basic platform usage scenarios.

Beyond the simple case, there are two more interesting usage scenarios, depending on whether modules act on flows in parallel or serially. In parallel processing (see figure 3(a)), flows are load-balanced, pushing different flows to different module hosts but processing each of them equally. In this case identical processing modules run on multiple module hosts (in the figure, hosts A and C). The controller sets up the switch's flow table so that a flow gets sent to either of the module hosts, thus load balancing the traffic; an algorithm such as hash-based Equal Cost Multi-path (ECMP), which is supported by many switches, could be used to accomplish this. To avoid reordering, all the packets from one flow must be processed by one module host. Flows could also be distributed unevenly based on the capabilities of the module hosts or their current load. Parallel modules are useful for quite a number of CPU-intensive network processing tasks, including intrusion detection, spam over Internet telephony and Denial-of-Service attack filtering, and monitoring and deep packet inspection.

In serial processing or pipelining (see figure 3(b)), the operations performed on flows are split across several module hosts and done one at a time. One example of an application for this is VPN termination, where one host could be used to perform the expensive encryption operation before another takes care of the tunneling. Serial processing is essential when each packet must be processed first by one module, then by another. Serial processing would also be useful if one of the hosts had dedicated hardware to perform an expensive operation at line rate, or if a module host did not have enough interfaces to carry out a particular function, such as acting as a router.

Combinations of serial and parallel are certainly possible, as well as heterogeneous parallel processing. In this case, different flows are processed on different modules hosts, but they are also processed by different modules. For example, traffic to a web server farm might be processed by server load



**Figure 4: Scenario: offloading to a separate module host for further processing.**

balancing modules, whereas traffic to a mail server might traverse a blacklist filter.

A more complex usage scenario is *flow splitting*, whereby a processing module is used to split a subset of traffic from a flow aggregate to another module for further processing (see figure 4). An application that fits rather well with this mechanism is intrusion detection: for example, module host A could be used to apply a quick, preliminary filter in order to separate out suspicious flows. Matching flows would then be sent to module host B for a more in-depth inspection, whereas those that do not match are sent back to the switch for immediate forwarding. It is worth pointing out that for simple filters, the actual splitting of flows could be done by the Openflow switch, thus off-loading some of the work away from the module hosts.

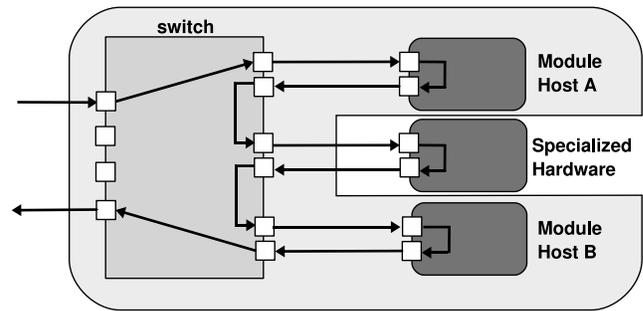
Finally, because a Flowstream platform consists of loosely-coupled hardware, it is possible to use it to encapsulate existing network middleboxes, such as a Bivio 7000 series [13] network appliance. For example, in figure 5, if a third-party, specialized hardware box cannot cope with the full traffic rate or needs some traffic excluded from it, it can be plugged into the flowstream switch and treated as a black box through which a subset of the traffic can be directed. Alternatively, it would only require a small amount of development work to port many commercial Linux-based middleboxes to a Flowstream architecture.

### 3.3 Module Migration

Flowstream architectures fit firmly into the trend of using arrays of cheap and potentially unreliable hardware, but providing robustness in software. To provide such robustness, we need to be able to migrate modules between hosts, both to manage changing load and to adapt to failures. All three of the mechanisms described for implementing modules can be migrated live between hosts:

- Today’s OS virtualization platforms can support live VM migration.
- Cluster computing platforms support live process migration.
- Click kernel forwarding paths can be reconfigured on the fly to include new elements.

This ability to live-migrate processing functions between hardware while simultaneously re-plumbing the switch’s flow table to match provides a powerful and flexible mechanism that can be for many purposes, including load-balancing and reducing costs by powering down under-utilized module hosts during quiet hours.



**Figure 5: Integration of specialized hardware in the platform.**

## 4. BENEFITS OF CONSOLIDATION

A Flowstream platform consolidates a number of middlebox systems into a single entity. For this to be worthwhile, we need to gain tangible benefits from the consolidation, benefits that make the whole greater than its parts; otherwise we are just shifting functions from one system to another. This consolidation has the following benefits:

- Increased tolerance of failures.
- Reduced equipment and maintenance costs.
- The ability to do dynamic reprovisioning.

**Increased tolerance of failures** can be achieved by allowing spare module hosts to take the place of a module host that has failed or is perceived to be about to fail. In a standard middlebox deployment each type of middlebox requires a spare system to be available in case of hardware failure. In the Flowstream architecture a smaller number of spare systems are required because module hosts are agnostic to the processing modules being run on them. Provided that the hardware profile of the failed or failing module host is exceeded by the remaining spare capacity, we can distribute the processing modules from the failing module host onto the spares. Module migration makes it possible to redistribute running modules, or as a last resort, to restart a module on a new host.

**Reduced equipment and maintenance costs** are achieved by separating the logical and physical systems and adopting Google’s model of using mass market commodity boxes. Expensive downtime is eliminated by migrating the processing modules to spare systems and then undertaking maintenance on an offline system.

**Dynamic reprovisioning** is an outcome of the load balancing scenario presented in section 3.2. The Flowstream architecture enables us, at fairly fine granularity, to increase or decrease the capability of any processing module by varying the allocation of flows and processing resources, with the possibility of shutting down whole systems during quiet periods and bringing them back up when load increases. Further, new features can easily be trialled by splitting or copying a small portion of the traffic to a new processing module without interrupting the live system.

## 5. RELATED WORK

NOX [7] relies on Openflow switches to provide a centralized programmatic interface in order to ease the management of enterprise networks. Flowstream is designed to

support the functionality of current and emerging middle-boxes by the consolidation of customized processing modules into an in-network flow-processing platform. Furthermore, virtualization allows for dynamic module migration reducing equipment and maintenance costs.

Pswitches [8] proposes the use of advanced commodity switches to control the paths of flows in data-center networks, and share expensive middleboxes. Ethane [2] introduces the use of flow-based switching in order to control and improve the security of enterprise networks. An Ethane switch is essentially an early Openflow-like technology. Our architecture takes the flexibility afforded by commodity switches further by building complex network processing functionality within PC clusters.

In [16] the authors propose combining a programmable controller with switches for traffic management purposes. Flowstream advocates the combination of commodity server and switching hardware to implement complex router applications beyond traffic management.

Other work [1] explores the scalability of software routers on general-purpose hardware, concentrating on issues such as the how much processing can be done per packet while maintaining line rate. The authors eventually propose a clustered software-router architecture that uses an interconnect of multiple servers in order to enhance scalability. While Flowstream can certainly be used as a router, it provides a general platform where more advanced, higher-layer processing can be done.

Finally, SuperCharging PlanetLab [20] decouples network nodes into a control/application plane running on commodity hardware and specialized network-processor hardware for packet forwarding. In contrast, Flowstream relies mostly on commodity hardware.

## 6. CONCLUSIONS

In this paper we introduced *Flowstream*, a new class of system architectures for building network flow processing platforms. These architectures are now possible thanks to the commoditization of x86 servers, switches and the availability of powerful open virtualization solutions.

We have outlined what a Flowstream platform looks like and discussed its benefits, including flexibility, scalability, fault tolerance and even the possibility of reducing energy costs by switching underused servers off. In addition, we covered some basic usage scenarios and showed how some of today's network applications would be run on such a platform. More importantly, we believe that the platform should be flexible enough to accommodate innovative applications as well.

We are currently undertaking an implementation and evaluation of a Flowstream-based system.

## 7. REFERENCES

- [1] Katerina Argyraki, Salman Abdul Baset, Byung-Gon Chun, Kevin Fall, Gianluca Iannaccone, Allan Knies, Eddie Kohler, Maziar Manesh, Sergiu Nedveschi, and Sylvia Ratnasamy. Can software routers scale? In *Proceedings of PRESTO'08*, Seattle, USA, August 2008.
- [2] Martin Casado, Michael Freedman, Justin Pettit, Nick McKeown, and Scott Shenker. Ethane: Taking control of the enterprise. In *Proceedings of SIGCOMM'07*, Kyoto, Japan, August 2007.
- [3] Open Flow Switch Consortium. Open flow switch. <http://www.openflowswitch.org>.
- [4] Norbert Egi, Adam Greenhalgh, Mark Handley, Mickael Hoerd, Felipe Huici, and Laurent Mathy. Fairness issues in software virtual routers. In *Proceedings of PRESTO'08*, Seattle, USA, August 2008.
- [5] Norbert Egi, Adam Greenhalgh, Mark Handley, Mickael Hoerd, Felipe Huici, and Laurent Mathy. Towards high performance virtual routers on commodity hardware. In *Proceedings of ACM CoNEXT 2008*, Madrid, Spain, December 2008.
- [6] Endace. Endace ninjabox network monitoring. <http://www.endace.com/ninjabox.html>.
- [7] Natasha Gude, Teemu Koponen, Justin Pettit, Ben Pfaff, Martin Casado, Nick McKeown, and Scott Shenker. Nox: Towards an operating system for networks. *ACM SIGCOMM Computer Communication Review*, 38(3):105–110, July 2008.
- [8] Dilip Joseph, Arsalan Tavakoli, and Ion Stoica. A policy-aware switching layer for data centers. In *Proceedings of ACM SIGCOMM 2008*, Seattle, USA, August 2008.
- [9] Eddie Kohler, Robert Morris, Benjie Chen, John Jahnotti, and M. Frans Kashiok. The click modular router. *ACM Transaction on Computer Systems*, 18(3):263–297, 2000.
- [10] Jiuxing Liu, A. Mamidala, V. Vishnu, and D.K. Panda. Evaluating infiniband performance with pci express. *Micro, IEEE*, 25(1):20–29, Jan.-Feb. 2005.
- [11] Check Point Software Technologies Ltd. Check point. <http://www.checkpoint.com/>.
- [12] Nick McKeown. Enterprise GENI Talk, October 2008.
- [13] Bivio Networks. Bivio 7000 series. <http://www.bivio.net/products/bivio7000.htm>.
- [14] F5 Networks. Big-ip product family. <http://www.f5.com/products/big-ip/>.
- [15] Packeteer. Packeteer products. <http://www.packeteer.com/products/>.
- [16] Hideyuki Shimonishi, Takashi Yoshikawa, and Atsushi Iwata. Off-the-path flow handling mechanism for high-speed and programmable traffic management. In *Proceedings of PRESTO'08*, Seattle, USA, August 2008.
- [17] Citrix Systems. Citrix NetScaler.
- [18] Citrix Systems. Citrix WANScaler.
- [19] Riverbed Technology. Riverbed. <http://www.riverbed.com/>.
- [20] Jon Turner, Patrick Crowley, John DeHart, Amy Freestone, Brandon Heller, Fred Kuhns, Sailesh Kumar, John Lockwood, Jing Lu, Michael Wilson, Charles Wiseman, and David Zar. Supercharging planetlab - a high performance, multi-application, overlay network platform. In *Proceedings of SIGCOMM'07*, Kyoto, Japan, August 2007.
- [21] Yi Wang, Eric Keller, Brian Biskeborn, Jacobus van der Merwe, and Jennifer Rexford. Virtual routers on the move: Live router migration as a network-management primitive. In *Proceedings of SIGCOMM'08*, Seattle, USA, August 2008.