

# Efficient Implementation of Dynamic Protocol Stacks in Linux

Ariane Keller  
ETH Zurich, Switzerland  
ariane.keller@tik.ee.ethz.ch

Daniel Borkmann  
ETH Zurich, Switzerland  
HTWK Leipzig, Germany  
borkmann@iogearbox.net

Wolfgang Mühlbauer  
ETH Zurich, Switzerland  
muehlbauer@tik.ee.ethz.ch

## ABSTRACT

TODO: rewrite abstract - beginning is copied from an old paper... Future network architectures aim at solving the shortcomings of the traditional, static Internet architecture. In order to provide optimal service they have to adapt their functionality to different networking situations. This can be achieved by dividing the networking functionality into modular blocks and combining them as required at runtime. In this paper we address the performance aspect of such architectures and we show that their performance is comparable with the performance of a standard Linux protocol stack.

## 1. INTRODUCTION

Some references that might be useful: [2] (ANA) and [6] (Click) and [3] (From protocol stack to protocol heap: role-based architecture) and [4] (PLUTARCH: an argument for network pluralism) and [1] (netgraph) and [7] (survey of next generation internet) and [5] (xKernel) and [8] (model for flexible high-performance communication subsystem). - should we explicitly say something on active networking or should we try to avoid it completely?

Placeholder: Lorem ipsum dolor sit amet, consectetur adipiscing elit. Pellentesque eu arcu ut est volutpat consequat sit amet dignissim enim. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Nunc magna purus, vehicula sit amet fringilla ac, interdum ut dui. Ut in magna tortor, vitae dignissim lorem. Praesent condimentum eros aliquam mi pharetra egestas. Sed sem tortor, iaculis non ornare consequat, auctor eget velit. Nam nibh nibh, ullamcorper vitae gravida non, rhoncus vitae mi. Nunc vestibulum suscipit justo in laoreet. Praesent ac porta ante. Integer sem urna, pretium sed dignissim id, laoreet sit amet sem. Cras ac risus nec nibh tempor gravida. Integer ac ligula sed orci luctus condimentum at quis dui. Etiam dignissim dignissim tellus, et dapibus elit venenatis nec. Sed hendrerit imperdiet lacinia. Sed enim purus, mattis vel ullamcorper vel, malesuada vitae turpis. Pellentesque nec lacus tortor, eget scelerisque lectus. Suspendisse lectus mauris, tempor eget

porttitor id, consectetur vel sem.

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Pellentesque eu arcu ut est volutpat consequat sit amet dignissim enim. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Nunc magna purus, vehicula sit amet fringilla ac, interdum ut dui. Ut in magna tortor, vitae dignissim lorem. Praesent condimentum eros aliquam mi pharetra egestas. Sed sem tortor, iaculis non ornare consequat, auctor eget velit. Nam nibh nibh, ullamcorper vitae gravida non, rhoncus vitae mi. Nunc vestibulum suscipit justo in laoreet. Praesent ac porta ante. Integer sem urna, pretium sed dignissim id, laoreet sit amet sem. Cras ac risus nec nibh tempor gravida. Integer ac ligula sed orci luctus condimentum at quis dui. Etiam dignissim dignissim tellus, et dapibus elit venenatis nec. Sed hendrerit imperdiet lacinia. Sed enim purus, mattis vel ullamcorper vel, malesuada vitae turpis. Pellentesque nec lacus tortor, eget scelerisque lectus. Suspendisse lectus mauris, tempor eget porttitor id, consectetur vel sem. Praesent ac porta ante. Integer sem urna, pretium sed dignissim id, laoreet sit amet sem. Cras ac risus nec nibh tempor gravida. Integer ac ligula sed orci luctus condimentum at quis dui. Etiam dignissim dignissim tellus, et dapibus elit venenatis nec. Sed hendrerit imperdiet lacinia. Sed enim purus, mattis vel ullamcorper vel, malesuada vitae turpis. Pellentesque nec lacus tortor, eget scelerisque lectus. Suspendisse lectus mauris, tempor eget porttitor id, consectetur vel sem.

## 2. THE BODY OF THE PAPER

Description of the base architecture see Figure 1.

- Environment: kernel space - kernel XXX
- Application Interface: socket API
- Hardware Interface: Vlink subsystem
- Configuration: what can be configured?, Command line interface
- Internals: state Management, Functional block notifier

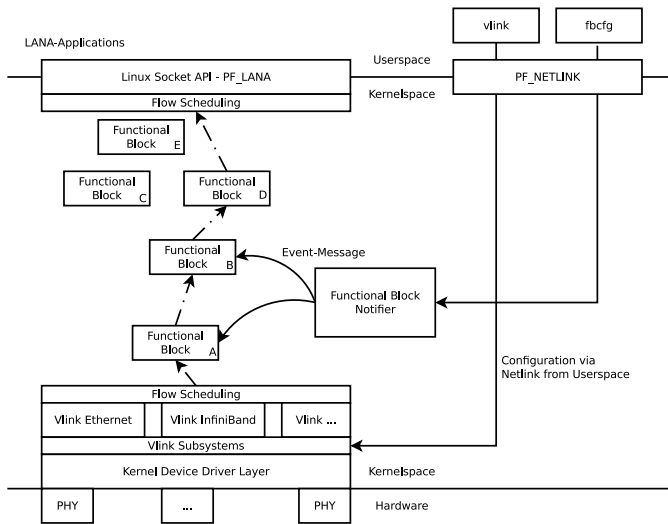
### 2.1 Improving the Performance

During the implementation of our framework we have evaluated different possibilities for the integration of our packet processing engine with the Linux kernel. We think the insights gained are interesting for other researcher that have to do fundamental changes on the Linux protocol stack and hence, we summarize them here. Our goal was to be able to process as many *minimum sized Ethernet frames* as the Linux kernel is able to process. In order to compare the performance of the Linux Kernel and the performance of our engine we have dropped all packets in the Linux Kernel

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

ANCS 2011 Brooklyn, New York, USA

Copyright 20XX ACM X-XXXXX-XX-X/XX/XX ...\$10.00.



**Figure 1: Architecture** - in order to save space we should try to come up with a picture that only uses 1 column.

protocol stack as soon as they were arrived (TODO: where exactly?). In our system the packets were processed by the fb\_eth functional block followed by two fb\_dummy functional blocks that were simply forwarding the packets. We can distinguish the following three approaches:

- On each CPU there exists one high priority thread that is responsible for processing LANA packets. This approach leads to a starvation of the interrupt handler (ksoftirqd) and hence the maximal achieved packet rate is only about half as what is achieved by the protocol stack of the Linux kernel. Also changing the priority of the LANA thread to normal only slightly increases the throughput.
- Instead of relying completely on the Scheduler of the Linux Kernel we control preemption and scheduling explicitly. This approach still exhibits scheduling overhead, but it increases the performance to about two thirds of the performance of the Linux Kernel.
- Instead of executing the LANA functions in a dedicated thread they are executed directly in the ksoftirqd function. With this approach approximately 95% of the performance of the Linux kernel is achieved.

The corresponding numbers are listed in Table 1.

mechanism	performance
dedicated kernel thread (high priority)	700000
dedicated kernel thread (normal priority)	750000
dedicated kernel thread (controlled scheduling)	900000
execution in ksoftirqd	130000
Linux kernel stack	138000

**Table 1: Performance evaluation (pps) of different approaches to receiving packets in the Linux kernel.** The packets are 64 Bytes long. The evaluation was done on a TODO: CPU/RAM/NETWORKCARD/KERNEL

## 2.2 Current State

- Base machinery working
- Implemented Functional blocks: fb\_eth, fb\_ethvlink, fb\_pflana, fb\_bpf, fb\_dummy

## 3. CONCLUSIONS AND FUTURE WORK

- network system implemented that is as performant as the linux stack but allows for flexibility (inclusion of new protocols and change at runtime)
- Automatic protocol stack setup
- Protocol Stack negotiation between different nodes
- Performance evaluation with specific scenarios in comparison to other approaches such as click or linux.

## 4. ACKNOWLEDGMENTS

The research leading to these results has received funding from the European Union Seventh Framework Programme under grant agreement n°257906.

## 5. REFERENCES

- [1] All About Netgraph.  
<http://people.freebsd.org/~Lijjulan/netgraph.html> (Aug 10).
- [2] G. Bouabene, C. Jelger, C. Tschudin, S. Schmid, A. Keller, and M. May. The autonomic network architecture (ANA). *Selected Areas in Communications, IEEE Journal on*, 28(1):4–14, Jan. 2010.
- [3] R. Braden, T. Faber, and M. Handley. From protocol stack to protocol heap: role-based architecture. *SIGCOMM Comput. Commun. Rev.*, 33(1):17–22, 2003.
- [4] J. Crowcroft, S. Hand, R. Mortier, T. Roscoe, and A. Warfield. Plutarch: an argument for network pluralism. In *Proceedings of ACM SIGCOMM FDNA Workshop*, August 2003. Karlsruhe, Germany.
- [5] N. C. Hutchinson and L. L. Peterson. The X-Kernel: An architecture for implementing network protocols. *IEEE Trans. Softw. Eng.*, 17(1):64–76, 1991.
- [6] E. Kohler, R. Morris, B. Chen, J. Jannotti, and M. Kaashoek. The click modular router. *ACM Trans. Comput. Syst.*, 18(3):263–297, 2000.
- [7] Subharthi Paul, Jianli Pan and Raj Jain. Architectures for the Future Networks and the Next Generation Internet: A Survey, 2009.  
<http://www.cse.wustl.edu/~jain/papers/i3survey.htm> (Oct 09).
- [8] M. Zitterbart, B. Stiller, and A. N. Tantawy. A model for flexible high-performance communication subsystems. *IEEE Journal on Selected Areas in Communications*, 11(4):507–518, May 1993.