

Huge pages impact analysis on high-speed network packet capture in the Linux kernel

Master's thesis

Computer and Information Engineering

Laura Trivelloni

Advisor: Prof. Marco Cesati

Co-Advisors: Eng. Emiliano Betti
Eng. Pierpaolo Santucci

April 15, 2019



Introduction





Introduction

Problem analysis

Proposed solution

Results

Conclusions

Packet capture

- traffic analysis
- network monitoring

Requirements

- high-speed
- low packet loss rate





Introduction

Problem analysis

Proposed solution

Results

Conclusions

Packet capture

- traffic analysis
- network monitoring

Requirements

- high-speed
- low packet loss rate

10–40 Gigabit Ethernet





Introduction

Problem analysis

Proposed solution

Results

Conclusions

Linux is an **open-source** operating system.

Frameworks for **high-performance** packet capture:

- **PFQ**
- **Intel DPDK**
- **PF_RING ZC**





Introduction

Problem analysis

Proposed solution

Results

Conclusions

PFQ I/O

- highly optimized for **multi-core** architecture
- RX and TX **line-rate on 10-Gbit links**
- user-kernel memory-mapped buffer backed with **Huge Pages**





Introduction

Problem analysis

Proposed solution

Results

Conclusions

Virtual memory:

- software-controlled set of memory addresses
- data **protection** and **control** among processes

Translations from virtual to physical address cost CPU cycles.





Introduction

Problem analysis

Proposed solution

Results

Conclusions

Virtual memory:

- software-controlled set of memory addresses
- data **protection** and **control** among processes

Translations from virtual to physical address cost CPU cycles.

→ cache for resolved addresses:

Translation Lookaside Buffer





Introduction

Problem analysis

Proposed solution

Results

Conclusions

Virtual memory:

- software-controlled set of memory addresses
- data **protection** and **control** among processes

Translations from virtual to physical address cost CPU cycles.

→ cache for resolved addresses:

Translation Lookaside Buffer

Not always enough to save CPU clocks:

- TLB miss event





Introduction

Problem analysis

Proposed solution

Results

Conclusions

How to reduce the CPU cycles spent for translations?

- fewer translations → fewer CPU cycles
- **wider address range** for each TLB entry
- the same number of TLB entry covers **more addresses**
- more likely TLB **hits** events





Introduction

Problem analysis

Proposed solution

Results

Conclusions

How to reduce the CPU cycles spent for translations?

- fewer translations → fewer CPU cycles
- **wider address range** for each TLB entry
- the same number of TLB entry covers **more addresses**
- more likely TLB **hits** events

Huge Pages
4K → 2M





Introduction

Problem analysis

Proposed solution

Results

Conclusions

Transparent Huge Pages

- **on-demand** allocation
- without user awareness

HugeTLB filesystem

- **pre-allocated** pages
- **application awareness**

Really available?

Not *transparent*!



Problem analysis



The AF_PACKET kernel subsystem



Introduction

Problem analysis

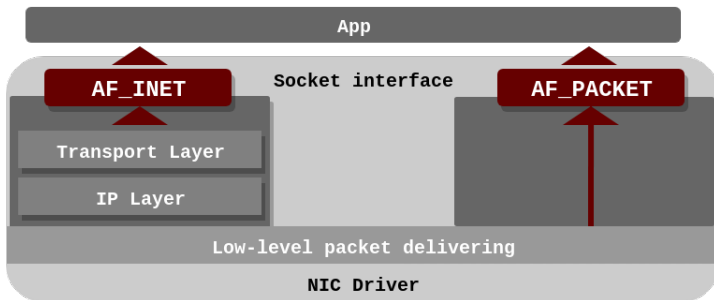
Proposed solution

Results

Conclusions

The AF_PACKET implements the networking domain of the **raw sockets**:

- low-level access to the network interface



Traffic analysis in Linux



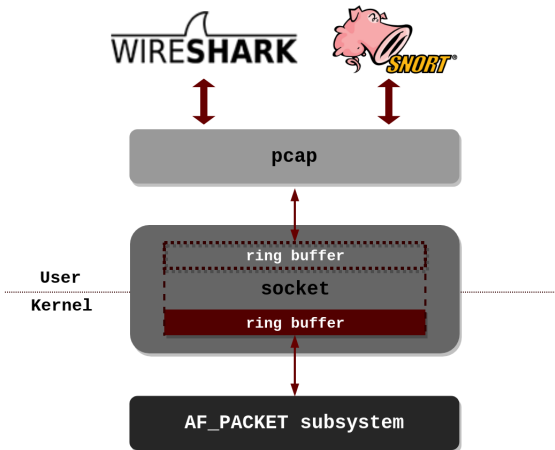
Introduction

Problem
analysis

Proposed
solution

Results

Conclusions



PFQ vs AF_PACKET



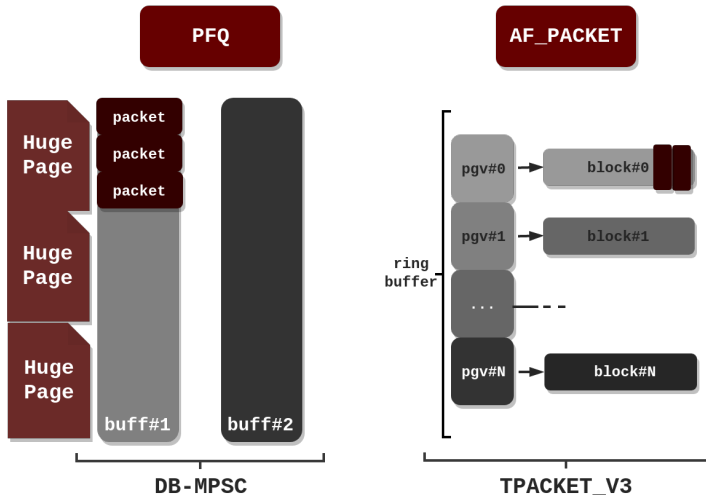
Introduction

Problem analysis

Proposed solution

Results

Conclusions





Introduction

Problem analysis

Proposed solution

Results

Conclusions

Premises:

- the Linux kernel **does not use Huge Pages** in the buffer reception
- PFQ that optimizes the packet capture **uses Huge Pages**





Introduction

Problem analysis

Proposed solution

Results

Conclusions

Premises:

- the Linux kernel **does not use Huge Pages** in the buffer reception
- PFQ that optimizes the packet capture **uses Huge Pages**

Our question

Can make sense introducing a buffer allocation backed with the Huge Pages in the AF_PACKET domain?



Proposed solution



Our solution



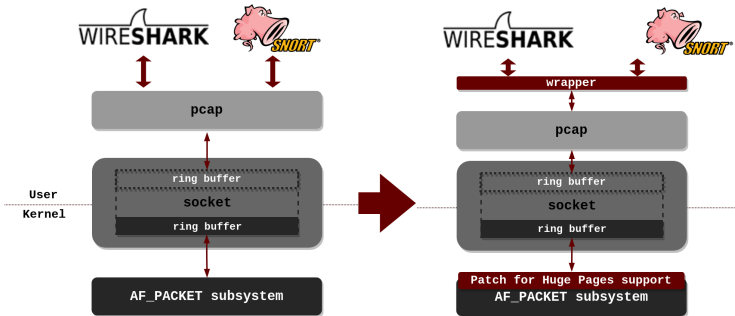
Introduction

Problem analysis

Proposed solution

Results

Conclusions





Introduction

Problem analysis

Proposed solution

Results

Conclusions

The solution is exported in user space as a socket configuration **option**.

Application awareness implies:

- easy **testing** activities





Introduction

Problem analysis

Proposed solution

Results

Conclusions

The solution is exported in user space as a socket configuration **option**.

Application awareness implies:

- easy **testing** activities
- to reach **compatibility** with existing libraries





Introduction

Problem analysis

Proposed solution

Results

Conclusions

The solution is exported in user space as a socket configuration **option**.

Application awareness implies:

- easy **testing** activities
- to reach **compatibility** with existing libraries
 - pcap





Introduction

Problem analysis

Proposed solution

Results

Conclusions

Attempts of block allocation in order of preference:

- 1 **physically contiguous direct** memory request
- 2 only **virtually contiguous** memory request
- 3 physically contiguous **forced** memory request



New block allocation policies



Introduction

Problem
analysis

Proposed
solution

Results

Conclusions

One allocation policy among the following can be forced:

- 1 memory backed with *pre-allocated huge pages*



New block allocation policies



Introduction

Problem
analysis

Proposed
solution

Results

Conclusions

One allocation policy among the following can be forced:

- 1 memory backed with *pre-allocated huge pages*
 - HugeTLB filesystem



New block allocation policies



Introduction

Problem
analysis

Proposed
solution

Results

Conclusions

One allocation policy among the following can be forced:

- 1 memory backed with *pre-allocated huge pages*
 - HugeTLB filesystem
- 2 memory backed with *on-demand huge pages*



New block allocation policies



Introduction

Problem
analysis

Proposed
solution

Results

Conclusions

One allocation policy among the following can be forced:

- 1 memory backed with *pre-allocated huge pages*
 - HugeTLB filesystem
- 2 memory backed with *on-demand huge pages*
 - Transparent Huge Pages



New block allocation policies



Introduction

Problem analysis

Proposed solution

Results

Conclusions

One allocation policy among the following can be forced:

- 1 memory backed with *pre-allocated huge pages*
 - HugeTLB filesystem
- 2 memory backed with *on-demand huge pages*
 - Transparent Huge Pages
- 3 physically contiguous direct memory request





Introduction

Problem analysis

Proposed solution

Results

Conclusions

One allocation policy among the following can be forced:

- 1 memory backed with *pre-allocated huge pages*
 - HugeTLB filesystem
- 2 memory backed with *on-demand huge pages*
 - Transparent Huge Pages
- 3 physically contiguous direct memory request
- 4 only virtually contiguous memory request



New block allocation policies



Introduction

Problem analysis

Proposed solution

Results

Conclusions

One allocation policy among the following can be forced:

- 1 memory backed with *pre-allocated huge pages*
 - HugeTLB filesystem
- 2 memory backed with *on-demand huge pages*
 - Transparent Huge Pages
- 3 physically contiguous direct memory request
- 4 only virtually contiguous memory request
- 5 physically contiguous forced memory request



New block allocation policies



Introduction

Problem analysis

Proposed solution

Results

Conclusions

One allocation policy among the following can be forced:

- 1 memory backed with *pre-allocated huge pages*
 - HugeTLB filesystem
- 2 memory backed with *on-demand huge pages*
 - Transparent Huge Pages
- 3 physically contiguous direct memory request
- 4 only virtually contiguous memory request
- 5 physically contiguous forced memory request

If not specified, the **original** allocation policy is performed.





Introduction

Problem
analysis

**Proposed
solution**

Results

Conclusions

Making cascaded modifications to a library is a costly operation.

- exploiting library **dynamic linking**





Introduction

Problem analysis

Proposed solution

Results

Conclusions

Making cascaded modifications to a library is a costly operation.

- exploiting library **dynamic linking**
- **alternative routines** corresponding to the original symbols at compile-time





Introduction

Problem analysis

Proposed solution

Results

Conclusions

Making cascaded modifications to a library is a costly operation.

- exploiting library **dynamic linking**
- **alternative routines** corresponding to the original symbols at compile-time
- using **environment variables** to transfer not supported information





Introduction

Problem analysis

Proposed solution

Results

Conclusions

Making cascaded modifications to a library is a costly operation.

- exploiting library **dynamic linking**
- **alternative routines** corresponding to the original symbols at compile-time
- using **environment variables** to transfer not supported information

No need to recompile the pcap library!



Results





Introduction

Problem
analysis

Proposed
solution

Results

Conclusions

- generated synthetic traffic





Introduction

Problem
analysis

Proposed
solution

Results

Conclusions

- generated synthetic traffic
 - **pktgen** kernel tool





Introduction

Problem analysis

Proposed solution

Results

Conclusions

- generated synthetic traffic
 - **pktgen** kernel tool
- traffic saved as a pcap file





Introduction

Problem analysis

Proposed solution

Results

Conclusions

- generated synthetic traffic
 - **pktgen** kernel tool
- traffic saved as a pcap file
 - **tcpdump** tool





Introduction

Problem analysis

Proposed solution

Results

Conclusions

- generated synthetic traffic
 - **pktgen** kernel tool
- traffic saved as a pcap file
 - **tcpdump** tool
- pcap trace sent using **tcpreplay** tool





Introduction

Problem analysis

Proposed solution

Results

Conclusions

- generated synthetic traffic
 - **pktgen** kernel tool
- traffic saved as a pcap file
 - **tcpdump** tool
- pcap trace sent using **tcpreplay** tool
- rate at 10 Gb/s





Introduction

Problem
analysis

Proposed
solution

Results

Conclusions

- **hash** computation on each received packet content





Introduction

Problem
analysis

Proposed
solution

Results

Conclusions

- **hash** computation on each received packet content
 - traffic analysis simulation





Introduction

Problem analysis

Proposed solution

Results

Conclusions

- **hash** computation on each received packet content
 - traffic analysis simulation
- variable **delay factor** every 100 packets





Introduction

Problem analysis

Proposed solution

Results

Conclusions

- **hash** computation on each received packet content
 - traffic analysis simulation
- variable **delay factor** every 100 packets
 - 1000 multiplications repeated for *delay factor* times





Introduction

Problem analysis

Proposed solution

Results

Conclusions

- **hash** computation on each received packet content
 - traffic analysis simulation
- variable **delay factor** every 100 packets
 - 1000 multiplications repeated for *delay factor* times
 - buffer overloading to simulate traffic analysis peaks





Introduction

Problem analysis

Proposed solution

Results

Conclusions

- **hash** computation on each received packet content
 - traffic analysis simulation
- variable **delay factor** every 100 packets
 - 1000 multiplications repeated for *delay factor* times
 - buffer overloading to simulate traffic analysis peaks
- statistics computed on 10 runs for each configuration



Test machines



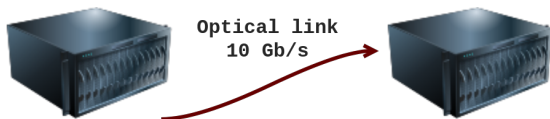
Introduction

Problem
analysis

Proposed
solution

Results

Conclusions



Packet generator

- CPU: Intel® Xeon E5
- NIC: Intel® 82599ES

Packet receiver

- CPU: Intel® Xeon E5
- NIC: Intel® X710





Introduction

Problem analysis

Proposed solution

Results

Conclusions

- **Performance Measurement Unit** provides hardware-based counters
- **perf** tool analyser to measure events:
 - **CPU cycles**
 - **TLB accesses**
 - **TLB misses**
- **Time Stamp Counter** register
 - **CPU cycles** counter in Intel x86



TLB load miss rates



Introduction

Problem analysis

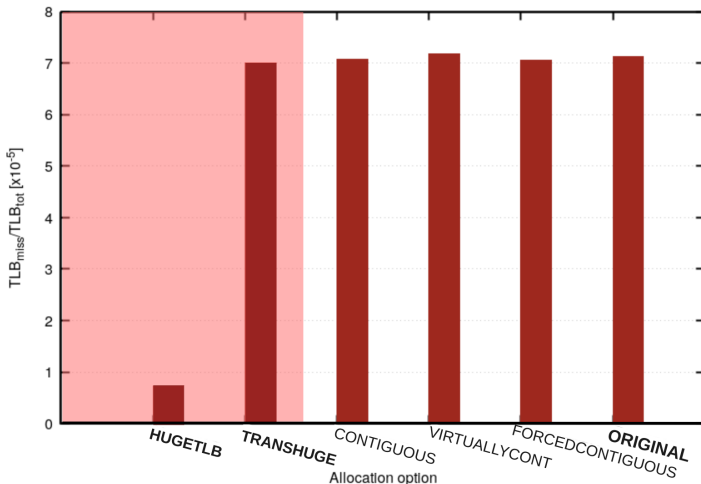
Proposed solution

Results

Conclusions



Average TLB load miss rates comparison



Packet loss



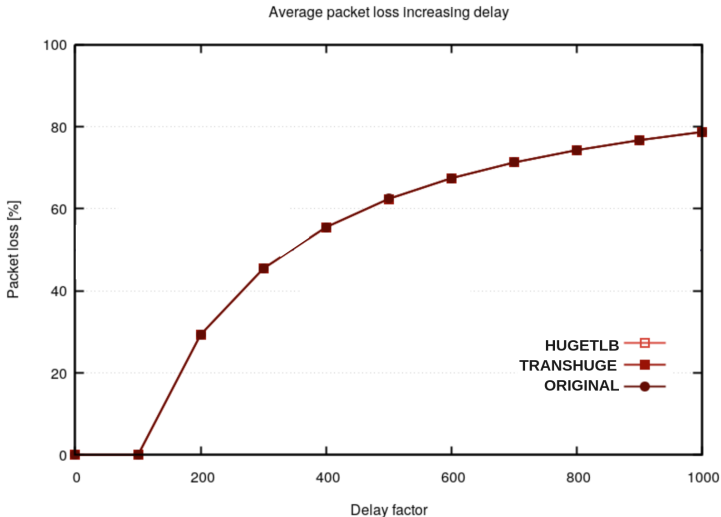
Introduction

Problem
analysis

Proposed
solution

Results

Conclusions



Conclusions





Introduction

Problem
analysis

Proposed
solution

Results

Conclusions

Impact analysis of two new allocation policies based on
huge pages:

- **compatibility** with widespread applications
- effective **improvement** in terms of TLB miss rates
- **no gain** in terms of packet loss rate



Thank You for listening!



Laura Trivelloni

