

# Improving network packet capture by using Linux kernel's huge pages

Master's thesis  
*Computer and Information Engineering*

Laura Trivelloni

**Advisor:** Prof. Marco Cesati

**Co-Advisors:** Ing. Emiliano Betti  
Ing. Pierpaolo Santucci

April 8, 2019





1. Introduction
2. Problem analysis
3. Proposed solutions
4. Results
5. Conclusions

# Introduction





## Introduction

## Problem analysis

## Proposed solutions

## Results

## Conclusions

### Internet of Things (IoT) phenomenon:

- multitude of type of connected devices
- sometimes sensitive data flow through the network





## Introduction

## Problem analysis

## Proposed solutions

## Results

## Conclusions

**Network Traffic Analysis** is a common activity aimed at:

- network monitoring
- recognition of traffic characteristics
- information gathering for fighting criminal or terroristic activities



# Problem analysis





## Introduction

## Problem analysis

## Proposed solutions

## Results

## Conclusions

The OS's **socket buffer** has a key role:

- to store temporarily incoming packets before to be analysed
- to mitigate:
  - peaks of traffic
  - large per-packet processing time





## Introduction

## Problem analysis

## Proposed solutions

## Results

## Conclusions

Improving the performance of the buffer management to face with:

- **increasing data rate** supported by the network cards
  - 10 Gigabit Ethernet
- **packet loss** due to processing delay on packet
  - more or less frequent delays based on the type of analysis







## Introduction

## Problem analysis

## Proposed solutions

## Results

## Conclusions

Low-level performance improvement using larger than normal memory pages → **huge pages**

- wider address range for each TLB entry
- the same number of TLB entry covers more addresses
- more likely TLB hits events



# AF\_PACKET kernel module



## Introduction

## Problem analysis

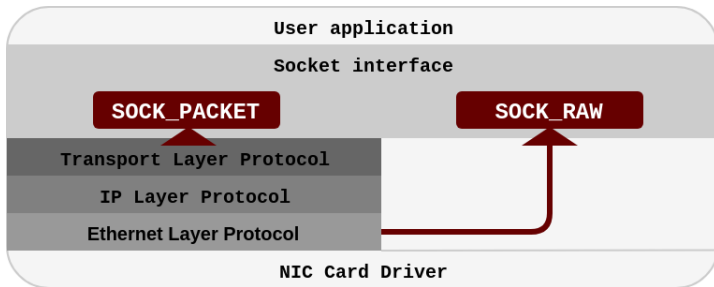
## Proposed solutions

## Results

## Conclusions

The AF\_PACKET module implements the networking domain of the **raw sockets**:

- direct polling on the Ethernet Layer
- to avoid TCP/IP stack analysis overhead



# Buffer implementation



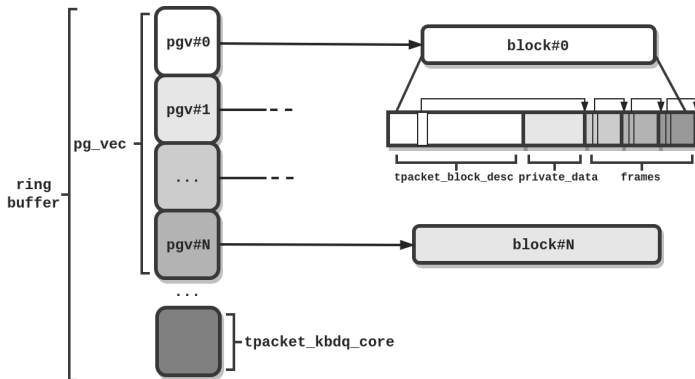
## Introduction

## Problem analysis

## Proposed solutions

## Results

## Conclusions



# Current buffer allocation policy



## Introduction

## Problem analysis

## Proposed solutions

## Results

## Conclusions

Attempts of block allocation in order of preference:

- 1 **physically contiguous** non-blocking memory reclaim
- 2 only **virtually contiguous** memory request
- 3 physically contiguous **repeated** memory reclaim



# Proposed solutions



# New block allocation policies



Introduction

Problem  
analysis

**Proposed  
solutions**

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  
solutions

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  
solutions

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  
solutions

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  
solutions

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 non-blocking memory reclaim



# New block allocation policies



Introduction

Problem  
analysis

Proposed  
solutions

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 non-blocking memory reclaim
- 4 only virtually contiguous memory request



# New block allocation policies



Introduction

Problem  
analysis

Proposed  
solutions

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 non-blocking memory reclaim
- 4 only virtually contiguous memory request
- 5 physically contiguous repeated memory reclaim





## Introduction

## Problem analysis

## Proposed solutions

## 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 non-blocking memory reclaim
- 4 only virtually contiguous memory request
- 5 physically contiguous repeated memory reclaim

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





## Introduction

## Problem analysis

## Proposed solutions

## 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 implementations



# THP solution



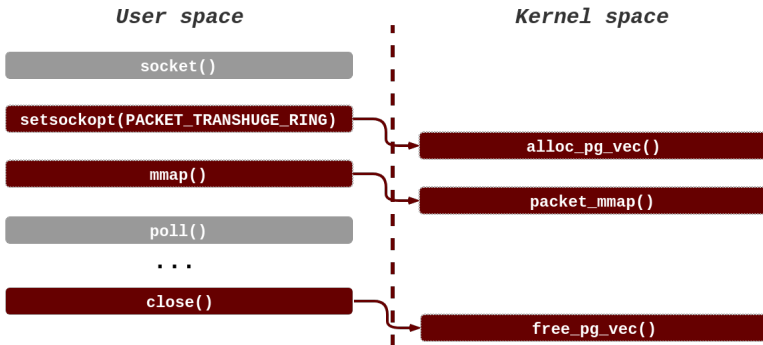
## Introduction

## Problem analysis

## Proposed solutions

## Results

## Conclusions



# HugeTLB filesystem solution



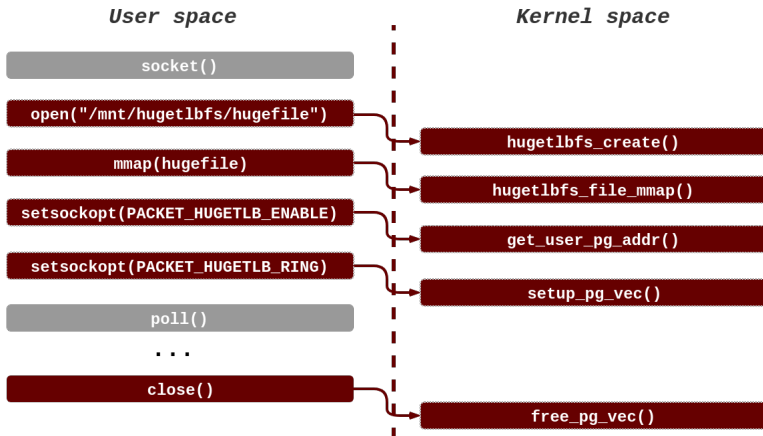
## Introduction

## Problem analysis

## Proposed solutions

## Results

## Conclusions







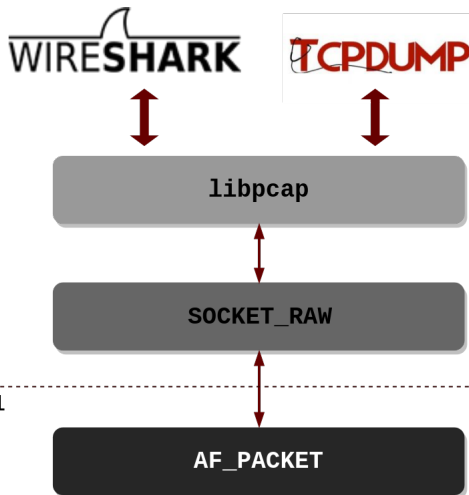
Introduction

Problem  
analysis

Proposed  
solutions

Results

Conclusions





## Introduction

## Problem analysis

## Proposed solutions

## Results

## Conclusions

Making cascaded modifications to a library is a costly operation. → **LD\_PRELOAD trick**

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



# Results





## Introduction

## Problem analysis

## Proposed solutions

## Results

## Conclusions

- generated synthetic randomized traffic
  - **pktgen** kernel tool
- traffic saved as a pcap trace
  - **tcpdump** tool
- pcap trace sent using **tcpreplay** tool
  - more instances to reach the desired rate
- rate at 10 Gb/s
  - **nload** tool



# TLB load miss rates



## Introduction

## Problem analysis

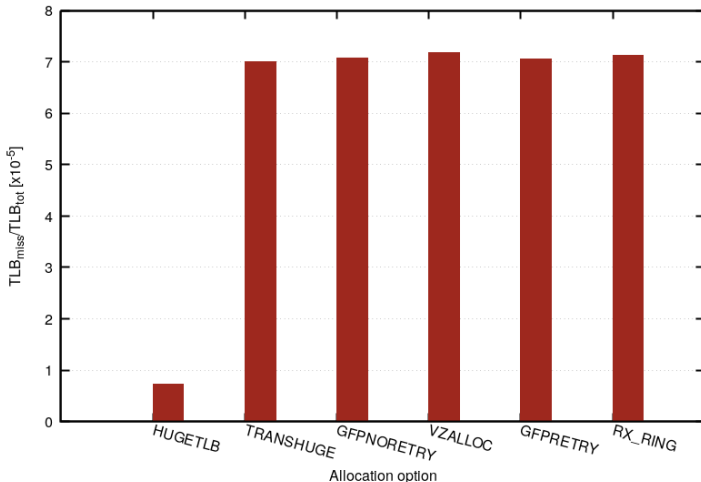
## Proposed solutions

## Results

## Conclusions



Average TLB load miss rates comparison



# TLB load miss rates - main policies



Introduction

Problem  
analysis

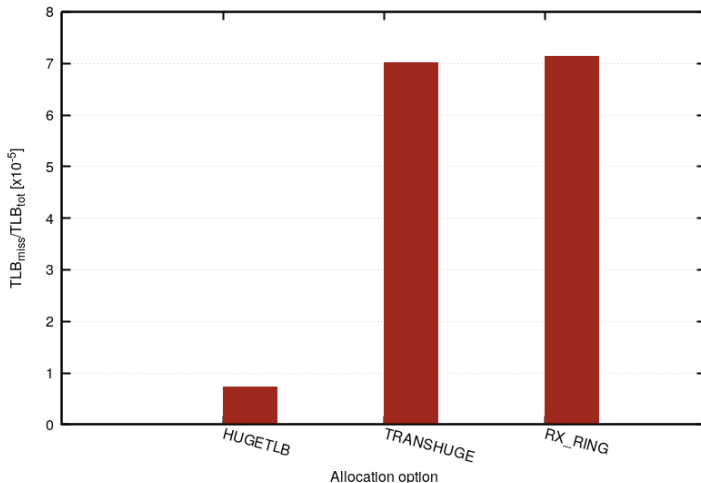
Proposed  
solutions

Results

Conclusions



Average TLB load miss rates comparison



# TLB store miss rates



## Introduction

## Problem analysis

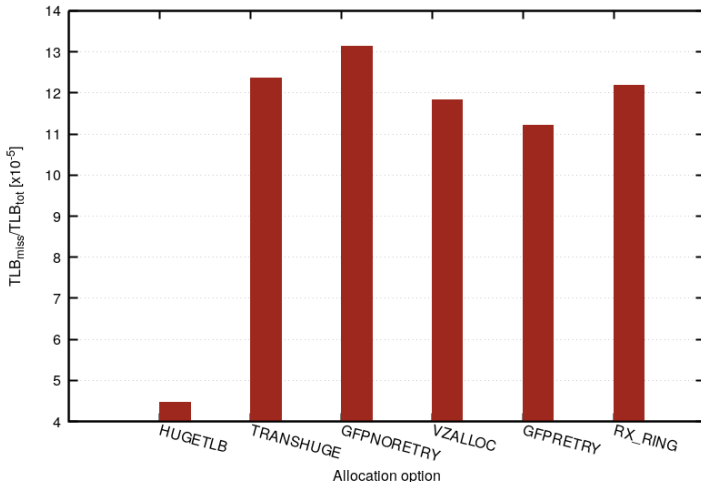
## Proposed solutions

## Results

## Conclusions



Average TLB store miss rates comparison



# TLB load miss rates - main policies



Introduction

Problem  
analysis

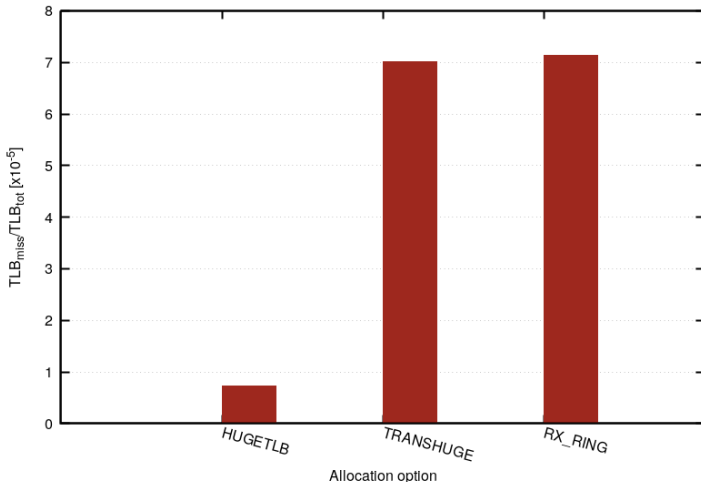
Proposed  
solutions

Results

Conclusions



Average TLB store miss rates comparison





# Packet loss



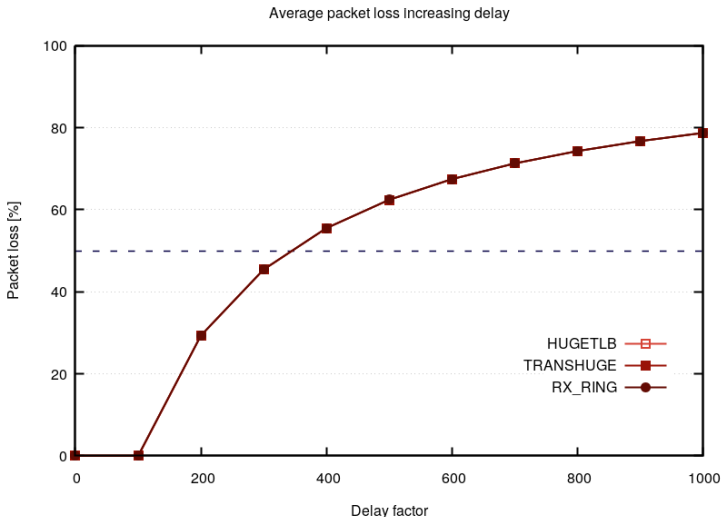
Introduction

Problem  
analysis

Proposed  
solutions

Results

Conclusions



# Conclusions





Introduction

Problem  
analysis

Proposed  
solutions

Results

**Conclusions**

In conclusion, ...



# Ideas for future works



## Introduction

## Problem analysis

## Proposed solutions

## Results

## Conclusions

- grouping all block allocation modes in only one option
- making allocation policies completely **transparent** to the user
  - all the allocation modes in order of preference (as original implementation)



# Thank You for listening!



Laura Trivelloni

