Introduction
Network stack
Packet ingress flow
Methods to capture packets

**University of Sao Paulo - USP**

# Network packet capture in Linux kernelspace
## An overview of the network stack in the Linux kernel

Beraldo Leal

beraldo@ime.usp.br
http://www.ime.usp.br/~beraldo/

Institute of Mathematics and Statistics - IME
University of Sao Paulo - USP

25th October 2011

Introduction
Network stack
Packet ingress flow
Methods to capture packets

University of Sao Paulo - USP

# Outline

Introduction

Network stack

Packet ingress flow

Methods to capture packets

Introduction
Network stack
Packet ingress flow
Methods to capture packets

**University of Sao Paulo - USP**

## Introduction

- Sniffers;
- Improvements in packet reception;
- Linux kernel network subsystem;

Introduction
Network stack
Packet ingress flow
Methods to capture packets

University of Sao Paulo - USP

# Sniffers

- tcpdump, wireshark, snort, etc;
- Using the well-known library libpcap;
- Not suitable for $> 10$ Gbps;
- Packet loss;

Introduction
Network stack
Packet ingress flow
Methods to capture packets

University of Sao Paulo - USP

## Improvements in packet reception

- Commodity hardware for packet capture;
  - 3COM
  - Intel
  - endace, ...
- Many Interruptions
- NEW API or NAPI (interruption coalescence)
- zero-copy
  - Direct Memory Access - DMA
  - mmap()

Introduction
Network stack
Packet ingress flow
Methods to capture packets

University of Sao Paulo - USP

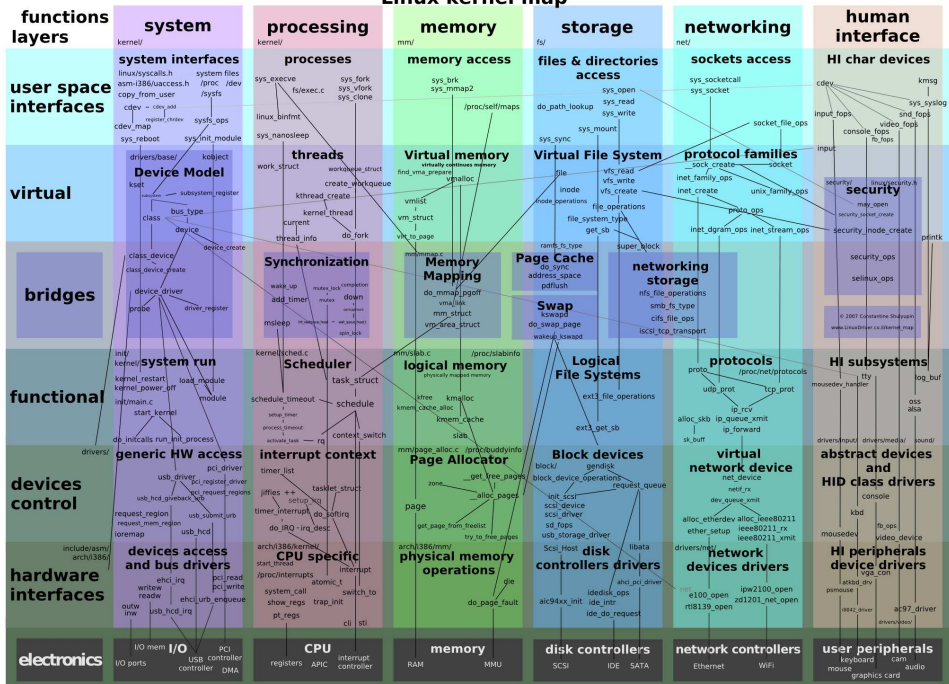# Linux kernel network subsystem

- Kernel number of files: 36.680 [1] [2]
- net/ number of files: 1.293 (  3.5% )
- drivers/net/ number of files: 1.935 (  5.27% )
- Kernel SLOC: 9.723.525
- net/ SLOC: 480.928 (  5% )
- drivers/net/ SLOC: 1.155.317 (  12% )

---

[1]kernel 3.0.0
[2]source: wc, find, cat, etc..

# Linux kernel map

| functions layers | system kernel/ | processing kernel/ | memory mm/ | storage fs/ | networking net/ | human interface |
|---|---|---|---|---|---|---|

**user space interfaces**

- **system interfaces**
  linux/syscalls.h · system files
  asm-i386/uaccess.h · /proc /dev
  copy_from_user · /sysfs
  cdev ← cdev_add
  linux/fs.h · register_chrdev
  cdev_map
  sys_reboot · sys_init_module
- **processes**
  sys_execve · sys_fork
  fs/exec.c · sys_vfork · sys_clone
  linux_binfmt
  sys_nanosleep
- **memory access**
  sys_brk
  sys_mmap2
  /proc/self/maps
- **files & directories access**
  sys_open
  sys_read
  sys_write
  do_path_lookup
  sys_mount
  sys_sync
- **sockets access**
  sys_socketcall
  sys_socket
  socket_file_ops
- **HI char devices**
  cdev · kmsg
  input_fops · sys_syslog
  console_fops · video_fops
  fb_fops
  input

**virtual**

- **Device Model**
  drivers/base/ · kobject
  subsystem_register
  class · bus_type
  device
  class_device
  class_device_create
  device driver
  probe · driver_register
  device_create
- **threads**
  work_struct · workqueue_struct
  create_workqueue
  kthread_create
  kernel_thread
  current
  thread_info · do_fork
- **Virtual memory**
  virtually continuous memory
  find_vma_prepare
  vmalloc
  vmlist
  vm_struct
  vm_to_page
- **Virtual File System**
  vfs_read · file
  vfs_write
  inode · vfs_create
  inode_operations · file_operations
  file_system_type · get_sb
- **protocol families**
  sock_create · socket
  inet_family_ops
  inet_create · unix_family_ops
  proto_ops
  inet_dgram_ops inet_stream_ops
- **security**
  security/ · security.h
  may_open
  security_file_ops
  security_inode_create
  security_ops
  selinux_ops
  printk

**bridges**

- **Synchronization**
  wake_up
  add_timer · mutex
  down
  mutex_lock · up
  msleep
  rt_mutex_lock · set_current_task
  spin_lock
- **Memory Mapping**
  mm/mmap.c
  do_mmap_pgoff
  vma_link
  mm_struct
  vm_area_struct
- **Page Cache**
  do_sync
  address_space
  pdflush
- **Swap**
  kswapd
  do_swap_page
  wakeup_kswapd
- **networking storage**
  nfs_file_operations
  smb_fs_type
  cifs_file_ops
  iscsi_tcp_transport
  ramfs_fs_type · super_block

© 2007 Constantine Shulyupin
www.LinuxDriver.co.il/kernel_map

**functional**

- **system run**
  init/ kernel/
  kernel_restart · load_module
  kernel_power_off
  init/main.c · module
  start_kernel
  do_initcalls run_init_process
- **Scheduler**
  kernel/sched.c
  task_struct
  schedule_timeout · schedule
  setup_timer
  context_switch
  activate_task · rq
- **logical memory**
  physically mapped memory
  kfree · kmalloc
  kmem_cache_alloc
  kmem_cache
  slab
- **Logical File Systems**
  ext3_file_operations
  ext3_get_sb
- **protocols**
  proto · /proc/net/protocols
  udp_prot · tcp_prot
  ip_rcv
  ip_queue_xmit
  ip_forward
  sock_skb
  sk_buff
- **HI subsystems**
  tty · log_buf
  mousedev_handler
  oss
  alsa
  drivers/input/ drivers/media/ sound/

**devices control**

- **generic HW access**
  drivers/
  pci_driver
  pci_register_driver · request_regions
  usb_hcd_giveback_urb · request_regions
  request_region · usb_submit_urb
  request_mem_region · usb_hcd
  ioremap
- **interrupt context**
  timer_list
  jiffies ++ · tasklet_struct
  setup_irq
  timer_interrupt · do_softirq
  do_IRQ · irq_desc
- **Page Allocator**
  block/
  block_device_operations
  zone · request_queue
  alloc_pages
  init_scsi
  scsi_device
  scsi_driver
  get_page_from_freelist
  try_to_free_pages
- **Block devices**
  gendisk
  sd_blank
  Scsi_Host · libata
  disk controllers drivers
- **virtual network device**
  net_device
  netif_rx
  dev_queue_xmit
  alloc_etherdev alloc_ieee80211
  ether_setup ieee80211_rx
  ieee80211_xmit
- **network devices drivers**
  drivers/net/
  e100_open ipw2100_open
  rtl8139_open zd1201_net_open
- **abstract devices and HID class drivers**
  console
  kbd
  mousedev · fb_ops
  video_device
- **HI peripherals device drivers**
  vga_con
  atkbd_drv
  psmouse
  i8042_drv · ac97_driver

**hardware interfaces**

- **devices access and bus drivers**
  include/asm/
  arch/i386/
  ehci_irq · pci_read
  readw · pci_write
  writew
  outw · ehci_hcd urb_enqueue
  inw · usb_hcd_irq
- **CPU specific**
  arch/i386/kernel/
  /proc/interrupts
  atomic_t
  system_call · switch_to
  interrupt · trap_init
  pt_regs
  cli · sti
- **physical memory operations**
  arch/i386/mm/
  gle
  do_page_fault
- **disk controllers drivers**
  idedisk_ops
  ide_init
  aic94xx_init · ide_do_request

**electronics**

| I/O | CPU | memory | disk controllers | network controllers | user peripherals |
|---|---|---|---|---|---|
| I/O ports · USB controller · PCI controller · DMA | registers · APIC · interrupt controller | RAM · MMU | SCSI · IDE · SATA | Ethernet · WiFi | keyboard · mouse · cam · graphics card |

Introduction
**Network stack**
Packet ingress flow
Methods to capture packets

**University of Sao Paulo - USP**

# Network stack

## L5: Application

http, ftp, ssh, telnet, ... (message)

## L4: Transport

tcp, udp, ... (segment)

## L3: Network

ipv4, ipv6, ... (datagram/packet)

## L1/2: Link / host-to-network

ethernet, token ring, ... (frame)

Introduction
**Network stack**
Packet ingress flow
Methods to capture packets

**University of Sao Paulo - USP**

## Important data structs:

- `net_device`
  - `include/linux/netdevice.h`
- `sk_buff`
  - `include/linux/skbuff.h`

Introduction
**Network stack**
Packet ingress flow
Methods to capture packets

**University of Sao Paulo - USP**

## Important data structs:

- `net_device` (include/linux/netdevice.h)
  - `unsigned int mtu`
  - `unsigned int flags`
  - `unsigned char dev_addr[MAX_ADDR_LEN]`
  - `int promiscuity`

Introduction
**Network stack**
Packet ingress flow
Methods to capture packets

University of Sao Paulo - USP

## Important data structs:

- `sk_buff (include/linux/skbuff.h)`
  - `struct sk_buff *next;`
  - `struct sk_buff *prev;`
  - `ktime_t tstamp;`
  - `struct net_device *dev;`
  - `unsigned int len;`
  - `unsigned int data_len;`
  - `__u16 mac_len;`
  - `__u8 pkt_type;`
  - `__be16 protocol;`
  - `sk_buff_data_t transport_header; (old h)`
  - `sk_buff_data_t network_header; (old nh)`
  - `sk_buff_data_t mac_header; (old mac)`

Introduction
**Network stack**
Packet ingress flow
Methods to capture packets

**University of Sao Paulo - USP**

## Important sk_buff routines

- `alloc_skb();`
- `dev_alloc_skb();`
- `kfree_skb();`
- `dev_kfree_skb();`
- `skb_clone();`
- `skb_network_header(skb);`
- `skb_transport_header(skb);`
- `skb_mac_header(skb);`

Introduction
Network stack
Packet ingress flow
Methods to capture packets

University of Sao Paulo - USP

## Packet ingress flow

- When working in interrupt driven model, the nic registers an interrupt handler;
- This interrupt handler will be called when a frame is received;
- Typically in the handler, we allocate `sk_buff` by calling `dev_alloc_skb()`;
- Copies data from nic's buffer to this struct just created;
- nic call generic reception routine `netif_rx()`;
- `netif_rx()` put frame in per cpu queue;
- if queue is full, drop!
- `net_rx_action()` decision based on `skb->protocol`;
- This function basically dequeues the frame and delivery a copy for every protocol handler;
  - `ptype_all` and `ptype_base` queues

Introduction
Network stack
**Packet ingress flow**
Methods to capture packets

**University of Sao Paulo - USP**

# Packet ingress flow

- `ip_v4_rcv()` will receive the ip datagram (if is a ipv4 packet);
- ip checksum, check ip headers, ....
- `ip_rcv_finish()` makes route decision (`ip_forward()` or `ip_local_delivery()`)
- `ip_local_delivery()` defrag fragmented packets, and call `ip_local_deliver_finish()`
- `ip_local_deliver_finish()` find protocol handler again;
- `tcp_v4_rcv()`, `udp_rcv()`, or other L4 protocol handler
- ...

**ip_local_deliver_finish()**
(net/ipv4/ip_input.c)
find protocol handler or
send icmp_dst_unreach

NF_IP_FORWARD

<continue>

**Layer 3
Network**

**ip_forward()**
(net/ipv4/ip_forward.c)
handle route alert;
send redirect if
necessary;
decrease TTL;
verify if frag is
possible (mtu)

NF_IP_LOCAL_IN

**ip_local_deliver()**
(net/ipv4/ip_input.c)
defrag fragmented
packets

NF_IP_PRE_ROUTING

**ip_rcv_finish()**
(net/ipv4/ip_input.c)
find route and handle
IP options

**ip_error()**
(net/ipv4/route.c)
routing error, send
icmp pkt

**ip_rcv()**
(net/ipv4/ip_input.c)
verify skb, IP headers
and IP checksum

**arp_rcv()**
(handle arp requests
and replies)

**packet_rcv()**
<tcpdump_process>
<dhcpd process>
<...>

<...>

**netif_rx()**
(net/core/dev.c)

input_queue
[cpu]

**net_rx_action()**
(net/core/dev.c)
decision based on
skb->protocol field

**Layer 1/2
Physical/Link**

**Network Drivers
(drivers/net/*)**

**Application**

userspace

kernelspace

**Socket Layer**
**(net/core/sock.c)**

**tcp_v4_do_rcv()**
(net/ipv4/tcp_ipv4.c)
check for socket state

**_tcp_v4_lookup()**
(net/ipv4/tcp_ipv4.c)
check for socket in
LISTEN, with dst_port

generate ICMP
error

**tcp_v4_rcv()**
(net/ipv4/tcp_ipv4.c)
check for tcp headers

**udp_rcv()**
(net/ipv4/udp.c)
check for udp headers

<...>

**Layer 4**
**Transport**

**ip_local_deliver_finish()**
(net/ipv4/ip_input.c)
find protocol handler or
send icmp_dst_unreach

NF_IP_FORWARD

<continue>

**ip_forward()**
(net/ipv4/ip_forward.c)
handle route alert;
send redirect if
necessary;
decrease TTL;
verify if frag is
possible (mtu)

**Layer 3**
**Network**

NF_IP_LOCAL_IN

**ip_local_deliver()**
(net/ipv4/ip_input.c)
defrag fragmented
packets

NF_IP_PRE_ROUTING

**ip_rcv_finish()**
(net/ipv4/ip_input.c)
find route and handle
IP options

**ip_error()**
(net/ipv4/route.c)
routing error, send
icmp pkt

Introduction
Network stack
Packet ingress flow
Methods to capture packets

**University of Sao Paulo - USP**

# Methods to capture packets

- protocol handler
  - register a function to handler packets with `dev_add_pack()`
- netfilter hooks
- userspace tools;
  - socket `AF_PACKET`, libpcap, ...

Introduction
Network stack
Packet ingress flow
Methods to capture packets

University of Sao Paulo - USP
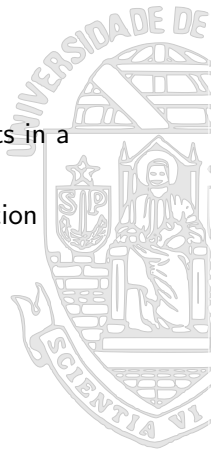
```
1    struct packet_type my_proto;
2
3    int my_packet_rcv(struct sk_buff *skb, struct net_device *dev,
4                      struct packet_type *pt, struct net_device *orig_dev) {
5
6        printk(KERN_ERR "+ 1!\n");
7
8        kfree_skb(skb);
9        return 0;
10   }
11
12   static int hello_init(void) {
13        printk("<1> Hello world!\n");
14
15        my_proto.type = htons(ETH_P_ALL);
16        my_proto.dev  = NULL;
17        my_proto.func = my_packet_rcv;
18
19        dev_add_pack(&my_proto);
20        return 0;
21   }
22
23   static void hello_exit(void) {
24        dev_remove_pack(&my_proto);
25        printk("<1> Bye, cruel world\n");
26   }
27   module_init(hello_init);
28   module_exit(hello_exit);
```

Introduction
Network stack
Packet ingress flow
Methods to capture packets

**University of Sao Paulo - USP**

```
1   int my_packet_rcv(struct sk_buff *skb, struct net_device *dev, struct packet_type *pt, struct net_device
        *orig_dev)
2   {
3       switch (skb->pkt_type) {
4               case PACKET_HOST:
5                   printk("PACKET_HOST - ");
6                   break;
7               case PACKET_BROADCAST:
8                   printk("PACKET_BROADCAST - ");
9                   break;
10              case PACKET_MULTICAST:
11                  printk("PACKET_MULTICAST - ");
12                  break;
13              case PACKET_OTHERHOST:
14                  printk("PACKET_OTHERHOST - ");
15                  break;
16              case PACKET_OUTGOING:
17                  printk("PACKET_OUTGOING - ");
18                  break;
19              case PACKET_LOOPBACK:
20                  printk("PACKET_LOOPBACK - ");
21                  break;
22              case PACKET_FASTROUTE:
23                  printk("PACKET_FASTROUTE - ");
24                  break;
25      }
26      printk("%s 0x%.4X 0x%.4X \n", skb->dev->name, ntohs(skb->protocol), ip_hdr(skb)->protocol)
27
28      kfree_skb(skb);
29      return 0;
30  }
```

Introduction
Network stack
Packet ingress flow
Methods to capture packets

University of Sao Paulo - USP
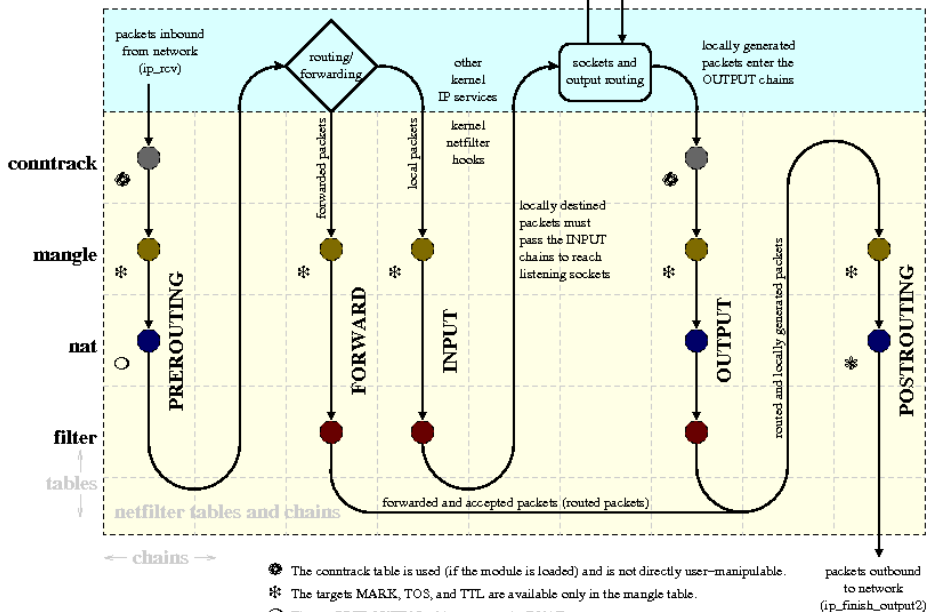
# Netfilter hooks

- iptables = userspace;
- netfilter = kernelspace;
- Netfilter is merely a series of hooks in various points in a protocol stack;
- packet filtering, network address [and port] translation (NA[P]T) and other packet mangling;
- www.netfilter.org

# Netfilter Packet Traversal

http://linux-ip.net/nf/nfk-traversal.png

Martin A. Brown, martin@linux-ip.net

applications

packets inbound from network (ip_rcv)

routing/forwarding

other kernel IP services

sockets and output routing

locally generated packets enter the OUTPUT chains

**conntrack**

**mangle**

**nat**

**filter**

forwarded packets

local packets

kernel netfilter hooks

locally destined packets must pass the INPUT chains to reach listening sockets

routed and locally generated packets

PREROUTING

FORWARD

INPUT

OUTPUT

POSTROUTING

tables

netfilter tables and chains

← chains →

forwarded and accepted packets (routed packets)

❀ The conntrack table is used (if the module is loaded) and is not directly user-manipulable.

❆ The targets MARK, TOS, and TTL are available only in the mangle table.

○ The nat PREROUTING table supports the DNAT target.

❆ The nat POSTROUTING table supports SNAT and MASQUERADE targets.

cf. http://www.docum.org/qos/kptd/

cf. http://open-source.arkoon.net/kernel/kernel_net.png

packets outbound to network (ip_finish_output2)

# Linux Kernel 2.4
# Packet handling

**Network Drivers (drivers/net/\*)**

raise NET_RX_SOFTIRQ

**netif_rx**
(net/core/dev.c)

**net_rx_action**
(net/core/dev.c)
[NET_RX softirq]

input queue
[cpu]

**ip_forward_finish**
(net/ipv4/ip_forward.c)
o handle IP options
o fragment packet

**ip_finish_output2**
(net/ipv4/ip_output.c)
o calls hh output or
o dst output routine

**dev_queue_xmit**
(net/core/dev.c)

NF_IP_POST_ROUTING

**packet_rcv**
<tcpdump process>
<dhtcpd process>
<...>

**arp_rcv**
handle arp requests
and replies

**ip_rcv**
(net/ipv4/ip_input.c)
verify skb, IP header
and IP checksum

NF_IP_PRE_ROUTING

**ip_local_deliver**
(net/ipv4/ip_input.c)
defrag fragmented packet

NF_IP_LOCAL_IN

**ip_local_deliver_finish**
(net/ipv4/ip_input.c)
find protocol handler or
send icmp_dest_unreach

**ip_rcv_finish**
(net/ipv4/ip_input.c)
o find input route if unknown
o handle IP options

NF_IP_FORWARD

**ip_forward**
(net/ipv4/ip_forward.c)
o handle router alert
o verify TTL
o verify strict routing
o send redirect if needed
o decrease TTL
o verify that frag is possible (mtu)

**ip_error**
(net/ipv4/route.c)
routing error : send icmp pkt

**ip_finish_output**
(net/ipv4/ip_output.c)

**ip_output**
(net/ipv4/ip_output.c)

NF_IP_LOCAL_OUT

**ip_build_xmit**
**ip_build_xmit_slow**
**ip_build_and_send_pkt**
**ip_queue_xmit**
(net/ipv4/ip_output.c)
o create and build IP packet
o find output route
o ...

**icmp_send**
(net/ipv4/icmp.c)
send an ICMP message

**Local IP Services**

Hervé Lefort – Arkoon Network Security – Feb 2002

Introduction
Network stack
Packet ingress flow
Methods to capture packets

University of Sao Paulo - USP

## References

- `br.kernelnewbies.org/node/150` has many links

Introduction
Network stack
Packet ingress flow
Methods to capture packets

University of Sao Paulo - USP

# Thankyou! Question?

Introduction
Network stack
Packet ingress flow
Methods to capture packets

**University of Sao Paulo - USP**

# Network packet capture in Linux kernelspace
## An overview of the network stack in the Linux kernel

Beraldo Leal

beraldo@ime.usp.br
http://www.ime.usp.br/~beraldo/

Institute of Mathematics and Statistics - IME
University of Sao Paulo - USP

25th October 2011