

# 11

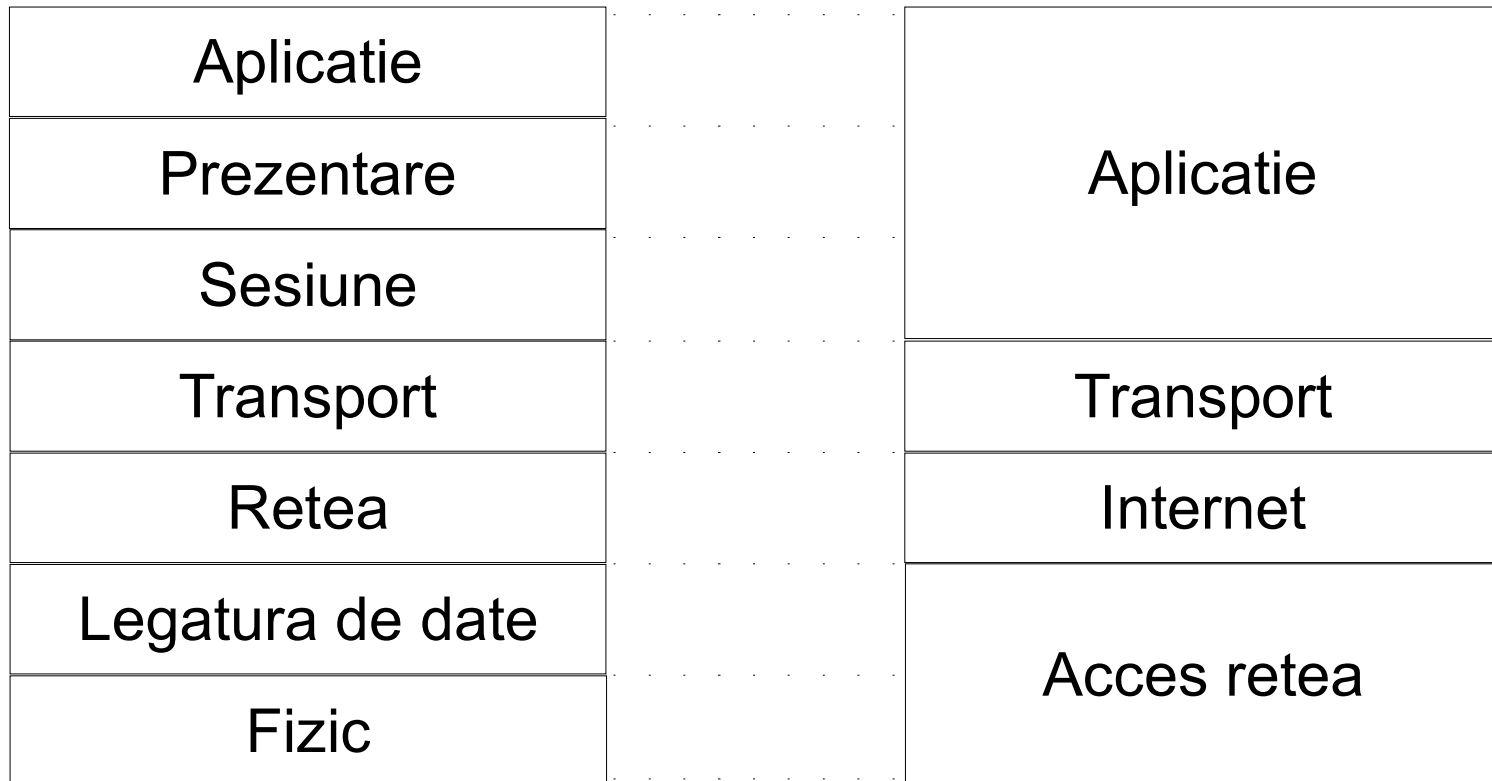
## Gestiunea rețelei - Linux

14 mai 2009

- Dați un exemplu în care tabelele de pagini nu sunt suficiente pentru a trata un page fault, situație în care sistemul de operare trebuie să acceseze descriptorul spațiului de adresă al procesului.
- Fie  $a$  un fișier,  $b$  un hard-link către  $a$  și  $c$  un sof-link către  $a$ . Câte inode-uri sunt asociate cu cele trei fișiere?
- Ce operații sunt efectuate la trunchierea unui fișier într-un sistem de fișiere (asupra inode-ului, dentry-ului, superbloc-ului, vectorului ce descrie blocurile libere, etc.)

- Implementarea socketilor în Linux
- Rutarea în Linux
- Interfața cu hardware-ul
- Optimizări

- Retelele sunt controlate prin intermediul unor protocoale, de obicei grupate într-o stivă de protocoale
- Avantajele unei stive de protocoale: modularitate
  - Reduce complexitate
  - Mai ușor de înțeles
  - Ușor de modificat



Berkely Socket Interface

Transport layer

TCP

UDP

Network layer

IP

Routing

NetFilter

Data link layer

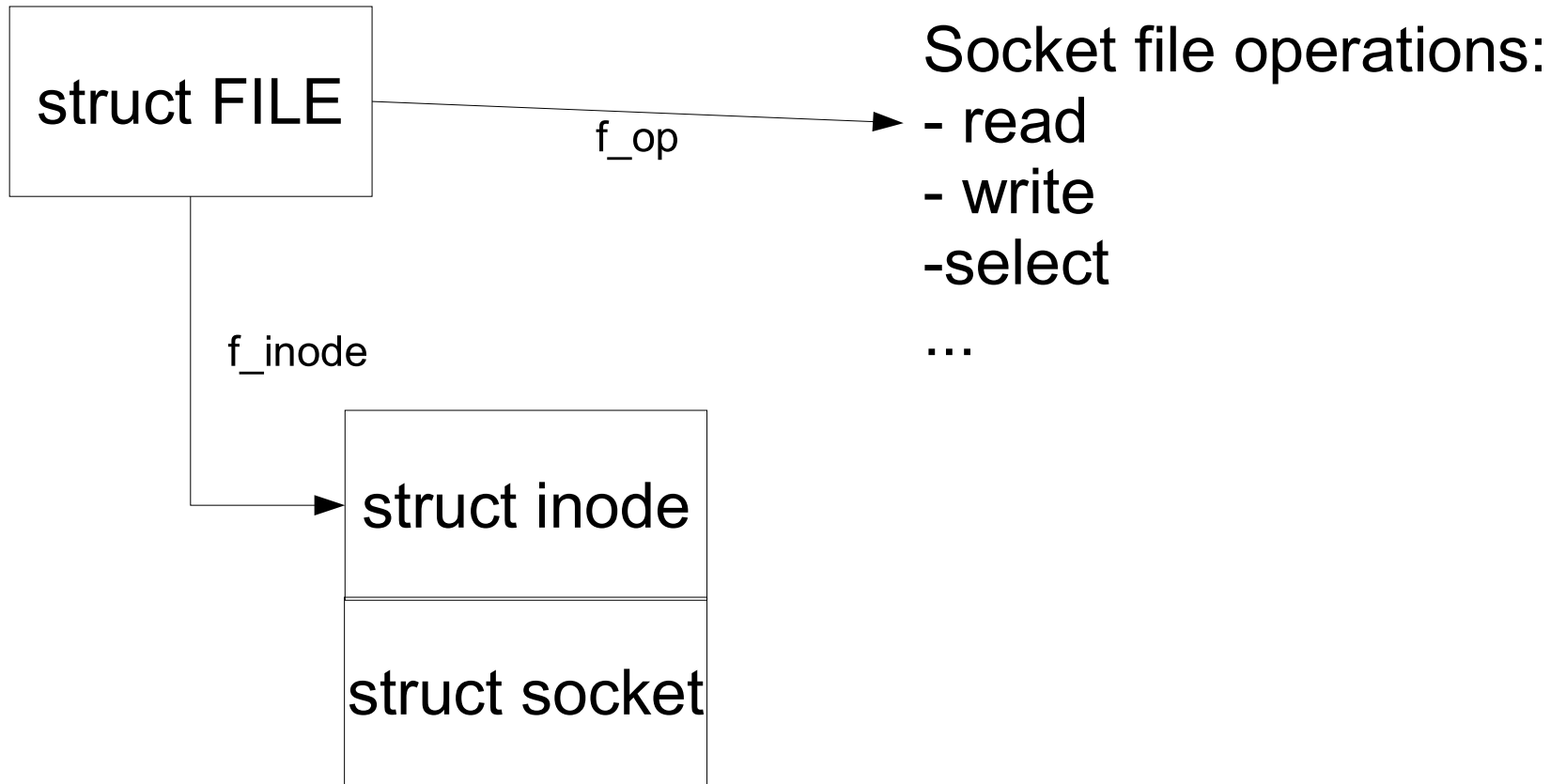
ETH

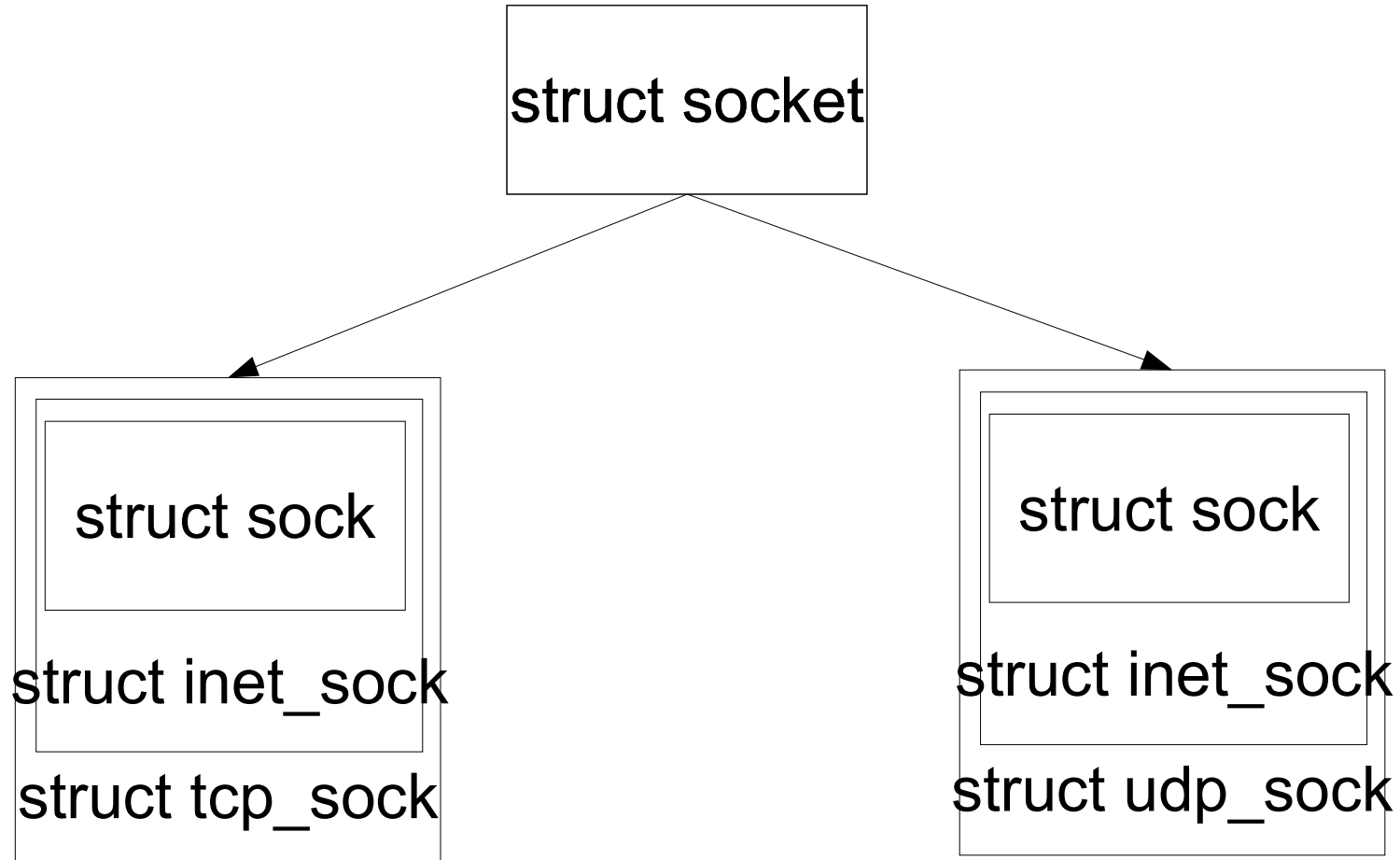
ARP

BRIDGING

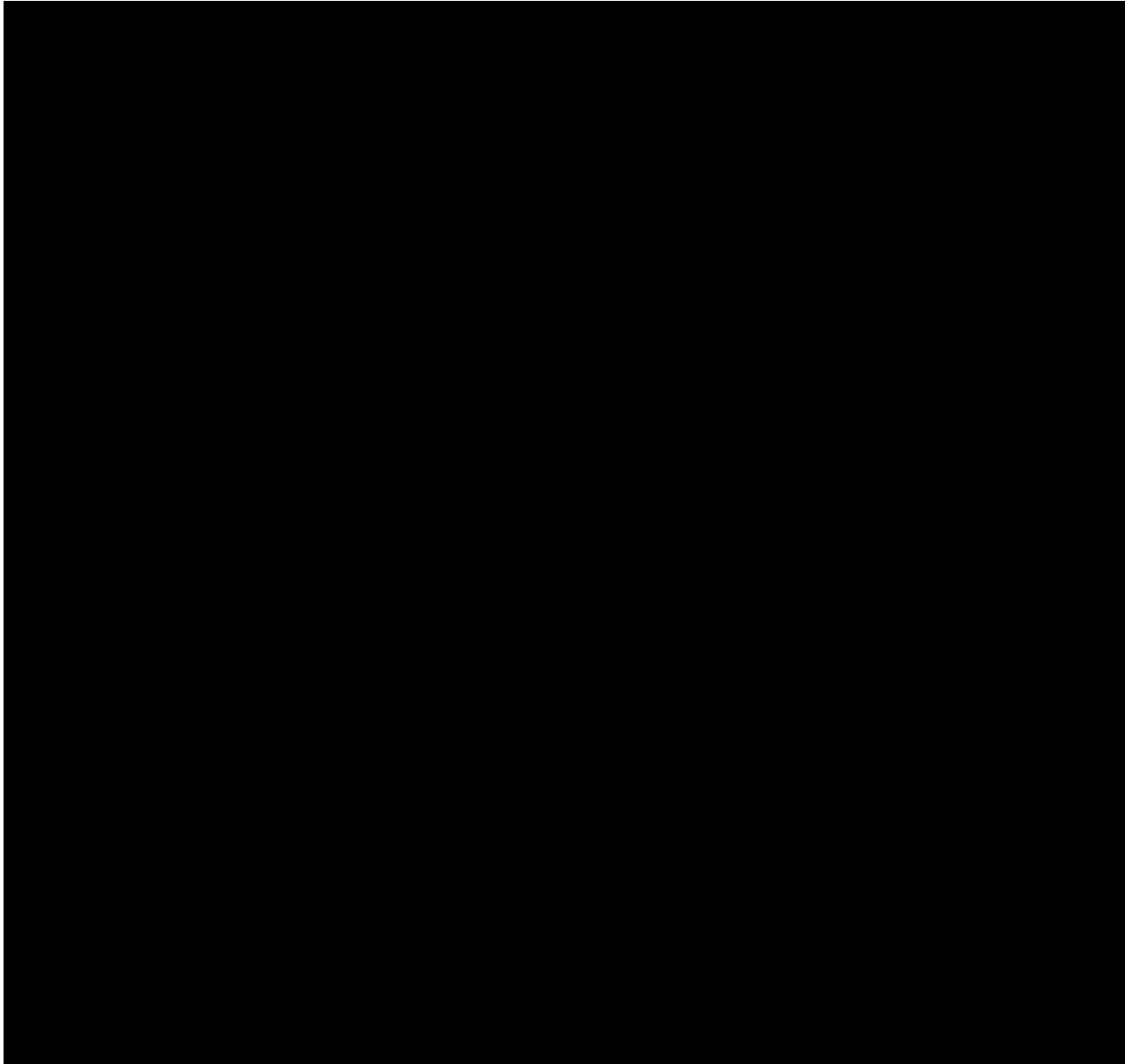
Queing discipline

Network device drivers









```
[tavi@tropaila tavi]$ /sbin/route
```

```
Kernel IP routing table
```

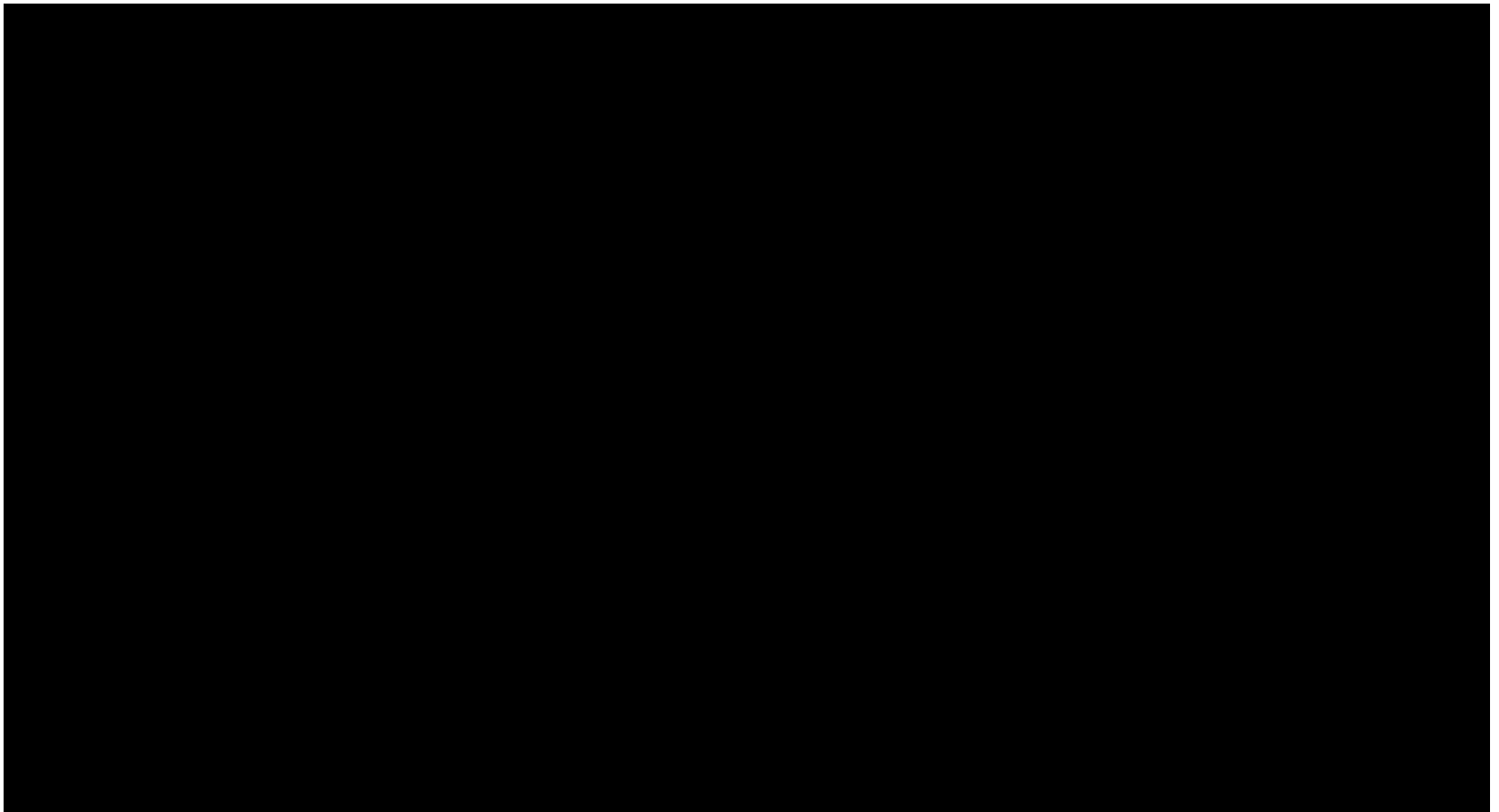
Destination	Gateway	Genmask	Flags	Metric	Ref	Use	Iface
192.168.0.0	*	255.255.255.0	U	0	0	0	eth0
172.16.1.0	192.168.58.1	255.255.255.0	UG	0	0	0	vmnet1
172.16.125.0	*	255.255.255.0	U	0	0	0	vmnet8
192.168.58.0	*	255.255.255.0	U	0	0	0	vmnet1
127.0.0.0	*	255.0.0.0	U	0	0	0	lo

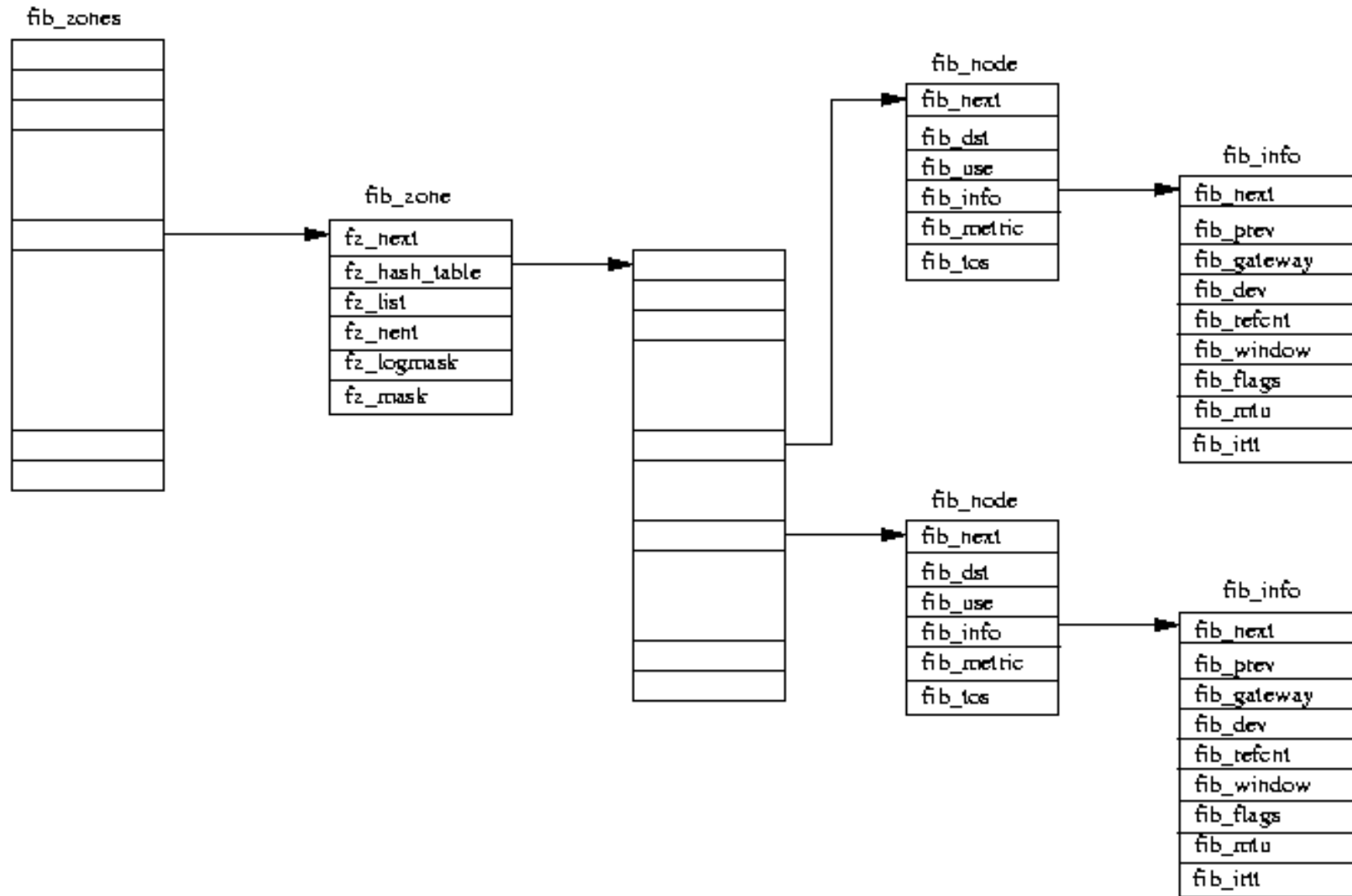
- Contine informatii necesare rutarii unui pachet pentru o retea cunoscuta
  - Adresa retelei
  - Masca retelei
  - Interfata sau next hop-ul pe unde trebuie trimis pachetul
  - Metrica

- Pentru pachetele primite se verifica daca adresa destinatie coincide cu vreo adresa locala
  - Daca da, se trimite pachetul local
- Daca nu, se ruteaza pachetul
  - Se determina
    - Interfata pe care trebuie trimis pachetul
    - Adresa de nivel 2 a next hop-ului
    - Adresa de nivel 3 a next hop-ului
  - daca este necesar se fragmenteaza / defragmenteaza pachetul
  - se suprascrive header-ul de nivel 2 si se transmite pachetul de interfata determinata

- Procesul de rutare este initiat din bottom-half handler-ul asociat gestiunii retelei
- Pentru determinarea informatiilor de rutare se consulta in ordine
  - Cache-ul de rute
  - Forwarding Information Database (FIDB)

- Structura de date ce contine informatii necesare rutarii
- O tabela cu 32 de intrari, cate una pentru fiecare masca de retea posibila
- O intrare in tabela contine un pointer catre o lista de structuri ce descriu rute cu aceasi masca de retea

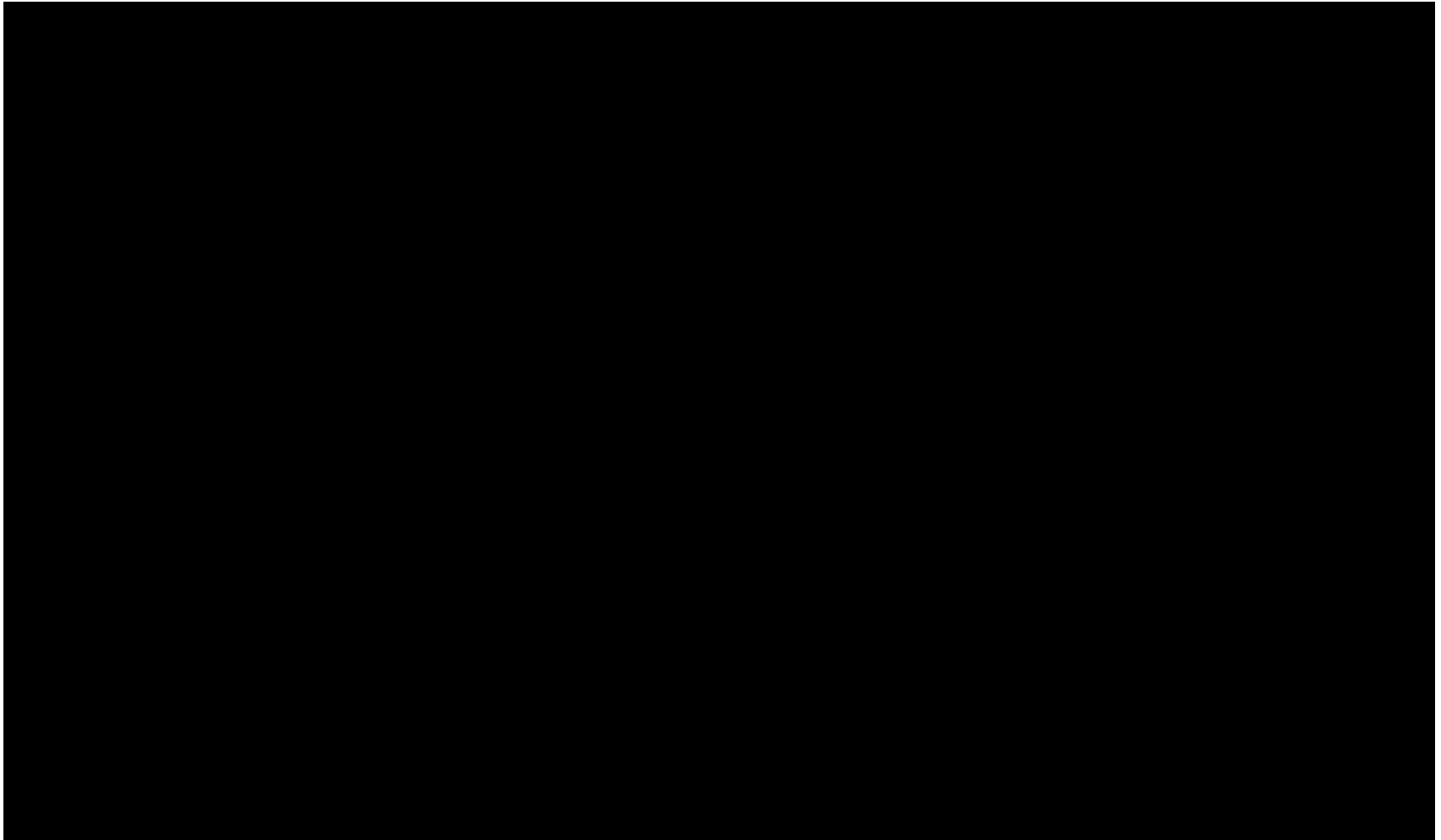




- Se consulta FIB incepand cu zona asociata celei mai specifice masti (/32)
- Se inspecteaza hash table-ul fiecarei zone
  - Pentru a mari viteza de cautare se folosesc contoare pentru a indica cat de des a fost folosita o ruta; in cadrul unei liste din hash table nodurile se tin ordonate descrescator dupa acest contor
- Daca nu se face match in hash table-ul unei zone, se trece la urmatoarea
- Daca nu se face match pe nici o ruta se face drop pe pachet



- Rutarea fiecarui pachet prin FIB este extrem de lenta atunci cand tabela de rutare nu este foarte mica
- Pentru a mari viteza se foloseste un cache de rute
  - Un hash table cu 256 de intrari
  - Fiecare intrare contine o lista cu noduri ce descriu
    - un triplet <adresa sursa, adresa destinatie, TOS>
    - informatii necesare de rutere al acestui tip de pachet: adresele de nivel 2 si 3, interfata de iesire
  - Nodurile sunt ordonate descrescator dupa un contor ce specifica cat de des a fost folosit acesta



- Se calculeaza indexul pentru hash table
  - O combinatie intre adresa sursa, adresa destinatie si TOS
- Se parcurge lista asociata indexului
- Daca nu se gaseste nici o potrivire se ruteaza prin FIB
  - Dupa determinarea informatiilor necesare rutarii se introduce o intrare in cache-ul de rute
- Daca se gaseste o potrivire se actualizeaza contorul nodului, si se reordoneaza lista
- Intrarile din cache-ul de rute sunt invalidate periodic

- Infrastructura folosită în Linux pentru implementarea facilităților de firewall & NAT
- Bazat pe hook-uri:
  - NF\_IP\_PRE\_ROUTING
  - NF\_IP\_LOCAL\_IN
  - NF\_IP\_FORWARD
  - NF\_IP\_LOCAL\_OUT
  - NF\_IP\_POST\_ROUTING
  - NF\_IP\_NUMHOOKS



```

struct sk_buff {
    struct sk_buff      *next;
    struct sk_buff      *prev;

    struct sock         *sk;
    ktime_t             timestamp;
    struct net_device   *dev;
    char                cb[48];

    unsigned int        len,
                       data_len;
    __u16               mac_len,
                       hdr_len;

    void                (*destructor)(struct sk_buff *skb);

    sk_buff_data_t      transport_header;
    sk_buff_data_t      network_header;
    sk_buff_data_t      mac_header;
    sk_buff_data_t      tail;
    sk_buff_data_t      end;

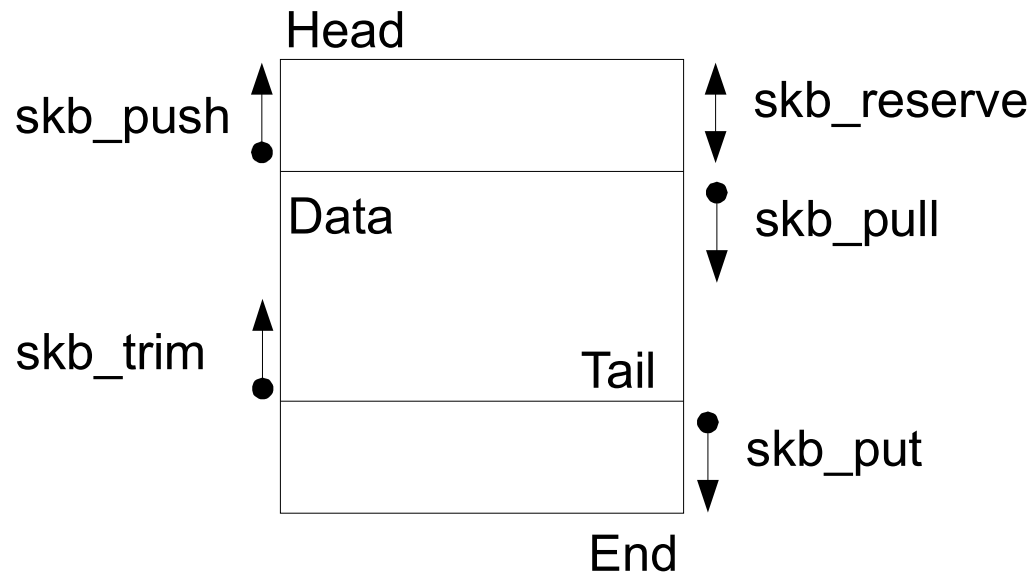
    unsigned char       *head,
                       *data;
    unsigned int        truesize;
    atomic_t            users;
};

```

```
unsigned char *skb_transport_header(const struct sk_buff *skb);  
void skb_reset_transport_header(struct sk_buff *skb);  
void skb_set_transport_header(struct sk_buff *skb, const int offset);  
unsigned char *skb_network_header(const struct sk_buff *skb);  
void skb_reset_network_header(struct sk_buff *skb);  
void skb_set_network_header(struct sk_buff *skb, const int offset);  
unsigned char *skb_mac_header(const struct sk_buff *skb);  
int skb_mac_header_was_set(const struct sk_buff *skb);  
void skb_reset_mac_header(struct sk_buff *skb);  
void skb_set_mac_header(struct sk_buff *skb, const int offset);
```

```

/* reserve head room */
void skb_reserve(struct sk_buff *skb, int len);
/* add data to the end */
unsigned char *skb_put(struct sk_buff *skb, unsigned int len);
/* add data to the top */
unsigned char *skb_push(struct sk_buff *skb, unsigned int len);
/* discard data at the top */
unsigned char *skb_pull(struct sk_buff *skb, unsigned int len);
/* discard data at the end */
unsigned char *skb_trim(struct sk_buff *skb, unsigned int len);
  
```





- Operații

- `int hard_start_xmit(struct sk_buff *skb, struct net_device *dev);`
- `int open(struct net_device *dev);`
- `int stop(struct net_device *dev);`
- `int do_ioctl(struct net_device *dev, struct ifreq *ifr, int cmd);`
- `int change_mtu(struct net_device *dev, int new_mtu);`
- `int create(struct sk_buff *skb, struct net_device *dev, unsigned short type, const void *daddr, const void *saddr, unsigned len);`
- `int rebuild(struct sk_buff *skb);`

- **Recepția pachetelor se face din intrerupere**

- `int netif_rx(struct sk_buff *skb);`
- `int netif_receive_skb(struct sk_buff skb); /* NAPI */`

- Flags

- IFF\_UP, IFF\_BROADCAST, IFF\_RUNNING, etc.

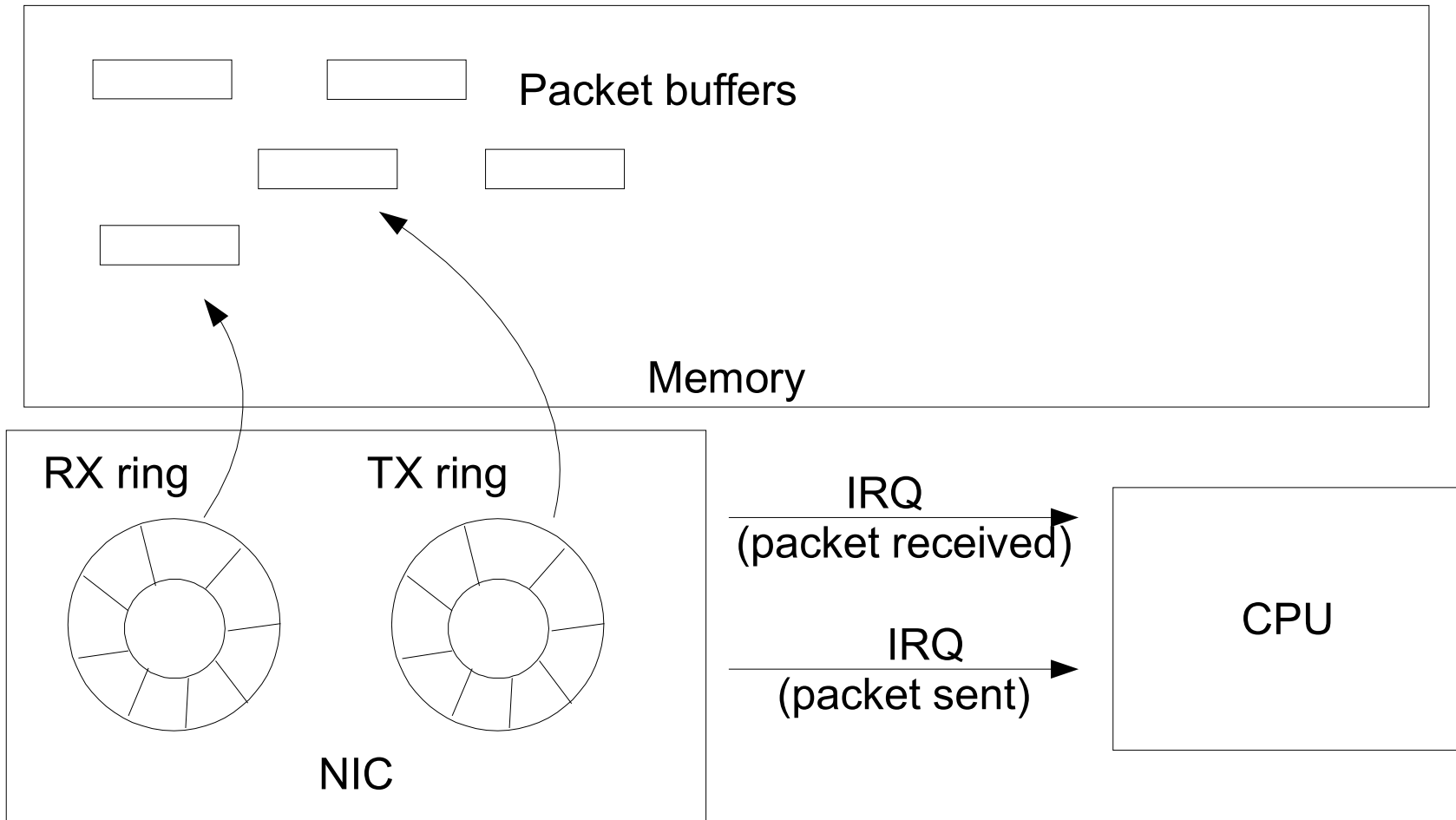
- Features

- NETIF\_F\_SG, NETIF\_F\_CSUM, NETIF\_F\_GSO, NETIF\_F\_TSO, etc.

```
dev_ixianet = alloc_netdev(sizeof(struct net_device_stats),
                          "ixint0",
                          ixnam_llm_probe);

if(dev_ixianet)
{
    SET_MODULE_OWNER(dev_ixianet);

    ret = register_netdev(dev_ixianet);
}
else
{
    ret = -ENOMEM;
}
```



- RX ring – umplut de SO cu buffere pentru pachetele recepționate
- TX ring – actualizat de SO pentru fiecare pachet de trimis
- Transferul pachetelor din/în memorie se face de către placa de rețea
- Procesorul este întrerupt doar cand
  - S-a recepționat un pachet
  - Placa de rețea confirmă trimiterea unui pachet

- Se face în context softirq
- În context IRQ
  - La RX: se face doar extragerea din ringul RX si plasarea pachetului în coada de procesare a socket-ului
  - La TX: se face (eventual) doar eliberarea memoriei pentru un pachet trimis

- Se așteaptă un număr de pachete sau un timeout înainte de a trimite întreruperea
- Se reduce costul procesării pe RX atunci când pachetele vin în burst
- Crește latența și jitter-ul

- Același scop ca și interrupt coalescing, dar fără a influența jitter-ul sau latența
- Nu necesită suport hardware
- În momentul generării unei întreruperi
  - Se dezactivează întreruperea plăcii de rețea
  - Se procesează toate pachetele
  - Abia apoi se activează întreruperea la loc

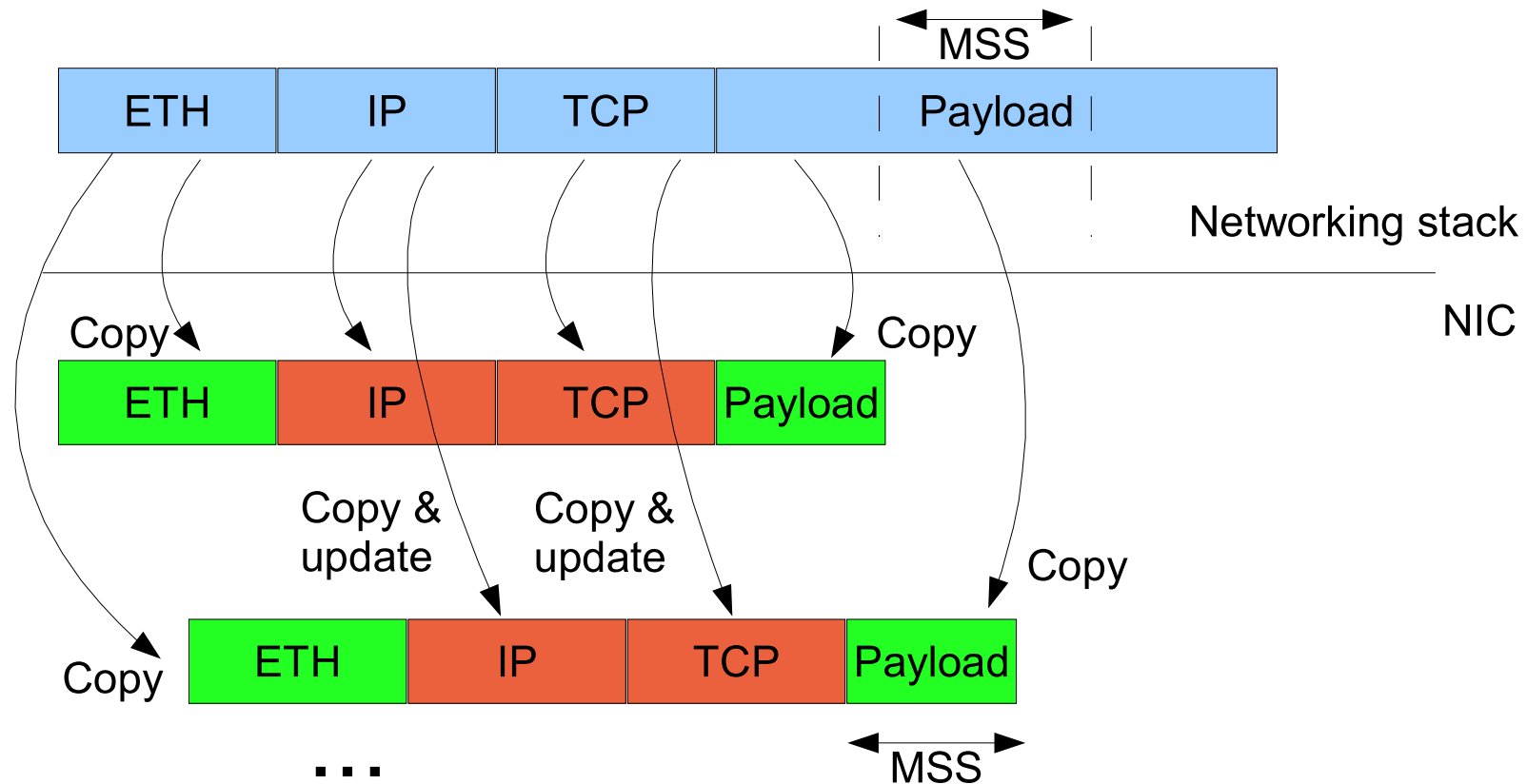
- Un pachet poate fi format din zone ne-contigue de memorie
- Descriptorii de TX/RX mai au un bit: end-of-packet
- Necesari pentru zero-copying



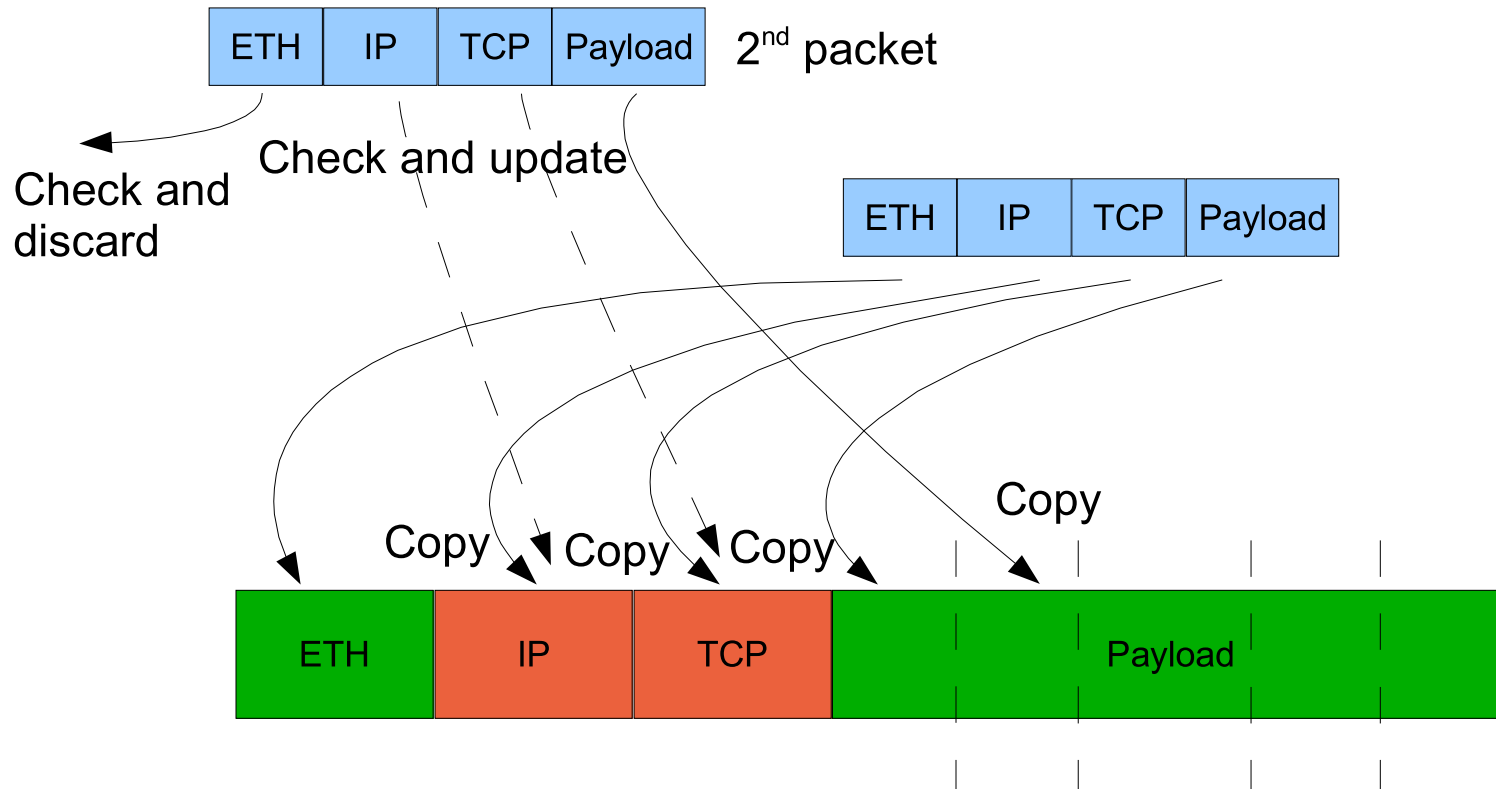
- Checksum-ul pe pachetele TCP/UDP se efectuează de către placa de rețea
- Placa de rețea Primește pseudo-header checksum-ul și poziția checksum-ului în header-ul TCP
- Necesari pentru zero-copying

- Full offload - Implementarea stivei TCP/IP direct în placa de rețea
- Probleme:
  - Scalabilitate
  - Securitate
  - Conformanță
- Performanța stivei de rețea este direct proporțională cu numărul de pachete procesate pe secundă (doar headerele sunt procesate)

- Exemplu: dacă un end-point poate procesa 60K pachete pe secundă, atunci pentru următoarele dimensiuni de pachete obținem:
  - 1538 -> 738Mbps
  - 2038 -> 978Mbps
  - 9038 -> 4.3Gbps
  - 20738 -> 9.9Gbps
- Stateless offload
  - Stiva folosește pachete mari
  - TSO (Transmit/TCP Segmentation Offload): hardware-ul primește pachete mari și le sparge în pachete mai mici pe care apoi le trimite pe fir
  - LRO (Large Receive Offload): hardware-ul agreghează pachetele mici în pachete mari pe care le trimite apoi la stivă



- În hardware – de către placa de rețea
- În software – înainte de trimiterea către device driver (GSO – generic segmentation offload)



- Dacă se recepționează pachete ale aceluiași segment în ordine ele se vor agrega și către stiva se va trimite un pachet mare (până la 64KB)
  - Agregarea se poate face în hardware sau software
  - Pentru a reduce latența, în momentul în care se iesi din bulca NAPI se face flush

- ```
long splice(int fd_in, loff_t *off_in, int fd_out,  
           loff_t *off_out, size_t len, unsigned int flags);
```
- Se mută date de la fd\_in la fd\_out, fără a se face copierea dacă este posibil
- Unul din file descriptori trebuie să fie neapărat un pipe
- Offset-ul pentru pipe trebuie să fie NULL
- Linux 2.6.25 oferă support complet pentru fișiere, socketi și memorie (via vmsplice)

- Controller pentru copierea din memorie în memorie via DMA
- Folosit la copierea din/în user-space a datelor socketilor

?

