

## Systemy Operacyjne 2

### Obsługa sieci w Linuksie

Arkadiusz Chrobot

Katedra Systemów Informatycznych

20 czerwca 2024

1 / 17

## Plan

Wprowadzenie

Stos TCP/IP

Sterowniki urządzeń sieciowych

Filtr sieciowy

2 / 17

## Wprowadzenie

Unix jest jednym z pierwszych systemów operacyjnych, które oferowały implementację komunikacji sieciowej. Współcześnie większość serwerów w Internecie pracuje pod kontrolą Linuksa, systemu operacyjnego, który jest kompatybilny z Uniksem. Ten wykład jest krótkim przeglądem budowy i działania podsystemu jądra Linuksa odpowiedzialnego za obsługę sieci. Tematyka ta jest złożona, dlatego tylko najważniejsze zagadnienia są tutaj zaprezentowane. Treść wykładu jest podzielona na trzy części:

- ▶ przetwarzanie pakietu na poziomie jądra,
- ▶ sterowniki urządzeń sieciowych,
- ▶ implementacja filtra sieciowego.

3 / 17

## Stos TCP/IP

Część jądra odpowiedzialna za przetwarzanie nadchodzących i wychodzących pakietów (ang. *packets*) sieciowych jest nazywana stosem TCP/IP. Rysunek 1 przedstawia jej schemat<sup>1</sup>. Podsystem jądra Linuksa odpowiedzialny za komunikację sieciową składa się z trzech części, które odpowiadają trzem warstwom modelu ISO/OSI — warstwie łącza danych, warstwie sieciowej i warstwie transportowej. Aby wysłać dane przez sieć proces użytkownika wywołuje odpowiednie wywołanie systemowe, które aktywuje metodę `write()` obiektu pliku związanego z gniazdem sieciowym tego procesu. W zależności od użytego protokołu transportowego ta metoda wywołuje funkcję jądra `tcp_sendmsg()` lub `udp_sendmsg()`. Są one odpowiedzialne za utworzenie nagłówka odpowiedniego protokołu.

<sup>1</sup>Wykład został przygotowany na podstawie strony: <https://docs.kernel.org/networking/index.html> i książki Williama Stallinsa „Systemy operacyjne”, PWN, Warszawa, 2009

4 / 17

Notatki

Notatki

Notatki



## Sterowniki urządzeń sieciowych

Główną strukturą danych używaną przez sterownik urządzenia sieciowego jest struktura typu `struct net_device`. Najważniejszymi polami tej struktury są: `mtu` — określa maksymalny rozmiar pakietu jaki urządzenie może obsłużyć, `flags` — określa stan urządzenia, `dev_addr` — wskaźnik na adres MAC, `promiscuity` — licznik, który przechowuje liczbę żądań przełączenia urządzenia w tryb bezładny (ang. *promiscuous mode*), `ip_ptr` — wskazuje na nagłówek przechowujący dane specyficzne dla protokołu IP w wersji 4, `netdev_ops` — wskazuje na strukturę wskaźników na funkcje, które wykonują takie operacje, jak np. wysyłanie pakietów, `rx_handler` — wskazuje na procedurę obsługi przerwania nadajnika.

Wcześniejsze implementacje sterowników urządzeń sieciowych wymagały, aby to urządzenia potwierdzało odebranie każdego pakietu przerwaniem. Prowadziło to do przeciążenia systemu w przypadku dużego ruchu sieciowego. W serii 2.6 jądra dodano nowe API sterowników urządzeń sieciowych, nazwane NAPI.

9 / 17

Notatki

---

---

---

---

---

---

---

---

---

---

## Sterowniki urządzeń sieciowych

NAPI umożliwia sterownikowi przełączenie urządzenia sieciowego w tryb przeglądania (ang. *polling*) pozwalający mu na zgromadzenie odpowiedniej liczby przychodzących pakietów, które są przetwarzane później przez jądro. To redukuje liczbę zgłaszanych przez urządzenie przerw i w konsekwencji zmniejsza obciążenie jądra ich obsługą. Niektóre z przychodzących pakietów mogą nawet być odrzucone zanim zostaną przekazane jądro do przetwarzania. To rozwiązanie nazywa się *dławieniem pakietów* (ang. *packet throttling*). Aby użycie NAPI było możliwe wymagane jest wsparcie sprzętowe w postaci tzw. cyklicznego bufora dla transmisji DMA (ang. *DMA ring*) lub odpowiednio duża wolna przestrzeń w RAM, aby można było ją przydzielić dla buforów dla tych transmisji.

10 / 17

Notatki

---

---

---

---

---

---

---

---

---

---

## Filtr sieciowy

Filtr sieciowy (ang. *netfilter*) jest zbiorem wskaźników na funkcje, nazywanych *uchwyty* (ang. *hooks*), rozlokowanych w strategicznych miejscach stosu TCP/IP. Te wskaźniki mogą być użyte do zaimplementowania zapory sieciowej (ang. *firewall*) lub takich rozwiązań jak NAT (ang. *Network Address Translation*). Funkcje wskazywane przez uchwyty są zazwyczaj definiowane w modułach jądra<sup>2</sup>. W podsystemie sieciowym jądra znajduje się pięć takich uchwytów:

**NF\_IP\_PRE\_ROUTING** funkcje związane z tym uchwytem są wywoływane, gdy odbierany jest pakiet,

**NF\_IP\_LOCAL\_IN** funkcje związane z tym uchwytem przeprowadzają wstępne przetwarzanie pakietów dostarczonych do urządzenia sieciowego (ang. *host*),

**NF\_IP\_FORWARD** funkcje związane z tym uchwytem przeprowadzają przetwarzanie pakietów, które powinny być przekazane do innych urządzeń sieciowych,

<sup>2</sup><http://www.paulkiddie.com/creating-a-netfilter-kernel-module-which-filters-udp-packets>

11 / 17

Notatki

---

---

---

---

---

---

---

---

---

---

## Filtr sieciowy

**NF\_IP\_POST\_ROUTING** funkcje związane z tym uchwytem przetwarzają pakiety z ustalonymi trasami, które mają zostać wysłane,

**NF\_IP\_LOCAL\_OUT** funkcje związane z tym uchwytem przetwarzają pakiety, które zostały wysłane lokalnie.

Każda funkcja związana z dowolnym z tych uchwytów może przeprowadzać dowolną operację na pakiecie, jaka jest konieczna, ale musi ostatecznie zwrócić jedną z następujących wartości:

**NF\_ACCEPT** pakiet zaakceptowany do dalszego przetwarzania,

**NF\_DROP** pakiet został odrzucony,

**NF\_REPEAT** należy ponownie wywołać funkcję dla tego pakietu,

**NF\_STOLEN** funkcja „wykrada” pakiet, co oznacza, że będzie on przetwarzany w inny sposób niż pozostałe pakiety,

12 / 17

Notatki

---

---

---

---

---

---

---

---

---

---

## Filtr sieciowy

`NF_QUEUE` pakiet jest umieszczony w kolejce, skąd zostanie przekazany do przestrzeni użytkownika,

`NF_STOP` przetwarzanie pakietu zostało wstrzymane.

Pojedyncza funkcja związana z uchwycem jest reprezentowana przez strukturę typu `struct nf_hook_ops`. Definicja tego typu jest zamieszczona w Listingu nr 1. Pole `list` pozwala przechowywać te struktury w liście. Pole `hook` jest wskaźnikiem na funkcję przetwarzającą pakiet. Pole `dev` jest wskaźnikiem na strukturę reprezentującą urządzenie sieciowe. Składowa `priv` jest wskaźnikiem na obszar pamięci przechowujący dane prywatne funkcji przetwarzającej pakiety. Pole `pf` przechowuje identyfikator rodziny protokołów, której pakiety będą przetwarzane przez funkcję. Pole `hooknum` przechowuje numer uchwytu, a `priority` priorytet funkcji, który określa porządek w jakim funkcje związane z uchwycem są wykonywane (np. stała `NF_IP_PRI_FIRST` oznacza najwyższy priorytet).

13 / 17

Notatki

---

---

---

---

---

---

---

---

## Filtr sieciowy

```
1 struct nf_hook_ops
2 {
3     struct list_head    list;
4     nf_hookfn           *hook;
5     struct net_device   *dev;
6     void                *priv;
7     u_int8_t            pf;
8     unsigned int        hooknum;
9     int                 priority;
10 };
```

Listing 1: Definicja typu `struct nf_hook_ops`

14 / 17

Notatki

---

---

---

---

---

---

---

---

## Filtr sieciowy

Struktury typu `struct nf_hook_ops` są rejestrowane przy użyciu funkcji `nf_register_net_hook()`, a wyrejestrowywane przy pomocy funkcji `nf_unregister_net_hook()`. Typ wartości zwracanej przez funkcję przetwarzającą pakiet to `unsigned int`<sup>3</sup>. Przyjmuje ona trzy argumenty: adres obszaru pamięci przechowującego jej prywatne dane (jest przekazywany przez parametr typu `void *`), adres bufora pakietu (bufor jest typu `struct sk_buff`) i adres struktury, która przechowuje stan uchwytu. Ta struktura jest typu `struct nf_hook_state`.

<sup>3</sup>Na wcześniejszych slajdach opisano możliwe wartości do zwrócenia.

15 / 17

Notatki

---

---

---

---

---

---

---

---

## Pytania

?

Notatki

---

---

---

---

---

---

---

---

16 / 17

Dziękuję Państwu za uwagę!

Notatki

---

---

---

---

---

---

---

---

Notatki

---

---

---

---

---

---

---

---

Notatki

---

---

---

---

---

---

---

---

Notatki

---

---

---

---

---

---

---

---