

Vypracované otázky k MSZ pro rok 2022

Specializace NNET

1. května 2022

Vladimír Dušek, xdusek27

Specializace Počítačové sítě – NNET

1. Architektura superskalárních procesorů a algoritmy zpracování instrukcí mimo pořadí, predikce skoků.
2. Paměťová konzistence a předbíhání operací čtení a zápisu, podpora virtuálního adresového prostoru.
3. Datový paralelismus SIMD, HW implementace a SW podpora.
4. Architektury se sdílenou pamětí UMA a NUMA, zajištění lokality dat.
5. Problém koherence pamětí cache na systémech se sdílenou pamětí, protokol MSI.
6. Paralelní zpracování v OpenMP: Smyčky, sekce a tasky a synchronizační prostředky.
7. Pravděpodobnost, podmíněná pravděpodobnost, nezávislost.
8. Náhodná proměnná, typy náhodné proměnná, funkční a číselné charakteristiky, významná rozdělení pravděpodobnosti.
9. Bodové a intervalové odhady parametrů, testování hypotéz o parametrech.
10. Vícevýběrové testy, testy o rozdělení, testy dobré shody.
11. Regresní analýza.
12. Markovské řetězce a základní techniky pro jejich analýzu.
13. Randomizované algoritmy (Monte Carlo a Las Vegas algoritmy).
14. Problém generalizace strojového učení a přístup k jeho řešení (trénovací, validační a testovací sada, regularizace, předtrénování, multi-task learning, augmentace dat, dropout, ...)
15. Generativní modely a diskriminativní přístup ke klasifikaci (gaussovský klasifikátor, logistická regrese, ...)
16. Neuronové sítě a jejich trénování (metoda gradientního sestupu, účelová (loss) funkce, výpočetní graf, aktivační funkce, zápis pomocí maticového násobení, ...)
17. Neuronové sítě pro strukturovaná data (konvoluční a rekurentní sítě, motivace, základní vlastnosti, použití)
18. Prohledávání stavového prostoru (informované a neinformované metody, lokální prohledávání, prohledávání v nejistém prostředí, hraní her, CSP úlohy)
19. Klasifikace formálních jazyků (Chomského hierarchie), vlastnosti formálních jazyků a jejich rozhodnutelnost.
20. Konečné automaty (jazyky přijímané KA, varianty KA, minimalizace KA, Mihill-Nerodova věta).
21. Regulární množiny, regulární výrazy a rovnice nad regulárními výrazy.
22. Zásobníkové automaty (jazyky přijímané ZA, varianty ZA).
23. Turingovy stroje (jazyky přijímané TS, varianty TS, lineárně omezené automaty, vyčíslitelné funkce).
24. Nerozhodnutelnost (problém zastavení TS, princip diagonalizace a redukce, Postův korespondenční problém).
25. Časová a paměťová složitost (třídy složitosti, úplnost, SAT problém).
26. Postrelační a rozšířené relační databáze (objektový a objektově relační databázový model – struktura a operace; podpora práce s XML a JSON dokumenty v databázích).
27. NoSQL databáze (porovnání relačních a NoSQL; CAP věta a ACID/BASE principy; typy NoSQL databází; dotazování v NoSQL databázích; agregace dat pomocí Map-Reduce a agregační pipeline).
28. Získávání znalostí z dat (pojem znalost; typické zdroje dat; základní úlohy získávání znalostí; analytické projekty a proces získávání znalostí z dat).

29. Porozumění datům (důvod a cíl; popisné charakteristiky dat a vizualizační techniky; korelační analýza).
30. Prostorové DB (problematika mapování prostoru, ukládání, indexace; využití).
31. Indexace (nejen) v prostorových DB (kD-Tree a Grid File (a jejich varianty), R-Tree).
32. Lambda kalkul (definice všech pojmu, operací...).
33. Práce v lambda kalkulu (demonstrace reprezentace čísel a pravdivostních hodnot a operací nad nimi).
34. Haskell – lazy evaluation (typy v jazyce včetně akcí, uživatelské typy, význam typových tříd, demonstrace lazy evaluation).
35. Prolog – způsob vyhodnocení (základní princip, unifikace, chování vestavěných predikátů, operátor řezu – vhodné a nevhodné užití).
36. Prolog – změna DB/programu za běhu (demonstrace na prohledávání stavového prostoru, práce se seznamy).
37. Model PRAM, suma prefixů a její aplikace.
38. Distribuované a paralelní algoritmy – algoritmy nad seznamy, stromy a grafy.
39. Interakce mezi procesy a typické problémy paralelismu (synchronizační a komunikační mechanismy).
40. Distribuované a paralelní algoritmy – předávání zpráv a knihovny pro paralelní zpracování (MPI).
41. Distribuovaný broadcast, synchronizace v distribuovaných systémech.
42. Klasifikace a vlastnosti paralelních a distribuovaných architektur, základní typy jejich topologií.
43. Distribuované a paralelní algoritmy – algoritmy řazení, select, algoritmy vyhledávání.
44. Bezdrátové lokální sítě (Wifi, Bluetooth).
45. Hledání minimální kostry obyčejného grafu (pojmy, stromy a kostry, Kruskalův algoritmus, Primův algoritmus).
46. Hledání nejkratších cest ze zdrojového uzlu do všech ostatních uzel grafu (Bellman-Fordův algoritmus, Dijkstrův algoritmus).
47. Klasifikace algoritmů volby koordinátora, algoritmus Bully a jeho složitost.
48. Podmínky konsistentního globálního stavu distribuovaného systému.
49. Principy distribuovaného zpracování MapReduce, průběh a jednotlivé operace distribuovaného výpočtu pomocí MapReduce, jeho implementace v Apache Hadoop a Apache Spark.
50. Symetrická kryptografie. Vlastnosti, vlastnosti bezpečného algoritmu, délka klíče, útok silou, příklady symetrických algoritmů, Feistelovy šifry, DES, režimy činnosti, proudové šifry.
51. Asymetrická kryptografie, vlastnosti, způsoby použití, poskytované bezpečnostní funkce, elektronický podpis a jeho vlastnosti, hybridní kryptografie, algoritmus RSA, generování klíčů, šifrování, dešifrování.
52. Hašovací funkce, klíčovaný haš a MAC a jejich použití a vlastnosti.
53. Správa klíčů v asymetrické kryptografii (certifikáty X.509).
54. Základní architektury přepínačů, algoritmy pro plánování, řešení blokování, vícestupňové přepínací sítě.
55. Základní funkce směrovače, zpracování paketů ve směrovači, typy přepínání a architektur.
56. Metody pro výpočet směrování v sítích (Bellman-Ford, Dijkstra, Path vector, DUAL).
57. Řízení toku dat (flow-control) a prevence zahlcení (congestion-control) na transportní vrstvě (MP-TCP, QUIC, SCTP, DCCP).
58. Metody detekce síťových incidentů (signatury, statistické metody) a nástroje (IDS/IPS).
59. Sítě Peer-to-Peer: vlastnosti, chování, způsoby směrování. Strukturované a nestrukturované sítě.

60. Události v JavaScriptu (smyčka událostí, asynchronní programování, klientské události, obsluha událostí)

61. Přenos a distribuce webových dat (URI, protokol HTTP, proxy HTTP, CDN, XHR)

62. Bezpečnost webových aplikací (SOP, XSS, CSRF, bezpečnostní hlavičky HTTP)

Obsah

| | | |
|----|---|-----|
| 1 | BMS – Bezdrátové lokální sítě (Wifi, Bluetooth). | 5 |
| 2 | GAL – Hledání minimální kostry obyčejného grafu (pojmy, stromy a kostry, Kruskalův algoritmus, Primův algoritmus). | 6 |
| 3 | GAL – Hledání nejkratších cest ze zdrojového uzlu do všech ostatních uzel grafu (Bellman-Fordův algoritmus, Dijkstrův algoritmus). | 16 |
| 4 | PDI – Klasifikace algoritmů volby koordinátora, algoritmus Bully a jeho složitost. | 25 |
| 5 | PDI – Podmínky konsistentního globálního stavu distribuovaného systému. | 31 |
| 6 | PDI – Principy distribuovaného zpracování MapReduce, průběh a jednotlivé operace distribuovaného výpočtu pomocí MapReduce, jeho implementace v Apache Hadoop a Apache Spark. | 35 |
| 7 | KRY – Symetrická kryptografie. Vlastnosti, vlastnosti bezpečného algoritmu, délka klíče, útok silou, příklady symetrických algoritmů, Feistelovy šifry, DES, režimy činnosti, proudové šifry. | 42 |
| 8 | KRY – Asymetrická kryptografie, vlastnosti, způsoby použití, poskytované bezpečnostní funkce, elektronický podpis a jeho vlastnosti, hybridní kryptografie, algoritmus RSA, generování klíčů, šifrování, dešifrování. | 56 |
| 9 | KRY – Hašovací funkce, klíčovaný haš a MAC a jejich použití a vlastnosti. | 62 |
| 10 | KRY – Správa klíčů v asymetrické kryptografii (certifikáty X.509). | 69 |
| 11 | PDS – Základní architektury přepínačů, algoritmy pro plánování, řešení blokování, vícestupňové přepínací sítě. | 75 |
| 12 | PDS – Základní funkce směrovače, zpracování paketů ve směrovači, typy přepínání a architektur. | 92 |
| 13 | PDS – Metody pro výpočet směrování v sítích (Bellman-Ford, Dijkstra, Path vector, DUAL). | 108 |
| 14 | PDS – Řízení toku dat (flow-control) a prevence zahlcení (congestion-control) na transportní vrstvě (MP-TCP, QUIC, SCTP, DCCP). | 115 |

Kapitola 1

BMS – Bezdrátové lokální sítě (Wifi, Bluetooth).

1.1 Metadata

- Předmět: Bezdrátové a mobilní sítě (BMS)
- Přednáška:
 - [[todo]]
- Záznam:
 - [[todo]]

1.2 Úvod a kontext

[[todo]]

Kapitola 2

GAL – Hledání minimální kostry obyčejného grafu (pojmy, stromy a kostry, Kruskalův algoritmus, Primův algoritmus).

2.1 Metadata

- Předmět: Grafové algoritmy (GAL)
- Přednáška:
 - 5) Stromy, minimální kostry, Jarníkův a Borůvkův algoritmus.
 - 6) Růst minimální kostry, algoritmy Kruskala a Prima.
- Záznam:
 - 2020-10-22
 - 2020-10-29

2.2 Úvod a kontext

Orientovaný graf Orientovaný graf je dvojice $G = (V, E)$, kde V je konečná množina uzlů a $E \subseteq V \times V$ je množina hran.

Neorientovaný graf Neorientovaný graf je dvojice $G = (V, E)$, kde V je konečná množina uzlů a $E \subseteq \binom{V}{2}$ je množina hran. (Hrana je tedy dvouprvková množina, avšak běžně se držíme stejného značení jako u orientovaných grafů a používáme dvojici.)

Ohodnocený graf Ohodnocený graf je takový graf, jehož každá hrana má přiřazenou nějakou hodnotu, typicky definovanou pomocí váhové funkce $w : E \mapsto \mathbb{R}$.

Podgraf Graf $G' = (V', E')$ je podgraf grafu $G = (V, E)$ jestliže $V' \subseteq V$ a $E' \subseteq E$.

Sled Posloupnost uzlů $\langle v_0, v_1, \dots, v_k \rangle$, kde $(v_{i-1}, v_i) \in E$ pro $i = 1, \dots, k$ se nazývá sled délky k z v_0 do v_k .

Uzavřený sled Sled $\langle v_0, v_1, \dots, v_k \rangle$ se nazývá uzavřený, pokud existuje hrana (v_0, v_k) .

Dosažitelnost Pokud existuje sled s z uzlu u do uzlu v , říkáme, že v je dosažitelný z u sledem s , značeno $u \xrightarrow{s} v$.

Tah Tah je sled ve kterém se neopakují hrany.

Cesta Cesta je sled ve kterém se neopakují uzly.

Souvislý graf Neorientovaný graf se nazývá souvislý, pokud mezi libovolnými dvěma uzly existuje cesta.

Kružnice Uzavřená cesta se nazývá kružnice.

Cyklus Orientovaná kružnice se nazývá cyklus (první a poslední uzel je shodný).

Prostý graf Orientovaný graf bez cyklů se nazývá prostý.

Acyklický graf Graf je bez cyklů, resp. kružnic, se nazývá acyklický.

Strom Graf, který je souvislý a acyklický, se nazývá strom.

Kostra Strom, který tvoří podgraf souvislého grafu na množině všech jeho vrcholů, se nazývá kostra (*spanning tree*).

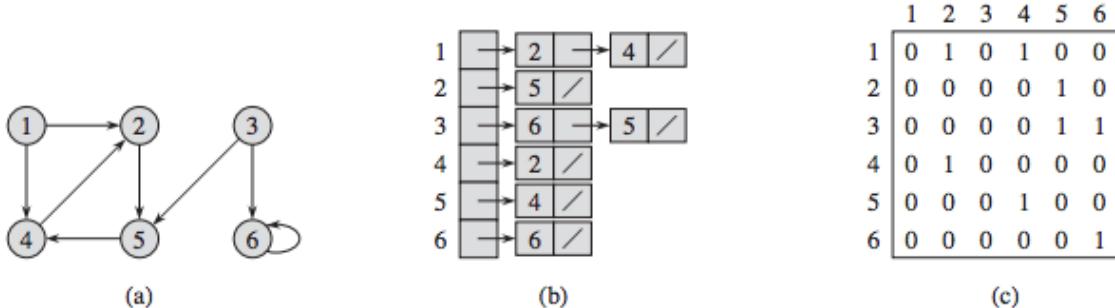
Minimální kostra Necht' $G = (V, E)$ je souvislý neorientovaný graf s váhou funkcí $w : E \mapsto \mathbb{R}$. Minimální kostra (*MST, minimum spanning tree*) je strom $G' = (V, E')$, kde $E' \subseteq E$ a

$$w(E') = \sum_{(u,v) \in T} w(u, v)$$

je minimální ze všech možných alternativních kostér.

Seznam sousedů Seznam sousedů (*Adj, adjacency list*) je reprezentace grafu v paměti. Jde o preferovanou variantu pro řídké grafy – kde $m \ll n^2$. Pro každý uzel máme definovaný seznam jeho sousedů.

Matice sousednosti Matice sousednosti (*adjacency matrix*) je reprezentace grafu v paměti. Jde o preferovanou variantu pro husté grafy – kde m je skoro n^2 .



Obrázek 2.1: Příklad reprezentace grafu pomocí seznamu sousedů a matice sousednosti.

2.3 Generický algoritmus

Hledání minimální kostry je problém, který lze řešit algoritmy, které spadají do kategorie tzv. hladových (*greedy*) deterministických algoritmů. Spočívají v tom, že průběžně odhadují kostru přidáváním dalších hran a nikdy se nemusejí vracet (neprovádí se *backtracking*). Generický algoritmus tvoří jakousi základní kostru pro další, už konkrétní, algoritmy.

Řez Necht' $G = (V, E)$ je graf. Řez grafu G je dvojice $(S, V - S)$, kde $\emptyset \subseteq S \subseteq V$.

Křížení Hrana $(u, v) \in E$ kříží řez $(S, V - S)$, pokud jeden její konec je v S a druhý v $V - S$.

Respektování Necht' $A \subseteq E$ je množina hran. Řez $(S, V - S)$ respektuje množinu hran A , pokud žádná hrana v A nekříží řez $(S, V - S)$.

Lehkost Necht' $(S, V - S)$ je řez a B je množina hran, která ho kříží. Hrana z množiny B s nejmenší hodnotou se nazývá lehká.

Bezpečnost Necht' $G = (V, E)$ je souvislý neorientovaný graf s reálnou váhovou funkcí w . Necht' $A \subseteq E$ je součástí nějaké minimální kostry G . Necht' $(S, V - S)$ je řez, který respektuje A . Necht' (u, v) je lehká hrana křížící $(S, V - S)$. Pak hrana (u, v) je bezpečná pro A .

```

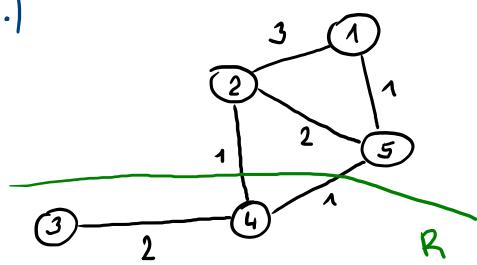
1 def generic_mst(G):
2     # G je graf
3     A = {}# A je mnozina hran rozpracovane minimalni kostry
4     while netvorí_kostru(A, G):
5         for hrana in G.E:
6             if je_bezpecna(A, hrana):
7                 A += {hrana}
8

```

Výpis 2.1: Generický algoritmus. Před každou iterací algoritmu je množina A podmnožinou nějaké minimální kostry. Hrana $(u, v) \in E$ je bezpečná pro A , pokud $A \cup \{(u, v)\}$ je podmnožinou nějaké minimální kostry.

Př. $G = (V, E)$ $\psi: E \rightarrow \mathbb{R}$ (definované obrazem)

1.)



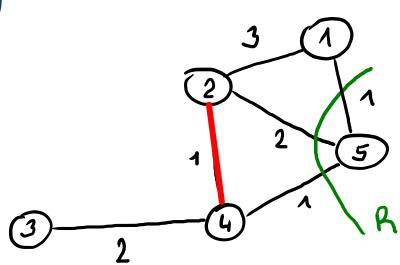
$$A = \{\emptyset\}$$

$$R = (\{\{1, 2, 5\}, \{3, 4\}\})$$

$$LH = \{(2, 4), (4, 5)\}$$

$$A \leftarrow A \cup \{(2, 4)\}$$

2.)



$$A = \{(2, 4)\}$$

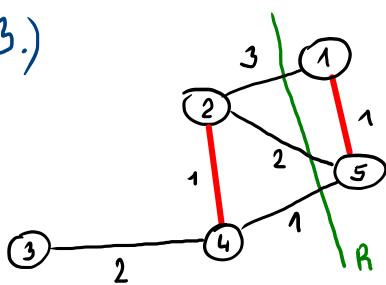
$$R = (\{\{1, 2, 3, 4\}, \{5\}\})$$

$$LH = \{(1, 5), (4, 5)\}$$

$$A \leftarrow A \cup \{(1, 5)\}$$

Obrázek 2.2: Příklad, část 1.

3.)



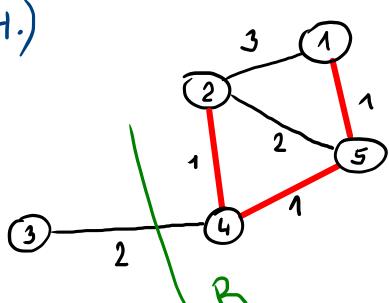
$$A = \{(1, 5), (2, 4)\}$$

$$R = (\{\{1, 5\}, \{2, 3, 4\}\})$$

$$LH = \{(4, 5)\}$$

$$A \leftarrow A \cup \{(4, 5)\}$$

4.)



$$A = \{(1, 5), (2, 4), (4, 5)\}$$

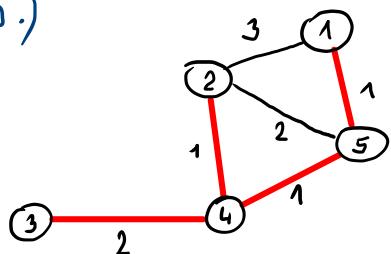
$$R = (\{\{1, 2, 4, 5\}, \{3\}\})$$

$$LH = \{(3, 4)\}$$

$$A \leftarrow A \cup \{(3, 4)\}$$

Obrázek 2.3: Příklad, část 2.

5.)



$$A = \{(3,4), (2,4), (5,4), (1,5)\}$$

A je minimální kostra

*(minimál už nemůže učítat
věz, když by vyspetoval A)*

Obrázek 2.4: Příklad, část 3.

2.4 Kruskalův algoritmus

Kruskalův a Primův algoritmus se liší v tom, jakým způsobem vybírají bezpečnou hranu. Kruskalův algoritmus nahlíží na A jako na les a hledá hranu s nejmenším ohodnocením, která spojuje stromy v lese. Na konci je A jeden strom.

```

1 def kruskal_mst(G):
2     # G je graf
3
4     # inicializace, kazdy uzel je ve sve mnozine
5     A = {}# A je mnozina hran rozpracovane minimalni kostry
6     for v in G.V:
7         make_set(v)
8
9     # seradit vzestupne podle w
10    E = sort(G.E, G.w)
11
12    for (u, v) in E:
13        if find_set(u) != find_set(v):
14            A += {(u, v)}
15            union(u, v)
16
17    return A

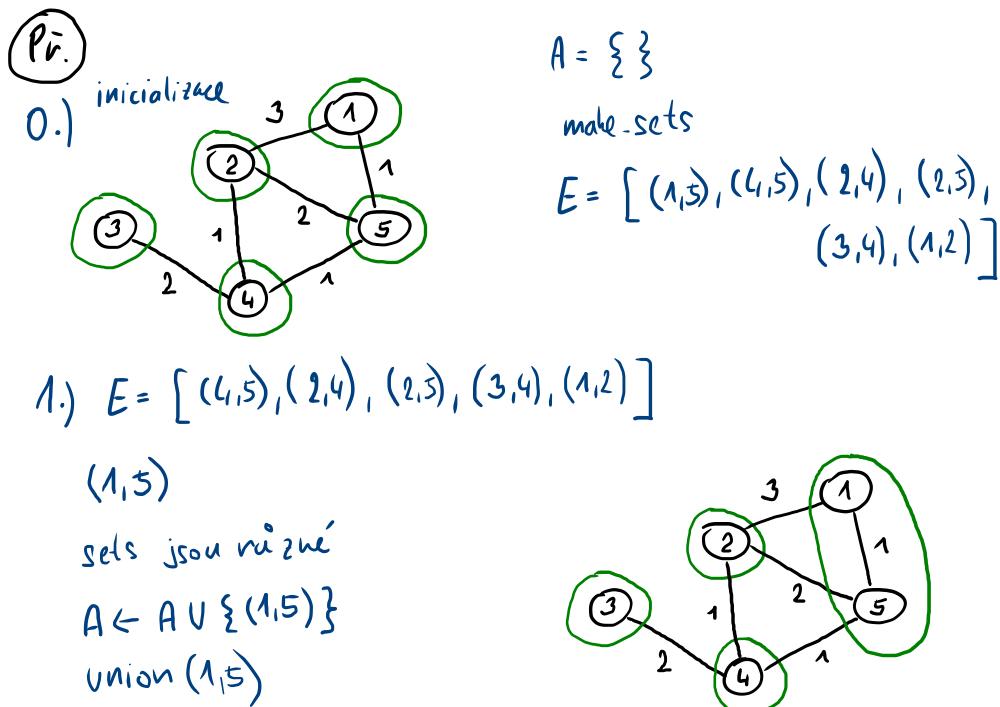
```

Výpis 2.2: Kruskalův algoritmus. Funkce `make_set(v)` vytvoří množinu obsahující v , `find_set(v)` vrátí reprezentanta množiny ve které se nachází v , `union(u, v)` sjednotí dvě množiny obsahující u a v .

2.4.1 Složitost

- Řádek 5 – $O(1)$
- Řádek 6-7 – n -krát složitost *make_set* (n je počet uzlů).
- Řádek 10 – $O(m \cdot \log(m))$ (m je počet hran).
- Řádky 12-15 – Závisí na implementaci *find_set* a *union*.
 - Při implementaci seznamem s heuristickou celkem: $O(m + n \cdot \log(n))$.
 - Při stromové implementaci s váhami a zkratkami celkem: $O((m + n) \cdot \alpha(n))$.
Kde α je velmi pomalu rostoucí funkce ($\alpha \leq 4$).
- Pro souvislý graf platí $m > n$. Proto množinové operace stojí $O(m \cdot \alpha(n))$. Jelikož $\alpha(n) = O(\log(n)) = O(\log(m))$, tak celková složitost je $O(m \cdot \log(m))$.
- Dále platí $m < n^2$, pak $\log(m) = O(\log(n))$, proto celkem: $O(m \cdot \log(n))$.

2.4.2 Příklad



Obrázek 2.5: Příklad, část 1.

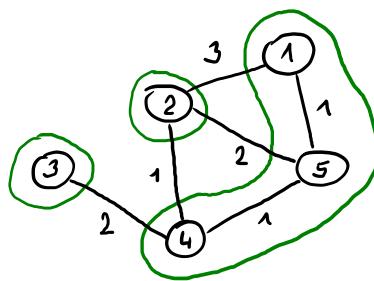
$$2.) E = [(2,4), (2,5), (3,4), (1,2)]$$

(4,5)

sets jsou různé

$$A \leftarrow A \cup \{(4,5)\}$$

union (4,5)



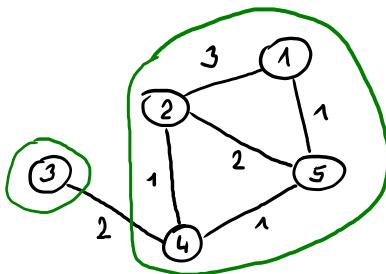
$$3.) E = [(2,5), (3,4), (1,2)]$$

(2,4)

sets jsou různé

$$A \leftarrow A \cup \{(2,4)\}$$

union (2,4)

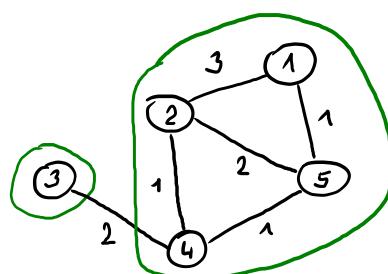


Obrázek 2.6: Příklad, část 2.

$$4.) E = [(3,4), (1,2)]$$

(2,5)

sets nejsou různé



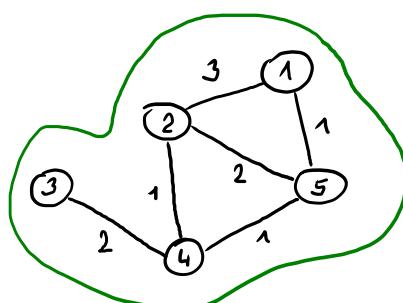
$$4.) E = [(1,2)]$$

(3,4)

sets jsou různé

$$A \leftarrow A \cup \{(3,4)\}$$

union (3,4)



$$5.) E = [] , (1,2) , \text{ sets nejsou různé}$$

Obrázek 2.7: Příklad, část 3.

2.5 Primův-Jarníkův algoritmus

Primův algoritmus buduje tzv. A strom. Má zadaný určitý uzel, ze kterého hledá nejbližší další uzel, který by připojil. A pak další a další.

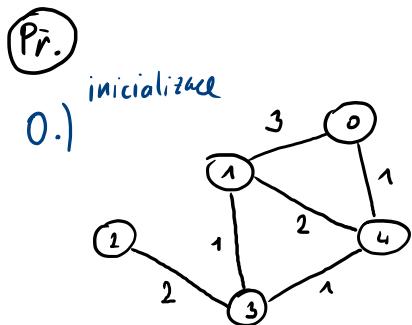
```
1 def prim_mst(G, r):
2     # G je graf
3     # r je výchozí uzel
4
5     for u in G.V:
6         key[u] = INF # pole cen prechodu, kolik stojí prechod do vrcholu na indexu
7         pi[u] = NULL # pole predchudcu, kdo je predchudce vrcholu na indexu
8
9     key[r] = 0
10    Q = Queue(G.V) # prioritní fronta uzlu
11
12    while not Q.empty():
13        u = Q.extract_min(key) # vrati prvek z Q s nejmensi hodnotou v key
14
15        # pro vsechny sousedy uzlu u (Adj je seznam sousedu)
16        for v in Adj[u]:
17            # pokud je levnejsi cesta a jeste to není prozkoumaný uzel
18            if v in Q and w(u, v) < key(v):
19                pi[v] = u
20                key[v] = w(u, v)
21                Q.decrease_key(key) # aktualizace prioritni fronty
22
23    return pi
```

Výpis 2.3: Primův algoritmus.

2.5.1 Složitost

- Řádky 5-10 – $O(n)$ za použití binární haldy (n je počet uzlů).
- Řádky 12-13 – While cyklus se provede n -krát a protože $extract_min$ stojí $O(\log(n))$, tak je celková složitost $O(n \cdot \log(n))$.
- Řádek 16 – For cyklus se provede $O(m)$ krát, protože délka všech seznamů sousedů je dohromady $2m$ (m je počet hran).
- Řádek 18-20 – $O(1)$.
- Řádek 21 – $O(\log(n))$.
- Jelikož $m > n$, tak celkem $O(n \cdot \log(n) + m \cdot \log(n)) = O(m \cdot \log(n))$.

2.5.2 Příklad



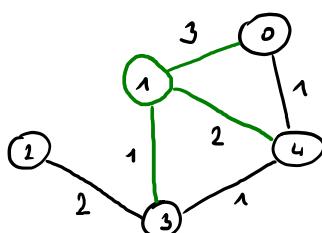
$$n = 1$$

$$\text{key} = [\infty, 0, \infty, \infty, \infty]$$

$$\pi = [\text{NULL}, \text{NULL}, \text{NULL}, \text{NULL}, \text{NULL}]$$

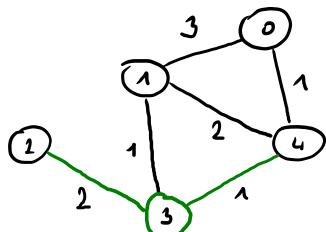
$$Q = [0, 1, 2, 3, 4]$$

1.) $u = 1$
 $Q = [0, 1, 3, 4]$
 $v \in \{0, 3, 4\}$
 $\text{key} = [3, 0, \infty, 1, 2]$
 $\pi = [1, \text{NULL}, \text{NULL}, 1, 1]$

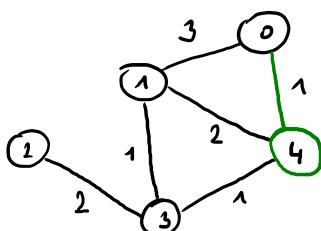


Obrázek 2.8: Příklad, část 1.

2.) $u = 3$
 $Q = [0, 1, 4]$
 $v \in \{1, 2, 4\}$
 $\text{key} = [3, 0, 2, 1, 1]$
 $\pi = [1, \text{NULL}, 3, 1, 3]$



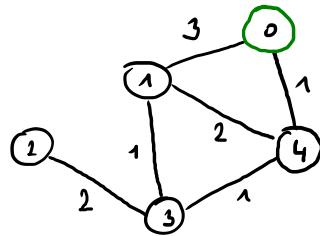
3.) $u = 4$
 $Q = [0, 1]$
 $v \in \{0, 1, 3\}$
 $\text{key} = [1, 0, 2, 1, 1]$
 $\pi = [4, \text{NULL}, 3, 1, 3]$



Obrázek 2.9: Příklad, část 2.

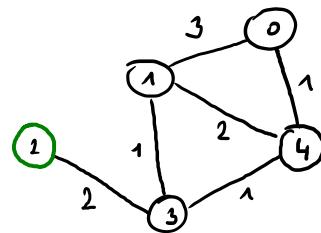
4.) $u = 0$
 $Q = [1]$
 $v \in \{1, 4\}$

$key = [1, 0, 2, 1, 1]$
 $\pi = [4, \text{NULL}, 3, 1, 3]$



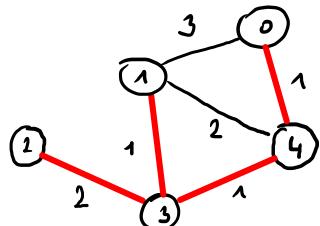
5.) $u = 2$
 $Q = []$
 $v \in \{3\}$

$key = [1, 0, 2, 1, 1]$
 $\pi = [4, \text{NULL}, 3, 1, 3]$



Obrázek 2.10: Příklad, část 3.

6.) $key = [1, 0, 2, 1, 1]$
 $\pi = [4, \text{NULL}, 3, 1, 3]$



Obrázek 2.11: Příklad, část 4.

Kapitola 3

GAL – Hledání nejkratších cest ze zdrojového uzlu do všech ostatních uzlů grafu (Bellman-Fordův algoritmus, Dijkstrův algoritmus).

3.1 Metadata

- Předmět: Grafové algoritmy (GAL)
- Přednáška:
 - 7) Nejkratší cesty z jednoho vrcholu, Bellman-Fordův algoritmus, nejkratší cesta z jednoho vrcholu v orientovaných acyklických grafech.
 - 8) Dijkstrův algoritmus. Nejkratší cesty ze všech vrcholů.
- Záznam:
 - 2020-11-05

3.2 Úvod a kontext

Viz. „Úvod a kontext“ v předchozích otázkách z tohoto předmětu.

Cena cesty Necht' $G = (V, E)$ je ohodnocený graf s váhovou funkcí $w : E \mapsto \mathbb{R}$. Cena cesty $p = \langle v_0, v_1, \dots, v_k \rangle$ je suma

$$w(p) = \sum_{i=0}^k w(v_i, v_{i+1})$$

Cena nejkratší cesty Cena nejkratší cesty z u do v je

$$\delta(u, v) = \begin{cases} \min(\{w(p) : u \xrightarrow{p} v\}) \\ \infty \text{ pokud cesta neexistuje} \end{cases}$$

Nejkratší cesta Nejkratší cesta z u do v je pak libovolná cesta p taková, že $w(p) = \delta(u, v)$.

Cena cesty se záporným cyklem Pokud na cestě z u do v existuje záporný cyklus (cyklus jehož celková cena je záporná), pak $\delta(u, v) = -\infty$.

Záporné ohodnocení hran Pokud na cestě z u do v neexistuje záporný cyklus, tak algoritmy pracují dobře i se záporným ohodnocením hran.

Reprezentace cesty Cestu reprezentujeme pomocí pole předchůdců π .

Hledání nejkratších cest ze všech uzlů do jednoho Tento problém lze řešit stejnými algoritmy. Graf se transponuje (převrácení orientace hran), provede se algoritmus pro problém „hledání nejkratších cest ze jednoho uzlu do všech ostatních uzlů“ a poté se transponuje zpět.

Reprezentace nejkratší cesty Nejkratší cestu grafu $G = (V, E)$ reprezentujeme pomocí pole předchůdců π , kde $\pi[v]$ označuje předchůdce uzlu $v \in V$ na nejkratší cestě. Podgraf předchůdců pak je $G_\pi = (V_\pi, E_\pi)$, $V_\pi = \{v \in V : \pi[v] \neq \text{NULL}\} \cup \{s\}$, $E_\pi = \{(\pi[v], v) \in E : v \in V_\pi - \{s\}\}$. V okamžiku dokončení algoritmu výpočtu nejkratších cest je G_π strom nejkratších cest. Tj. kořenový strom obsahující nejkratší cesty ze zdroje s do všech ostatních uzlů.

3.3 Pomocné funkce

Představené algoritmy pracují z důvodu efektivity se sledy a nikoliv s cestami (bylo by nutné stále kontrolovat, zda nebyla porušena podmínka cesty), ačkoliv je problém nazývá hledání nejkratší cesty.

```

1 def initialize_single_source(G, s):
2     # G je graf
3     # s je vychozí uzel
4     for v in G.V:
5         d[v] = INF # d je pole vzdalenosti
6         pi[v] = NULL # pi je pole predchudcu
7         d[s] = 0

```

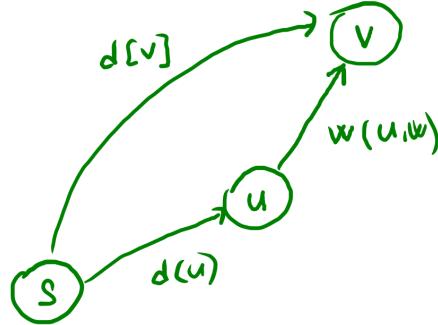
Výpis 3.1: Pomocná inicializační funkce. Složitost je $\Theta(n)$, kde n je počet uzlů.

```

1 def relax(u, v, w):
2     # u a v jsou uzly grafu
3     # w je vahova funkce
4     if d[v] > d[u] + w(u, v):
5         d[v] = d[u] + w(u, v)
6         pi[v] = u

```

Výpis 3.2: Pomocná funkce *relax*. Složitost je $O(1)$.



Obrázek 3.1: Ukázka činnosti funkce *relax*.

3.4 Bellman-Fordův algoritmus

Slouží pro řešení v obecných grafech, mohou obsahovat cykly a záporné hrany. Záporné cykly je však nutné detektovat a vrátit specifickou hodnotu. V podstatě se jedná o *brute force* algoritmus, provede se relaxace $n - 1$ -krát pro každou hranu.

```

1 def bellman_ford(G, s, w):
2     # G je graf
3     # s je vychozi uzel
4     # w je vahova funkce
5
6     # faze inicializace
7     initialize_single_source(G, s)
8     n = len(G.V) # pocet uzlu
9
10    # faze relaxace: provedeni (n-1) * m relaxaci (m je pocet hran)
11    for _ in range(0, n-1):
12        for u, v in G.E:
13            relax(u, v, w)
14
15    # faze detekce zaporneho cyklu
16    for u, v in G.E:
17        if d[u] > d[v] + w(u, v):
18            return NULL
19
20    return pi

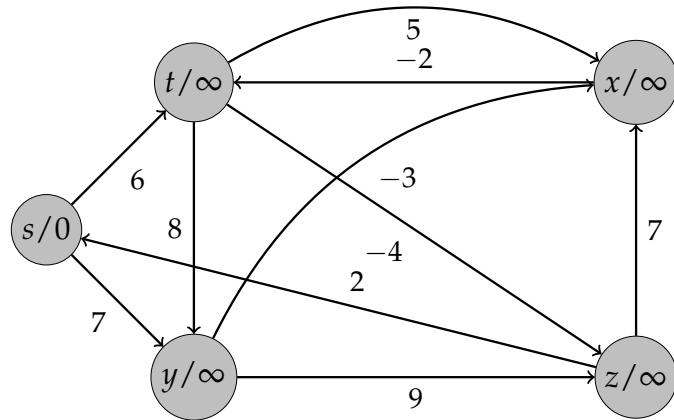
```

Výpis 3.3: Algoritmus Bellman-Ford. Proč $n - 1$ iterací? Protože mezi libovolnými dvěma uzly v grafu, existuje cesta o maximálním počtu hran $n - 1$.

3.4.1 Složitost

- Řádek 7, 8 – $\Theta(1)$.
- Řádky 11, 12, 13 – $(n - 1) \cdot \Theta(m) = \Theta(n \cdot m)$, kde n je počet uzelů a m je počet hran grafu.
- Řádek 16, 17, 18 – $\Theta(m)$.
- Celkem $\Theta(n \cdot m)$.

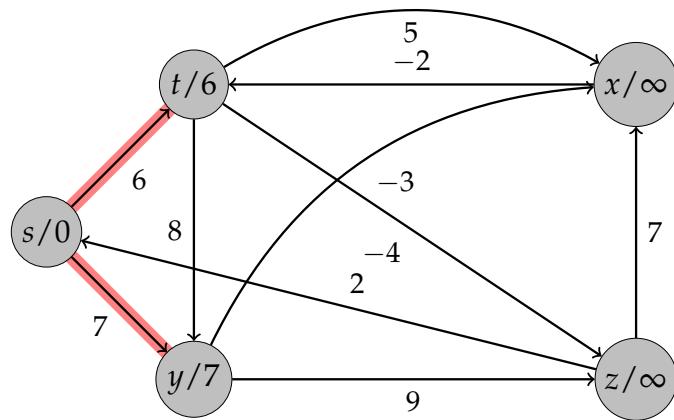
3.4.2 Příklad



Obrázek: Práce algoritmu Bellman-Ford.

- ▶ Pokud $(u, v) \in E$ je označená, pak $\pi[v] = u$
- ▶ Hrany se relaxují v tomto pořadí:
 $(t, x), (t, y), (t, z), (x, t), (y, x), (y, z), (z, x), (z, s), (s, t), (s, y)$.

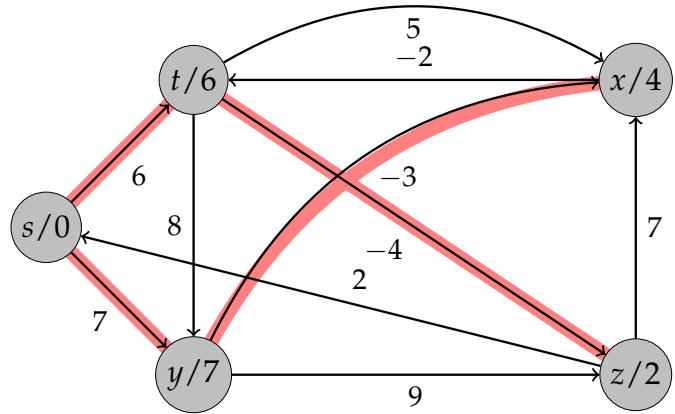
Obrázek 3.2: Příklad, část 1.



Obrázek: Práce algoritmu Bellman-Ford.

- ▶ Pokud $(u, v) \in E$ je označená, pak $\pi[v] = u$
- ▶ Hrany se relaxují v tomto pořadí:
 $(t, x), (t, y), (t, z), (x, t), (y, x), (y, z), (z, x), (z, s), (s, t), (s, y)$.

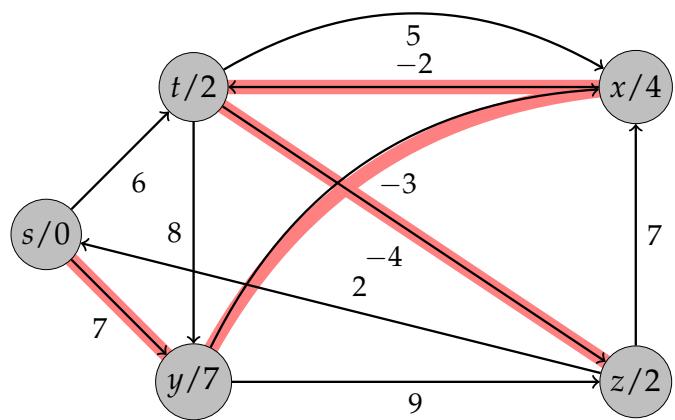
Obrázek 3.3: Příklad, část 2.



Obrázek: Práce algoritmu Bellman-Ford.

- ▶ Pokud $(u, v) \in E$ je označená, pak $\pi[v] = u$
- ▶ Hrany se relaxují v tomto pořadí:
 $(t, x), (t, y), (t, z), (x, t), (y, x), (y, z), (z, x), (z, s), (s, t), (s, y)$.

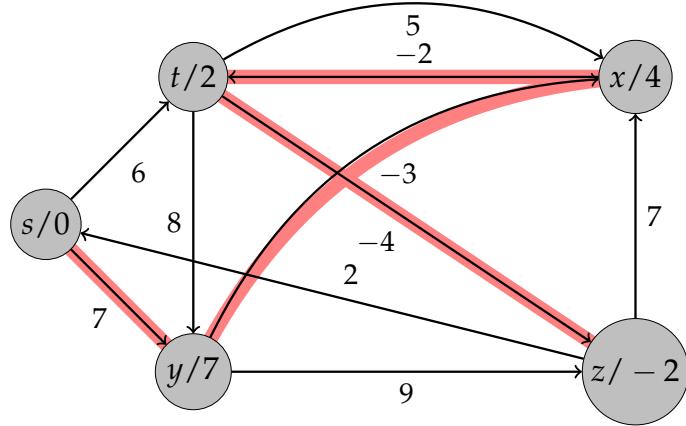
Obrázek 3.4: Příklad, část 3.



Obrázek: Práce algoritmu Bellman-Ford.

- ▶ Pokud $(u, v) \in E$ je označená, pak $\pi[v] = u$
- ▶ Hrany se relaxují v tomto pořadí:
 $(t, x), (t, y), (t, z), (x, t), (y, x), (y, z), (z, x), (z, s), (s, t), (s, y)$.

Obrázek 3.5: Příklad, část 4.



Obrázek: Práce algoritmu Bellman-Ford.

- ▶ Pokud $(u, v) \in E$ je označená, pak $\pi[v] = u$
- ▶ Hrany se relaxují v tomto pořadí:
 $(t, x), (t, y), (t, z), (x, t), (y, x), (y, z), (z, x), (z, s), (s, t), (s, y)$

Obrázek 3.6: Příklad, část 5.

3.5 Dijkstrův algoritmus

Slouží pro řešení v acyklických grafech bez záporných hran. Pro takto omezený problém existují rychlejší algoritmy než pro problém v obecných grafech.

```

1 def dijkstra(G, s, w):
2     # G je graf
3     # s je výchozí uzel
4     # w je vahová funkce
5
6     # fáze inicializace
7     initialize_single_source(G, s)
8     Q = Queue(G.V) # prioritní fronta uzlu
9     S = {}# množina uzlu, která už byla prozkoumána
10
11    # fáze relaxace
12    while not Q.empty():
13        u = Q.extract_min(d) # vrátí prvek z Q s nejménší hodnotou v d
14        S += {u}
15        # pro všechny sousedy uzlu u (Adj je seznam sousedů)
16        for v in Adj[u]:
17            relax(u, v, w)
18
19        Q.decrease_key(d) # aktualizace prioritní fronty
20
21    return d, pi

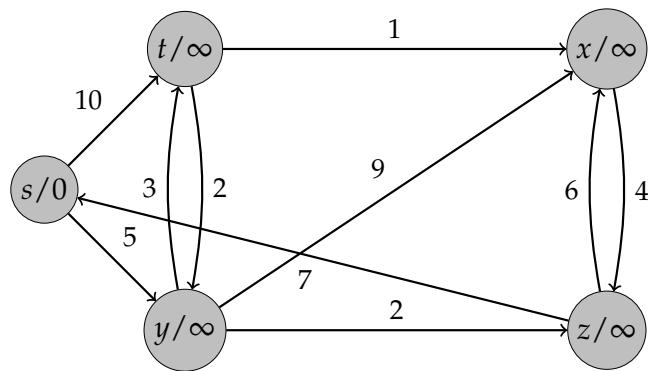
```

Výpis 3.4: Algoritmus Dijkstra.

3.5.1 Složitost

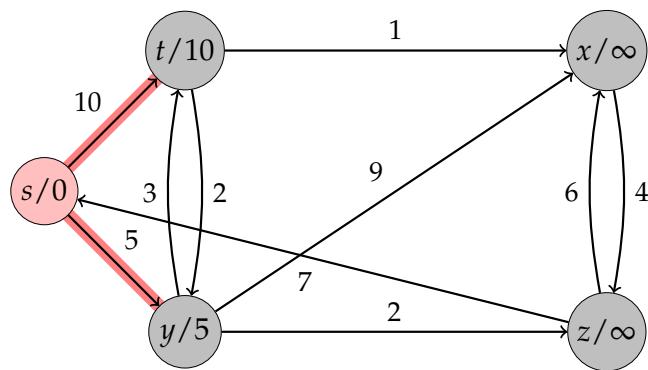
- Předpokládejme implementaci prioritní fronty pomocí pole.
- Řádek 8, 18 – $O(1)$.
- Řádek 11 – While cyklus se provede n -krát, kde n je počet uzlů.
- Řádek 12 – $O(n)$, najítí minima v poli uzlů. Celkově (s cyklem) $O(n^2)$.
- Řádek 16 – $O(m)$, pro všechny hrany. Celkově (s cyklem) $O(m \cdot n)$.
- Celkem $O(n^2 + m) = O(n^2)$.
- Pro řídké grafy lze využít implementaci fronty pomocí binární haldy a získat tak $O(m \cdot \log(n))$.
- Při implementaci fronty pomocí Fibonacciho haldy dostaneme časovou složitost $O(n \cdot \log(n) + m)$.

3.5.2 Příklad



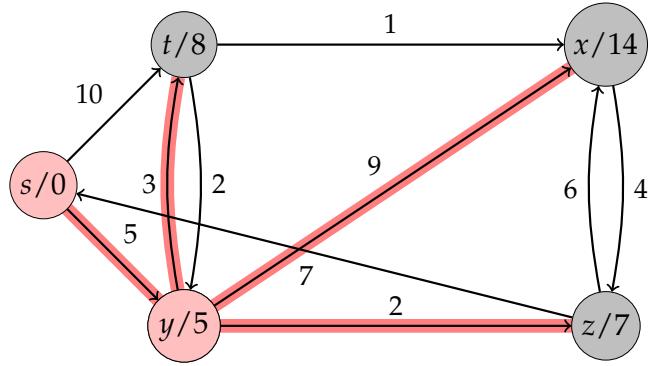
Obrázek: Práce Dijkstrova algoritmu. Označené uzly značí uzly z množiny S .

Obrázek 3.7: Příklad, část 1.



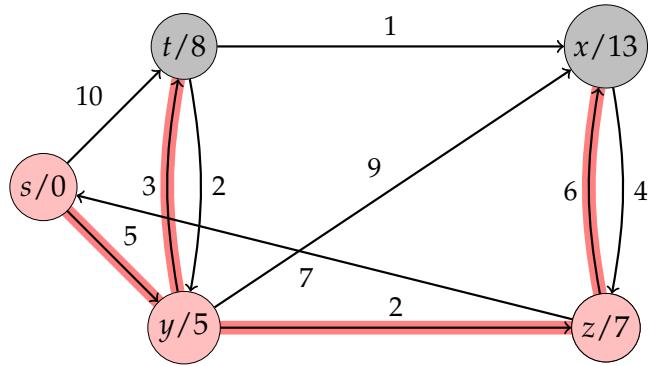
Obrázek: Práce Dijkstrova algoritmu. Označené uzly značí uzly z množiny S .

Obrázek 3.8: Příklad, část 2.



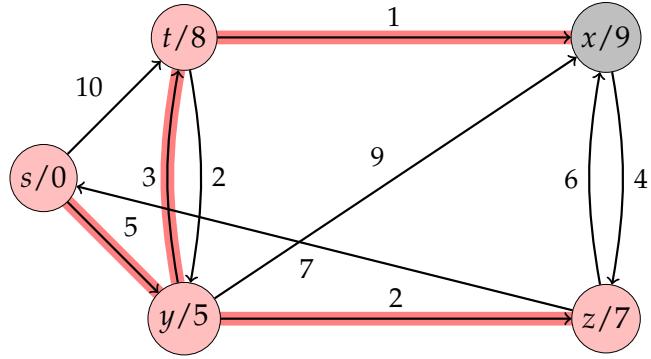
Obrázek: Práce Dijkstrova algoritmu. Označené uzly značí uzly z množiny S .

Obrázek 3.9: Příklad, část 3.



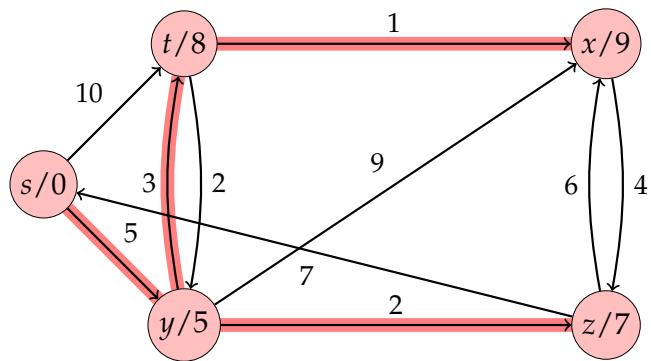
Obrázek: Práce Dijkstrova algoritmu. Označené uzly značí uzly z množiny S .

Obrázek 3.10: Příklad, část 4.



Obrázek: Práce Dijkstrova algoritmu. Označené uzly značí uzly z množiny S .

Obrázek 3.11: Příklad, část 5.



Obrázek: Práce Dijkstrova algoritmu. Označené uzly značí uzly z množiny S .

Obrázek 3.12: Příklad, část 6.

Kapitola 4

PDI – Klasifikace algoritmů volby koordinátora, algoritmus Bully a jeho složitost.

4.1 Metadata

- Předmět: Prostředí distribuovaných aplikací (PDI)
- Přednáška:
 - 7) Synchronizace
- Záznam:
 - 2020-11-02

4.2 Úvod a kontext

- Mějme množinu procesů v rámci distribuovaného systému. Řešíme problém nalezení shody na nějaké věci (synchronizační problém). Problém můžeme rozdělit na dvě situace:
 - **Problém volby koordinátora** – Výběr jednoho z procesů, který bude vedoucím procesem (koordinátor). Tento proces pak může vykonat určitou činnost nebo může sloužit ostatním procesům k realizaci význačné role v systému.
 - **Problém vzájemného vyloučení** – Předpokládejme, že konkrétní zdroj může v daném okamžiku používat pouze jeden proces. Tento problém se běžně vyskytuje ve víceprocesorových systémech, ale také v distribuovaných systémech.
- Synchronizační problémy lze v rámci operačních systémů nebo multiprocesorových systémů řešit pomocí provádění atomických operací, sdílené paměti apod. (je pro ně podpora v rámci operačního systému nebo hardwaru). V distribuovaných systémech nic takového není z principu možné a proto se synchronizační problémy řeší pomocí zasílání zprav, resp. algoritmicky.

4.3 Problém volby koordinátora

- Předpokládáme:

- Každý proces má unikátní ID.
 - Procesy neznají stav (běžící, neběžící) dalších procesů.
 - Každý proces zná ID dalších procesů (záleží na topologii).
- Cíl:
 - Dosáhnutí shody mezi všemi procesy na procesu, který je koordinátor.
 - Kritérium výběru koordinátora může být různé. Např. na základě proces ID (proces s největším ID se stane koordinátorem).

4.4 Bully algoritmus

Pro topologii každý s každým – každý proces může komunikovat s každým dalším procesem. Používá tři druhy zpráv: ELECTION, OK, COORDINATOR.

4.4.1 Postup

- Proces P, který má podezření, že chybí koordinátor, může zahájit volby.
 1. Proces P odešle zprávu ELECTION všem procesům s větším ID.
 2. Pokud nikdo neodpoví, P vyhrává volby a stává se koordinátorem.
 3. Pokud některý z procesů s větším ID odpoví (zpráva OK), tak přebírá řízení a práce P je ukončena.
 4. Pokud P obdrží zprávu ELECTION od procesů s menším ID, pošle jim odpověď OK na zablokování procesů.
- Nakonec zůstane pouze P (nový koordinátor), který o tom informuje ostatní zasláním zprávy COORDINATOR.
- Pokud se proces probudí nebo je restartován, první akcí je vyvolání voleb.

4.4.2 Složitost

Složitost z hlediska počtu zpráv.

Nejhorší případ (iniciátor s nejmenším ID):

- $(n - 1)$ iterací
- $2(n - 1)$ zpráv ELECTION a OK pro každou iteraci
- $(n - 1)$ zpráv COORDINATOR
- Celkem: $(n - 1) \times 2(n - 1) + (n - 1) \approx n^2$

Nejlepší případ (iniciátor s největším ID):

- $(n - 1)$ zpráv COORDINATOR
- Celkem: $(n - 1)$

4.4.3 Příklad

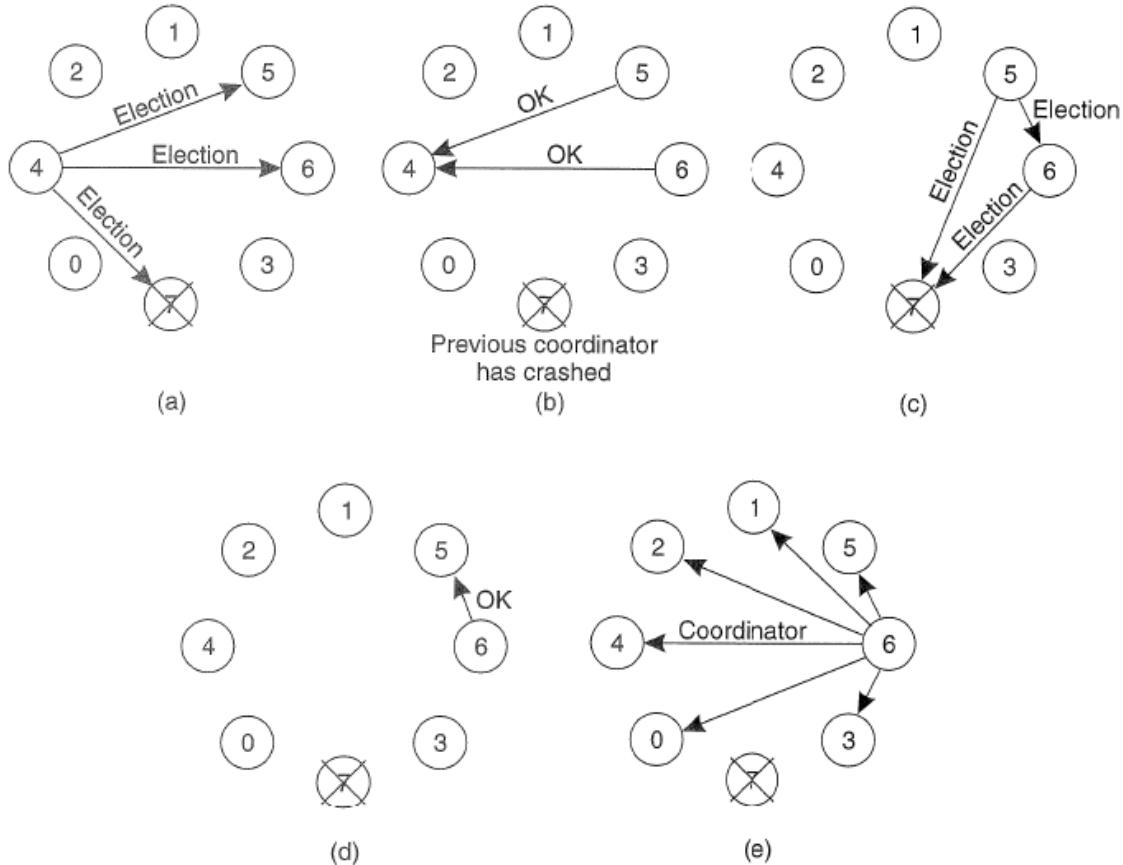


Figure 5-11. The bully election algorithm. (a) Process 4 holds an election. (b) Processes 5 and 6 respond, telling 4 to stop. (c) Now 5 and 6 each hold an election. (d) Process 6 tells 5 to stop. (e) Process 6 wins and tells everyone.

Obrázek 4.1: Příklad činnosti Bully algoritmu.

4.5 Ring Algoritmus

Pro kruhovou topologii – procesy jsou uspořádané do kruhu podle svého proces ID. Každý proces musí vědět nejenom o svém následovníkovi, ale také o jeho následníkovi, který funguje jako „záloha“, v případě že by se přímý následník stal nedostupný. Používá dva druhy zpráv: ELECTION, COORDINATOR.

4.5.1 Postup

- Proces P, který má podezření, že chybí koordinátor, může zahájit volby.
 1. Zašle zprávu ELECTION obsahující jeho ID dalšímu procesu (pokud další proces nereaguje, proces P zašle stejnou zprávu dalšímu v kruhu).
 2. Každý člen topologie přijme zprávu ELECTION, přidá do ní své ID a přepošle zprávu dalšímu procesu.

- Když se zpráva vrátí k procesu P, je zpráva převedena na zprávu COORDINATOR a poslána následujícímu procesu v topologii, aby bylo možné nahlásit:
 1. Novým koordinátorem se stává proces s nejvyšším ID.
 2. Členové sítě jsou stále aktivní.
- Po síti může obíhat více zpráv zároveň.

4.5.2 Složitost

Složitost z hlediska počtu zpráv.

Vždy $2n \approx n$ zpráv. Jedno kolečko „oběhne“ zpráva ELECTION a druhé zpráva COORDINATOR.

4.5.3 Příklad

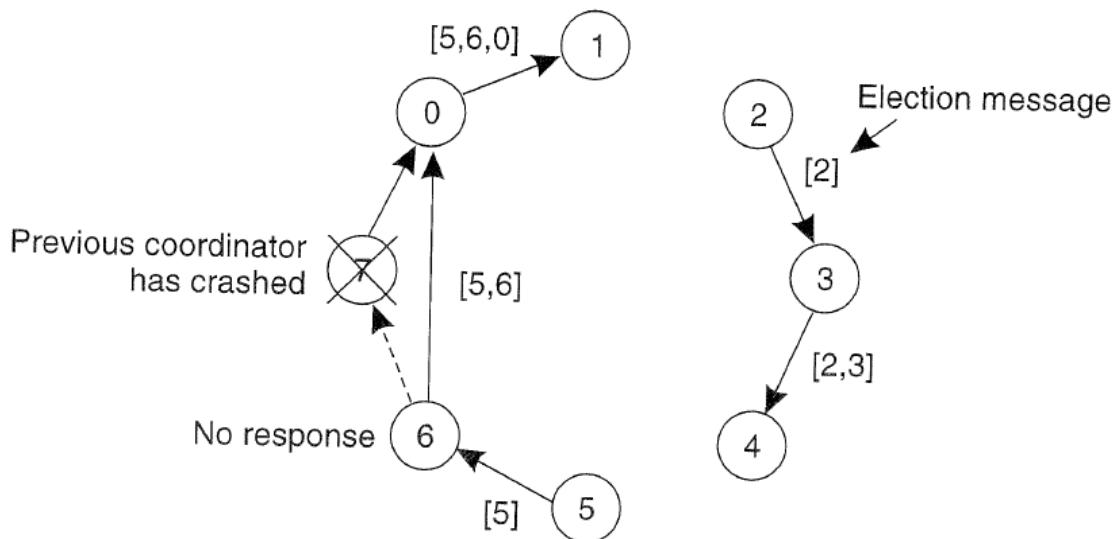


Figure 5-12. Election algorithm using a ring.

Obrázek 4.2: Příklad činnosti Ring algoritmu.

4.6 Algoritmus pro obecnou topologii

Předpokládáme, že nemáme ani kruhovou topologii ani spojení každý s každým. Např.: peer-to-peer sítě, sensorové sítě, ...

4.6.1 Postup

- V první iteraci se broadcastem posílá zpráva ELECTION.
- Každý uzel si uloží od kterého souseda dostal zprávu ELECTION jako první. Tím vzníká kostra grafu (*spanning tree*).

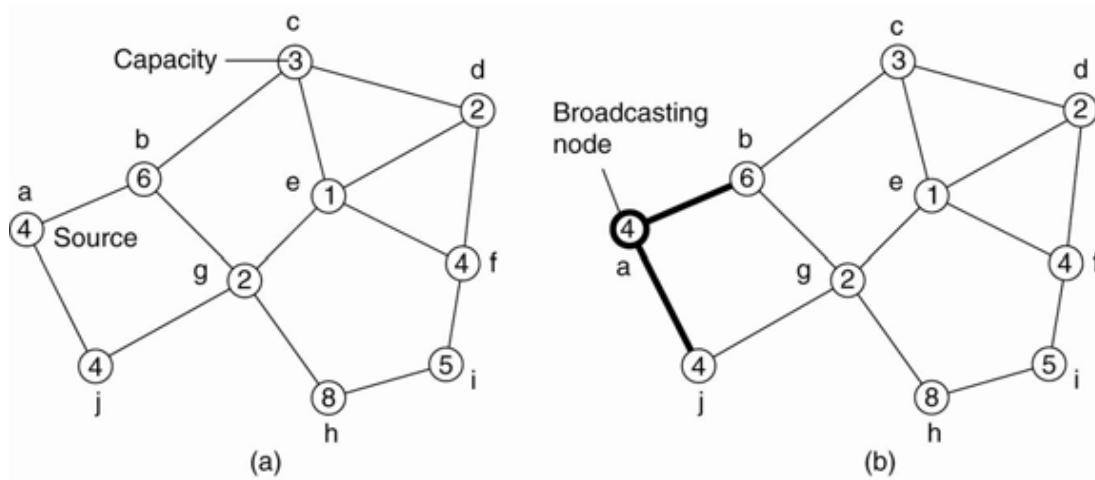
- Uložený soused je poté využijí pro zpětnou komunikaci. To znamená, že další komunikace už probíhá přes strom, nikoliv přes broadcast. Tím je ušetřena některé komunikace.

4.6.2 Složitost

Složitost z hlediska počtu zpráv.

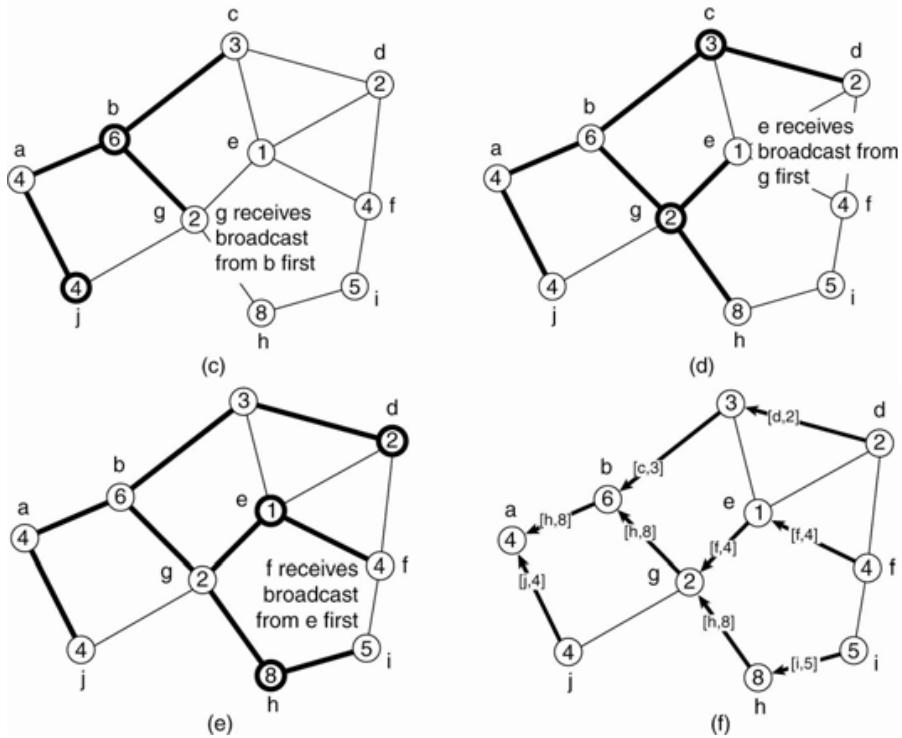
- Inicializační broadcast: počet hran grafu.
- Odpověď: počet hran kostry grafu.
- Result broadcast: počet hran kostry grafu.

4.6.3 Příklad



Node *a* initiates an election.

Obrázek 4.3: Příklad činnosti algoritmu pro obecnou topologii, část 1.



In the end, source a notes that h is the best leader and broadcasts this info to all nodes.

Obrázek 4.4: Příklad činnosti algoritmu pro obecnou topologii, část 2.

Kapitola 5

PDI – Podmínky konsistentního globálního stavu distribuovaného systému.

5.0.1 Metadata

- Předmět: Prostředí distribuovaných aplikací (PDI)
- Přednáška:
 - 4) Globální stav a snapshots
- Záznam:
 - 2020-10-12

5.1 Úvod a kontext

Distribuovaný systém Distribuovaný systém je množina procesů p_1, p_2, \dots, p_n , které jsou propojeny komunikačními kanály. V systému neexistuje žádná globální paměť ani globální hodiny. Procesy spolu komunikují pouze zasíláním zpráv skrze komunikačními kanály.

Komunikační kanál Komunikační kanál mezi procesy p_i a p_j značíme C_{ij} .

Událost Rozlišujeme tři typy událostí: interní událost procesu, zaslání zprávy a přijetí zprávy.

Zpráva Zpráva m_{ij} značí zprávu zaslanou procesem p_i procesu p_j . $send(m_{ij})$ značí odeslání zprávy a $recv(m_{ij})$ přijetí.

Stav procesu Lokální stav procesu p_i značíme LS_i . Lokální stav je definován jako sekvence všech událostí, o kterých proces p_i ví. Nechť e je libovolná událost, $e \in LS_i$ značí, že událost e patří do lokálního stavu procesu p_i , $e \notin LS_i$ značí, že událost e nepatří do lokálního stavu procesu p_i .

Stav komunikačního kanálu Stav komunikačního kanálu C_{ij} značíme SC_{ij} a je definován množinou zpráv, které obsahuje. Pro kanál C_{ij} můžeme definovat jeho stav na základě

lokálních stavů procesů LS_i a LS_j :

$$transit(LS_i, LS_j) = \{m_{ij} \mid send(m_{ij}) \in LS_i \wedge rec(m_{ij}) \notin LS_j\}$$

5.2 Model komunikace

- FIFO – Komukační kanál funguje jako fronta zpráv *first in, first out*. Kanál tedy zachovává pořadí zpráv sám o sobě.
- non-FIFO – Komunikační kanál se chová jako datová struktura množina, do které odesílatel vkládá zprávy a příjemce je odebírá v náhodném pořadí.
- Causal ordering (kauzální uspořádání) – Systém, který podporuje kauzální doručení zpráv splňuje následující vlastnost. Pro jakékoli dvě zprávy m_{ij} a m_{kj} platí, pokud $send(m_{ij}) \rightarrow send(m_{kj})$, pak i $recv(m_{ij}) \rightarrow recv(m_{kj})$.

5.3 Konzistentní globální stav

Globální stav Globální stav distribuovaného systému je kolekce lokálních stavů procesů a komunikačních kanálů.

$$GS = \left\{ \bigcup_i LS_i, \bigcup_{i,j} SC_{ij} \right\}$$

Časoprostorový diagram Diagram pro vizualizaci komunikace procesů v distribuovaném systému. Viz obrázek 5.1 a 5.2.

Konzistentní globální stav Konzistentní globální stav (*snapshot*) je stav systému v určitém časovém okamžiku. Lze si jej představit jako řez v časoprostorovém diagramu, který rozděluje diagram na dvě části: minulost a budoucnost. Aby byl řez (globální stav) konzistentní, tak pokud je doručení nějaké zprávy v minulosti, musí být v minulosti i její odeslání. Formálně jde o globální stav, který splňuje následující podmínky:

$$send(m_{ij}) \in LS_i \Rightarrow m_{ij} \in SC_{ij} \oplus recv(m_{ij}) \in LS_j$$

,

$$send(m_{ij}) \notin LS_i \Rightarrow m_{ij} \notin SC_{ij} \wedge recv(m_{ij}) \notin LS_j$$

K čemu je *snapshot* *Snapshot* lze využít např. pro tvorbu záloh systému nebo při zotavování systému po chybách.

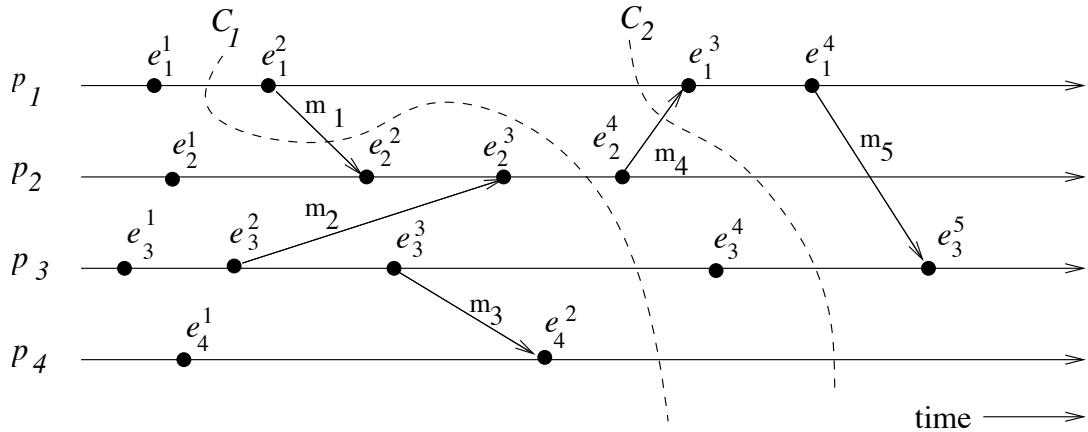
Jak lze *snapshot* vytvořit Absence globální sdílené paměti, globálních hodin a nepředvídatelná délka zpoždění v odesílání zpráv v distribuovaném systému činí problém vytváření snapshotů netriviálním. Způsob vytváření lze rozdělit do dvou kategorií: na základě algoritmů a na základě checkpointů.

Problémy při zaznamenávání snapshotu Jak rozlišit mezi zprávami, které mají být součástí snapshotu a které nikoliv?

- Zprávy, které jsou odeslány procesem před zaznamenáním svého snapshotu, jsou zaznamenány do stavu.
- Zprávy, které jsou odeslány procesem po zaznamenání svého snapshotu, nejsou zaznamenány do stavu.

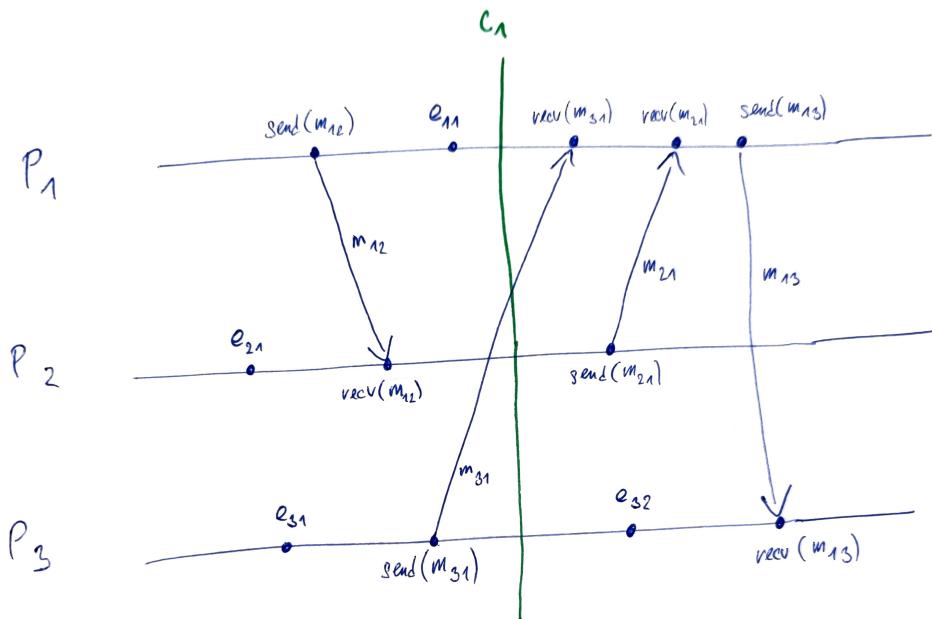
Jak rozpoznat okamžik, ve kterém má proces zaznamenat snapshot?

- Proces p_j musí zaznamenat svůj snapshot před zpracováním zprávy m_{ij} , která byla poslána procesem p_i po zaznamenání jeho snapshotu.



Obrázek 5.1: Příklad řezu v časoprostorovém diagramu. Řez C_1 je nekonzistentní, kvůli zprávě m_1 . Řez C_2 je konzistentní a zpráva m_4 je zachycena ve stavu kanálu Ch_{21} .

(48)



$$GS_{C_1} = \left\{ LS_{P_1} \cup LS_{P_2} \cup LS_{P_3} \cup SC_{12} \cup SC_{13} \cup SC_{23} \right\} \quad SC_{12} = \{\} \quad SC_{13} = \{m_{21}\} \quad SC_{23} = \{\}$$

$$LS_{P_1} = (\text{send}(m_{12}), e_{11}) \quad LS_{P_2} = (e_{21}, \text{recv}(m_{12})) \quad LS_{P_3} = (e_{31}, \text{send}(m_{13}))$$

Obrázek 5.2: Příklad konzistentního globální stavu formálně.

Kapitola 6

PDI – Principy distribuovaného zpracování MapReduce, průběh a jednotlivé operace distribuovaného výpočtu pomocí MapReduce, jeho implementace v Apache Hadoop a Apache Spark.

6.1 Metadata

- Předmět: Prostředí distribuovaných aplikací (PDI)
- Přednáška:
 - 9) Programovací model MapReduce a Apache Hadoop
 - 10) Distribuované souborové systémy
 - 11) Apache Spark
- Záznam:
 - 2020-11-16
 - 2020-11-23

6.2 Úvod a kontext

OLTP OLTP (*Online Transactional Processing*, provozní databáze, systémy pro online zpracování transakcí) jsou standardní databázové systémy s pevnou strukturou dat definovanou pomocí databázového schématu. Jsou navrženy a optimalizovány pro chod provozních aplikací s primárním cílem zajistit rychlý a souběžný přístup k datům. To vyžaduje transakční zpracování, řízení souběžnosti a techniky obnovy (rollback), které zaručují konzistence dat. Díky těmto vlastnostem mají OLTP databáze špatný výkon při provádění složitých dotazů, které potřebují spojit mnoho relačních tabulek dohromady nebo agregovat velké objemy dat. Kromě toho obsahují typicky podrobná data a neobsahují historická data, která jsou při datové analýze potřeba.

OLAP OLAP (*Online Analytical Processing*, online analytické zpracování) je databázové paradigma specificky zaměřené na dotazy, zejména na analytické dotazy. Používají se zde jiné techniky indexování a optimalizace dotazů. Normalizace není pro toto paradigma žádoucí, protože rozděluje databázi na mnoho tabulek. Složité dotazy v takovém případě vyžadují rekonstrukci dat a s tím spojený vysoký počet spojování tabulek. Pracuje se s tzv. multidimenzionálními kostkami, avšak v pozadí jsou stále relační databáze.

NoSQL Potřeba ukládat proudy dat (zpracovávané v reálném čase bez možnosti poza-stavení), obrázky, multimédia, velké JSON soubory, ..., vedla ke vzniku NoSQL databází. NoSQL databáze používají jiné prostředky než tabulková schémata tradiční relační data-báze. Často jde o „hloupé“, nestrukturované uložiště klíč-hodnota.

BigData Velká, nestrukturovaná (různorodá), rychle rostoucí data, která není možné uložit ani zpracovávat běžnými přístupy (na jednom uzlu, jedním uzlem). Produkují je např.: IoT senzory, sociální sítě, chatovací aplikace, webové vyhledávače, ... Pro jejich zpracování je nutné využít distribuované systémy (pro uložení i zpracování).

Distribuované zpracování dat Distribuované zpracování dat je zpracování velkých dat (*big data*) pomocí distribuovaných systémů. To s sebou přináší problémy. Jak zaručit vhodnou distribuci dat a výpočtu mezi uzly? Jak řešit nespolehlivost a výpadky uzlů? Jak a kam zajistit doručení výsledků výpočtu? ...

6.3 MapReduce

Algoritmy pro indexování webových stránek (Page Rank) přestávaly být udržitelné, bylo potřeba zvýšit jejich škálovatelnost. Google vydal příspěvek „MapReduce: Simplified Data Processing on Large Clusters“, kde bylo představeno paradigma MapReduce. Jde o para-digma distribuovaného výpočtu založené na funkcích *map* a *reduce* z funkcionálního pro-gramování.

Map Funkce *map* má ve funkcionálním programování 2 vstupní parametry a vrací seznam hodnot. První parametr je unární operátor (nebo funkce fungující jako unární operátor) a druhý je seznam hodnot. Výstupní seznam je spočítán jako aplikace unárního operátoru na vstupní seznam. Příklad:

$$\text{map}(\text{square}, [1, 2, 3, 4]) = [1, 4, 9, 16]$$

. V paradigmata MapReduce *map* vrací data jako seznam dvojic klíč-hodnota, přesněji:

$$\text{map}((\text{key}, \text{value})) \rightarrow [(\text{key}, \text{value})]$$

Reduce Funkce *reduce* má ve funkcionálním programování 2 vstupní parametry a vrací jednu hodnotu. První parametr je binární operátor (nebo funkce fungující jako binární operátor) a druhý je seznam hodnot. Výstupní hodnota je spočítána jako postupná aplikace binárního operátoru na všechny hodnoty ve vstupním seznamu. Příklad:

$$\text{reduce}(+, [1, 4, 9, 16]) = 30$$

. V paradigmata MapReduce *reduce* bere na vstupu klíč a seznam hodnot a vrací opět seznam dvojic klíč-hodnota, přesněji:

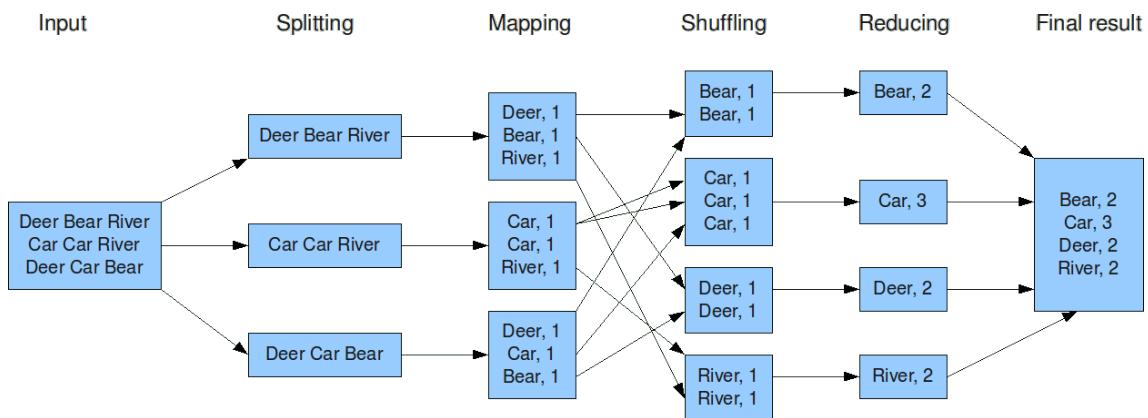
$$\text{reduce}(\text{key}, [\text{value}]) \rightarrow [(\text{key}, \text{value})]$$

```

1 def map(input_key: str, input_value: str) -> list[tuple[str, int]]:
2     # input_key - document name
3     # input_value - document content (etc. line)
4     result = []
5     for word in input_value.split(' '):
6         result.append((word, 1))
7     return result
8
9 def reduce(input_key: str, input_value: list[int]) -> tuple[str, int]:
10    result = 0
11    for val in input_value:
12        result += value
13    return (input_key, result)

```

Výpis 6.1: Příklad implementace funkcí *map* a *reduce* v paradigmatu MapReduce pro počítání četnosti slov ve vstupu v Pythonu.

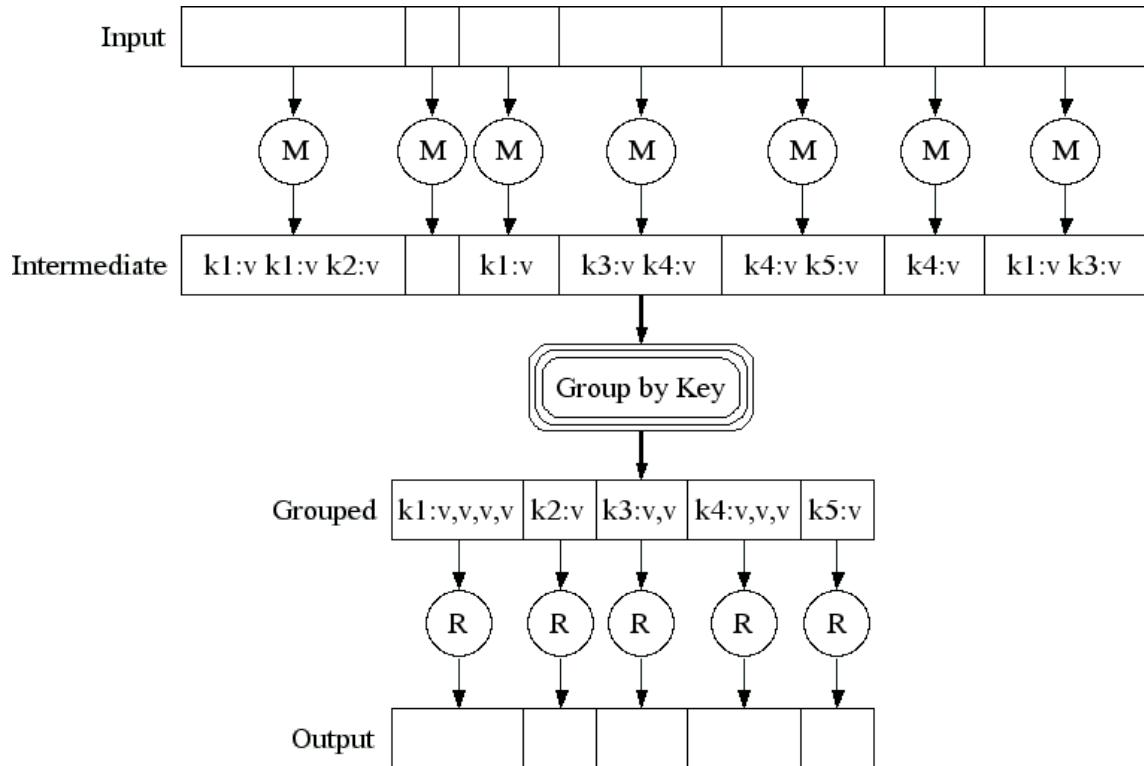


Obrázek 6.1: Úloha počítání četnosti slov v paradigmatu MapReduce v diagramu.

Průběh MapReduce Celý MapReduce probíhá v několika krocích, viz obrázek 6.1.

1. Input – Přípravený vstup pro distribuovaný výpočet (např. soubory ve virtuálním distribuovaném souborovém systému, viz dále HDFS).
2. Splitting – Rozdělení vstupu na části, které budou přiděleny jednotlivým uzelům. Může být výchozí (např. rozdělení textového souboru po řádcích) nebo definováno uživatelem.
3. Mapping – Každý uzel aplikuje funkci *map* na svoji přidělenou část. Uživatel definuje jak má funkce *map* vypadat.
4. Shuffling (také Grouping, Partitioning, Comparing)– Výpočetní uzly si mezi sebou vyměňí hodnoty, které spočítaly, na základě klíče. Tento krok zařizuje platforma pro distribuovaný výpočet sama o sobě, typicky na základě hashů klíčů. Tento krok je většinou *bottleneck*.

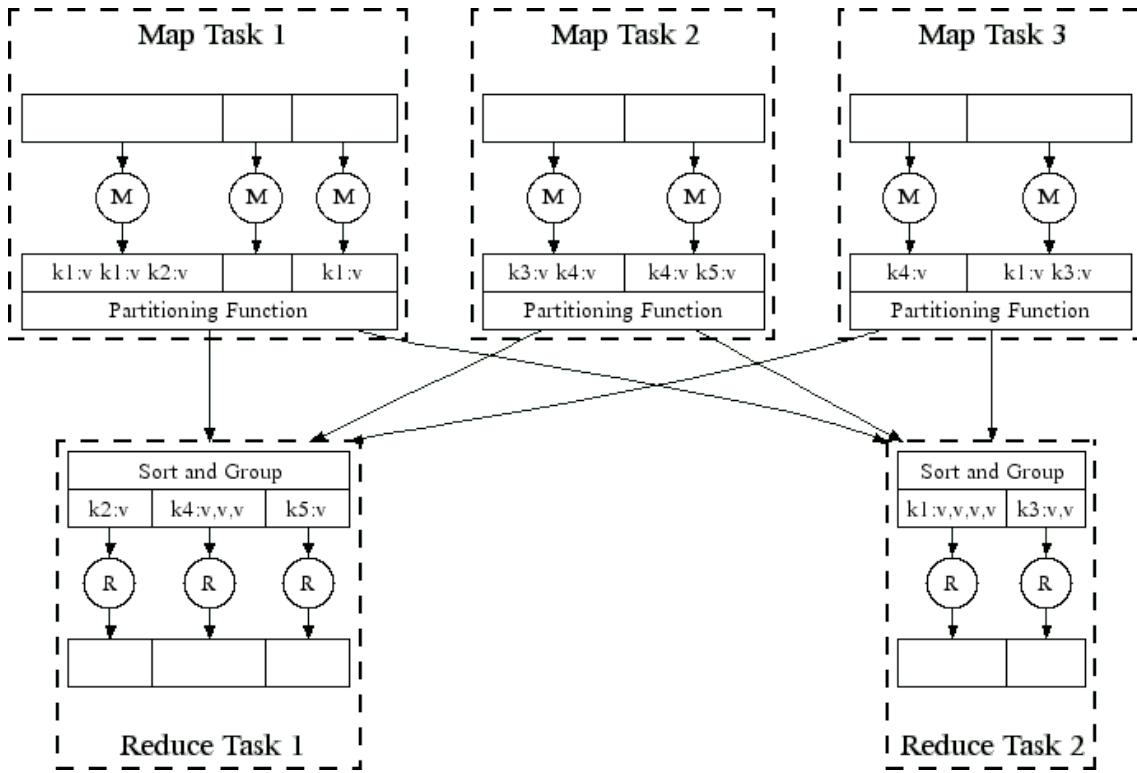
5. Reducing – Každý uzel zapojený do tohoto kroku (často je v tomto kroku potřeba méně uzlů, než v kroku mapping) aplikuje funkci *reduce* na svoji přidělenou část. Uživatel definuje jak má funkce *reduce* vypadat.
6. Final Result – Finální výsledek (např. zapsán do do virtuálního distribuovaného souborového systému, viz dále HDFS).



Obrázek 6.2: Výpočet MapReduce v obecném schématu.

Combiner Optimalizační krok, jde o „jakési“ provedení operace *reduce* už ve fázi *map* (každým uzlem). Tím je snížen počet mezivýsledků ve fázi Shuffling. Typicky funkce *combine* je stejná jako *reduce*.

Virtuální distribuovaný souborový systém Pro realizaci distribuovaného výpočtu je rovněž potřeba distribuovaný souborový systém (DFS). Ten je typicky realizován jako virtuální souborový systém nad jednotlivými souborovými systémy uzlů. Např.: GFS – Google File System, HDFS – Hadoop File System (viz dále). DFS obsahuje data samotná (*data nodes*) a metadata o tom, která data jsou na jakých uzlech (*name nodes*).



Obrázek 6.3: Výpočet MapReduce v obecném schématu a rozdělením práce na jednotlivé uzly (uzel je typicky víceprocesorový).

6.4 Apache Hadoop

Apache Hadoop je *open-source* implementace MapReduce paradigmatu vyvíjená Apache Software Foundation. Jde o implementaci v Java, ta je vhodná, jelikož díky JVM (Java Virtual Machine) je spouštění uživateli definovaných funkcí *map* a *reduce* snadné.

Hadoop MapReduce – Implementace MapReduce paradigmata. Data jsou čtena a ukládána na HDFS (včetně mezivýsledků). To znamená, můžeme pracovat v podstatě neomezenými daty, ale ukládání a načítání výpočet zpomalují.¹

HDFS HDFS (*Hadoop Distribute File System*) je virtuální distribuovaný souborový systém. Standardní soubor je rozdělen na datové bloky které jsou distribuovány na různé datové uzly. Architektura HDFS se skládá ze dvou typů uzlů – Name Node a Data Node. Name Node obsahuje alokační tabulkou pro souborový systém. Ví které datové bloky patří kterému souboru a kde jsou uloženy. Obsahuje další metadata jako názvy souborů, cesty, ... Data Node obsahuje datové bloky. Typicky redundancy a replikace, počítá se s možným selháním uzlů. Pro **čtení dat** se klient zeptá Name Nodu na konkrétní soubor v HDFS. Name Node vrátí metadata o souboru, na jakých Data Nodech se vyskytuje. Klient požádá příslušné Data Nody, ty mu pošlou data, která se na klientovi „poskládají“ do výsledného souboru. Pro **zápis dat** se klient zeptá Name Nodu, kam by měl zapisovat. Klient zapíše na příslušný Data Node. Data Node poté vyřeší replikace s dalšími uzly.

¹ Nebylo přednášeno podrobněji, pravděpodobně stačí princip obecného MapReduce, který byl vysvětlen v předchozí sekci.

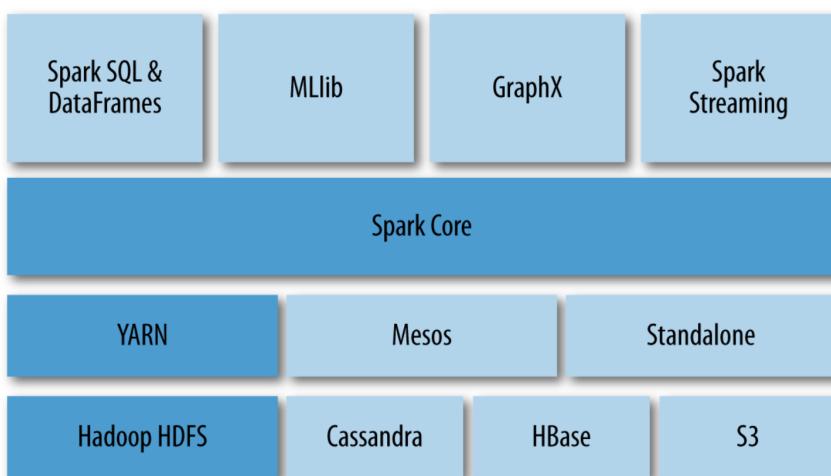
Hadoop YARN Hadoop YARN je plánovač (*scheduler*). Plánuje výpočet tak, aby proběhl co nejlepším způsobem na konkrétní distribuované architektuře. Plánovač má obecné obecné rozhraní a Hadoop YARN lze nahradit za jiný.

Hadoop Common Hadoop Common jsou další knihovny a ovladače pro klienty.

Další nástroje Nad Apache Hadoop existuje mnoho dalších nástrojů. Apache Pig pro *high level* programování map-reduce úloh. Apache Hive pro dolování dat nad Apache Hadoop. Apache HBase jako distribuovaná databáze nad Apache Hadoop, ...

6.5 Apache Spark

Apache Spark je *open-source* nástroj pro distribuované zpracování rozsáhlých dat vyvíjený Apache Software Foundation. Hlavní cíl je zvýšení rychlosti. Spark na to jde přesunutím co nejvíce výpočtů do operační paměti jednotlivých uzlů a tím pádem zminimalizovat počet zápisů a čtení z DFS (snaha odstranit *bottleneck* v kroku shuffling u Hadoopu). Tím ale vzniká jiný problém, a sice výpadek uzlu znamená, že data jsou ztraceny.



Obrázek 6.4: Architektura Apache Spark. Hlavní je Spark Core, zbytek funguje na systému pluginů a může používat HDFS, Hadoop YARN a Hadoop Common.

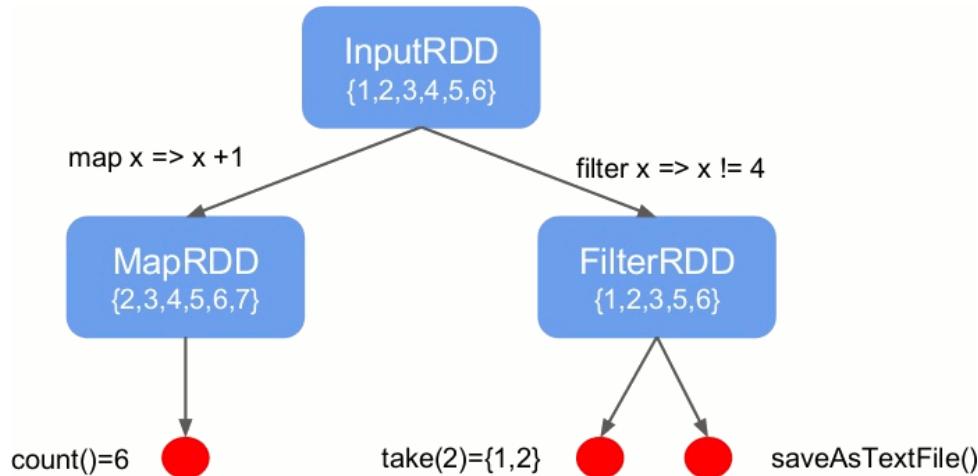
Resilient Distributed Dataset Resilient Distributed Dataset (RDD) je základní datová struktura Sparku. Jedná se o typované kolekce n-tic, které jsou neměnné (*read only*). Vstup je transformován na RDD a každá operace je pak transformace jednoho RDD na jiné.

$$RDD_1 \rightarrow map() \rightarrow RDD_2 \rightarrow reduce() \rightarrow RDD_3$$

Strategie vyhodnocování Spark uplatňuje strategii vyhodnocování *lazy evaluation*. Vyhodnocování výrazu je odloženo až do doby, dokud není potřeba jeho hodnota. Zabraňuje opakování vyhodnocování. Je vyhodnocována pouze ta část, která je potřeba. RDD funguje jako abstraktní datová struktura, nemusí obsahovat data uvnitř, ale pouze předpis jak data získat a získá je, až když jsou potřeba.

Struktura výpočtu Struktura výpočtu odpovídá orientovanému acyklickému grafu (DAG, *Directed Acyclic Graph*). Uzly jsou RDD a hrany jsou transformace. DAG je znám i dalším uzlům, takže pokud nastane výpadek uzlu a výpočet je ztracen, může být uzel snadno zastoupen.

Klíčové vlastnosti Klíčové vlastnosti Sparku jsou *lazy evaluation*, *in-mememory* a *parallel computing*.



Obrázek 6.5: Příklad výpočtu v Apache Spark.

Kapitola 7

KRY – Symetrická kryptografie. Vlastnosti, vlastnosti bezpečného algoritmu, délka klíče, útok silou, příklady symetrických algoritmů, Feistelovy šifry, DES, režimy činnosti, proudové šifry.

7.1 Metadata

- Předmět: Kryptografie (KRY)
- Přednáška:
 - 3) Symetrická kryptografie. Vlastnosti, vlastnosti bezpečného algoritmu, délka klíče, útok silou.
 - 4) Příklady symetrických algoritmů, Feistelovy šifry, DES, struktura, činnost, slabiny, režimy činnosti.
 - 5) Typické aplikace symetrické kryptografie.
- Záznam:
 - 2021-02-22
 - 2021-03-01
 - 2021-03-08

7.2 Úvod a kontext

Kryptografie Kryptografie (šifrování) je věda o metodách utajování smyslu zpráv převodem do podoby, která je čitelná jen se speciální znalostí.

Kryptoanalýza Kryptoanalýza je věda zabývající se metodami získávání obsahu šifrovaných informací bez přístupu k tajným informacím, které jsou za normálních okolností potřeba, tzn. především k tajnému klíči.

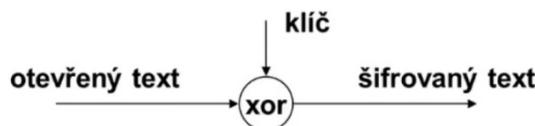
Kryptologie Jeden výraz pro kryptografiu a kryptoanalýzu.

Caesarova šifra Princip Caesarovy šifry je založen na tom, že všechna písmena zprávy jsou během šifrování zaměněna za písmeno, které se abecedně nachází o pevně určený počet míst dále (tj. posun je pevně zvolen). Caesarova šifra spadá do kategorie substitučních šifer (stejný znak je při více výskytech vždy zašifrován na stejný znak).

Vigenerova šifra Rozšíření Caesarovy šifry, klíč je delší než 1 znak. Klíč je řetězec, který reprezentuje posuny. V případě že vstup je delší než klíč, je klíč perioricky opakován. Vigenerova šifra spadá do kategorie polyalfabetických substitučních šifer (stejný znak může být při více výskytech zašifrován na jiný znak).

Vernamova šifra (*One Time Pad*) Vernamova šifra spadá do kategorie polyalfabetických substitučních šifer a je i dnes nerozluštitelná pokud:

- klíč je delší než vstupní text,
- klíč se nepoužije opakováně,
- klíč je náhodný.



Obrázek 7.1: Vernamova šifra.

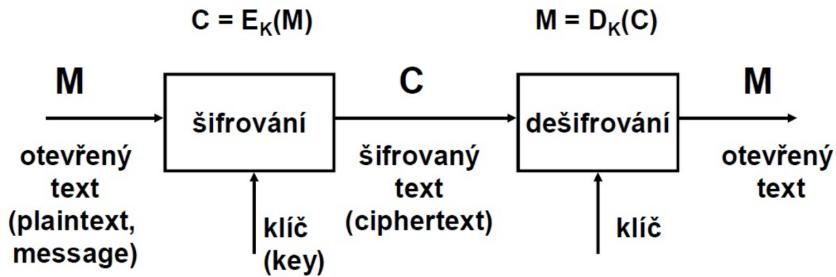
Autoklíč (*autokey*) Šifrování klíčem a když vstupní text je delší než klíč, tak se pokračuje šifrováním otevřeným nebo šifrovaným textem. Lze použít u Vigenerovy nebo Vernamovy šifry.

Symetrická kryptografie Algoritmy používají k šifrování i dešifrování stejný klíč. Výhodou symetrických šifer je jejich nízká výpočetní náročnost. Asymetrické šifry mohou být i stotisíckrát pomalejší. Nevýhodou je nutnost sdílení tajného klíče, takže jedna strana musí klíč vygenerovat a potom ho bezpečným způsobem předat druhé straně.

Typy útoků

- Ciphertext Only Attack (COA) – Útočník zná pouze zašifrovaný text a snaží se zjistit klíč nebo otevřený text. Nejčastější případ.
- Known Plaintext Attack (KPA) – Útočník zná zašifrovaný text a otevřený text a snaží se zjistit klíč.
- Chosen Plaintext Attack (CPA) – Útočník zná to co v KPA a navíc si text může zvolit.

Útok silou Pří útoku silou (*brute force*) zkouší útočník všechny teoreticky možné klíče, dokud nenajde ten správný.



Obrázek 7.2: Princip kryptografie, podle typu klíčů dělíme na symetrickou (tajný klíč) a asymetrickou (veřejný klíč, soukromý klíč).

Bezpečný algoritmus V moderní kryptografii je nepřijatelné utajování algoritmů (*security by obscurity*) – předpokládáme, že útočník zná šifrovací algoritmus. Bezpečnost musí záviset pouze na utajení klíče (Kerckhoffuv princip, *security by design*). Symetrický algoritmus je považován za bezpečný, pokud neexistuje rychlejší útok než útok silou.

Délka klíče Dnes je považováno 80 bitů a více za dostatečné. Typicky se délka zaokrouhuje na mocninu 2 (typicky 128b). Klíče symetrických algoritmů jsou kratší než asymetrických. Konkrétně: DES – 56b, 3DES – 112, AES – variabilní.

Využití Symetrická kryptografie je vhodná pro šifrování většího objemu dat. Narozdíl od asymetrické, která je pro tento účel příliš pomalá. Proto např. HTTPS využívá asymetrickou kryptografií pro výměnu symetrických klíčů a poté symetrickou kryptografií pro šifrování provozu.

Vlastnosti moderní kryptografie Symetrická kryptografie zaručuje všechny následující, kromě nepopiratelnosti – více entit má k dispozici klíč.

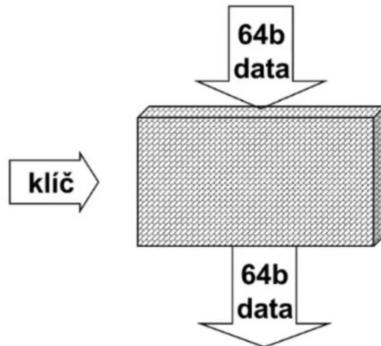
- Důvernost – Utajení informace. Bez znalosti klíče, není možné data číst.
- Autentizace – Prokázání, že zprávu skutečně poslal odesíatel a nikoliv útočník, který se za odesílatele vydává.
- Integrita – Prokázání, že nikdo nemohl data po cestě od odesílatele k příjemci změnit. Ochrana proti neoprávněné, neodhalené modifikaci zprávy.
- Nepopiratelnost – Pokud odesíatel data poslal, nemůže tuto skutečnost popřít.

7.3 Blokové šifry

Blokové šifry šifrují data po blocích pevně stanovené délky (64b, 128b, 256b, ...). Pokud je dat více, rozdělí se na více bloků, přičemž do zbylého místa v posledním je umístěno zarovnání *padding* (informace o délce zarovnání může být obsažena v posledním bytu). Příklady blokových šifer:

- Feistelova šifra (spíše princip)
- Data Encryption Standard (DES)
- Triple Data Encryption Algorithm (3DES)
- International Data Encryption Algorithm (IDEA)

- Blowfish
- Tiny Encryption Algorithm (TEA)
- Advanced Encryption Standard (AES)



Obrázek 7.3: Princip blokových šifer.

7.3.1 Feistelova šifra

Feistelova šifra (Feistelův princip) je koncept šifrování, který konkrétní algoritmy využívají. Jedná se o substituční-permutační síť. Vstupní blok je rozdělen na dvě poloviny L a R , výpočet výstupu pak vypadá následovně.

$$L_i = R_{i-1} \quad (7.1)$$

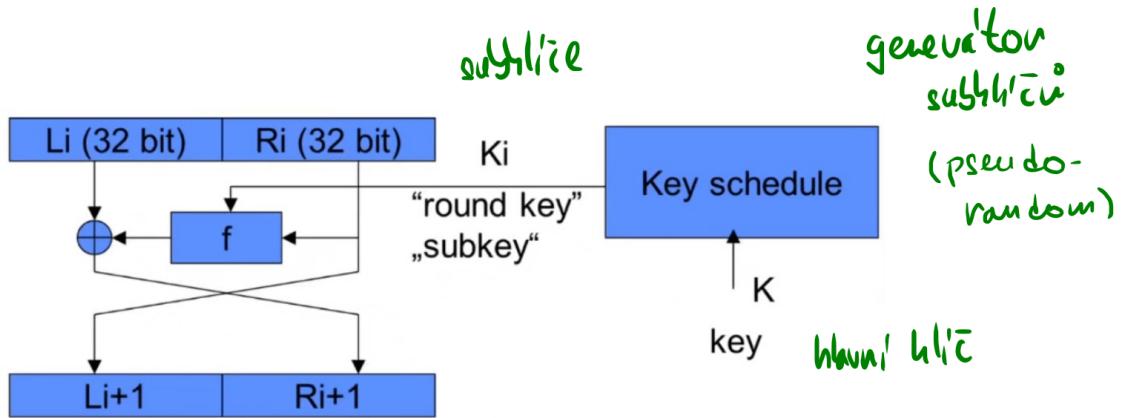
$$R_i = L_{i-1} \oplus F(R_{i-1}, K_i) \quad (7.2)$$

Funkce F F je funkce, na kterou Feistelova šifra neklade žádné požadavky. Jednotlivé algoritmy, využívající Festelovu šifru, funkci samy definují. Požadavky na funkci F , aby algoritmus byl bezpečný:

- skrytí vlastností zprávy;
- skrytí vlastností zprávy.

Subklíč K je tzv. subklíč, který je generován typicky nějakým pseudonáhodným generátorem na základě inicializačního klíče (hlavní).

Dešifrování Dešifrování se provádí stejným způsobem, pouze pořadí subklíčů je opačné.



Obrázek 7.4: Jeden krok opakování (Feistelův krok) vizuálně.

Příklad

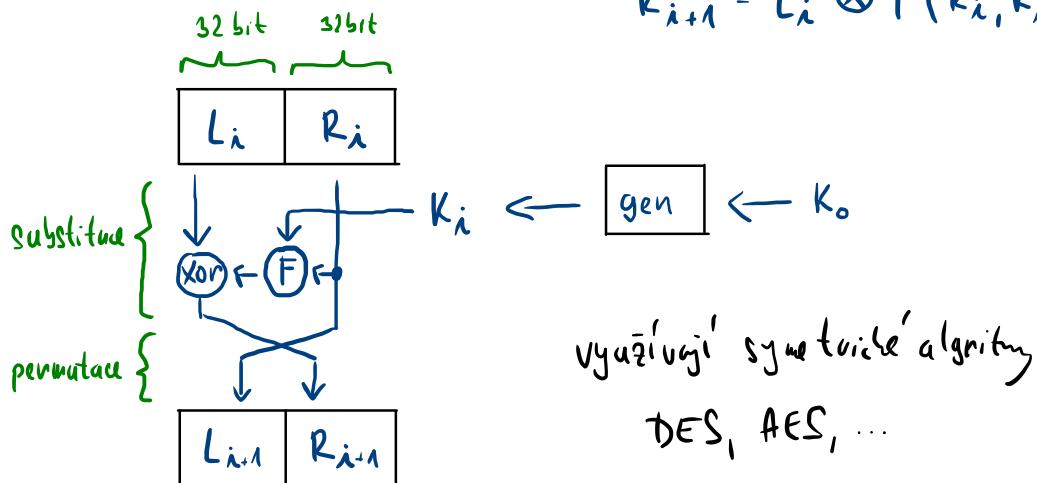
9. Nakreslit a popísať Feistelovu šifru, napísat algoritmus ktorý to používa.

- Princip kľúčových symetrických šífer

- Vstup rozdelen na kľúče o 64 bitoch
↳ vstup rozšíren

$$L_{i+1} = R_i$$

$$R_{i+1} = L_i \otimes F(R_i, K_i)$$



Obrázek 7.5: Feistelova šifra – příklad a rekapitulace.

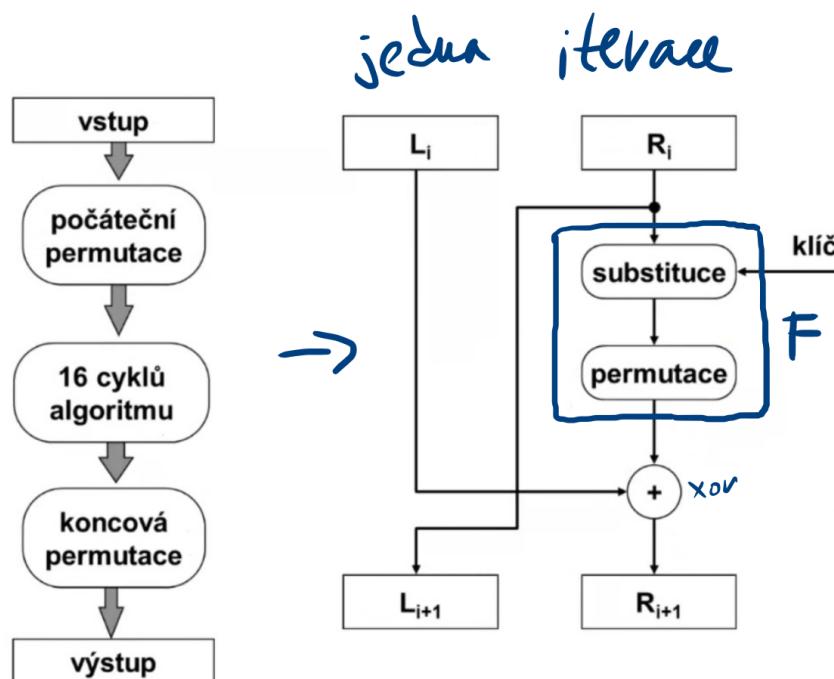
7.3.2 Data Encryption Standard (DES)

DES byl první algoritmus s veřejnou specifikací (*security by design*). Využívá princip Feistelovy šifry – 16 kol. Dodatečně přidává na začátek a konec permutaci navíc. Klíč je dlouhý 64b (resp. 56 významových bitů a 8 paritních).

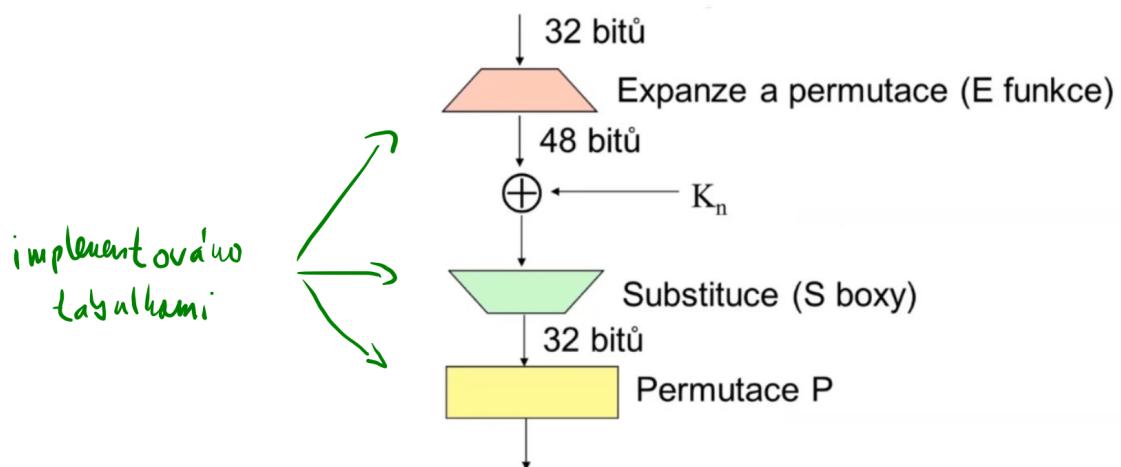
Slabiny

- 56 bitový klíč je příliš krátky a je možný útok silou.

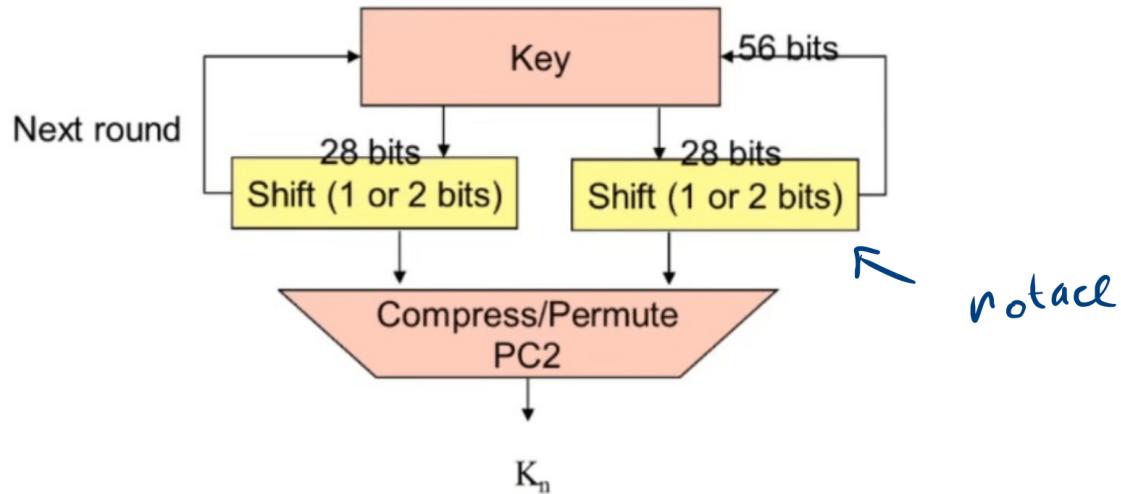
- Rozdílná velikost bloku a klíče (zvláštnost).
- Existence slabých a poloslabých klíčů.
- Není jasné proč zrovna 16 iterací a zda je to dostatečné.



Obrázek 7.6: DES – Schéma fungování algoritmu.



Obrázek 7.7: DES – funkce F.



Obrázek 7.8: DES – generování subklíčů.

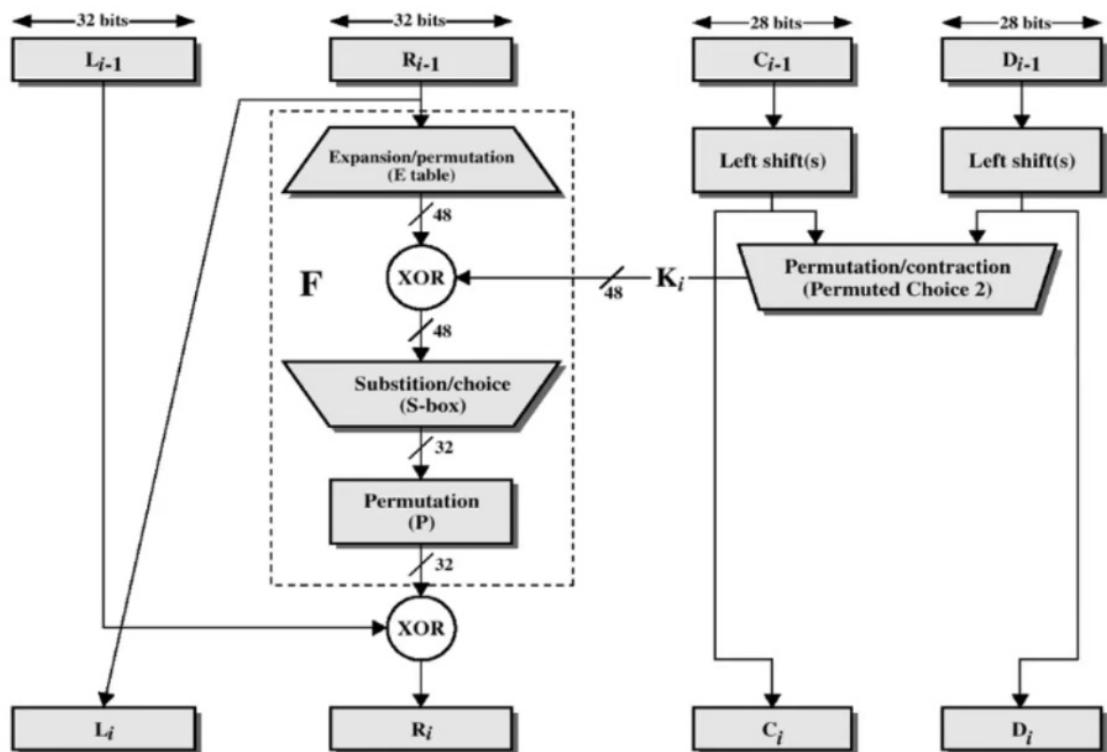
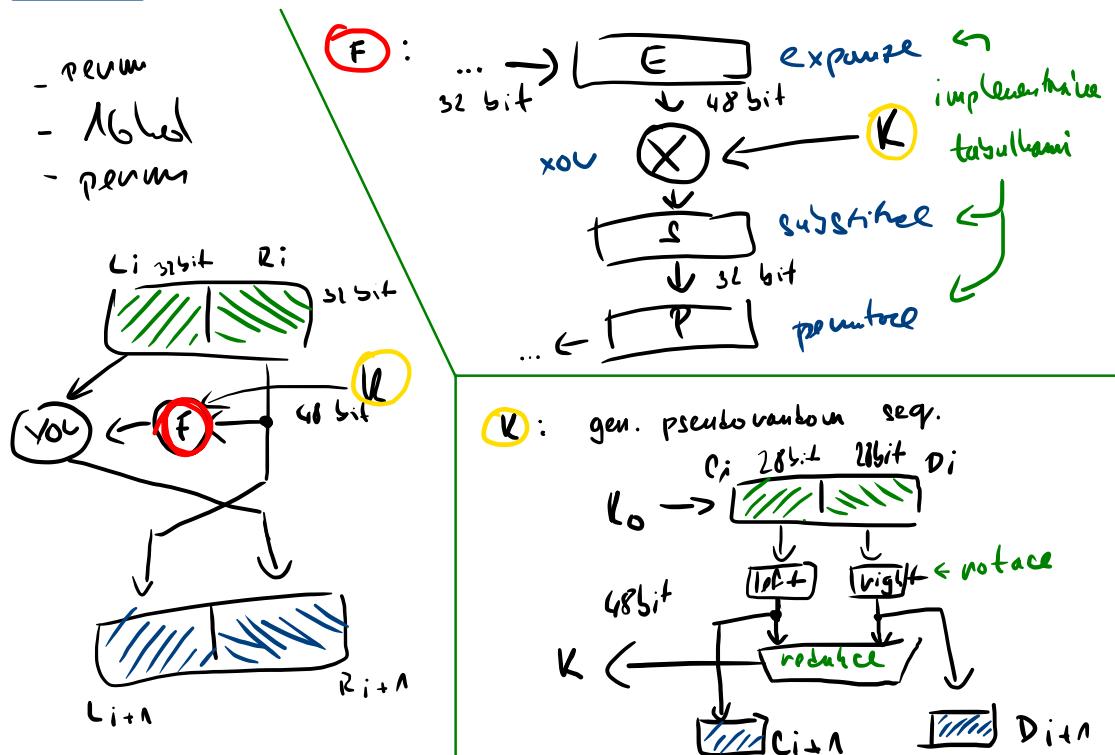


Figure 2.4 Single Round of DES Algorithm

Obrázek 7.9: DES – jedno kolo algoritmu. *Left shift* je ve skutečnosti bitová rotace.

Příklad

3. Nakreslete schema DES, včetně key scheduler, popište všechny 3 části funkce F, jaký mají účel a napишите proč je treba pochybovat o bezpečnosti DESu



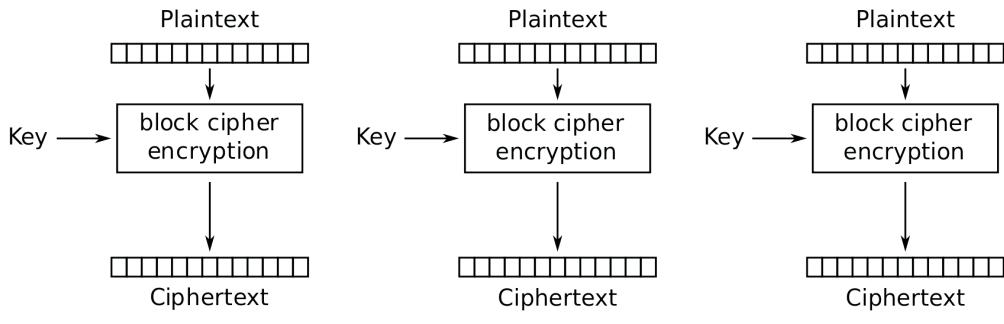
Obrázek 7.10: DES – příklad a rekapitulace.

7.4 Provozní režimy činnosti blokových šifer

Jak použít blokové šifry abychom byli schopni šifrovat data delší než jeden blok?

7.4.1 Electronic Code Book (ECB)

ECB („kódová kniha“) je výchozí *naivní* režim. Bloková šifra se při něm přímo aplikuje nezávisle na jednotlivé bloky, tedy při daném klíči odpovídá stejnému bloku otevřeného textu stejný blok šifrového textu. To má nežádoucí důsledky z hlediska bezpečnosti, v datech zůstane původní struktura, např. šifrovaný obrázek je rozpoznatelný.



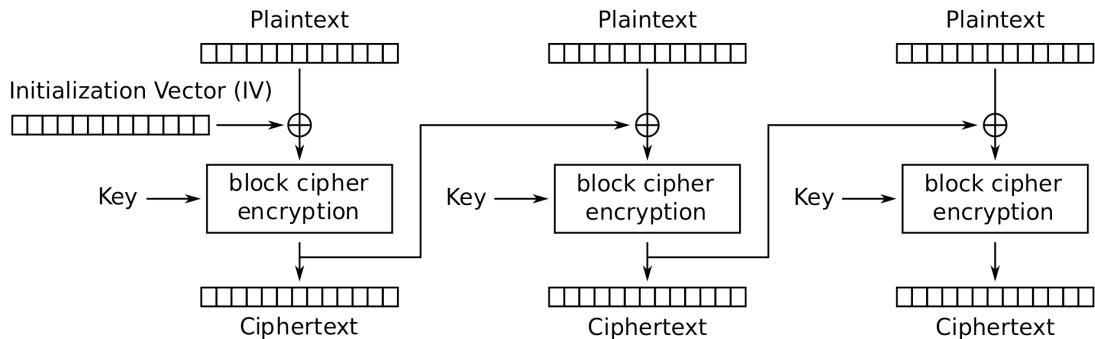
Electronic Codebook (ECB) mode encryption

Obrázek 7.11: Ukázka režimu ECB.

7.4.2 Cipher Block Chaining (CBC)

V režimu CBC („řetězení šifrových bloků“) je každý blok před šifrováním xorován zašifrovaným předchozím blokem a první blok je xorován inicializačním vektorem. Tento režim je široce používán. Nevýhody plynou ze zřetězené závislosti (šifrovaný blok závisí na všech předcházejících): Šifrování nelze paralelizovat a při poškození šifrového bloku nelze dešifrovat ani blok přímo následující. Dešifrování paralelizovat lze.

$$\begin{aligned} C_i &= E_K(P_i \oplus C_{i-1}) \\ P_i &= D_K(C_i) \oplus C_{i-1} \\ C_0 &= IV \end{aligned} \quad (7.3)$$



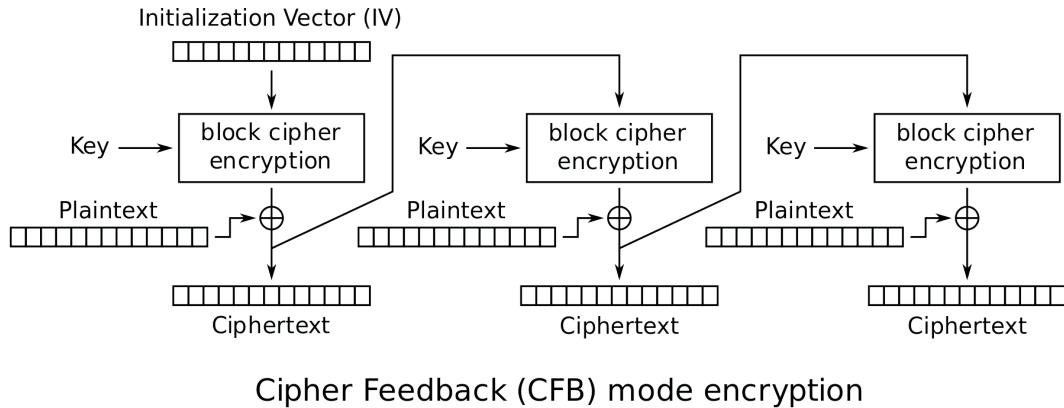
Cipher Block Chaining (CBC) mode encryption

Obrázek 7.12: Ukázka režimu CBC.

7.4.3 Cipher Feedback (CFB)

Režim CFB (šifrová zpětná vazba) se liší oproti CBC v prohození pořadí operací xor a šifrování – nejprve se zašifruje předchozí šifrovaný blok (resp. inicializační vektor) a výsledek se xoruje s otevřeným blokem. Toto prohození má významné implementační dopady: díky symetrii operace XOR vypadá dešifrovací funkce obdobně jako šifrovací. Šifruje pomaleji než CBC. Vstup není nutné zarovnávat. Plynou stejné nevýhody jako pro CBC.

$$\begin{aligned}
 C_i &= E_K(C_{i-1}) \oplus P_i \\
 P_i &= E_K(C_{i-1}) \oplus C_i \\
 C_0 &= IV
 \end{aligned} \tag{7.4}$$



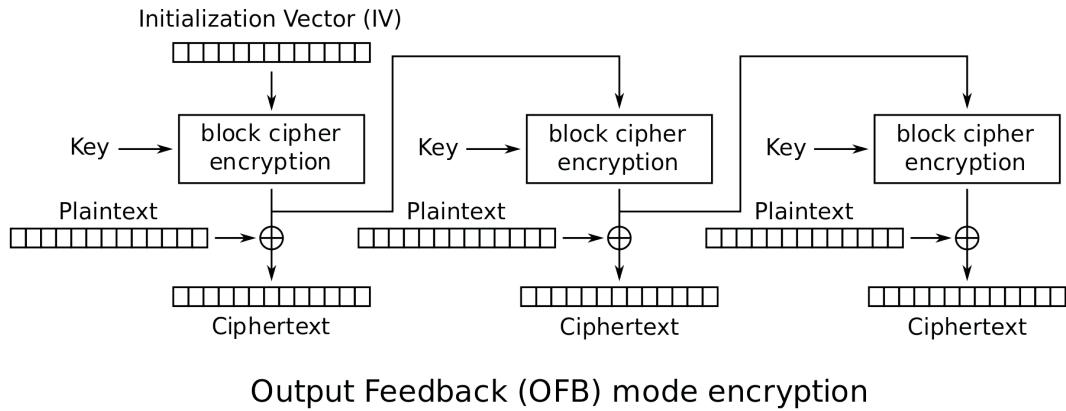
Obrázek 7.13: Ukázka režimu CFB.

7.4.4 Output Feedback (OFB)

Režim OFB (*výstupní zpětná vazba*) se liší od CFB pouze v tom, kde bere zpětnou vazbu. Šifrování probíhá pouhým xorováním otevřeného bloku s heslem, které je v každém kroku zašifrováno použitou blokovou šifrou. První blok hesla je získán zašifrováním inicializačního vektoru. Režim převádí blokovou šifru na synchronní proudovou šifru.

Slabina Celý blok encryption je pouze generátor pseudonáhdon posloupnosti (je nezávislá na otevřeném nebo šifrovaném textu). To umožňuje Known Plaintext Attack. Z toho plyne, že jedním klíčem není bezpečné šifrovat více než jednu zprávu.

$$\begin{aligned}
 C_i &= E_K(C_{i-1}) \\
 P_i &= P_i \oplus C_i \\
 C_0 &= IV
 \end{aligned} \tag{7.5}$$

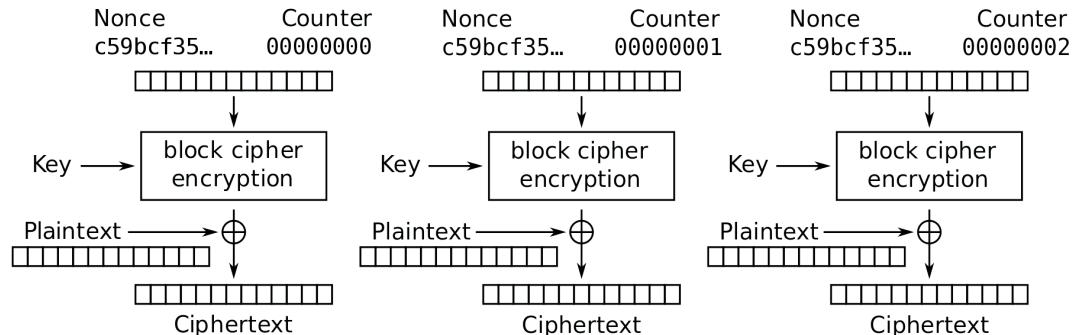


Obrázek 7.14: Ukázka režimu OFB.

7.4.5 Counter (CTR)

Režim CTR („čítačový režim“) převádí stejně jako OFB blokovou šifru na synchronní proudovou. Heslo, se kterým se blok otevřeného textu xoruje, je však získáno zašifrováním čítače, který se každou iteraci zvětšuje o pevně danou hodnotu, zpravidla o 1. Obsah čítače je opět před šifrováním nastaven inicializačním vektorem. Každý blok je šifrován nezávisle na ostatních, díky tomu je možné paralelizovat.

$$\begin{aligned} CTR_i &= CTR_{i-1} + 1 \\ P_i &= P_i \oplus E_k(CTR_i) \\ CTR_0 &= IV \end{aligned} \quad (7.6)$$



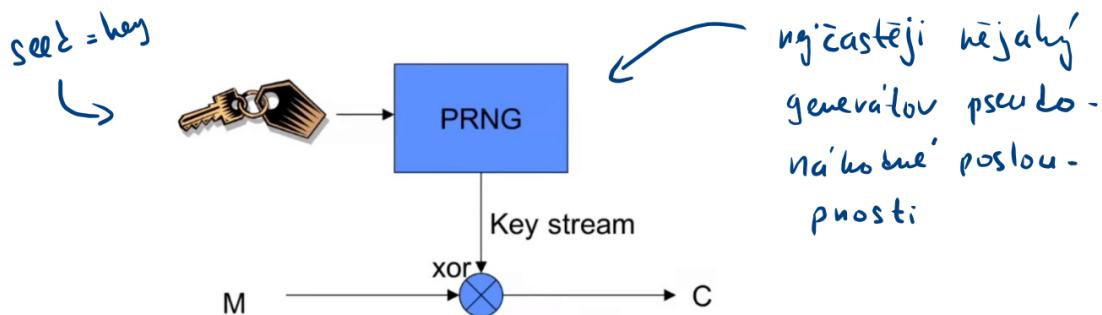
Obrázek 7.15: Ukázka režimu CTR.

7.5 Proudové šifry

Proudové šifry šifrují data jako *proud* (stream), nejčastěji po jednotlivých bytech. Dešifrování vždy probíhá stejným způsobem. Proudové šifry jsou rychlejší než blokové šifry a pro implementaci potřebují jednodušší hardware.

Problémy

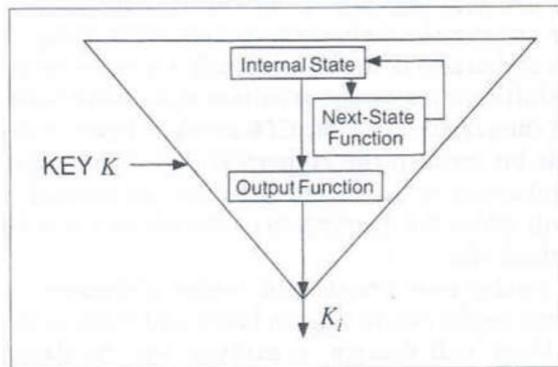
- Nezajišťují samy o sobě integritu.
- Na rozdíl od blokových šifer jsou náchylnější ke kryptoanalytickým útokům, pokud jsou nevhodně implementovány (počáteční stav nesmí být použit opakován) – „problém s inicializačním vektorem“.



Obrázek 7.16: Princip proudových šifer.

7.5.1 Rozdelení proudových šifer

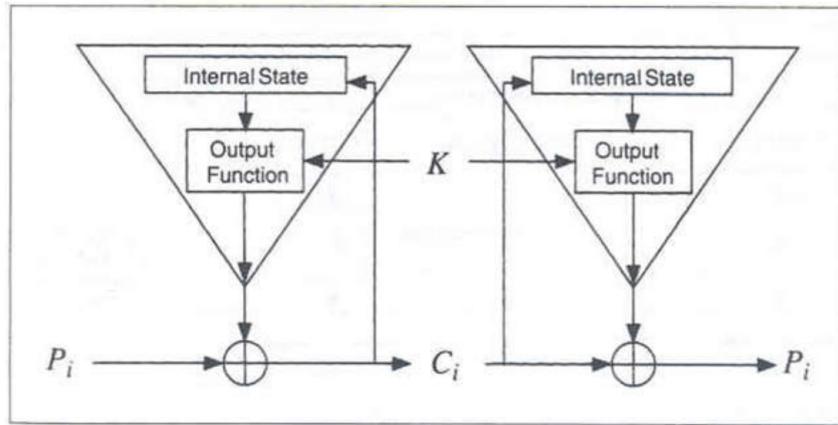
Synchronní proudové šifry Proud pseudonáhodných čísel *key stream* je generován nezávisle na vstupním textu nebo zašifrované zprávě. Poté dochází ke kombinaci vygenerovaných čísel se vstupujícím textem (k zakódování) nebo se šifrovaným textem (k dekódování). Nejběžnější formou kombinace keystreamu a vstupního textu je použití operace XOR. Např.: Vernamova šifra, DES v režimu OFB. Pokud se průběhu dešifrování něco ztratí, je konec.



Obrázek 7.17: Princip synchronní proudové šifry.

Samosynchronizující proudové šifry Proud pseudonáhodných čísel *key stream* závisí na pevném počtu předcházejících bytů šifrovaného (nebo otevřeného) textu. To znamená, že se šifra dokáže po chybě sama *zotavit* (resynchronize)¹. Např.: Vigenere Autokey, DES v režimu CFB.

¹Dnes se nesnášíme dešifrovat poškozená data, pokud nastane chyba v přenosu, vyžádáme si data znova.



Obrázek 7.18: Princip samosynchronizující proudové šifry.

7.5.2 Generátory PRNG

Generátory PRNG (*pseudo-random number generator*) generují pseudo-náhodnou posloupnost (*key stream*) z malého klíče (*seed*).

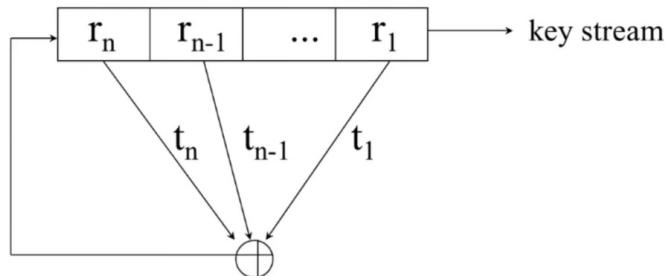
Blokové šifry v režimu OFB Blokové šifry v režimu OFB jsou pomalé.

Linear Feedback Shift Registers (LFSR) LFSR (posuvný registr s lineární zpětnou vazbou) je posuvný registr, jehož výstup je lineárně závislý na jeho předchozích výstupech a stavu. Mějme

- posuvný registr $R = (r_1, r_2, \dots, r_n)$,
- sekvenci zpětných vazeb $T = (t_1, t_2, \dots, t_n)$.

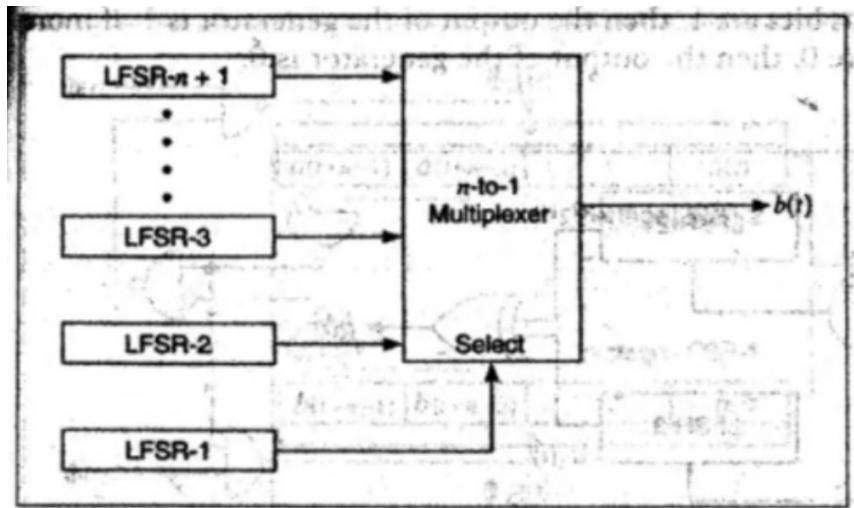
Alternativně lze zapsat polynomem:

$$T(x) = x^n + x^{n-1} + \dots + x + 1$$



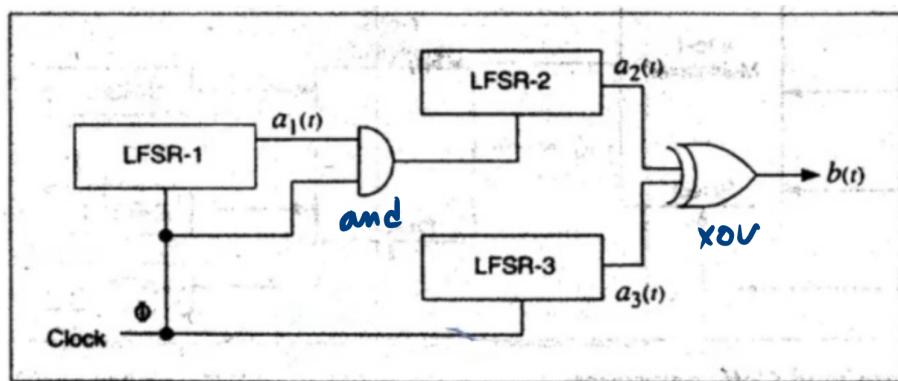
Obrázek 7.19: Příklad LFSR.

Geffe generátor Geffe generátor (kombinovaný generátor) je využití 2 a více LFSR propojených multiplexorem.



Obrázek 7.20: Příklad Geffe generátoru.

Stop and Go generátor Stop and Go generátor je několik LSFR s různým zdrojem hodin.



Obrázek 7.21: Příklad Stop and Go generátoru.

Kapitola 8

**KRY – Asymetrická kryptografie,
vlastnosti, způsoby použití,
poskytované bezpečnostní funkce,
elektronický podpis a jeho vlastnosti,
hybridní kryptografie, algoritmus
RSA, generování klíčů, šifrování,
dešifrování.**

8.1 Metadata

- Předmět: Kryptografie (KRY)
- Přednáška:
 - 6) Asymetrická kryptografie, vlastnosti, způsoby použití, poskytované bezpečnostní funkce.
 - 7) Elektronický podpis a jeho vlastnosti, hybridní kryptografie.
 - 8) Příklady asymetrických algoritmů, RSA.
- Záznam:
 - 2021-03-08
 - 2021-03-22

8.2 Úvod a kontext

Viz. „Úvod a kontext“ v předchozích otázkách z tohoto předmětu.

Asymetrická kryptografie

- V asymetrické kryptografii se používají páry klíčů (soukromý a veřejný). Soukromý je používán k dešifrování, resp. vytvoření digitálního podpisu. Veřejný je používán k šifrování, resp. ověření digitálního podpisu.

- Každý uživatel generuje svůj pár klíčů. Veřejný klíč je zveřejněn (znají ho všichni), soukromý je držen v tajnosti (zná ho pouze vlastník).
- Všechny asymetrické algoritmy jsou blokové.
- Asymetrické algoritmy jsou pomalejší než symetrické.

Způsoby použití Asymetrická kryptografie lze využít k:

- šifrování,
- digitálnímu podepisování,
- pro výměnu symetrického klíče (*key exchange*).

Vlastnosti Vlastnosti symetrické a asymetrické kryptografie¹.

| | Důvěrnost | Autentizace | Integrita | Nepopiratelnost |
|----------------------------|-----------|-------------|-----------|-----------------|
| Symetrická | ano | ? | ? | ne |
| Asymetrická - šifrování | ano | ? | ? | ne |
| Asymetrická - podepisování | ne | ano | ano | ano |
| Asymetrická - kombinace | ano | ano | ano | ano |

```

1 # Odesilatel (A):
2 msg = encrypt(msg, SK_A) # nechť msg je zprava k odeslání
3 msg = encrypt(msg, PK_B)
4 send(msg_2)
5
6 # Příjemce (B):
7 msg = receive()
8 msg = decrypt(msg, SK_B)
9 msg = decrypt(msg, PK_A)
```

Výpis 8.1: Kombinace klíčů obou stran u asymetrické kryptografie. Pořadí operací může být i opačné.

Digitální podpis Vytvoření digitálního podpisu konkrétních dat pomocí soukromého klíče podepisatele. Každý kdo zná veřejný klíč podepisatele, může pravost podpisu ověřit. Digitální podpis zajišťuje autentizaci, integritu a nepopiratelnost.

Algoritmy Algoritmy asymetrické kryptografie se nedají *vymyslet*, musí se objevit. Jsou založeny na těžkých matematických problémech.

- Problém batohu (*knapsack problem*) – MH (Merkle-Hellman)
- Faktorizace čísel – RSA (Rivest-Shamir-Adleman)
- Diskrétní logaritmus – DSA (Digital Signature Algorithm), DH (Diffie-Hellman)
- Eliptické křivky – ECDSA, ECDH

¹Otazníky – částečně, za předpokladu, ...

Problém batohu Problém batohu je NP-úplný problém kombinatorické optimalizace. Necht' x_1, x_2, \dots, x_n je množina objektů, každý objekt má svoji cenu v_i a svoji hmotnost w_i , dále mějme batoh, který má kapacitu W . Cílem je vybrat takovou množinu objektů, jejichž hmotnost je menší nebo rovna W a má nejvyšší možnou cenu². Formálně chceme maximalizovat sumu

$$\sum_{i=1}^n v_i \cdot x_i$$

, při splnění

$$\sum_{i=1}^n w_i \cdot x_i \leq W$$

, kde $x_i \in \{0, 1\}$.

Faktorizace čísel Faktorizace čísel označuje problém rozložení čísla na součin menších čísel, v nejčastější podobě pak rozklad celého čísla na součin prvočísel.

Diskrétní logaritmus Necht' p, g, k, Y jsou přirozená čísla, pro něž platí $Y \equiv g^k \pmod{p}$. Potom každé číslo k odpovídající uvedené rovnici nazveme diskrétní logaritmus o základu g z Y vzhledem k modulu p . Tato definice nedefinuje číslo k jednoznačně, proto se někdy upravuje tak, že ze všech možných diskrétních logaritmů ve smyslu předchozí definice se vybere ten nejmenší.

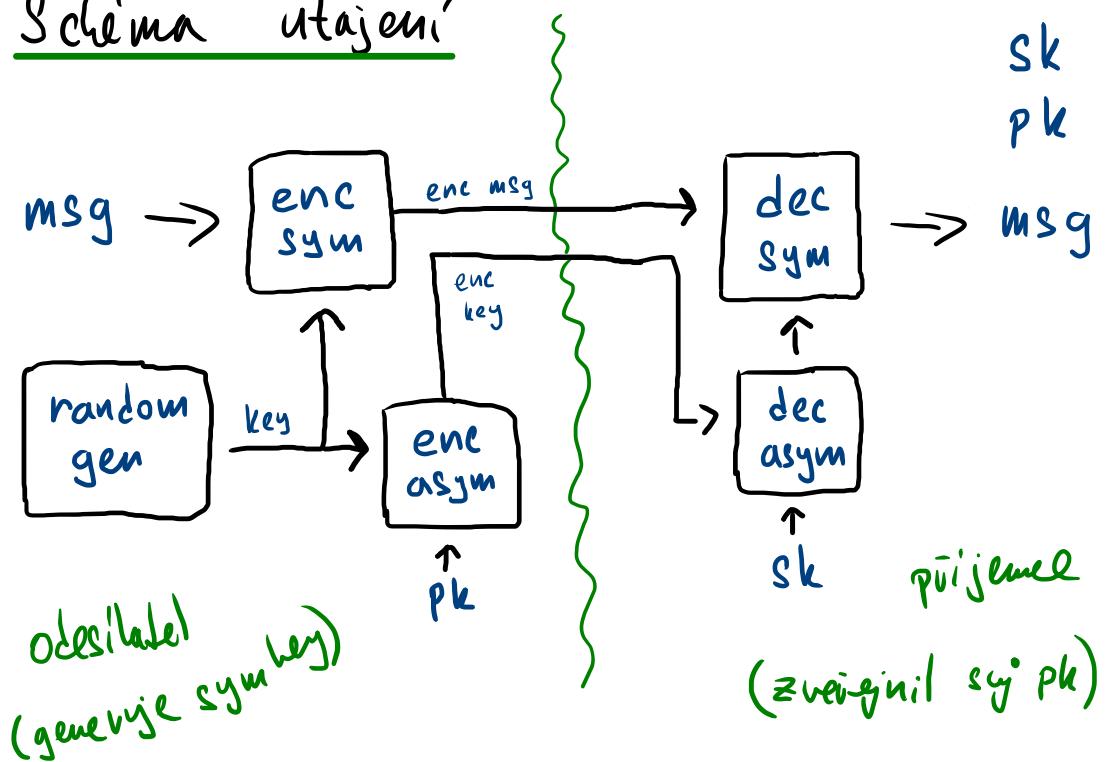
Eliptické křivky Jedná se o matematický aparát, na kterém aplikujeme různé algoritmy (DSA, DH).

8.3 Hybridní kryptografie

Hybridní kryptografie je kombinace symetrické a asymetrické kryptografie, ve které jsou využity přednosti obou (symetrická – rychlá, ale potřeba stejný klíč; asymetrická – pomalá, ale dva klíče). Asymetrická je využita pro bezpečné zaslání symetrického klíče.

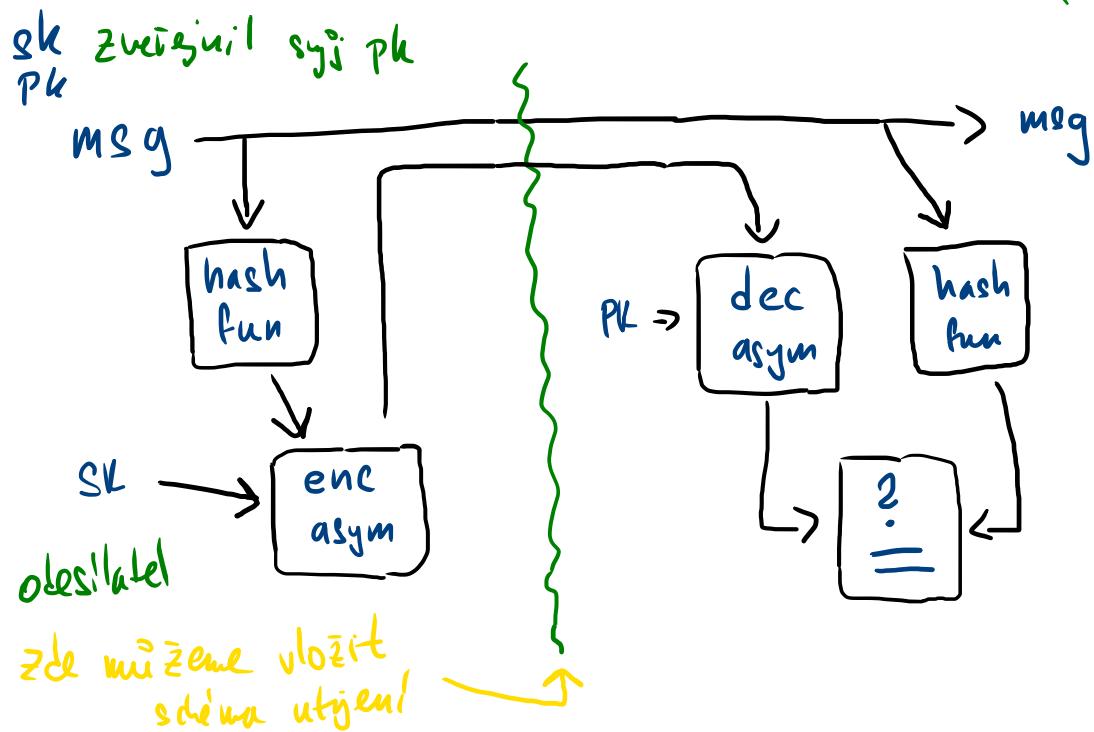
²Problém má více obdobných variant.

Schéma utajení



Obrázek 8.1: Schéma utajení hybridní kryptografie.

Schéma digitálního podpisu



Obrázek 8.2: Schéma digitálního podpisu hybride kryptografie.

8.4 RSA

Algoritmus RSA (Rivest-Shamir-Adleman) lze použít jak pro šifrování dat pro digitální podepisování. Je založen na problému faktorizace velkých čísel.

Klíče Klíče se skládají z:

- p, q – dvě náhodná soukromá prvočísla,
- n – veřejný modul ($n = p \cdot q$),
- e – veřejný exponent ($e < \Phi(n)$ \wedge $GCD(\Phi(n), e) = 1$), typicky 3 nebo $2^{16} + 1^3$,
- d – soukromý exponent,
- musí platit vztah: $e \cdot d \bmod \Phi(n) = 1$.

Veřejný klíč $PK = (n, e)$, soukromý klíč $SK = (n, d)$.

Postup generování Postup generování klíčů:

1. vygenerovat prvočísla p a q ,
2. spočítat modul $n = p \cdot q$,
3. spočítat $\Phi(n) = (p - 1) \cdot (q - 1)$,

³ GCD – největší společný dělitel

4. zvolit veřejný exponent $e < \Phi(n) \wedge GCD(\Phi(n), e) = 1$,
5. spočítat soukromý exponent d tak, že platí $e \cdot d \bmod \Phi(n) = 1$.

Šifrování a dešifrování Mějme zprávu m reprezentovanou jako celé číslo a zašifrovanou zprávu c reprezentovanou také jako celé číslo. Digitální podpis se vytváří stejným způsobem, pouze se prohodí exponenty.

$$c = m^e \bmod n \quad (8.1)$$

$$m = c^d \bmod n \quad (8.2)$$

Útoky a slabiny Pokud útočník rozloží číslo n na činitele p a q , tak může dopočítat soukromý klíč. Pokud útočník uhádně hodnotu $(p-1) \cdot (q-1)$, tak může dopočítat soukromý klíč. Šifrování malých čísel je zranitelné, proto se používá „předzpracování“ – zarovnání na X bitů (2048).

Příklad

1. RSA $n = 143$, $p = 11$, $q = 13$, $e = 7$. Vypočítat d , napsat VK, PK a zašifrovat číslo 9.

$$p = 11$$

$$q = 13$$

$$n = p q = 143$$

$$e = 7$$

$$\Phi(n) = (p-1)(q-1)$$

$$\Phi(143) = 10 \cdot 12 = 120$$

$$e \cdot d \bmod \Phi(n) = 1$$

$$7d \bmod 120 = 1$$

$$d = 103$$

$$VK = (n, e) = (143, 7)$$

$$SK = (n, d) = (143, 103)$$

$$m = 9$$

$$120x + 7y = 1$$

$$gcd(120, 7)$$

$$\begin{aligned} 120 &= 17(7) + 1 \\ 1 &= 120 - 17(7) \end{aligned}$$

$$\begin{aligned} 1 &= 120 - 17(7) \\ -17(7) &= 1 \\ -17 + 120 &= 103 \end{aligned}$$

$$C = m^e \bmod n = 9^7 \bmod 143 = 4782969 \bmod 143 = \underline{\underline{48}}$$

$$m = C^d \bmod n = 48^{103} \bmod 143 = \underline{\underline{9}}$$

Obrázek 8.3: Příklad RSA.

Kapitola 9

KRY – Hašovací funkce, klíčovaný haš a MAC a jejich použití a vlastnosti.

9.1 Metadata

- Předmět: Kryptografie (KRY)
- Přednáška:
 - KRY04_Asym_MNG.pdf
- Záznam:
 - 2021-03-22

9.2 Úvod a kontext

Viz. „Úvod a kontext“ v předchozích otázkách z tohoto předmětu.

Hashovací funkce Hashovací funkce je funkce (resp. algoritmus) pro převod vstupních dat do (relativně) malého čísla. Výstup hashovací funkce se označuje otisk, *fingerprint*, *digest* či *hash*. Jsou jednosměrné a odolné proti kolizím (viz vlastnosti).

Obecné vlastnosti Hashovací funkce by měla:

- Být aplikovatelná na argument o libovolné velikosti.
- Mít výstup konstantní délky.
- Dokázat spočítat výstup rychle.

Neklíčované hashovací funkce Hashovací funkce má pouze jeden argument – data. Např. MD2, MD4, MD5, SHS, SHA1, SHA2, SHA3.

$$f(data) \rightarrow hash$$

Klíčované hashovací funkce Hashovací funkce má dva argumenty – data a klíč. Také se jim někdy říká MAC (*message authentication code*).

$$f(data, key) \rightarrow hash$$

Lavinový efekt Lavinový efekt *Avalanche Effect* je žádoucí vlastností kryptografických algoritmů, typicky blokových šifer a kryptografických hašovacích funkcí, kdy se při nepatrné změně vstupu (například převrácení jednoho bitu) výrazně změní výstup (např. převrátí se polovina výstupních bitů). V případě kvalitních blokových šifer by taková malá změna klíče nebo otevřeného textu měla způsobit drastickou změnu šifrového textu.

9.3 Kryptografická odolnost hashovacích funkcí

Vlastnosti z hlediska odolnosti Hashovací funkce by z hlediska kryptografické odolnosti měly splňovat:

- *First preimage resistance* – Pro konkrétní y je výpočetně nezvládnutelné najít takové x , aby platilo $h(x) = y$. Útočník má k dispozici konkrétní hash, a snaží se pro něho nalézt zprávu.
- *Second preimage resistance* – Pro konkrétní x je výpočetně nezvládnutelné najít takové x' , aby platilo $h(x) = h(x')$. Útočník má k dispozici konkrétní zprávu (nemůže si ji zvolit), ke které se snaží nalézt jinou zprávu, která bude mít stejný hash.
- *Collision resistance* – Je výpočetně nezvládnutelné najít libovolnou dvojici x, x' takovou, aby platilo $x \neq x'$ a $h(x) = h(x')$. Útočník si může zvolit libovolnou zprávu, ke které se snaží nalézt jinou zprávu, která bude mít stejný hash. Pokud platí *collision resistance*, tak platí i *second preimage resistance*.

Narozeninový problém V teorii pravděpodobnosti je narozeninový problém úloha vypočítat minimální početnost skupiny lidí, ve které je alespoň 50% pravděpodobnost nalezení dvojice se stejným datem narození. Narozeninovým paradoxem je pak označována skutečnost, že tento počet (23) je mnohem menší než intuitivní odhad. Výsledek je intuitivnější, když uvážíme, že porovnání narozenin bude provedeno mezi všemi možnými dvojicemi jedinců. Při počtu 23 jedinců je třeba uvažovat $(23 \cdot 22)/2 = 253$ dvojic, což je více než polovina počtu dnů v roce (182, 5).

Jednodušší je nejprve spočítat jev opačný $\bar{p}(n)$, tedy pravděpodobnost, že všech n narozenin je rozdílných. Pro $n > 365$ je 1, jinak:

$$\begin{aligned} \bar{p}(n) &= 1 \cdot \left(1 - \frac{1}{365}\right) \cdot \left(1 - \frac{2}{365}\right) \cdots \left(1 - \frac{n-1}{365}\right) = \\ &= \frac{365 \cdot 364 \cdots (365 - n + 1)}{365^n} = \\ &= \frac{365!}{365^n (365 - n)!} \end{aligned} \tag{9.1}$$

$$p(n) = 1 - \bar{p}(n) \tag{9.2}$$

Narozeninový útok Mějme hashovací funkci, která má n bitový výstup (celkový počet možných hashů je 2^n). Útočník vytvoří dokument „přátelská dohoda“ a přibližně $2^{n/2}$ sémanticky ekvivalentní verzí (úprava bílých znaků, úprava pořadí celků, jiné formulace, ...). Podobně vytvoří dokument „nepřátelská dohoda“ a přibližně $2^{n/2}$ sémanticky ekvivalentní verzí. S pravděpodobností 0,5 existuje verze „přátelské dohody“ a „nepřátelské dohody“, které mají stejný hash. Pokud takové verze existují, útočník dá oběti podepsat „přátelskou dohodu“ \Rightarrow existuje validní podpis „nepřátelské dohody“.

Bezpečnostní cíle OWHF Bezpečnostní cíle OWHF (*one way hash function*). Některé protokoly nevyžadují bezkoliznost, proto má smysl řešit i tento případ.

- Vyžadované vlastnosti: *first preimage resistance* a *second preimage resistance*
- Cíl útočníka: vytvořit *first preimage* nebo *second preimage* (oba úkoly jsou stejně těžké)
- Složitost: $O(2^n)$ (n je počet bitů hashe)
- Požadovaná délka: $n \geq 80$

Bezpečnostní cíle CRHF Bezpečnostní cíle CRHF (*collision resistance hash function*).

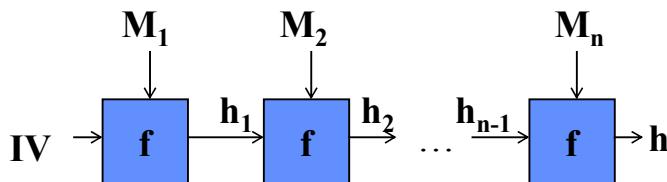
- Vyžadované vlastnosti: *collision resistance*
- Cíl útočníka: vytvořit kolizi
- Složitost: $O(2^{n/2})$ (n je počet bitů hashe) (kvůli narozeninovému útoku)
- Požadovaná délka: $n \geq 160$

Bezpečnostní cíle MAC Bezpečnostní cíle MAC (*message authentication code*).

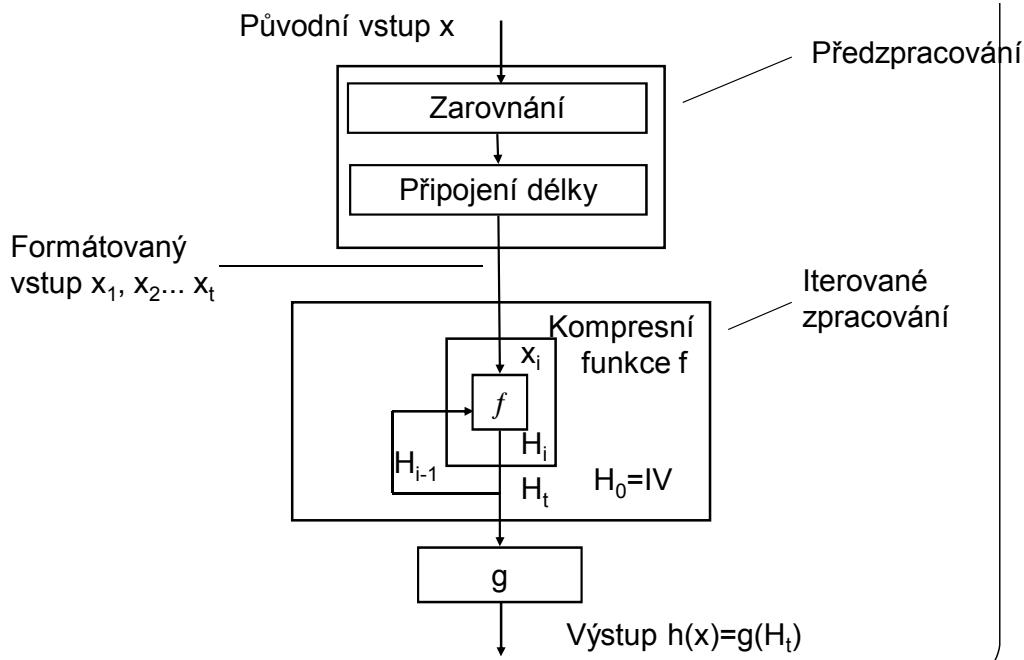
- Vyžadované vlastnosti: *computation resistance, key non-recovery*
- Cíl útočníka (útočník si může vybrat):
 - Vytvořit nový hash, který bude odpovídat nové zprávě
 - Nalézt klíč
- Složitost (n je počet bitů hashe, t je počet bitů klíče):
 - Vytvořit nový hash: $O(\max(2^{-n}, 2^{-t}))$
 - Nalézt klíč: $O(2^n)$
- Požadovaná délka: $n \geq 64 \wedge t \geq 64$

9.4 Hashovací funkce neklíčované

Nejčastější způsoby sestrojení hashovací funkce neklíčované jsou založené na principu iterace.



Obrázek 9.1: Schéma iterativní neklíčované hashovací funkce. Zpráva je rozdělena na n částí. f je tzv. kompresní funkce. IV je inicializační vektor, resp. konstanta. h_1 až h_{n-1} jsou mezivýsledky („mezihashe“) a h je výsledný hash.



Obrázek 9.2: Podrobnější schéma iterativní neklíčované hashovací funkce.

Jednotlivé kroky hashovací funkce:

- Předzpracování – Vstupní data jsou rozdělena na bloky o stejné délce. Je provedeno zarovnání posledního bloku. Je připojena informace o délce zprávy.
- Iterativní zpracování – V iteracích se postupně „přihashovávají“ vstupní bloky. Zpětná vazba pomocí stavové proměnné. Uvnitř kompresní funkce, která z delšího vstupu udělá kratší výstup.
- Postzpracování – Volitelný krok.

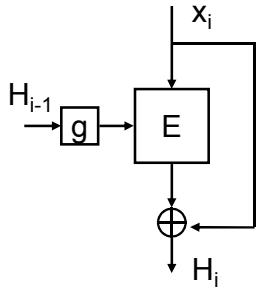
Merkelova meta-metoda Necht' f je kompresní funkce odolná proti kolizím. Hashovací funkce h na principu iterace využívající kompresní funkci f je rovněž odolná proti kolizím.

Merkel-Damgardovo zesílení Pokud je do vstupu hashovací funkce vložena délka zprávy, tak je zajištěno, že žádná zpráva není prefixem jiné zprávy.

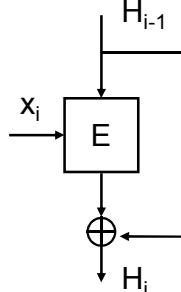
Zarovnání Nejednoznačné zarovnání (*ambiguous padding*) – připoj ke zprávě tolik bitů, aby délka zprávy byla násobkem délky bloku. Jednoznačné zarovnání (*unambiguous padding*) – připoj ke zprávě jeden bit a poté proved' nejednoznačné zarovnání.

9.4.1 Hashovací funkce s využitím blokových šifer

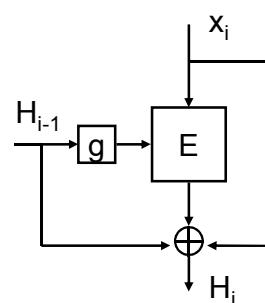
Alternativně lze využít pro konstrukci hashovacích funkcí blokové šifry. Avšak blokové šifry byly navrhovány pro jiný režim činnosti, kterém útočník nezná klíč (a není schopen ho ovlivnit), zná pouze šifrovaný text (ten je schopen ovlivnit). V tomto případě útočník může přímo ovlivňovat hodnoty klíče.



Matyas-Meyer-Oseas



Davies-Meyer

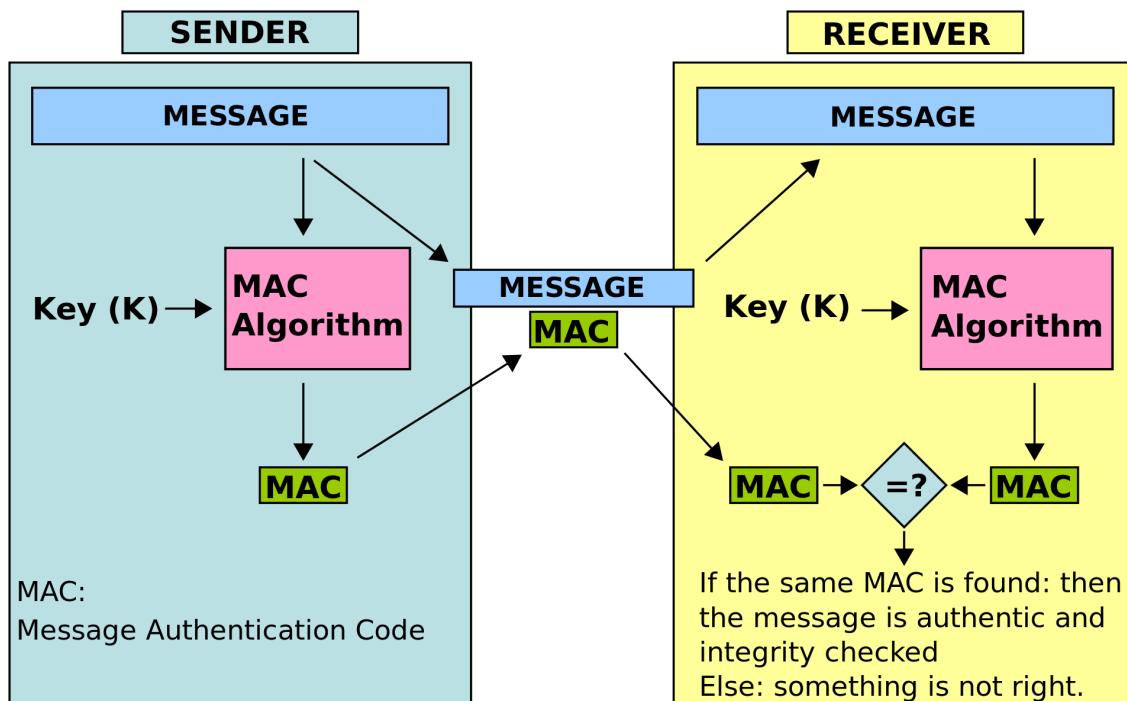


Miyaguchi-Preneel

Obrázek 9.3: Ukázka několika možných způsobů využití blokových šifer pro konstrukci kompresní funkce. S využitím iteračního způsobu lze zobecnit pro celou hashovací funkci.

9.5 MAC (*message authentication code*)

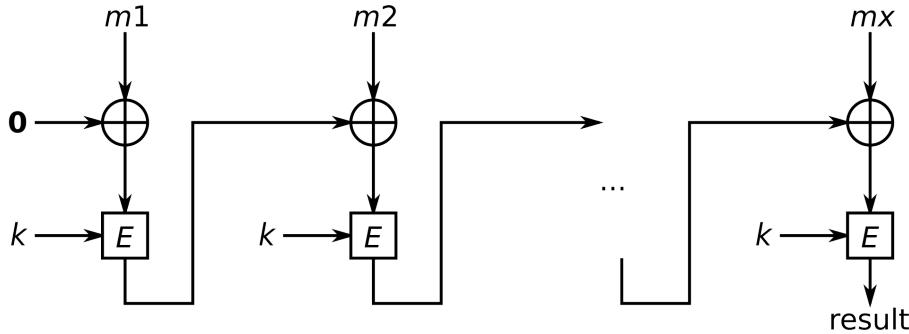
- Rodina hashovacích funkcí h_k , které jsou parametřitelné klíčem k .
- Vlastnosti (stejné jako u obecných hashovacích funkcí, pouze rozšířené o klíč):
 - Výstup $h_k(x)$ lze spočítat rychle, pokud je znám klíč k .
 - Jsou výpočteně bezpečné – při znalosti dvojice $(x, h_k(x))$ je výpočteně nemožné spočítat novou dvojici $(x', h_k(x'))$, pro $x \neq x'$, pokud není znám klíč.
- Využití: zajištění autentizace a integrity (nepopiratelnost zajistit nedokáže).



Obrázek 9.4: Schéma použití MAC. Pokud je stejný MAC výpočítán na straně příjemce, tak má jistotu, že zpráva nebyla po cestě změněna a že zprávu poslal skutečně odesílatel.

9.5.1 Sestrojení MAC pomocí blokové šifry v CBC

Pro sestrojení MAC hashovací funkce je využita symetrická bloková šifra v režimu CBC (*cipher block chaining*, šifrová zpětná vazba). Rozdíl oproti CBC šifrování spočívá v tom, že mezivýsledky se zahazují a pracuje se až s posledním blokem. Z něho se vezme určitý počet posledních bitů (podle požadované délky hashe – 32, 48, 64) a ten tvoří výsledný hash (MAC).



Obrázek 9.5: Ukázka sestrojení MAC funkce pomocí symetrické blokové šifry v režimu CBC.

Proč stačí výrazně menší délka? Protože klíč. Útočník sice může vyzkoušet všechny možné hashe, ale bez znalosti klíče, nezjistí, který je ten správný.

9.5.2 Sestrojení MAC pomocí neklíčované hashovací funkce

Pro sestrojení MAC hashovací funkce je využita neklíčovaná hashovací funkce. Klíč je připojen ke zprávě a je použita standardní hashovací funkce.

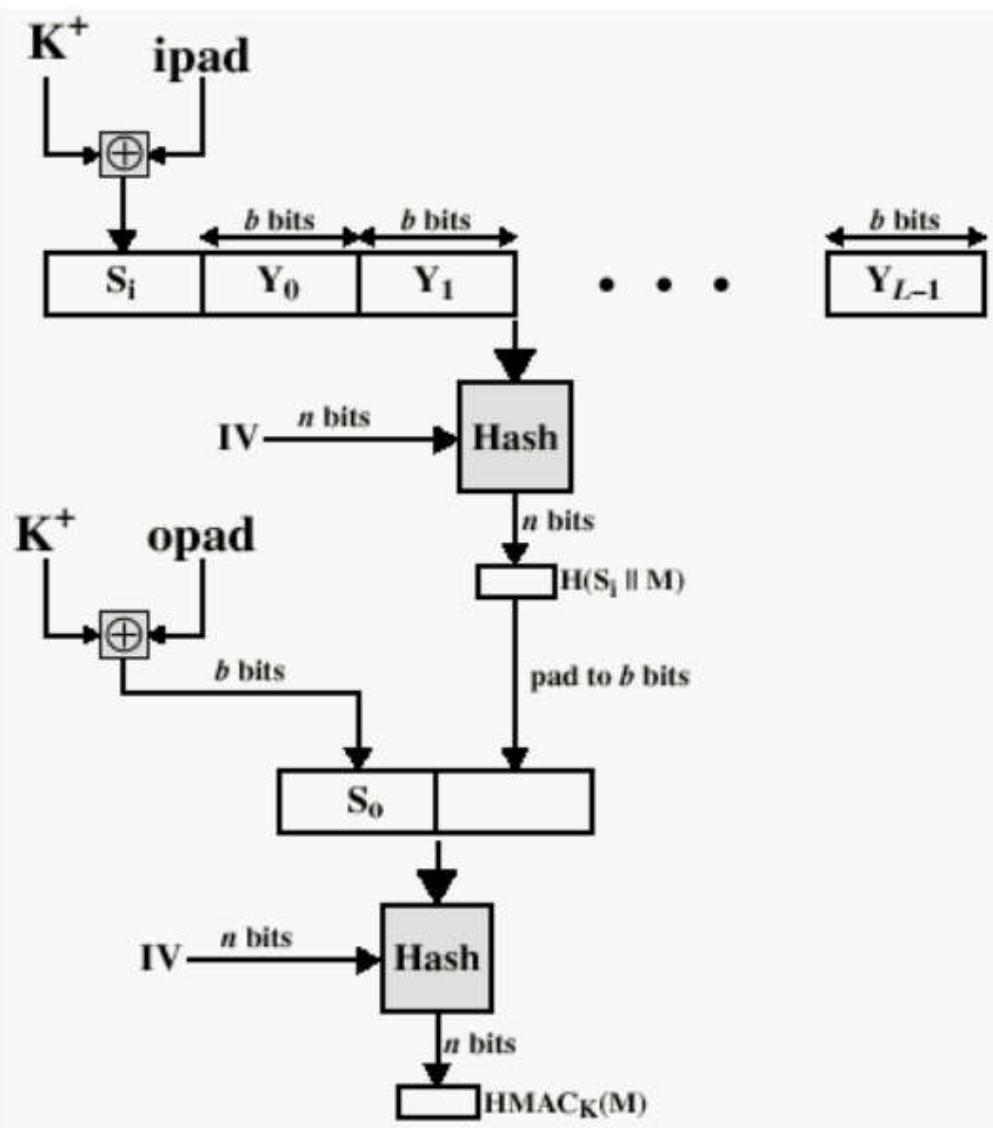
Secret prefix Klíč je přidán na začátek zprávy. Formálně: $H = h(k||x)$, kde H je výsledný hash (MAC), h je hashovací funkce, k je klíč a x je zpráva. Útočník může libovolně „při-hashovávat“ další bloky bez znalosti klíče a tím vytvářet nové validní hashe – $h(k||x||y)$, kde y je útočníkova zpráva \Rightarrow nepřijatelný způsob.

Secret suffix Klíč je přidán na konec zprávy. Formálně: $H = h(x||k)$. Útočník, který může zvolit x , může také vytvořit x' , pro které $h(x) = h(x')$ se složitostí $O(2^{n/2})$, kde n je délka hashe, bez ohledu na délku klíče k (narozeninový útok) \Rightarrow nepřijatelný způsob (útok který je neovlivnitelný délkou klíče).

Enveloping Klíč je přidán na začátek i na konec zprávy. Formálně: $H = h(k||p||x||k)$, kde p je zarovnání. Přijatelný způsob. Základ pro algoritmus HMAC.

9.5.3 HMAC (*hash function MAC*)

HMAC (*hash function MAC*) je do dnes používaný algoritmus. Specifikuje použití metody enveloping, ale ne, která hashovací funkce se použije.



Obrázek 9.6: Schéma HMAC; **ipad** a **opad** jsou vstupní/výstupní konstanty, které slouží k zarovnání; Y_i jsou bloky vstupní zprávy; IV je inicializační vektor.

Kapitola 10

KRY – Správa klíčů v asymetrické kryptografii (certifikáty X.509).

10.1 Metadata

- Předmět: Kryptografie (KRY)
- Přednáška:
 - KRY04_Asym_MNG.pdf
- Záznam:
 - 2021-03-29

10.2 Úvod a kontext

Viz. „Úvod a kontext“ v předchozích otázkách z tohoto předmětu.

Problém se zveřejňováním veřejných klíčů Jak můžu vědět, že publikovaný veřejný klíč patří opravdu entitě, které patřít má? Je potřeba zajistit autenticitu (pravost) veřejných klíčů – Vytvořit spolehlivou vazbu mezi veřejným klíčem a jménem jeho vlastníka.

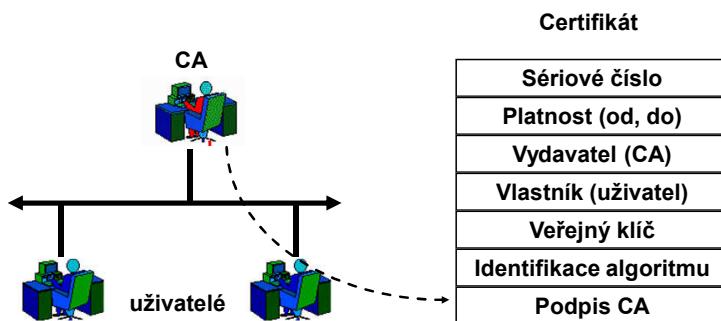
Systémy založené na veřejném klíči Systémy založené na veřejném klíči (PKI, *Public Key Infrastructure*) je označení infrastruktury správy a distribuce veřejných klíčů. PKI umožňuje pomocí přenosu důvěry používat cizí veřejné klíče a ověřovat jimi elektronické podpisy bez nutnosti jejich individuální kontroly.

10.3 Správa klíčů v asymetrické kryptografii

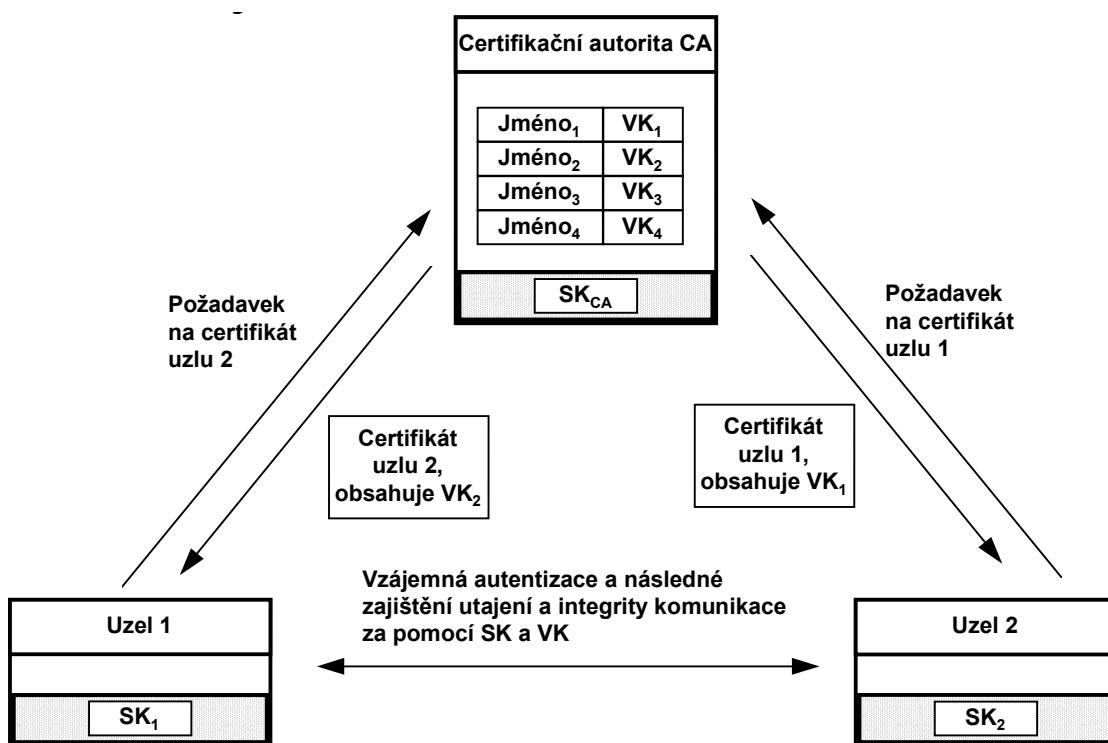
Certifikát Certifikace veřejného klíče. Nějaký prostředník (certifikační autorita), kterému důvěřujeme, se zaručuje, že konkrétní veřejný klíč, patří dané entitě.

Certifikační autorita Certifikační autorita (CA) je prostředník, který distribuuje certifikáty a které všichni důvěřují. CA negeneruje klíče uživatelům, ty si je generují samy.

Proces certifikace klíče CA podepíše veřejný klíč uživatele a jeho dalsí údaje (jméno, doba vydání, doba platnosti, ...) svým soukromým klíčem. Tyto podepsané údaje se nazývají certifikát.



Obrázek 10.1: Příklad certifikátu.



Obrázek 10.2: Příklad navázání bezpečné komunikace mezi dvěma entitami, které mají stejnou certifikační autoritu.

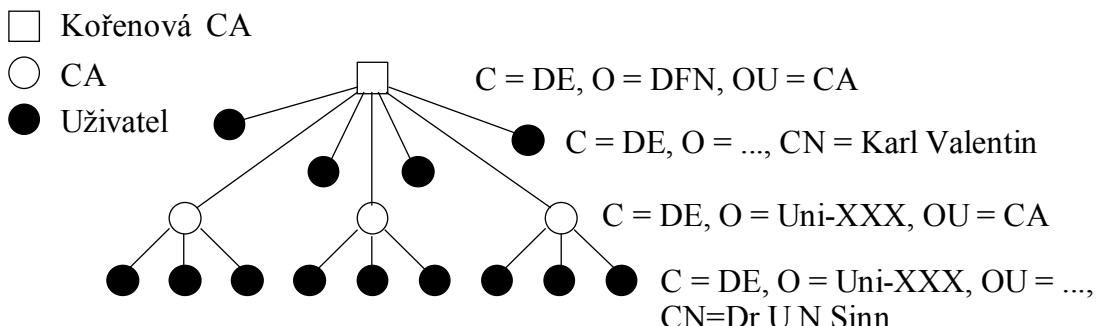
Popis navázání bezpečné komunikace (viz obrázek 10.2):

1. Uzel 1 si vygeneruje soukromý a veřejný klíč.
 2. Uzel 1 odešle veřejný klíč certifikační autoritě spolu se svým jménem (a dalšíma informacemi).
 3. CA vytvoří certifikát pro uzel 1 – svým soukromým klíčem podepíše veřejný klíč a jméno uzlu 1. CA odešle certifikát uzlu 1. CA odešle svůj veřejný klíč uzlu 1.
 4. Pokud uzel 2 chce také odesílat, provede také kroky 1-3.

5. Uzel 1 podepíše soubor a odešle ho uzlu 2 (soubor a podpis).
6. Uzel 2 si musí sehnat certifikát uzlu 1. Existují 3 způsoby jak to udělat.
 - Odesílatel zašle svůj certifikát společně se zprávou.
 - Příjemce si vyžádá certifikát odesílatele od certifikační autority.
 - Příjemce si vyžádá certifikát odesílatele od jiné služby (adresářové služby, LDAP).
7. Uzel 2 ověří podpis u certifikátu uzlu 1 veřejným klíčem certifikační autority.
8. Uzel 2 ověří podpis souboru pomocí veřejného klíče odesílatele (který je v certifikátu).

Strom certifikačních autorit Model s jednou globální CA je nemožný (příliš mnoho uživatelů, příliš velké vzdálosti, ...). Proto se používá strom certifikačních autorit. Veřejný klíč CA je certifikován jinou CA. CA nejvýše ve stromu se nazývá **kořenová certifikační autorita**.

- Certifikační autorita má svůj vlastní certifikát, který je podepsaný její certifikační autoritou.
- Koncový uživatel důvěruje stále pouze jedné entitě – kořenové certifikační autoritě, ale přibývá jedna úroveň ověřování navíc.
- Příjemce dostane zprávu s podpisem. Musí znát certifikát odesílatele (podepsaný CA), certifikát certifikační autority (podepsaný CA_{root}) a veřejný klíč kořenové CA¹.
- Úrovní certifikačních autorit může být více (nejčastěji 1-2).



Obrázek 10.3: Příklad stromu certifikačních autorit. C, O, OU je identifikátor entity.

Certifikační cesta Posloupnost certifikátů od certifikátu kořenové CA přes certifikáty dalších CA až k certifikátu komunikující protistrany.

Zneplatnění certifikátu Jak zrušit platnost certifikátu? Normálně se zruší sám, až skončí jeho platnost. Pokud je potřeba certifikát zneplatnit před jeho vypršením je třeba využít tzv. revokační seznam (CRL, *certificate revocation list*). Důvody zneplatnění certifikátu:

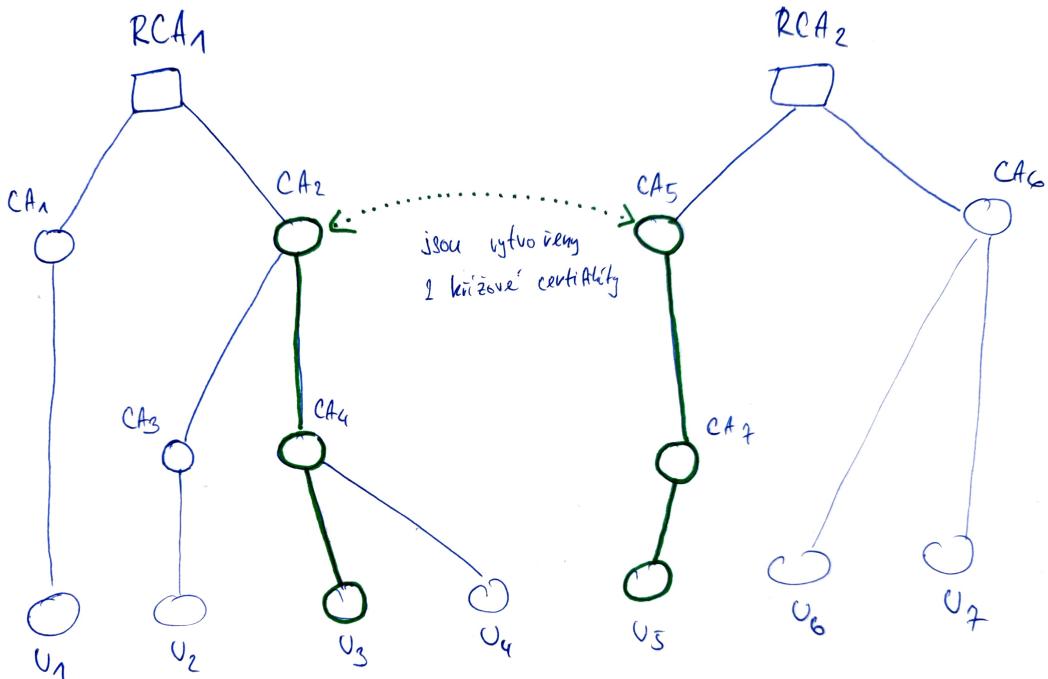
- soukromý klíč uživatele byl kompromitován,
- uživatel ztratil práva, která z certifikátu vyplývají (např. změna zaměstnavatele),
- soukromý klíč CA byl kompromitován (nikdy se nestalo).

¹Veřejný klíč kořenové certifikační autority se z praktických distribuuje ve formě „fiktivního certifikátu“

CRL CRL (*certificate revocation list*) je seznam zneplatněných certifikátů, takových, kterým ještě nevypršela platnost, ale je třeba je zneplatnit. CRL je podepsán CA, které ho spravuje a periodicky aktualizuje (může se zkracovat i růst). Jak se distribuuje:

- *Pull model* – Příjemce certifikátu si dle potřeby stáhne CRL od CA.
- *Push model* – CA pravidelně posílá CRL příjemcům certifikátu.

Křížový certifikát V případě, že spolu komunikují uživatelé, kteří nemají společnou kořenovou CA (jejich CA si nedůvěřují), je třeba využít tzv. křížový certifikát. Tím se komplikuje sestavování certifikační cesty, protože je třeba zvážit všechny křížové certifikáty, které jsou k dispozici.



Obrázek 10.4: Příklad křížového certifikátu. Uživatel U_3 chce navázat bezpečné spojení s uživatelem U_5 . Nemají společnou kořenovou CA, proto je třeba využít křížové certifikáty. CA_2 vytvoří křížový certifikát pro CA_5 a CA_5 vytvoří křížový certifikát pro CA_2 . Příklad: U_3 pošle podepsanou zprávu U_5 , jak bude vypadat certifikační cesta?
 $U_3 \leftarrow CA_4 \leftarrow CA_2 \leftarrow CA_5 \leftarrow RCA_2$.

10.4 Standard X.509

X.509 je standard pro systémy založené na veřejném klíči (PKI). Specifikuje formát certifikátů, formát CRL, parametry certifikátů, metody kontroly platnosti certifikátů, ...

```

1 Certificate ::= SIGNED SEQUENCE {
2     version [0] Version DEFAULT v1988,
3     serialNumber CertificateSerialNumber,
4     signature
5     AlgorithmIdentifier,
6     issuer
7     Name,
8     validity
9     Validity,
10    subject
11    Name,
12    subjectPublicKeyInfo SubjectPublicKeyInfo
13 }
14
15 Version ::= INTEGER {v1988(0) }
16
17 CertificateSerialNumber ::= INTEGER
18
19 Validity ::= SEQUENCE {notBefore UTCTime, notAfter UTCTime }
20
21 SubjectPublicKeyInfo ::= SEQUENCE {
22     algorithm
23     AlgorithmIdentifier,
24     subjectPublicKey
25     BIT STRING
26 }
27
28 AlgorithmIdentifier ::= SEQUENCE {
29     algorithm
30     OBJECT IDENTIFIER,
31     parameters
32     ANY DEFINED BY algorithm OPTIONAL
33 }

```

Výpis 10.1: Příklad definice certifikátu ve formátu X.509.

Význam položek Význam položek v definici certifikátu ve formátu X.509:

- Version – Standardně 0.
- Serial number – Sériové číslo certifikátu, spolu se jménem vydavatele jednoznačně identifikuje certifikát.
- Issuer – Jméno vydávající CA.
- Subject – Jméno vlastníka certifikátu.
- Validity – Doba platnosti certifikátu (`notBefore`, `notAfter`). Podpis je platný pouze pokud je datum podepsání v intervalu platnosti každého z certifikátů z certifikační cesty.
- SubjectPublicKeyInfo – Veřejný klíč vlastníka certifikátu a algoritmus, pro který je určen.
- Signature – Jakým algoritmem je certifikát podepsaný CA.

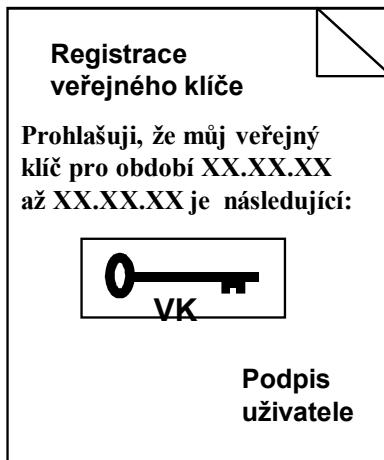
Prototypový certifikát Má strukturu certifikátu X.509. Uživatel si vygeneruje tzv. prototypový certifikát, který má standardní strukturu a vyplní informace na nějaké implicitní

hodnoty. Prototyp pošle CA spolu se svým veřejným klíčem, která certifikát dovyplní, podepíše a pošle zpět.

Registrační autorita Pokud chce uživatel vydat certifikát, kontaktuje tzv. registrační autoritu (nikoliv přímo CA).

Míra důvery v certifikát V praxi chceme více urovní důvěry, než pouze ostrý/žádný (např. chceme vytvořit testovací certifikát). To je řešeno jako rozšíření X.509 přidáním třídy certifikátu (*certification class*). Uživatel chce vydat certifikát od CA jisté třídy.

- Třída 1 – CA vůbec nekontroluje identitu žadatele. Lze jej získat anonymně. Používá se pro testovací certifikáty.
- Třída 2 – Identita žadatele musí být ověřena třetí stranou (notářsky ověřený formulář zasláný poštou).
- Třída 3 – Standardní certifikát. Žadatel musí osobně navštívit CA (resp. registrační autoritu). Osobní ověření totožnosti.
- Třída 4 – Stejně jako 3 a navíc je nutné prokázat oprávnění žadatele požadovat certifikát.



Obrázek 10.5: Příklad žádosti o certifikát.

Kapitola 11

PDS – Základní architektury přepínačů, algoritmy pro plánování, řešení blokování, vícestupňové přepínací sítě.

11.1 Metadata

- Předmět: Přenos dat, počítačové sítě a protokoly (PDS)
- Přednáška:
 - 04-switching.pdf
- Záznam:
 - 2021-03-05

11.2 Úvod a kontext

ISO/OSI model Referenční model ISO/OSI se používá jako názorný příklad řešení komunikace v počítačových sítích pomocí vrstevnatého modelu, kde jsou jednotlivé vrstvy nezávislé a snadno nahraditelné.

- **Aplikační vrstva** (L7, *application layer*)
 - Zajišťuje zpracování dat na nejvyšší úrovni (reprezentace dat, kódování, řízení dialogu, ...).
 - Tvořena procesy a aplikacemi, které komunikují po síti.
 - Bývá slučována s prezentační vrstvou (L6, prezentace dat a šifrování) a relační vrstvou (L5, koordinace a komunikace).
 - Příklad protokolů:
 - * Uživatelské – vykonávají služby přímo uživateli (Telnet, SSH, FTP, SMTP, HTTP, ...)
 - * Systémové – zajišťují sítové funkce (DNS, DHCP, SNMP, BOOTP, ...)
- **Transportní vrstva** (L4, *transport layer*)

- Rozděluje aplikační data (segmentace) na menší jednotky a zapouzdřuje je do segmentů (TCP) / datagramů (UDP).
- Vytváří logické spojení mezi procesy (přenáší data konkrétní aplikace ze zdrojového zařízení do aplikace na cílovém zařízení).
- Adresace: porty.
- Příklad protokolů: TCP, UDP, DCCP, SCTP, MP-TCP, QUIC

- **Sít'ová vrstva** (L3, *network layer*)

- Zapouzdřuje segmenty/datagramy do paketů.
- Řeší směrování.
- Adresace: IP adresa (logická adresa).
- Příklad protokolů: IPv4, IPv6, ARP, RARP, ICMP,IGMP

- **Linková vrstva** (L2, *data link layer*, vrstva síťového rozhraní, *network interface layer*)

- Zapouzdřuje pakety do rámčů.
- Zajišťuje *hop-by-hop* doručení.
- Adresace: MAC adresa (fyzická adresace).
- Příklad protokolů: Ethernet, Token Ring, FDDI, X.25, Frame Relay

- **Fyzická vrstva** (L1, *physical layer*)

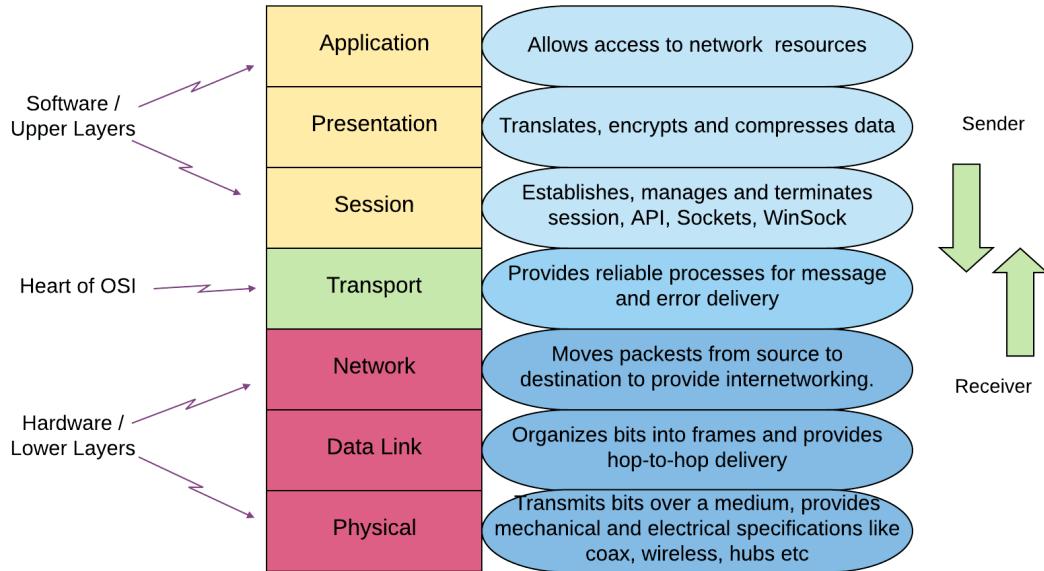
- Zajišťuje přenos bitů přes fyzické médium.

data, segment, datagram, paket, rámec, bit

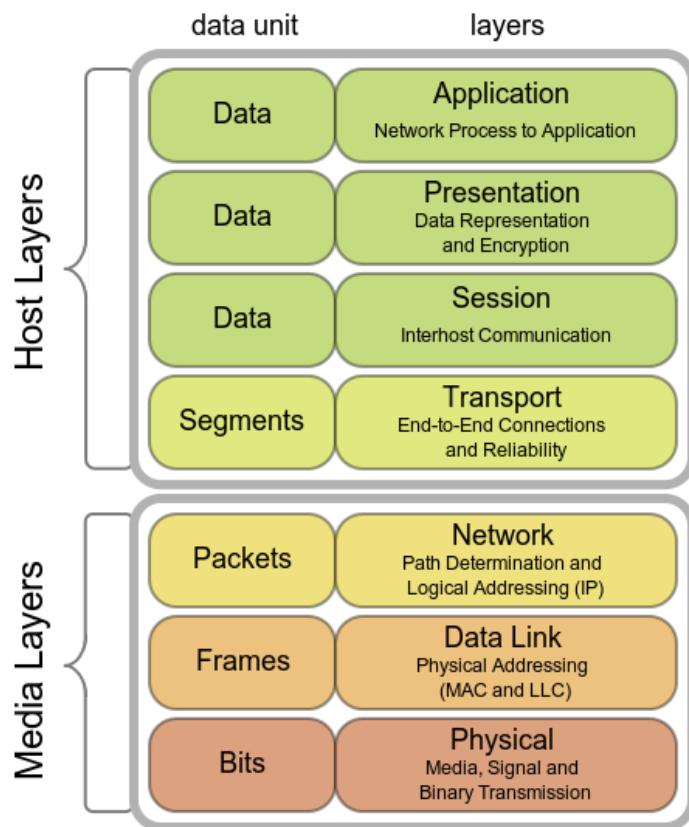
- Data – aplikační vrstva (L7)
- Segment – transportní vrstva (L4), TCP
- Datagram – transportní vrstva (L4), UDP
- Paket – sít'ová vrstva (L3)
- Rámec – linková vrstva (L2)
- Bit – fyzická vrstva (L1)

CAM tabulka CAM tabulka je datová struktura v přepínači, která uchovává informace o tom, za jakým fyzickým portem je zařízení s jakou MAC adresou. Jak se plní? Přepínač si ji plní automaticky. V momentě kdy k němu přijde paket, doplní si MAC adresu a port přes který paket přišel. Pokud přepínač neví za jakým portem je cílová MAC adresa, pošle to všem (*flooding*).

Propustnost Kolik je možné přenést dat za časovou jednotku.



Obrázek 11.1: Příklad OSI modelu z Linuxhint.



Obrázek 11.2: Příklad OSI modelu z Wiki.

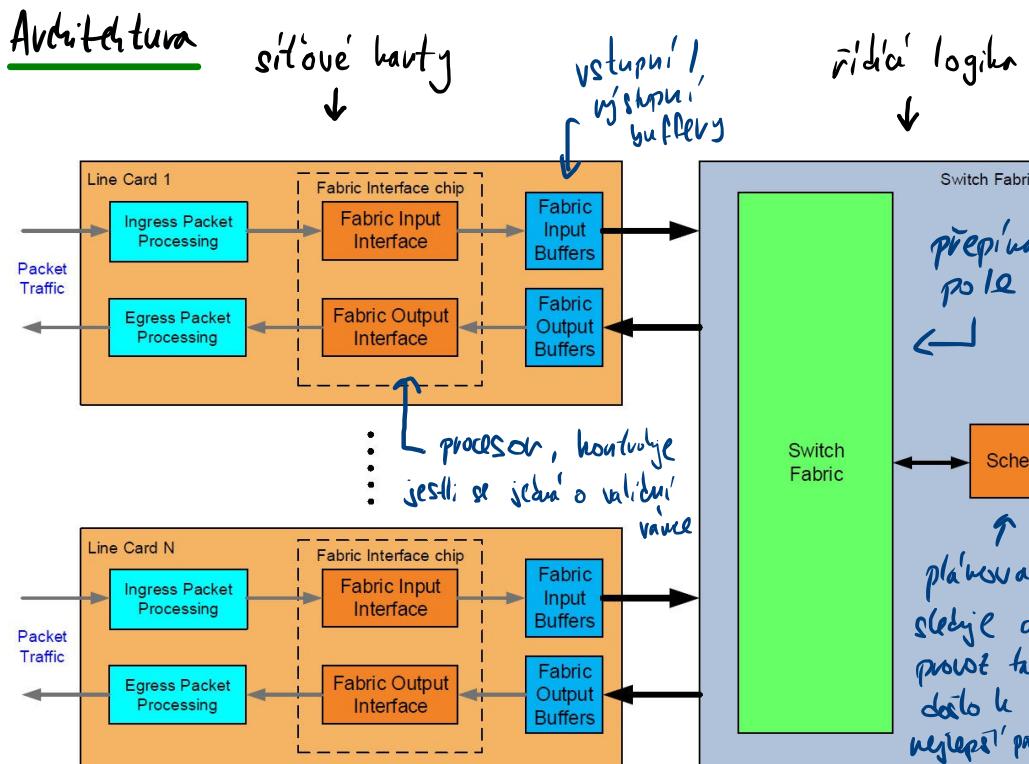
11.3 Obecná architektura přepínače

Sít'ový přepínač (*switch*) je aktivní prvek v počítačové síti, který propojuje jednotlivé prvky do hvězdicové topologie. Přepínač obsahuje sít'ové porty (až stovky), na něž se připojují sít'ová zařízení.

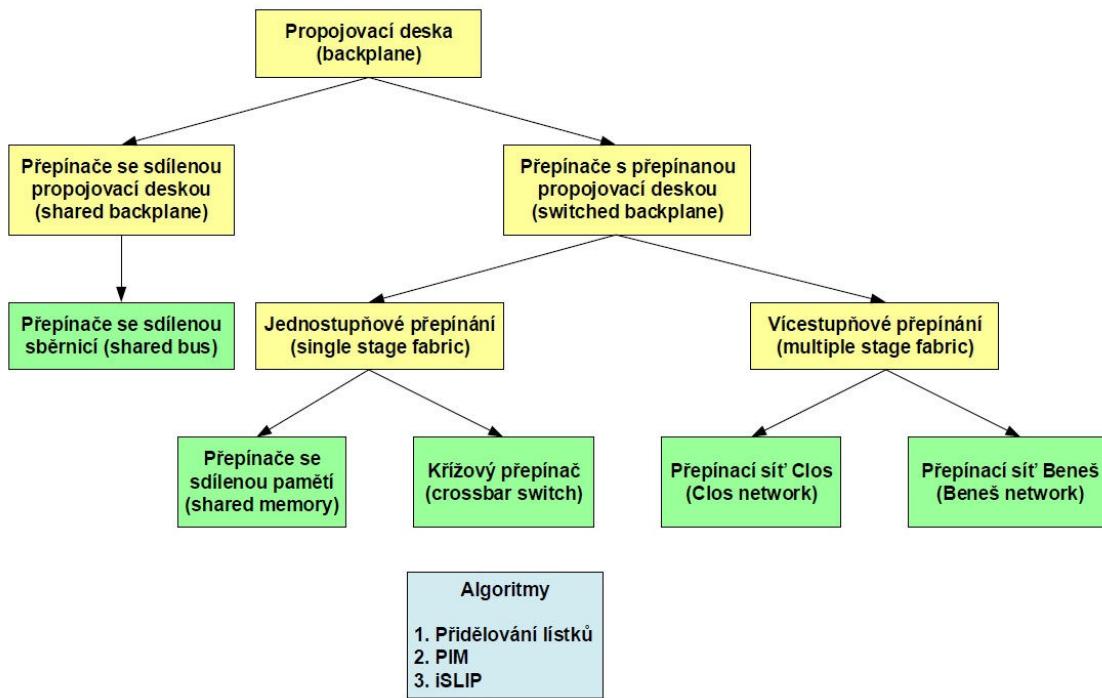
- Na jaké vrstvě OSI modelu pracuje? Pracuje s rámci (linková vrstva, L2).
- Na základě čeho provádí přepínání? Na základě cílové MAC adresy.
- Jaké typy přenosů přepíná? Na základě cílové MAC adresy rozslíšuje typ přenosu:
 - broadcast (samé 1),
 - multicast (speciální prefix pro IPv4 a IPv6),
 - unicast (cokoliv jiného).
- Co ovlivňuje rychlosť přepínání?
 - Hardware (typ paměti, rychlosť procesoru)
 - Logika přepínání
 - Šířka pásma cílového rozhraní

Požadavky na přenos

- Maximální využití sběrnice (požadavek na plánovač).
 - Aby došlo k maximálnímu přenosu dat přepínací logikou (co nejvíce přenosů v rámci taktu – „parallelizace“).
- Spravedlivé přidělování přenosového pásma (aby byly obsluhovány všechny porty).
- Zachování pořadí rámců (ne vždy je možné).



Obrázek 11.3: Obecná architektura přepínače.



Obrázek 11.4: Dělení přepínačů.

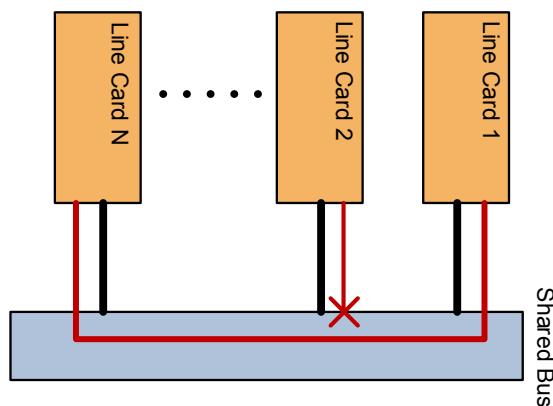
11.4 Přepínače se sdílenou propojovací deskou

11.4.1 Přepínače se sdílenou sběrnicí

- Může komunikovat pouze jeden port v daný čas, ostatní čekají (řízeno protokolem).
- Lze jednoduše realizovat broadcast a multicast.
- Mějme N portů (karet), rychlosť každé karty R bps, taktovací frekvenci r . Pak rychlosť sběrnice musí být $N \times R$ a šířka sběrnice

$$W = \frac{R \times N}{r}$$

- Pokud se přidají nové porty nebo se na nich zvýší rychlosť je potřeba zvýšit prospustnost sběrnice („něco za něco“).



Obrázek 11.5: Přepínač se sdílenou sběrnicí.

* Úloha č.1:

- Přepínač má 16 portů o rychlosti 100 Mb/s. Jaká je potřebná propustnost sběrnice a jaká musí být šířka sběrnice, aby zvládl požadovaný přenos? Uvažujte taktovací frekvenci sběrnice 40 MHz.

$$r = 40 \text{ MHz}$$

$$N = 16$$

$$R = 100 \text{ Mb/s}$$

propustnost - kolik dat za časovou jednotku

$$\hookrightarrow N \times R = 1,6 \text{ Gb/s}$$

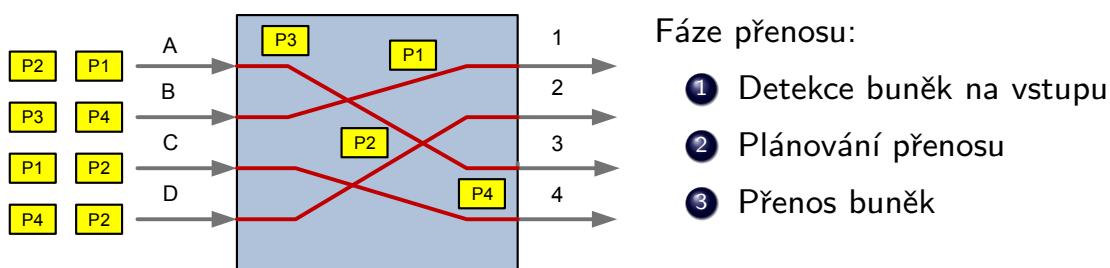
šířka sběrnice - kolik v jeden okamžik maximálně

$$\hookrightarrow \frac{N \times R}{r} = \frac{16 \cdot 100 \cdot 10^6}{40 \cdot 10^6} = 40 \text{ b}$$

Obrázek 11.6: Přepínač se sdílenou sběrnicí – příklad.

11.5 Přepínače s přepínanou propojovací deskou: jednostupňové přepínání

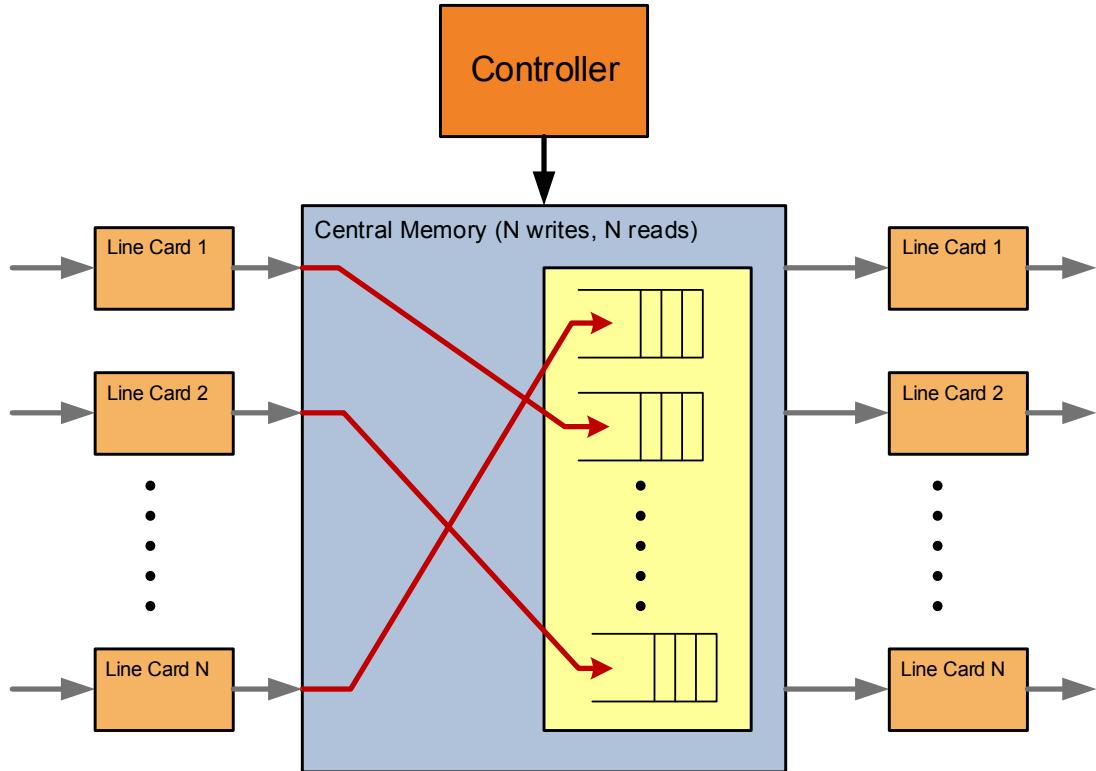
- Umožňuje paralelní přenos rámců.
- Je potřeba plánovač (*scheduler*), který plánuje přenos.
- Provedení přenosu: z bufferů vstupních karet se data přenesou do bufferů výstupních karet.



Obrázek 11.7: Činnost přepínače s přepínanou propojovací deskou.

11.5.1 Přepínače se sdílenou pamětí

- Centrální sdílená paměť, která obsahuje frontu pro každý výstupní buffer (každý port).
- Pro jeden přenos se musí 2x přistoupit do paměti (zápis a čtení).
- Zde narázíme na současné technologické limity paměti (rychlosť přístupu).



Obrázek 11.8: Činnost přepínače s přepínanou propojovací deskou a se sdílenou pamětí.
Obsahuje: N vstupních karet, řadič, sdílenou paměť s N frontami, N výstupních karet.

- Necht' R je rychlosť síťové kart, N je počet síťových karet, pak celková propustnosť přepínače je $BW = 2 \times N \times R$ (dva přístupy do paměti).
- Mějme data o velikosti C , pak doba přenosu je $t = C \div BW[s]$.

Úloha č.2:

- Jaká je potřebná doba přístupu do paměti u přepínače se sdílenou pamětí, který má 32 portů o rychlosti 1 Gb/s a velikost ukládané buňky je 40 bytů? Jak se tato doba změní pro rychlosť portů 10 Gb/s a 40 Gb/s?

$$C = 40 \text{ B} \quad N = 32 \quad R = 1 \text{ Gb/s}$$

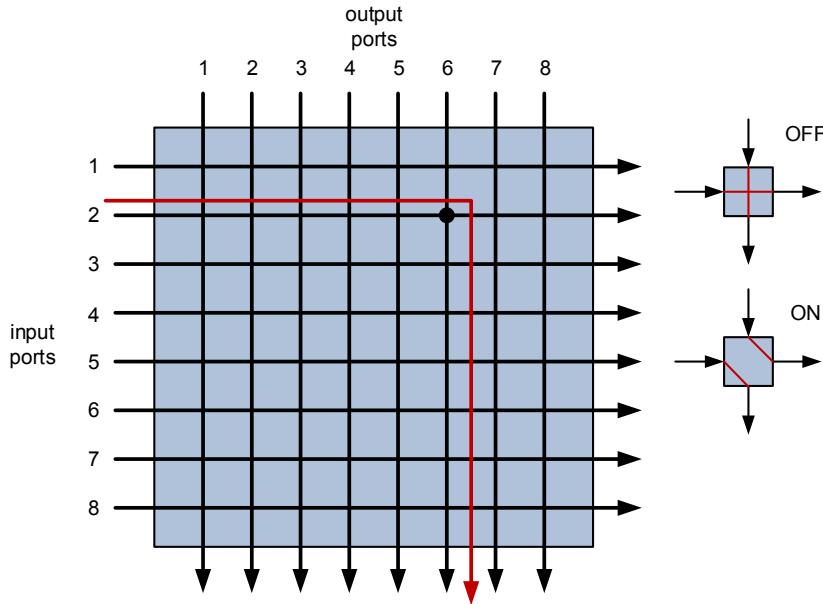
$$t = \frac{40 \cdot 8}{2 \cdot 32 \cdot 1 \cdot 10^9} = 5 \cdot 10^{-9} \text{ s} = 5 \text{ ns}$$

pro 10 $\rightarrow 0,5 \mu\text{s}$
 pro 40 $\rightarrow 0,125 \mu\text{s}$

Obrázek 11.9: Přepínač s přepínanou sběrnicí a se sdílenou pamětí – příklad.

11.5.2 Křížové přepínače (*crossbar*)

- Založený na jednom velkém propojovacím poli.
- Interně neblokující.
- Nativní podpora multicastu.
- Je potřeba N^2 propojení.
- Chybí redundance: pouze jedna cesta mezi vstupem a výstupem.



Obrázek 11.10: Přepínač s přepínánou sběrnicí a křížovými propoji.

11.6 Hledání spojů v křížovém přepínači

11.6.1 Problém párování

V kontextu přepínání v křížovém přepínači se jedná o problém párování v bipartitním grafu, kde množina vstupních portů je V_1 , množina výstupních portů je V_2 a hrany E jsou propoje vstupů a výstupů. Množina M se nazývá párování, pokud $M \subseteq E$, kde žádné dvě hrany z M , nemají společný vrchol.

Bipartitní graf Necht' $G = (V, E)$ je graf. Graf G je bipartitní, pokud platí $V = V_1 \cup V_2 \wedge V_1 \cap V_2 = \emptyset \wedge \forall e = \{u, v\}, e \in E : u \in V_1 \wedge v \in V_2$. Platí li navíc $E = V_1 \times V_2$, pak je graf úplný bipartitní.

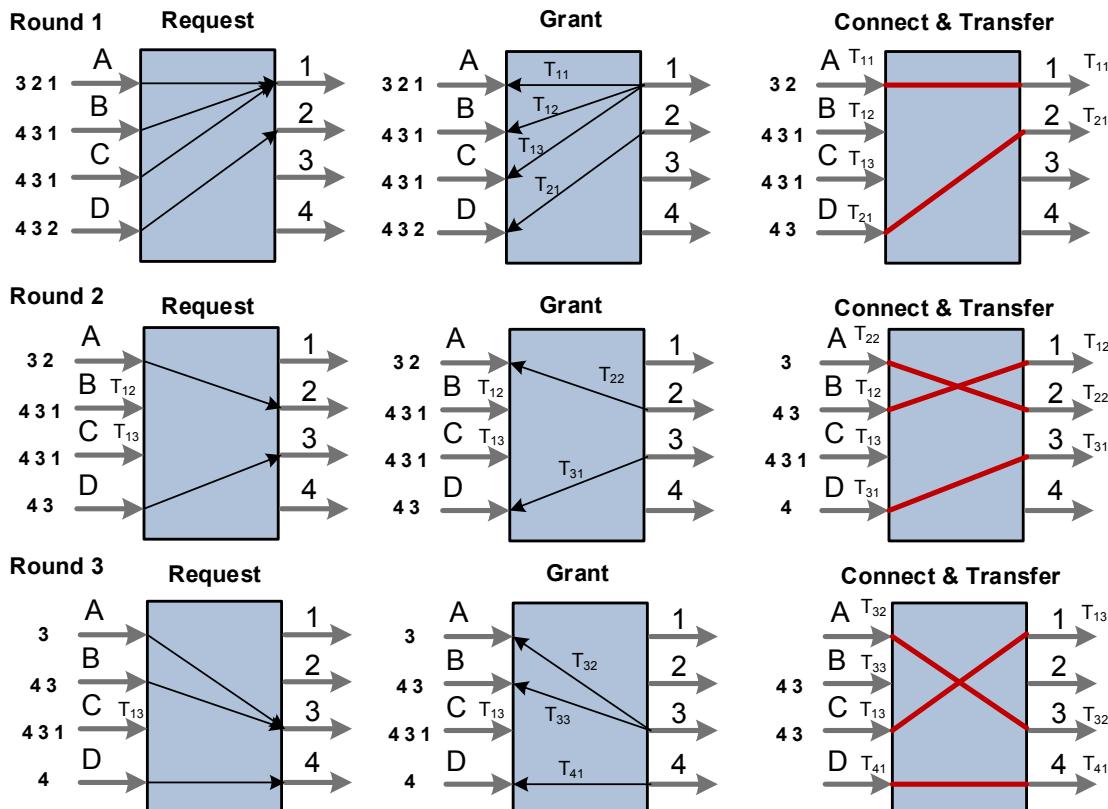
Největší párování Největší párování (*maximum matching*) je párování, které má největší počet hran (globální maximum). Složitost: $O(N^{5/2})$.

Maximální párování Maximální párování (*maximal matching*) je párování, kdy nelze přidat žádnou další hranu (lokální maximum). Složitost: $O(N + E)$.

Jelikož algoritmy pro hledání maximálního párování mají výrazně menší časovou složitost, volíme pro hledání spojů v křížovém přepínači ty.

11.6.2 Algoritmus přidělování lístků

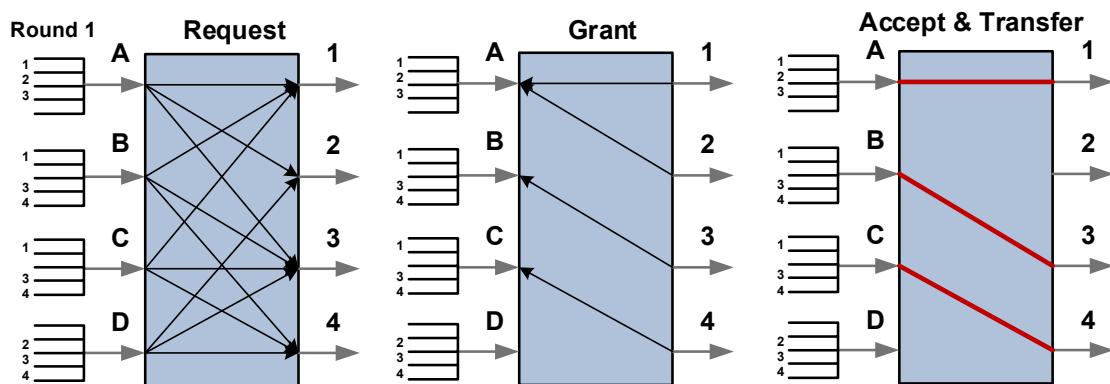
- Princip:
 - Výstupní port Q obsluhuje frontu požadavků na propojení.
 - Požadavek na vstupu P dostane od Q číslo T_{QX} pro obsloužení portu Q s pořadím X .
- Fáze:
 1. Žádost o lístek
 2. Přidělení lístku
 3. Propojení vstupů s výstupy a přenos
- Hodnocení:
 - Umožňuje přenos rámců s proměnlivou délkou
 - **Blokování na začátku fronty (head of line blocking)** – Situace, kdy více vstupů chce stejný výstup. Musí se zpracovat ve více kolech. Řešení: Virtuální výstupní fronty (každý vstupní port má frontu pro každý výstupní port). Tím může dojít k předbíhání rámců. Viz další algoritmy.



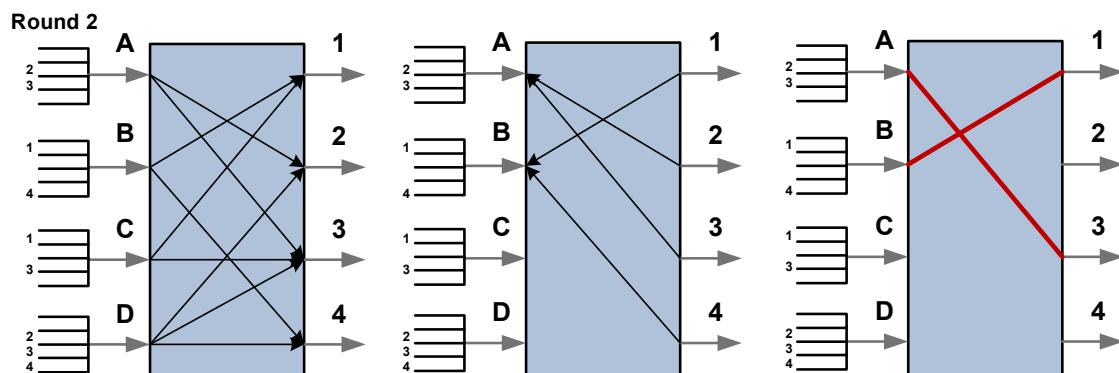
Obrázek 11.11: Příklad činnosti algoritmu přidělování lístků.

11.6.3 Algoritmus PIM (*Parallel Iterative Matching*)

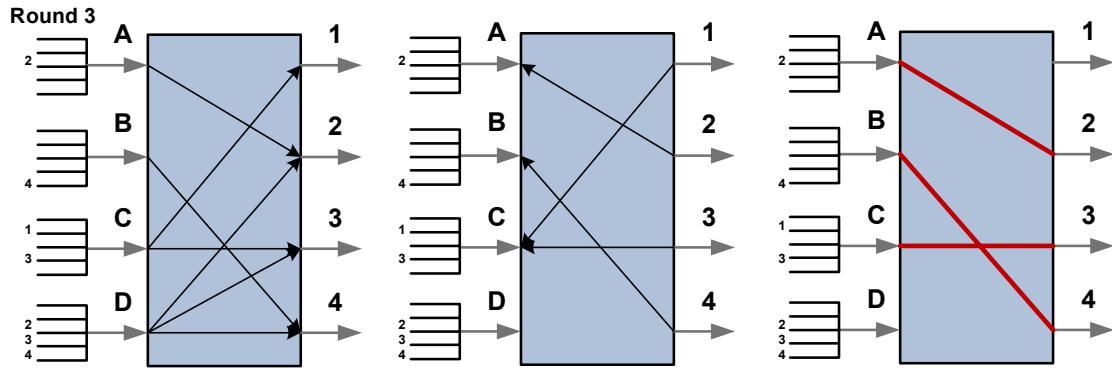
- Algoritmus přidělování lístků + virtuální výstupní fronty + náhodnost.
 - Posílají se požadavky pro všechny pakety ve frontě, nikoliv pouze pro první.
- Hledá maximální párování (pomocí náhodné volby, která zabrání vyhledování).
- Soutěžení o porty
 - Výstupní port – více žádostí naráz, jednu vyberu náhodně.
 - Vstupní port – povolení na více výstupů naráz, jeden vyberu náhodně.



Obrázek 11.12: Příklad činnosti algoritmu PIM, kolo 1.



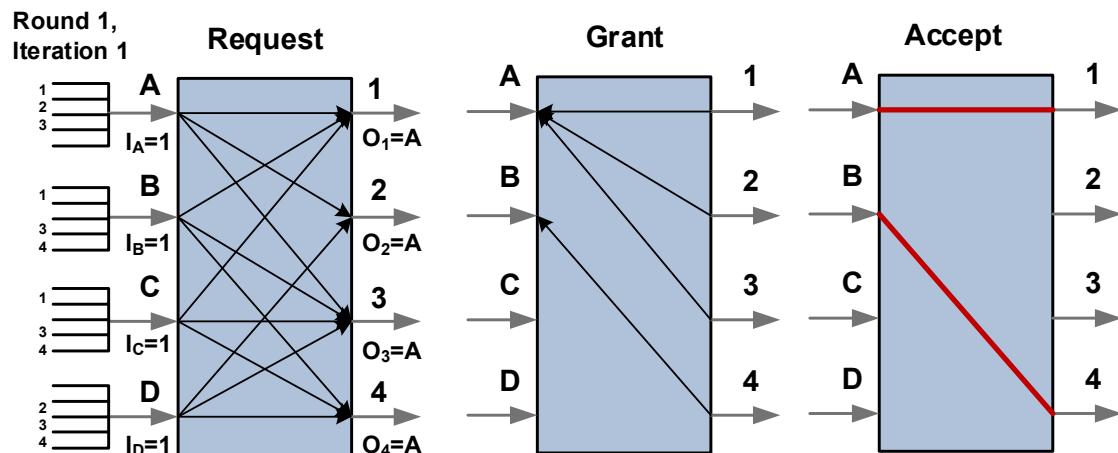
Obrázek 11.13: Příklad činnosti algoritmu PIM, kolo 2.



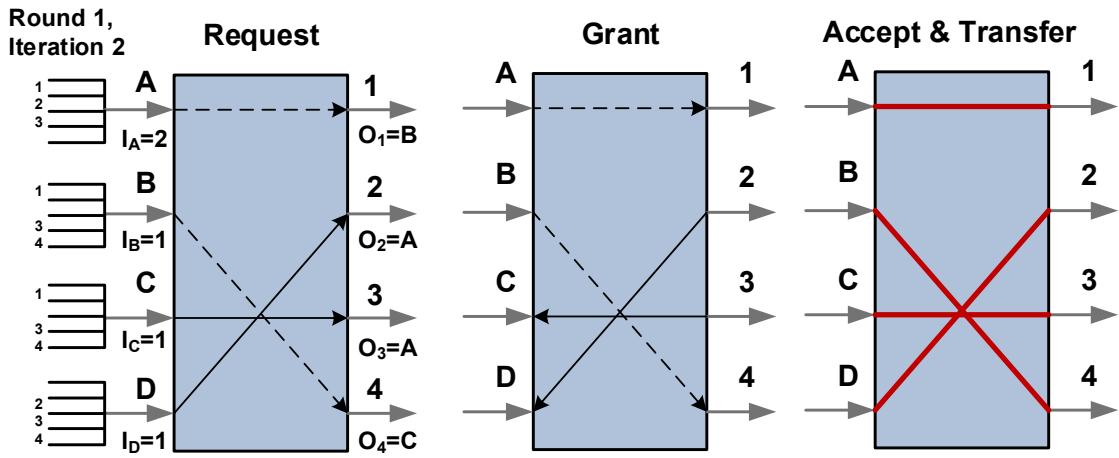
Obrázek 11.14: Příklad činnosti algoritmu PIM, kolo 3.

11.6.4 Algoritmus iSLIP

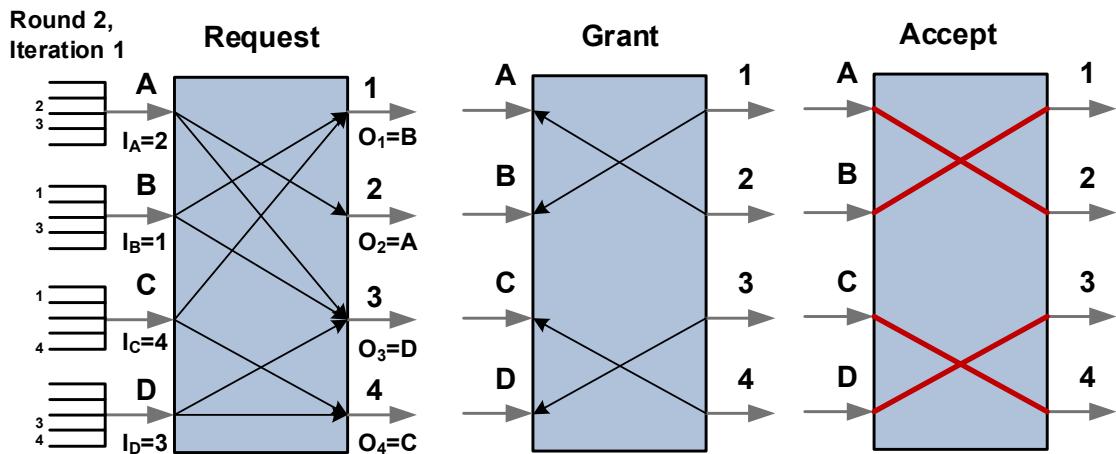
- Algoritmus přidělování lístků + virtuální výstupní fronty + iterace + ukazatele.
- Algoritmus „rotujících ukazatelů“ (vyhneme se náhodnému výběru).
 - Každý vstupní port má ukazatel (počítadlo), podle kterého se určuje priorita výstupu.
 - Každý výstupní port má ukazatel, podle kterého se určuje priorita vstupu.
 - Nová hodnota ukazatele je port na který se posílá / port od kterého se přijímá +1.
- Každé kolo má dvě iterace, resp. 6 fází – *request, grant, accept, request, grant, accept and transfer*.
 - Proč? Protože fáze *transfer* trvá nejdéle (optimalizace).



Obrázek 11.15: Příklad činnosti algoritmu iSLIP, kolo 1.



Obrázek 11.16: Příklad činnosti algoritmu iSLIP, kolo 2.



Obrázek 11.17: Příklad činnosti algoritmu iSLIP, kolo 3.

11.7 Přepínače s přepínánou propojovací deskou: vícestupňové přepínání

Vlastnosti jednostupňového křížového přepínače:

- S počtem portů roste kvadraticky velikost přepínacího pole.
- Vnitřní blokování (soupeření o porty).
- Blokování na začátku fronty (*head of line blocking*).
- Součástí přepínání je hledání maximálního párování.

Je možné zvětšit počet portů, aniž by se dramaticky rozšířila velikost přepínacího pole? Ano, pomocí vícestupňového přepínání (vstup → vnitřní přepínací obvody → výstup).

- Nemají vnitřní blokování – Nejsou potřeba plánovací algoritmy pro hledání maximálního párování.

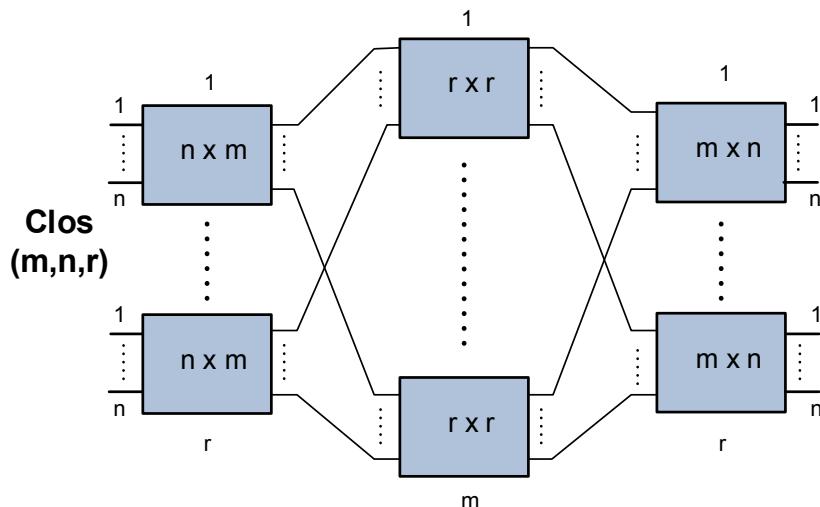
11.7.1 Přepínací síť Clos

$Clos(m, n, r)$

- r vstupních bloků s n vstupy (bloků = křížový přepínač)
- r výstupních bloků s n vstupy
- m vnitřních bloků s r vstupy

Closův teorém:

- Pokud $m \geq 2n - 1$, pak lze přidat nové propojení vstupu a výstupu bez přeskladání (sít' je neblokující).
 - Pokud všechny požadavky na vstupu jdou na jiné výstupy, pak pro jakoukoliv konfiguraci vstupů a výstupů je síť neblokující (pokud nejdou dva vstupy na jeden stejný výstup).

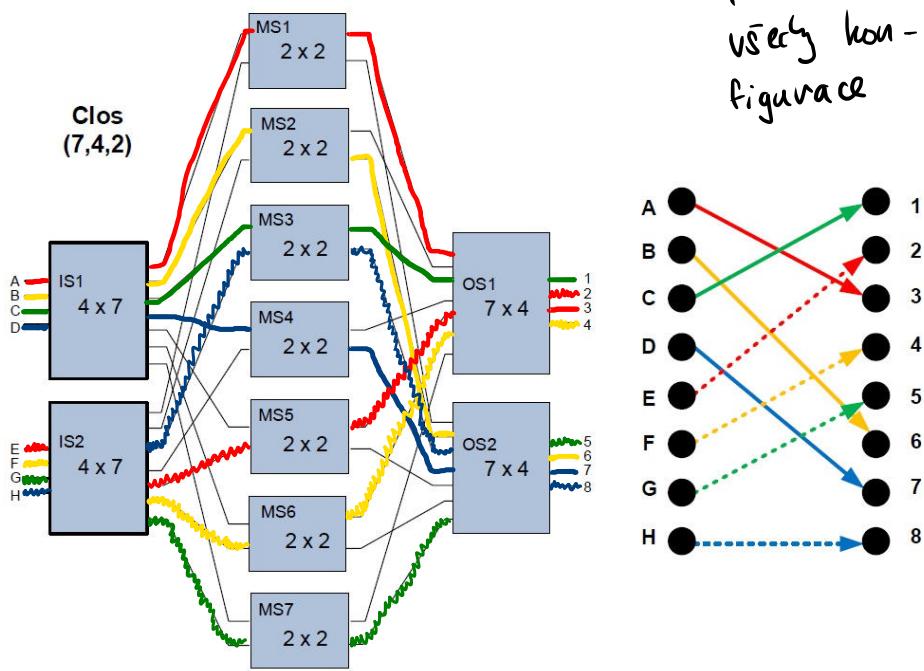


Obrázek 11.18: Obecné schéma Closovy sítě.

Příklad neblokující sítě Clos(7,4,2): osmiportový přepínač

- Platí Closův teorém, tj. $m \geq 2n - 1$

• Vždycky už každou propojení působí všechny konfigurace



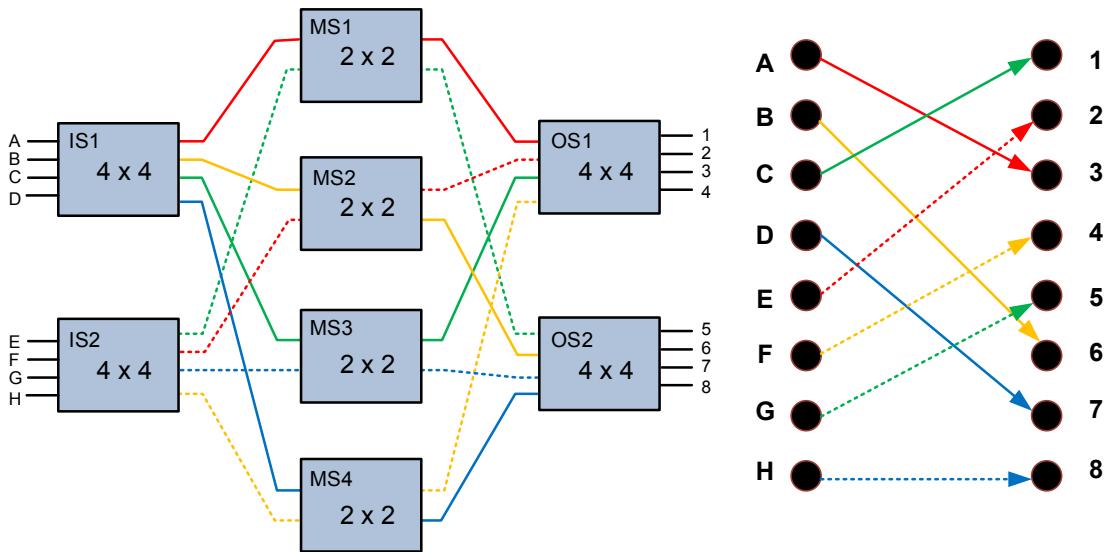
Obrázek 11.19: Příklad Closovy sítě.

11.7.2 Modifikovaná přepínací síť Clos

- Pokud $m \geq n$, pak je síť neblokující po přeskládání.
- Výpočetní složitost přeskládání je $O(N \log D)$, kde D je počet barev.
- Proč? Při plnění closova teorému je potřeba hodně bloků. Při této podmínce stačí výrazně méně a složitost výpočtu přeskládání je přijatelná.

Příklad sítě Clos(4,4,2): osmiportový přepínač

- Příklad neblokující konfigurace po přeskládání.

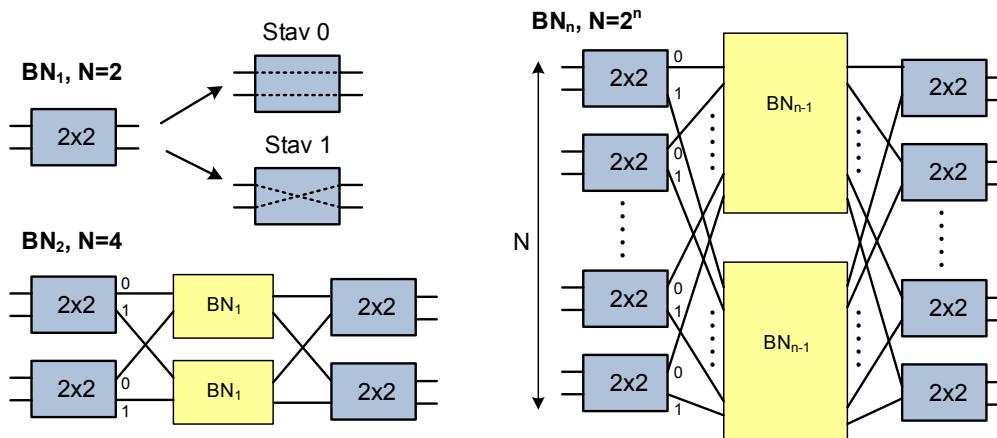


Obrázek 11.20: Příklad modifikované Closovy sítě.

11.7.3 Přepínací sít' Beneš

Benes(n) také $BN(n)$

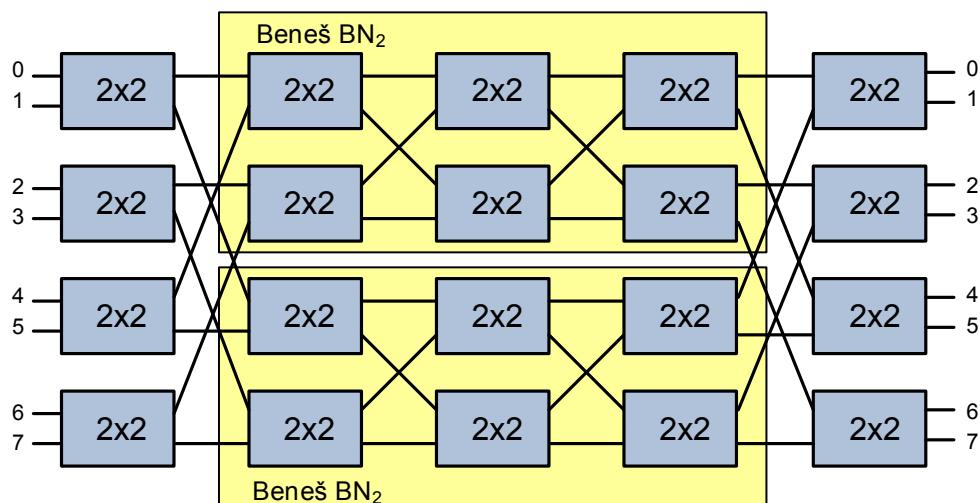
- Základem je modifikovaná síť $Clos(m = 2, n = 2, r = 1)$.
- Rekurzivní konstrukce sítě $Benes(n)$, kde n je stupeň rekurze.
- Vstupní a výstupní přepínače velikosti 2×2 .
- Celkový počet vstupních/výstupních portů $N = 2^n$.
- Prostřední část rekurzivních bloků $BN(n - 1)$.
- Neblokující po přeskládání – algoritmus hledání propojení (žádné soupeření o vnitřní uzly nebo výstupní porty).



Obrázek 11.21: Obecné schéma Benešovy sítě.

Příklad: přepínací síť Beneš BN₃

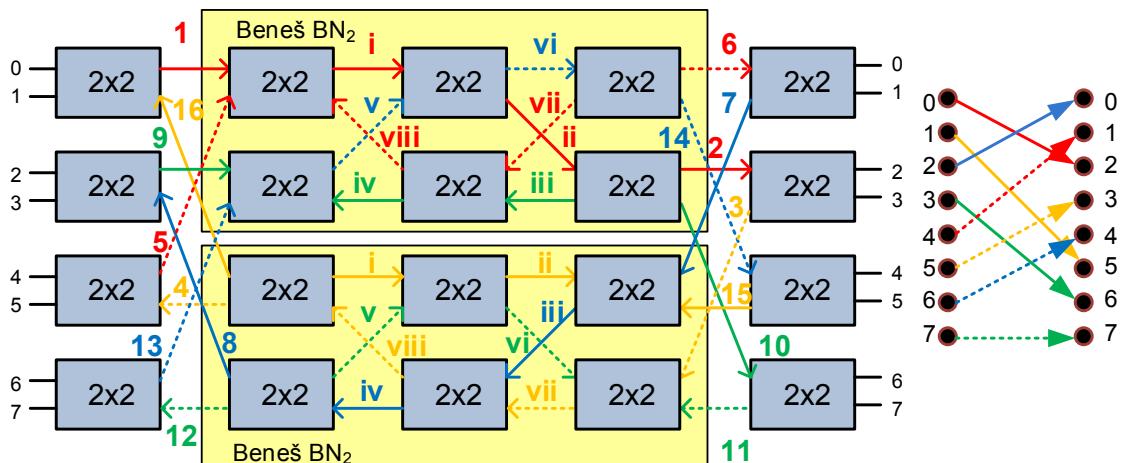
- $N = 8$ vstupů, 5 stupňů
- Použito např. u směrovače Cisco CSR-1.



Obrázek 11.22: Příklad Benešovy sítě.

Algoritmus hledání propojení

- Horní podsíť BN_i slouží pro dopředné směrování, spodní podsíť BN_i pro zpětné.
- Postupuje se od vstupů k výstupům a zpět, vždy pro sousední porty.
- Nejprve se propojí vstupní a výstupní bloky, poté rekursivně síť BN_{i-1}.



Obrázek 11.23: Příklad Benešovy sítě s přeskládáním.

Kapitola 12

PDS – Základní funkce směrovače, zpracování paketů ve směrovači, typy přepínání a architektur.

12.1 Metadata

- Předmět: Přenos dat, počítačové sítě a protokoly (PDS)
- Přednáška:
 - 05-routing.pdf
- Záznam:
 - 2021-03-12

12.2 Úvod a kontext

Viz. „Úvod a kontext“ v předchozích otázkách z tohoto předmětu.

NAT NAT (*Network Address Translation*) je metoda mapování IP adresního prostoru do jiného prostoru (typicky privátní adresy na veřejné adresy). Děje se tak úpravou hlaviček IP paketů během jejich přenosu přes směrovače (úprava zdrojové IP adresy a čísla portu). Směrovač si ukládá čtverice (WAN_IP : WAN_port, LAN_IP : LAN_port) aby mohl provádět i překlad zpět.

ACL ACL (*Access Control List*) je volitelná vrstva zabezpečení, která funguje jako brána firewall pro řízení provozu do jedné nebo více podsítí a z nich.

ARP ARP (*Address Resolution Protocol*) a RARP (*Reverse ARP*) je protokol, který komunikuje na síťové vrstvě (L3) a zajišťuje „překlad“ IP adres na MAC adresy a obráceně. Pouze pro IPv4, pro IPv6 je pro stejný účel využíván protokol ICMPv6 a zpráva *Neighbor Discovery*. Příklad využití: směrovač potřebuje získat MAC adresu next hopu (zná jeho IP adresu).

ICMP ICMP (*Internet Control Message Protocol*) je protokol, který komunikuje na síťové vrstvě (L3) a slouží pro řízení toku a detekce nedosažitelných uzlů.

MTU MTU (*Maximum Transmission Unit*) je největší velikost paketu, kterou lze v síti odeslat (přes výstupní rozhraní sítového prvku). Závisí na sítové kartě příslušného rozhraní.

Fragmentace paketů Uzel v síti (směrovač) dostane paket o velikosti n . Paket má přepo-
slat přes výstupní rozhraní do sítě, ve které je $MTU < n$. Aby paket mohl být odeslán,
musí být rozdělen (fragmentován) na více menších paketů (fragmenty) a odeslán po čás-
tech. Na straně příjemce pak musí nastat opačný proces – defragmentace.

MAC adresa MAC adresa je fyzická adresa zařízení, resp. sítové karty (identifikátor na L2). Na každý sítový port v přepínači/směrovači je jedna sítová karta (a teda i MAC adresa).

IP adresa IP adresa je logická adresa zařízení (identifikátor na L3). Pokud má zařízení více sítových karet, má typicky stále pouze jednu IP adresu. Přepínač nemá IP adresu vůbec, pracuje pouze na L2 vrstvě. Směrovač je výjimka a má 2 IP adresy, jednu pro komunikaci v lokální síti (LAN) a druhou pro internet (WAN).

12.3 Základní popis směrovače

Směrovač je sítové zařízení, které předává datové pakety mezi počítačovými sítěmi. Základní charakteristika:

- Pracuje s pakety (sítová vrstva, L3).
- Tvoří rozhraní mezi počítačovými sítěmi (provádí překlad NAT).
- Klasifikuje a filtruje pakety (firewall, ACL).
- Provádí fragmentaci a defragmentaci.

Činnost

1. Vypouzdření paketu z L2 (odebrání L2 hlavičky) a kontrola jestli je v pořádku (po-
mocí kontrolního součtu).
2. Vyhledání cesty kam se má paket směrovat a překlad adresy NAT (pomocí směrovací
tabulky).
3. Určení cílové MAC adresy na základě cílové IP adresy (pošle ARP dotaz).
4. Určení výstupního rozhraní.
5. Sestavení výsledného paketu podle výstupního rozhraní (zapouzdření do příslušné
L2 technologie – přidání L2 hlavičky).

Co ovlivňuje propustnost

- Rozbalení paketu.
- Vyhledání směrovací cesty.
- Překlad NAT.
- Vyhledání cílové MAC adresy.
- Zabalení paketu do správné technologie.

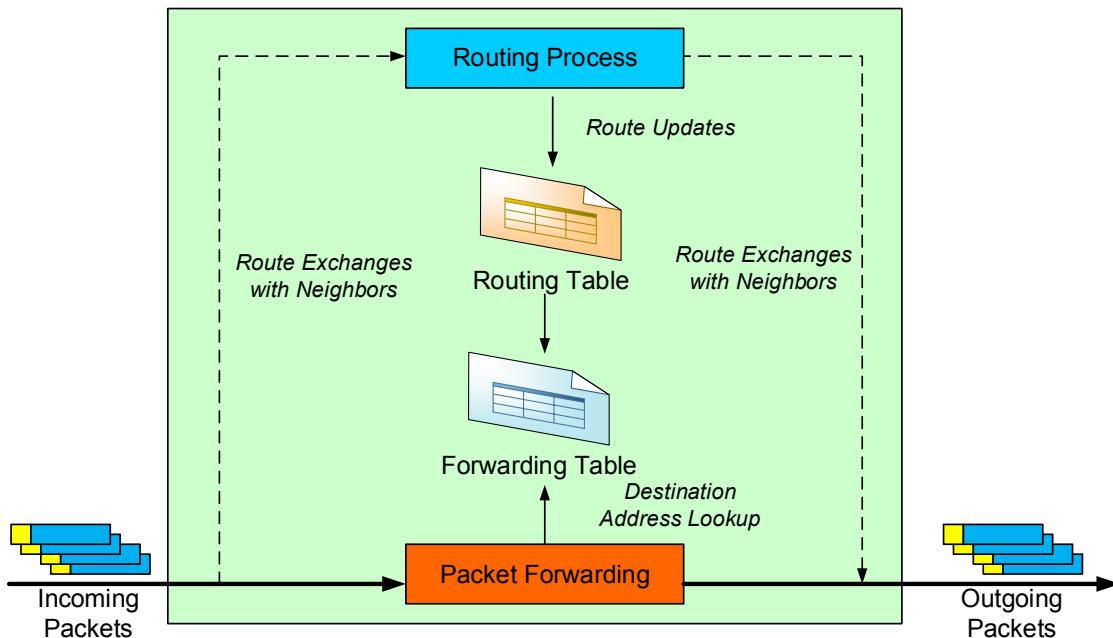
Typy

- Páteřní směrovače – součástí pátečních sítí mezi ISP.
- Hraniční směrovače – mezi zákaznickými sítěmi a ISP.
- Podnikové směrovače – propojení koncových systémů v podnikových (*enterprise*) sítích.
- Domácí směrovače.

12.4 Operace co směrovač vykonává

Pakety mohou být určeny buď přímo pro směrovač (aktualizace směrovací cesty přes směrovací protokoly – aktualizace směrovací tabulky) a nebo jinému uzlu v síti. V takovém případě směrovač paket přeposílá dalším uzlům v síti.

- **Směrování (routing)** – Určování cest paketů v prostředí počítačových sítí. Cílem je doručit paket adresátovi, pokud možno co nejefektivnější cestou. Směrování zajišťuje nejen směrovače, ale i koncové stanice (při vysílání).
- **Přepínání (forwarding, switching)** – Přepnutí paketu ze vstupního rozhraní na výstupní (viz architektura přepínačů).



Obrázek 12.1: Základní činnost směrovače.

Směrovací tabulka Směrovací tabulka (*routing table*) obsahuje informace nutné a dostačující pro směrování (ipprefix – IP prefix cílové sítě, **nexthop** – IP adresa dalšího uzlu, **metrika** – cena cesty, počet hopů). Informace do ní se získavají pomocí směrovacích protokolů a nebo staticky (administrátor provede manuálně). Směrovací tabulkou mají i koncové stanice (počítače), aby mohly posílat pakety do internetu.

| IP prefix | Next hop |
|--------------|---------------|
| 10.5.0.0/16 | 192.168.2.254 |
| 10.15.0.0/16 | 104.17.2.1 |
| 88.0.0.0/8 | 129.1.1.1 |
| | |

Obrázek 12.2: Příklad směrovací tabulky (bez metriky).

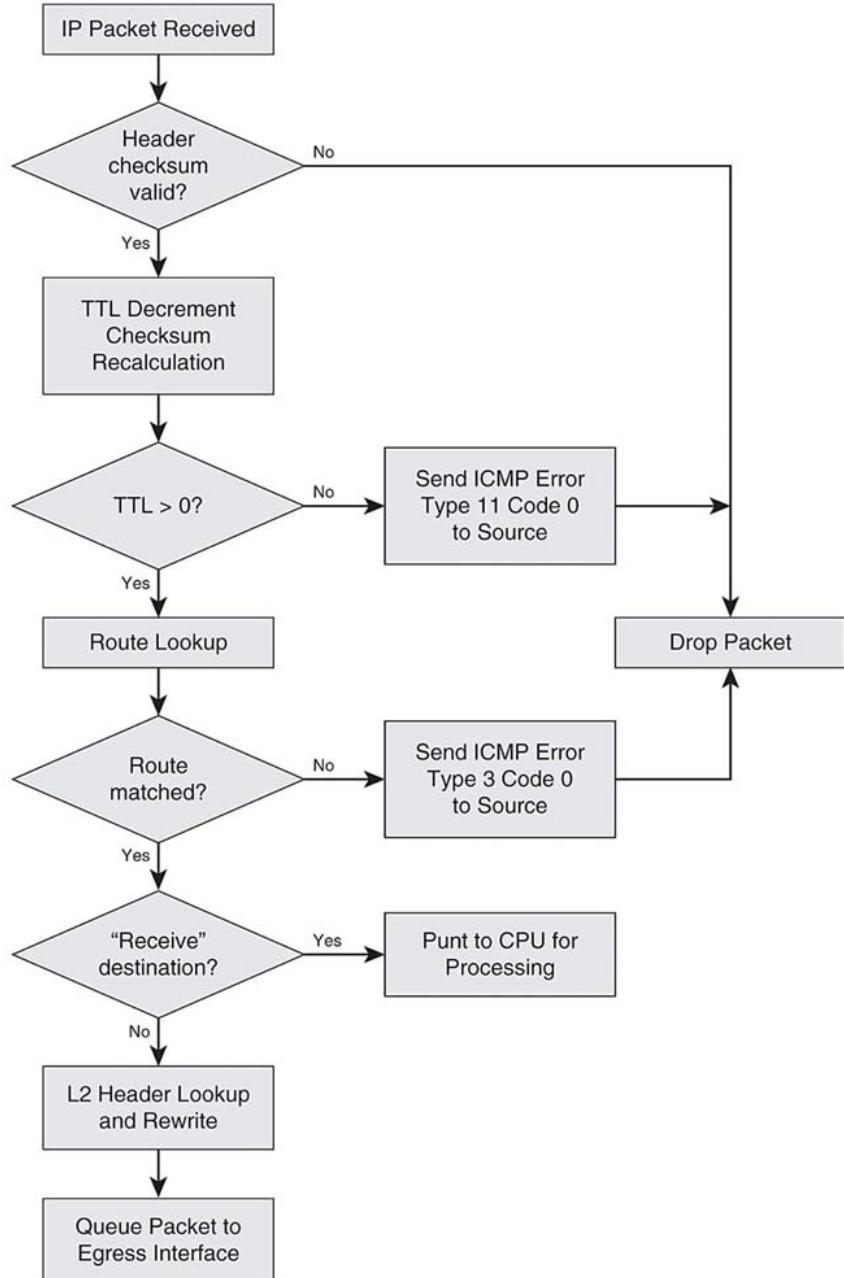
Přepínací tabulka Přepínací tabulka (FIB, *forwarding information base, forwarding table*) doplňuje směrovací tabulku o další informace, které jsou nutné pro sestavení cílového paketu (výstupní rozhraní, zdrojová MAC adresa, cílová MAC adresa). Informace do ní se získavají pomocí ARP protokolu a doplnění vlastních informací (*dst MAC address*).

| IP prefix | Interface | dst MAC address | src MAC address |
|--------------|-----------|-------------------|-------------------|
| 10.5.0.0/16 | eth0 | 00:0F:1F:CC:F3:06 | B8:6B:23:EA:FC:20 |
| 10.15.0.0/16 | eth1 | 01:12:11:A0:17:A0 | 00:16:17:E1:28:5F |
| 88.0.0.0/8 | eth2 | 01:3F:04:10:03:15 | 00:1B:63:8A:DB:1A |
| | | | |

Obrázek 12.3: Příklad přepínací tabulky.

Zpracování IP paketu (již vypouzdřeno z L2)

1. Validace hlavičky L3 (formát, verze, délka).
2. Kontrola hodnoty TTL a její dekrementace. Pokud je TTL 0, tak se paket zahodí a pošle ICMP zpráva.
3. Přepočítání kontrolního součtu.
4. Zpracování rozšířených voleb IP protokolu (timestamp, record route, strict source route).
5. Vyhledání cesty (next-hop) na základě cílové adresy (lokální doručení, unicast, multicast). Pokud se nepodaří vyhledat cesta, tak se paket zahodí a pošle ICMP zpráva.
6. Fragmentace a defragmentace IP paketů (pokud $MTU_{in} < MTU_{out}$).
7. Zpracování zpráv ICMP a IGMP.



Obrázek 12.4: Diagram zpracování paketu ve směrovači.

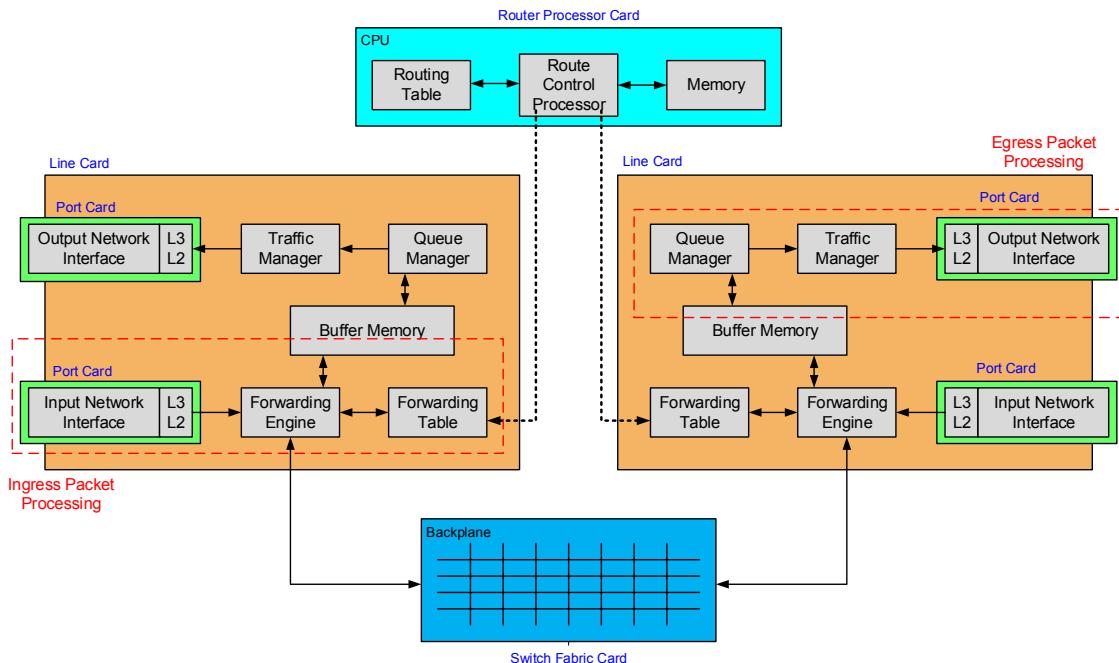
12.5 Architektura směrovače

Rozdělíme na fyzické části (z hlediska hardware) a funkční části (z hlediska komponent co něco vykonávají).

12.5.1 Fyzické části

- Procesorová část (*Router Processor Card*)
- Přepínací logika (*Switch Fabric Card*)

- Sít'ové karty (*Line Card*) – každá má vstupní a výstupní rozhraní (*Port Card*)



Obrázek 12.5: Fyzické části (procesory, paměti) směrovače.

12.5.2 Funkční části

Procesorová část procesor, paměť, směrovací tabulka

- Implementuje směrování na obecném CPU
- Zpracovává směrovací informace: aktualizace, udržuje sousedství
- Obsluhuje směrovací tabulkou
- Přenáší data do přepínací tabulky
- Zpracovává pakety, které nelze směrovat pomocí FIB
- Vytváří chybové zprávy ICMP

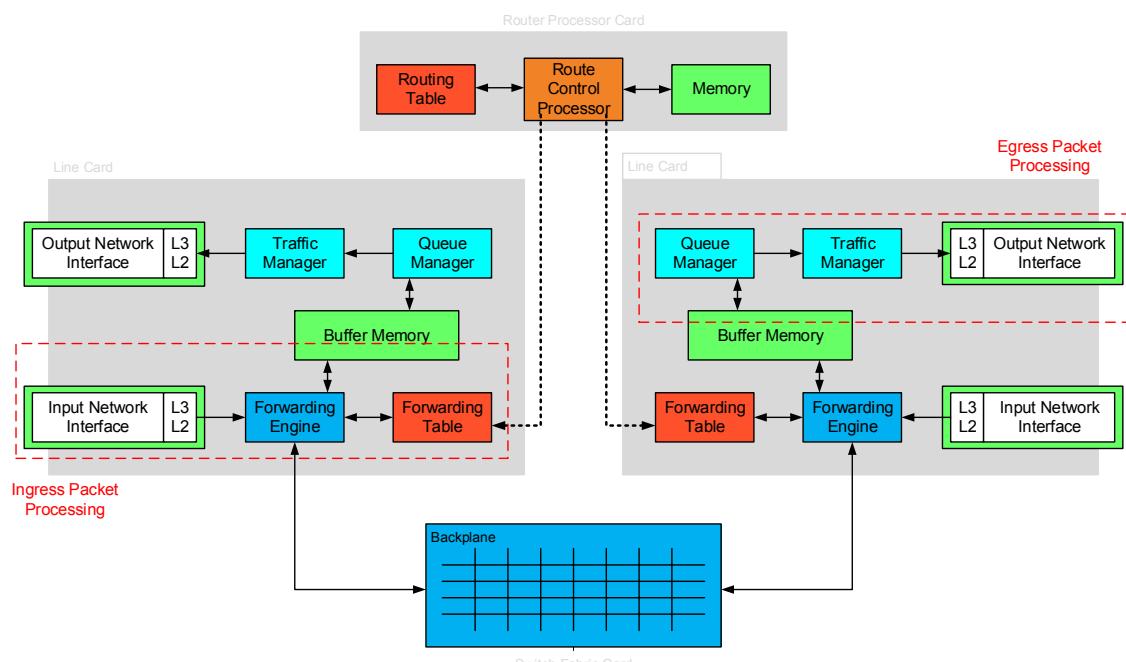
Přepínací logika propojovací deska (*Backplane*)

- Propojuje sít'ové rozhraní
- Vytváří sdílené (shared) či přepínané (switched) propojení
- Rychlosť přepínání odpovídá přenosovému pásmu všech rozhraní

Sít'ové karty

- Vstupní/Výstupní sít'ové rozhraní (*Input/Output Network Interface*)
 - Obsahuje více vstupních/výstupních portů
 - Odstraní zapouzdření L2
 - Předá hlavičku L3 přepínacímu modulu

- Uloží paket do paměti
- Zapouzdří odchozí pakety
- Řízení přepínání (*Forwarding Engine*)
 - Dostane hlavičku L3 ze síťového rozhraní
 - Určí výstupní rozhraní podle informace ve FIB
 - Provádí klasifikaci paketů pro podporu QoS na výstupu
- Správce front (*Queue Manager*)
 - Ukládá pakety do vyrovnávací paměti, pokud je výstupní port obsazen
 - Spravuje výstupní frontu → různé typy výstupních front
 - Při zaplnění fronty vybírá a zahazuje pakety podle definované politiky
- Správce provozu (*Traffic Manager*)
 - Prioritizuje a reguluje výstupní provoz podle požadavků QoS
 - Omezuje či ořezává výstupní provoz (shaping, policing)

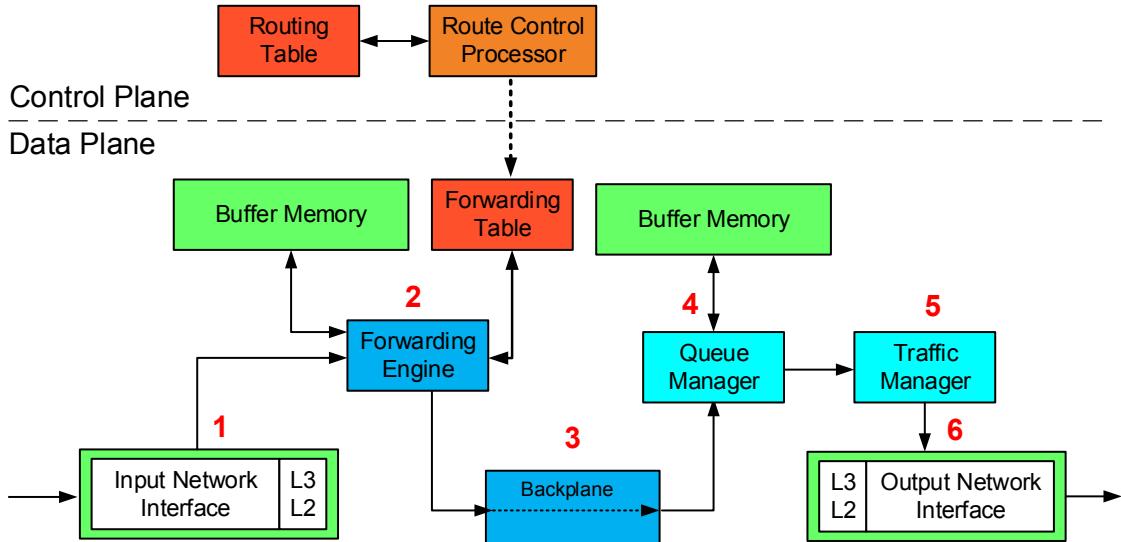


Obrázek 12.6: Funkční částí (procesy).

12.6 Zpracování paketu ve směrovači

Směrovač můžeme rozdělit na dvě části:

- *Data Plane* – Specifický hardware pro zpracování paketu (rychlé zpracování) – ASIC.
- *Control Plane* – Obecný procesor pro realizaci složitějších operací (pomalé zpracování).



Obrázek 12.7: Rozdělení přepínače na části *Data Plane* a *Control Plane* a zpracování rozděleno do 6. fází.

Kontext paketu Pomocná datová struktura, která reprezentuje paket uvnitř směrovače. Obsahuje vybrané informace z hlaviček. V rámci jednotlivých komponent směrovače se předává (pokud možno) tato struktura, pro větší efektivitu. Na počátku máme pouze částečný kontext, který se průběžně doplňuje. Na základě něho je pak sestaven výstupní paket.

Fáze zpracování Pokud některá z fází skončí neúspěšně, paket je zahozen.

1. Paket přijde na síťové rozhraní
 - Sít'ová karta zpracuje L2 rámcem, zkонтroluje FCS
 - Vytvoří kontext paketu: vloží L2 zdrojovou a cílovou adresu
 - Zpracování hlavičky L3: typ protokolu, kontrolní součet, TTL
 - Doplnění kontextu o informace L3: IP adresy, typ protokolu, DSCP + porty
2. Zpracování v přepínacím modulu
 - Vyhledání cesty v přepínací tabulce: next hop + výstupní rozhraní
 - Doplnění dalších informací do kontextu paketu
 - Paket uložen do vyrovnávací paměti → adresa vložena do kontextu
3. Přeposlání kontextu propojovací deskou
 - Paket i kontext jsou přeneseny na výstupní rozhraní
4. Zpracování správcem front
 - Podle priority v kontextu paketu je paket vložen do příslušné výstupní fronty
 - Obsluha fronty probíhá podle daného plánovacího algoritmu
5. Předání kontextu správci provozu
 - Kontrola omezení rychlosti (shaping) dle kontextu

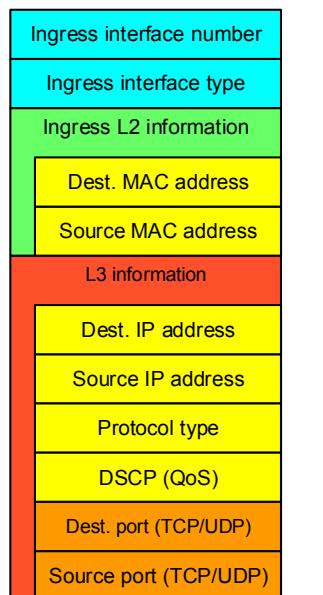
- Překročení rychlost: zahození či zpomalení

6. Výstupní síťové rozhraní

- L3: aktualizace TTL, přepočítání kontrolního součtu
- L2: přidání hlavičky, výpočet CRC
- Odeslání paketu (zapsání na výstupní médium)

Plný kontext paketu

- Kontext doplněn o výstupní informace
- Základní deska přenese paket i kontext na výstupní rozhraní
- Kontext obsahuje adresu uložení paketu v paměti
- Zpracování paketu předáno správci front



(a) Kontext paketu – částečný.



(b) Kontext paketu – úplný.

Obrázek 12.8: Kontext paketu (*ingress* – vstup, *egress* – výstup.)

Paket může být zahozen i na výstupní kartě. Proč?

- Vypršelo TTL
- Jsou plné fronty (politika RED, WRED)
- Filtrování na výstupu

Paket může být zpracován dvěma způsoby: rychlou cestou (přes *data plane*) a pomalou cestou (přes *control plane*).

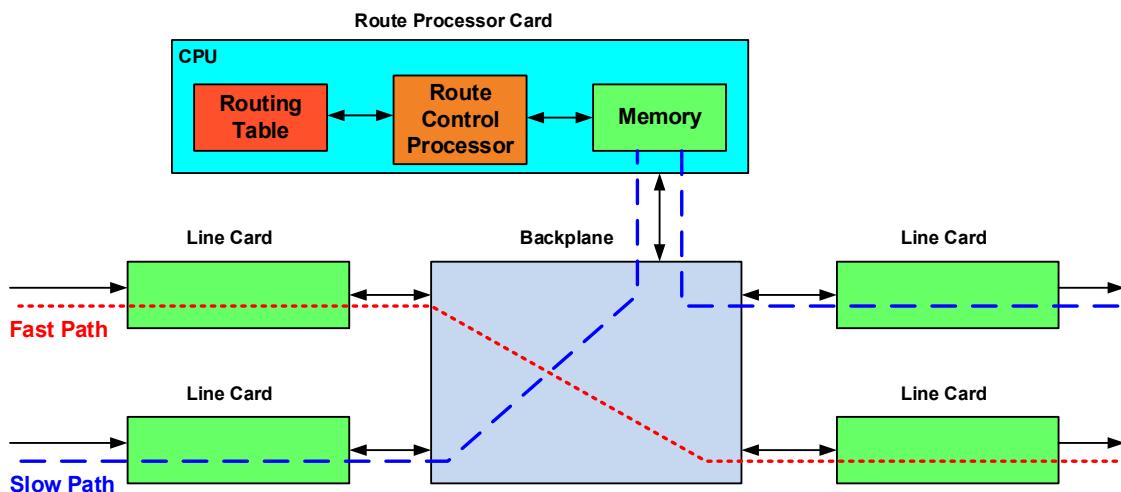
Rychlá cesta Paket je možné zpracovat pouze v hardwaru (*data plane*).

- Zpracování hlavičky IP (kontrola verze, délky paketu, snížení TTL, přepočítání kontrolního součtu)

- Přeposlání paketu ze vstupní na výstupní (lokální uložení, přeposlání na jeden port (unicast) či více portů (multicast))
- Klasifikace paketu na základě informací z hlavičky (optimalizované datové struktury pro rychlé uložení a vyhledání)
- Uložení do front, plánování (různé typy front, různé typy obsluhy)

Pomalá cesta Paket je částečně zpracován v hardwaru, ale většina zpracování musí proběhnout softwarově (*control plane*). Síťová karta nedokáže zpracovat paket sama a musí se využít centrální procesor.

- Zpracování ARP (zjištění výstupní adresy L2: první paket, ostatní pakety)
- Fragmentace a defragmentace
- Generování zpráv ICMP
- Zpracování směrovacích informací



Obrázek 12.9: Dva způsoby zpracování paketu ve směrovači – pomalou cestou a rychlou cestou.

12.7 Typy přepínání paketů

Přepínání paketů z jednoho rozhraní na druhé na základě směrovacích informací je jedna z nejdůležitějších funkcí směrovače. Proces přepínání paketů zahrnuje:

- Zjištění, zda cíl cesty paketu je dosažitelný.
- Vyhledání nejbližšího uzlu na cestě (next-hop) a určení výstupního rozhraní.
- Vyhledání informací pro vytvoření L2 hlavičky paketu na výstupu.

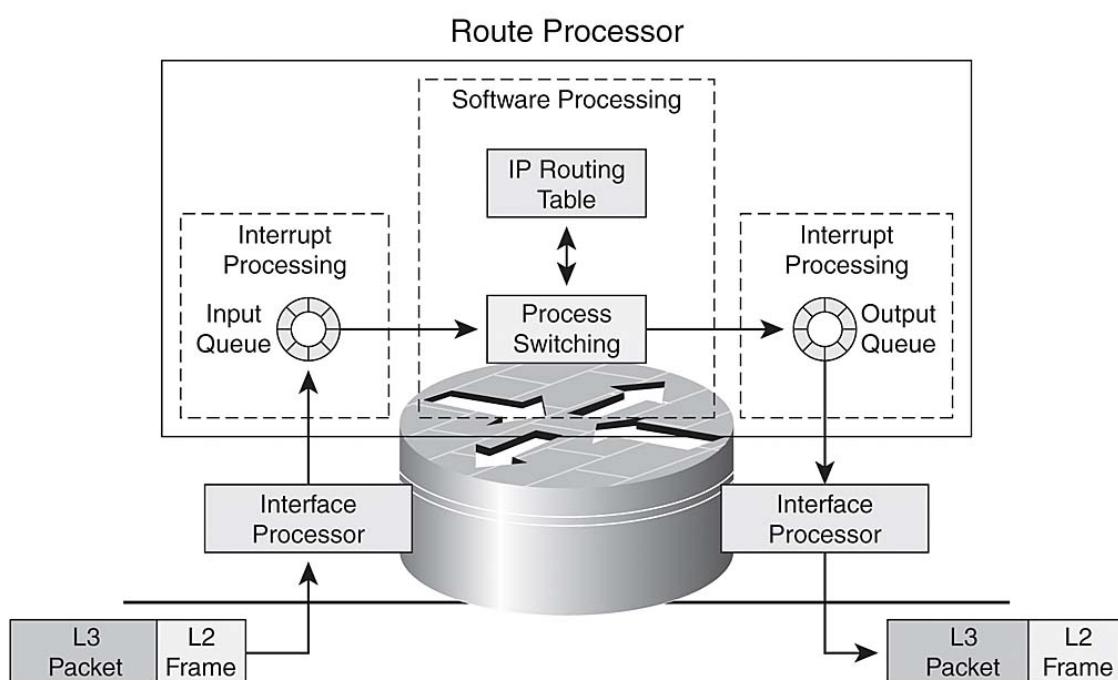
12.7.1 Softwarové přepínání (*Process Switching*)

- Přepínání pomocí centrálního procesoru s centrální pamětí.
 - Pomalé – přístup do paměti, obecný procesor (softwarové zpracování).
- Pro každý paket se hledá cesta ve směrovací tabulce a určuje se MAC adresa cíle.

- Směrovač využívá standardní mechanismus přepíná procesů v OS (přeřušení).

Fáze softwarového přepínání

1. I/O procesor detekuje paket na vstupním médiu. Přenese ho do vstupního bufferu.
2. I/O procesor vygeneruje přerušení. Během přerušení určí centrální procesor typ paketu a zkopíruje ho do centrální paměti.
3. Centrální plánovač zjistí, že ve vstupní frontě je paket. Naplánuje jeho další zpracování procesem.
4. Proces pro zpracování vyhledá ve směrovací tabulce další uzel (next hop) a výstupní rozhraní. V paměti ARP cache vyhledá MAC adresu dalšího uzlu.
5. Přepíše L2 hlavičku paketu a umístí paket do výstupní fronty na výstupním portu.
6. Paket vložen do fronty na výstupním portu.
7. I/O procesor detekuje paket ve vysílací frontě. Zapíše ho na síťové médium.



Obrázek 12.10: Softwarové přepínání.

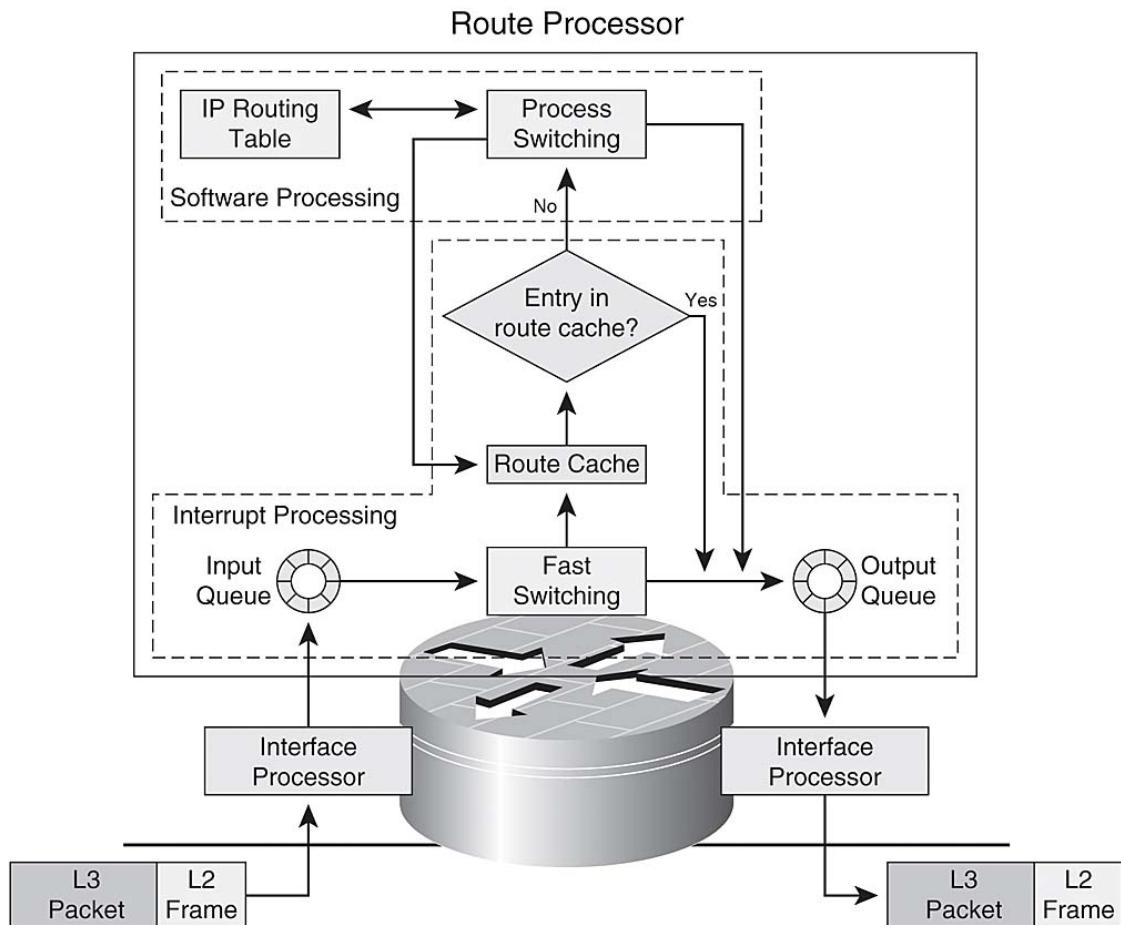
Optimalizace

- Nejdéle trvá: na základě cíle vyhledávání informací ve směrovací tabulce, najdu next hop, najdu příslušný záznam v ARP tabulce
- Můžeme využít cache (která si bude pamatovat poslední záznamy).

12.7.2 Rychlé přepínání (*Fast Switching*)

- Využívá vyrovnávací paměť *route cache* s předpočítanou L2 hlavičkou.
- První paket toku se přepíná softwarově, další pakety toku jdou rychlou cestou.

- Při přepínání paketu je vložen do paměti *route cache* nový záznam.
- Neexistuje synchronizace mezi směrovací tabulkou, ARP cache a *route cache*.
- Záznam v *route cache* se zneplatní při změně ARP cache či směrovací tabulky.
- Při zaplnění paměti nad určitou mez se začnou záznamy náhodně zahazovat.

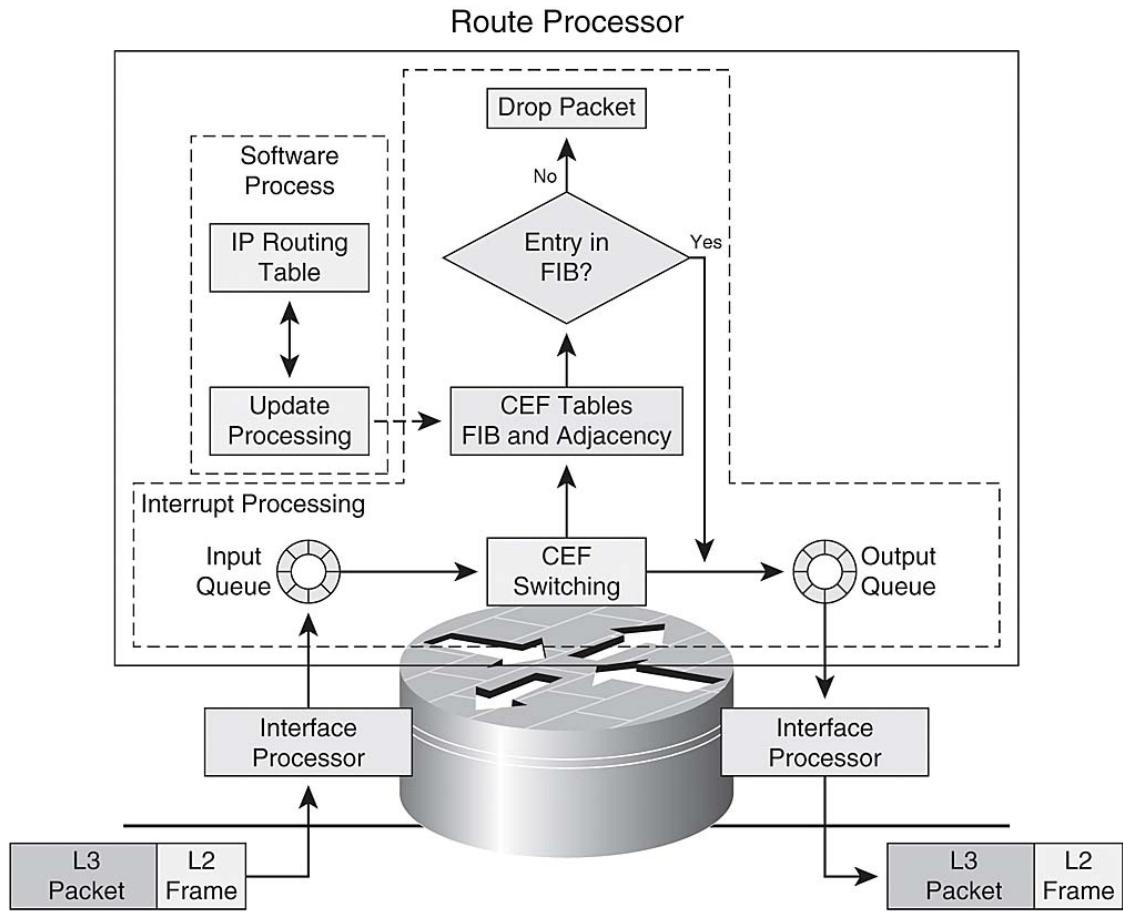


Obrázek 12.11: Rychlé přepínání.

12.7.3 Expresní přepínání CEF (Cisco Express Forwarding)

[[todo]]

- Tabulka CEF se předpočítá na základě směrovací tabulky a tabulky sousedů ještě před příchodem paketu → nedochází k softwarovému přepínání.
- Oddělení směrovacích informací od L2 dat → nedochází ke stárnutí záznamů při expiraci záznamu v tabulce ARP.
- Změny ve směrovací tabulce či tabulce ARP se okamžitě propagují do tabulky CEF.
- ARP tabulka se synchronizuje se záznamy v tabulce sousedů.



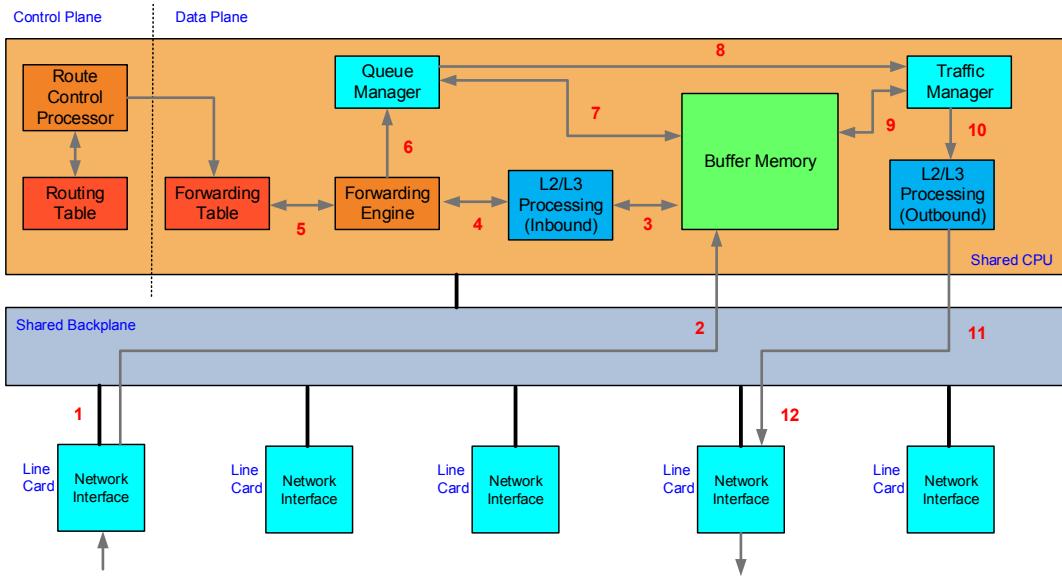
Obrázek 12.12: Expresní přepínání CEF.

12.8 Přehled architektur

Architektury směrovačů rozdělmě podle způsobu přepínání paketů.

12.8.1 Architektura se sdíleným procesorem

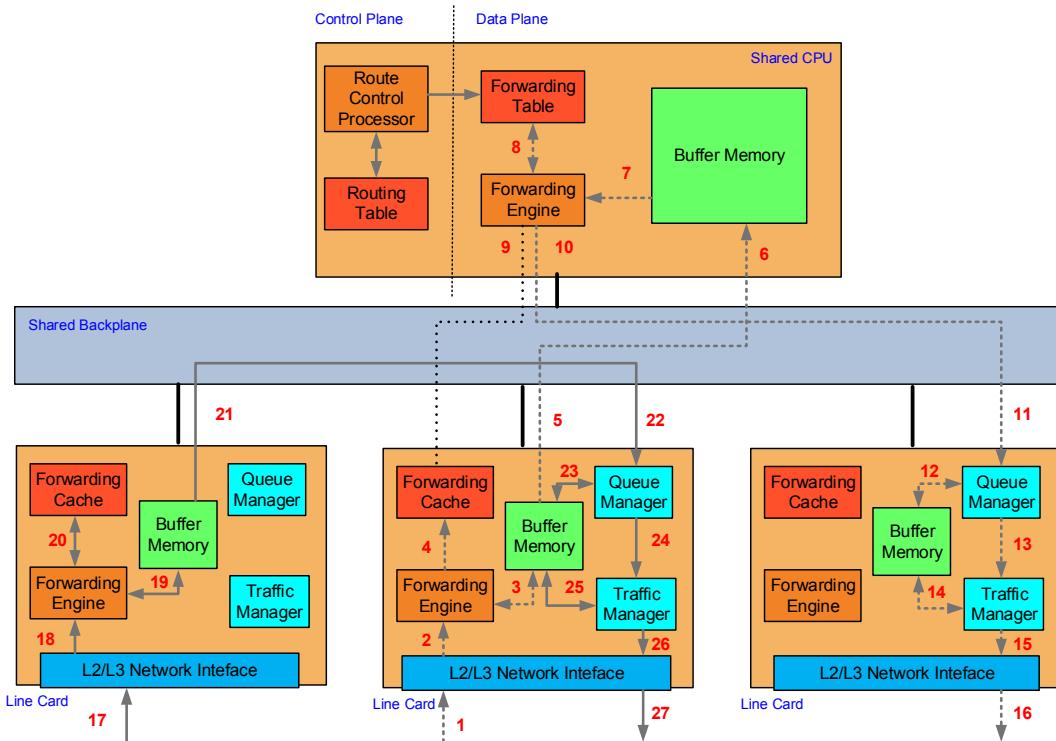
- Využívá softwarové přepínání → každý paket zpracován na CPU.
- Cykly CPU rozděleny mezi přenos, směrování a další operace.
- Sdílená sběrnice i procesor ⇒ levné, ale pomalé.



Obrázek 12.13: Architektura se sdíleným procesorem.

12.8.2 Architektura se sdíleným procesorem a pamětí cache na kartě

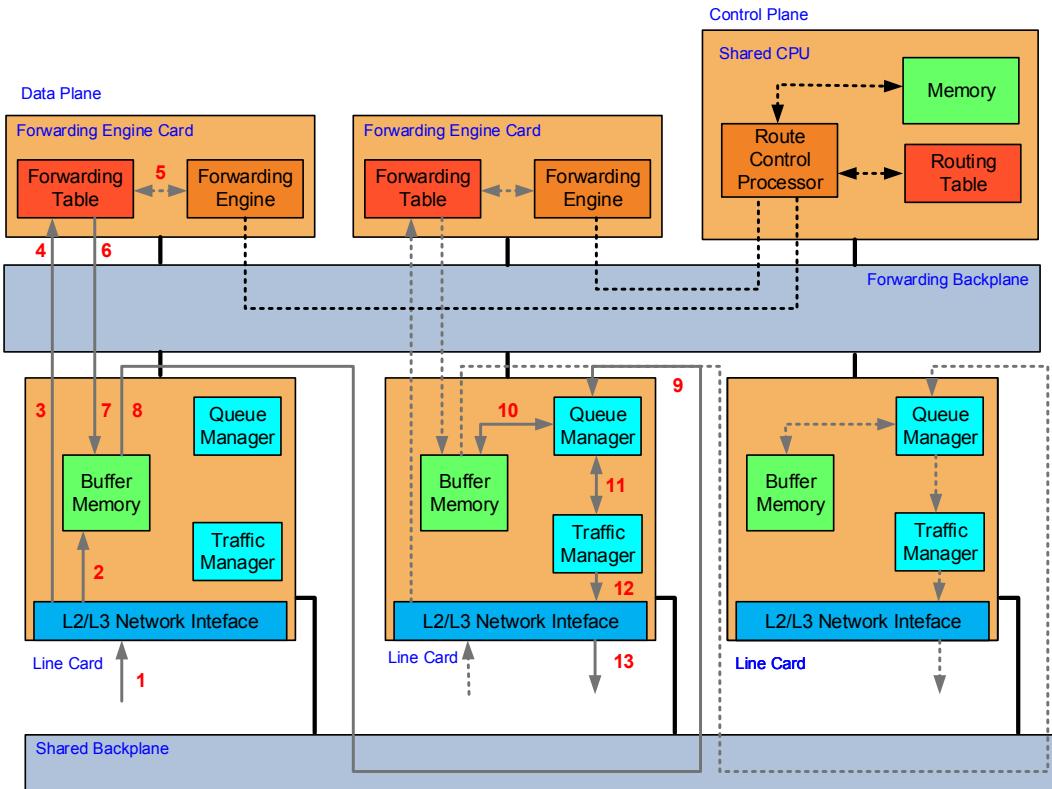
- Varianta s pamětí cache na kartě → synchronizace přepínacích tabulek.
- Sít'ová karta obsahuje FE pro zpracování hlaviček, paměť a přepínací tabulkou.
- Rychlé přepínání (Fast Switching): první paket vs. další pakety.



Obrázek 12.14: Architektura se sdíleným procesorem a pamětí cache na kartě.

12.8.3 Architektura s nezávislými moduly FE

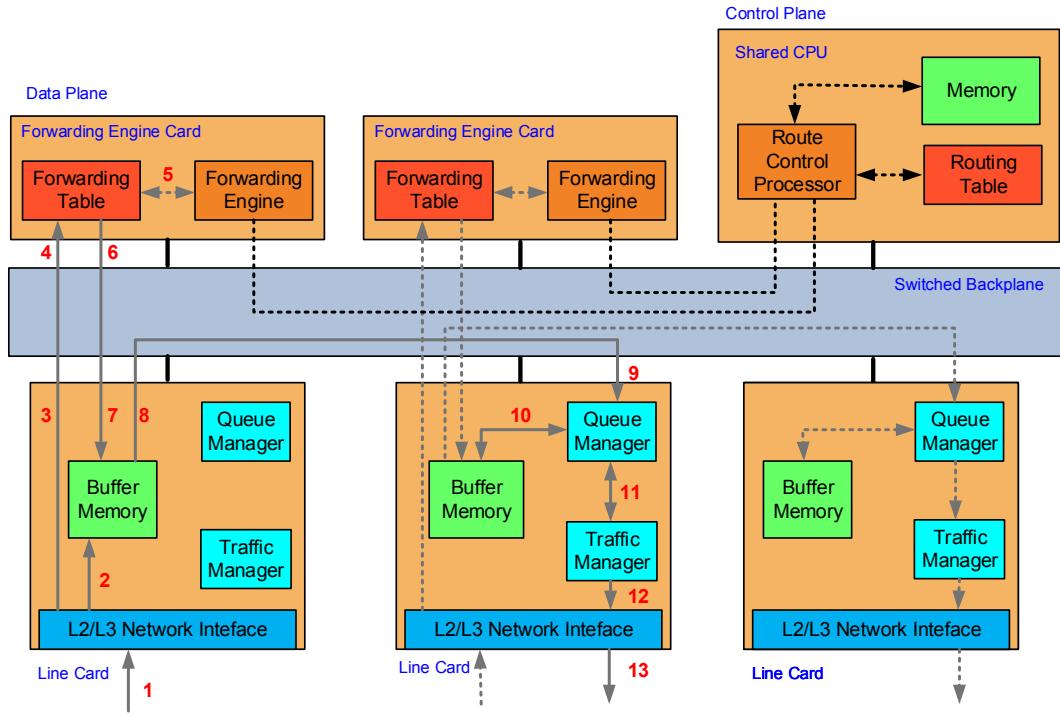
- Přepínací moduly FE implementovány na speciálních kartách.
- Paralelní zpracování paketů, dvě sběrnice (sdílená a přepínaná).



Obrázek 12.15: Architektura s nezávislými moduly FE.

12.8.4 Architektura s nezávislými moduly FE a přepínanou sběrnicí

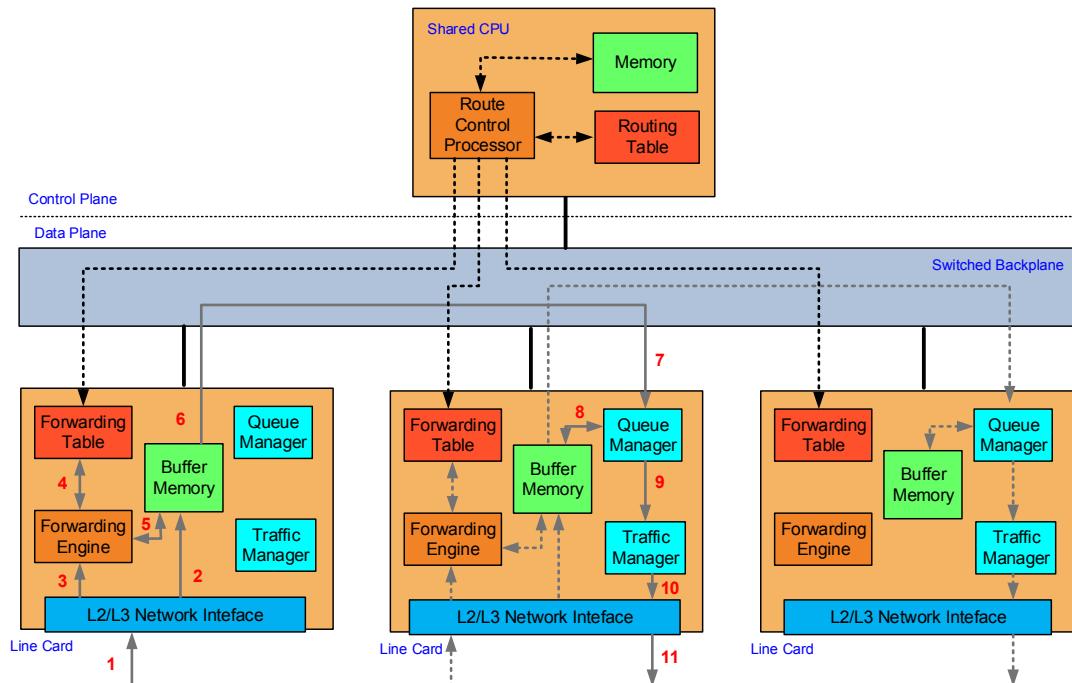
- Varianta s jednou přepínanou sběrnicí ⇒ vyšší propustnost.



Obrázek 12.16: Architektura s nezávislými moduly FE a přepínanou sběrnicí.

12.8.5 Distribuovaná architektura (Shared Nothing)

- Veškeré zpracování paketu přeneseno do síťového modulu.
- Oddělení procesu směrování a přepínání → využití technologie CEF.



Obrázek 12.17: Distribuovaná architektura (Shared Nothing).

Kapitola 13

PDS – Metody pro výpočet směrování v sítích (Bellman-Ford, Dijkstra, Path vector, DUAL).

13.1 Metadata

- Předmět: Přenos dat, počítačové sítě a protokoly (PDS)
- Přednáška:
 - 09-teorie-smerovani.pdf
- Záznam:
 - 2021-04-16

13.2 Úvod a kontext

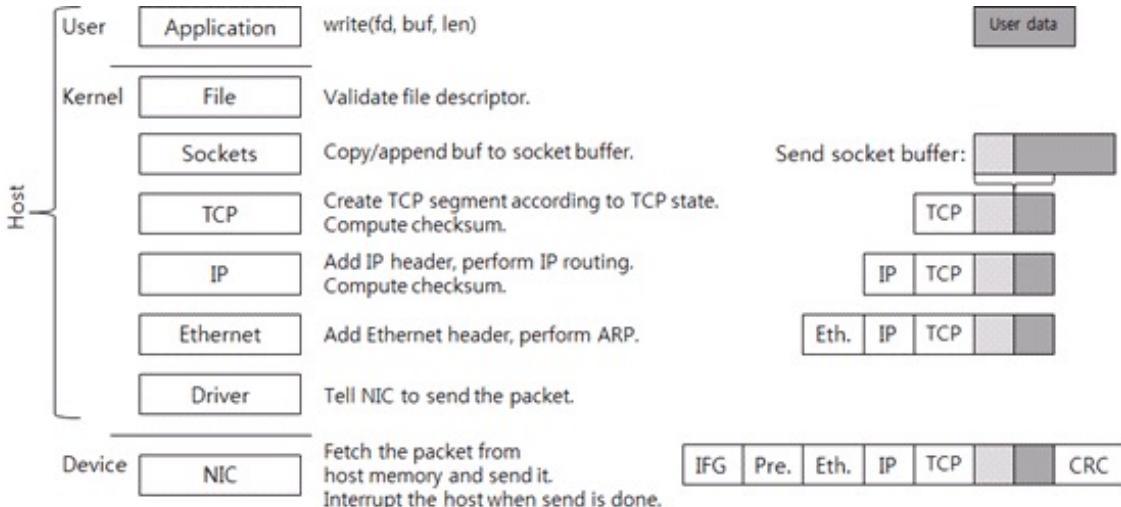
Viz. „Úvod a kontext“ v předchozích otázkách z tohoto předmětu.

Adresace

- Port (transportní vrstva, L4)
 - Identifikuje aplikaci v rámci zařízení.
 - Jak se mění při směrování paketu internetem: zůstává stejný s výjimkou zdrojového portu při překladu NAT.
 - Velikost: 16 bit
 - Prostor: plochý
- IPv4, IPv6 (sítová vrstva, L3)
 - Identifikuje uzel v rámci sítě (tzv. logická adresace).
 - Jak se mění při směrování paketu internetem: zůstává stejná s výjimkou zdrojové IP adresy při překladu NAT.
 - Velikost: 32 bit, 128 bit
 - Prostor: pseudohierarchie (A, B, C, D, E), pseudohierarchie (prefix + interface ID)

- MAC (linková vrstva, L2)
 - Identifikuje síťové rozhraní (sítovou kartu).
 - Jak se mění při směrování paketu internetem: mění se *hop-by-hop*.
 - Velikost: 32 bit, 128 bit
 - Prostor: pseudohierarchie (A, B, C, D, E), pseudohierarchie (prefix + interface ID)

Datový tok Identifikovaný čtveřicí (zdrojová IP, zdrojový port, cílová IP, cílový port).



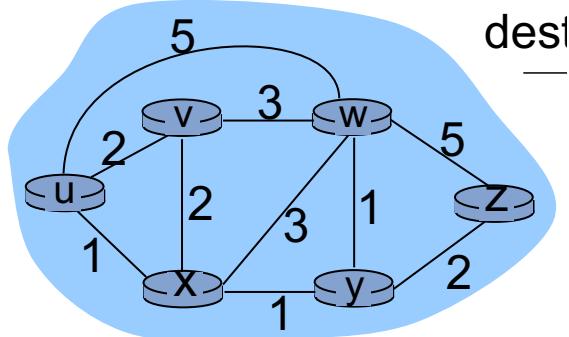
Obrázek 13.1: Operace na jednotlivých vrstvách ISO modelu.

13.3 Směrování

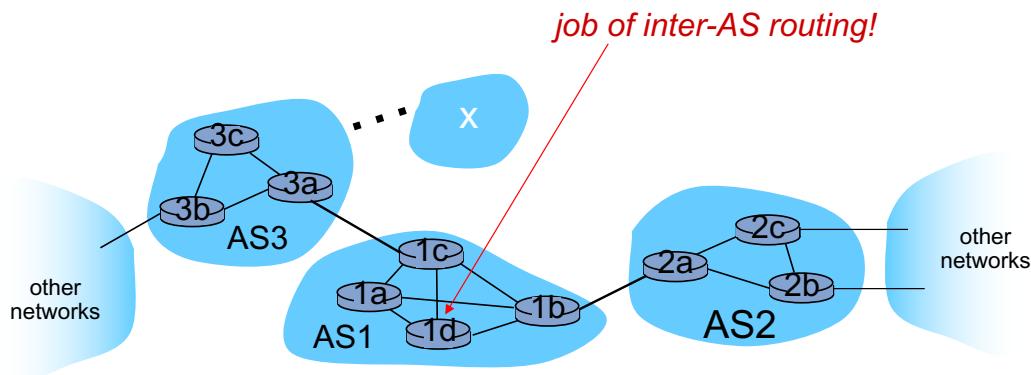
Z pohledu směrování můžeme počítačovou síť zjednodušit na neorientovaný, ohodnocený graf, kde uzly jsou směrovače a hrany jsou propojení mezi něma. Směrování je pak hledání nejkratší cesty mezi dvěma uzly, resp. ze zdrojového do všech.

Autonomní systém Není možné směrovat přes celý internet (příliš velký, příliš mnoho uzlů). Využíváme abstrakce, kdy se na podsíť díváme jako na uzel – tzv. autonomní systém. Tímto agregujeme více informací do jedné (přesné adresy uzlů, na vyšší úrovni reprezentujeme pouze adresou sítě). Internet je pak síť sítí.

Hierarchické směrování Hierarchické směrování znamená, že nesměrujeme mezi každým směrovačem na světě, ale využíváme autonomní systémy. Směrujeme tedy globálně mezi autonomníma systémama a pak uvnitř mezi směrovačema.



Obrázek 13.2: Příklad počítačové sítě zobrazené jako neorientovaný ohodnocený graf
 $G = (V, E)$, $V = \{u, v, w, x, y, z\}$, $E = \{(u, v), \dots\}$, $w = \{((u, v), 2), \dots\}$.



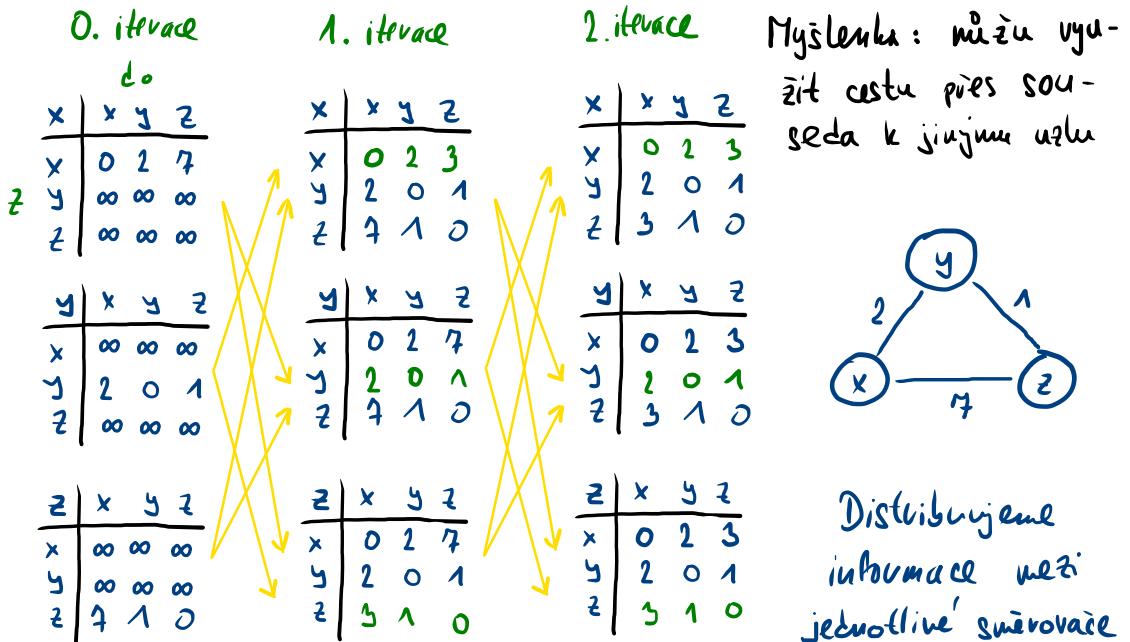
Obrázek 13.3: Příklad autonomních systémů (např. VUT, CESNET, UPC, ...).

13.4 Distance Vector přístup

- Distribuované směrovací informace – Každý uzel, zná pouze omezenou část topologie sítě. Konkrétně má informace pouze co sám zná a informace od svých sousedů.
- Používá se pro směrování uvnitř autonomních systémů.
- *Single-metric*.

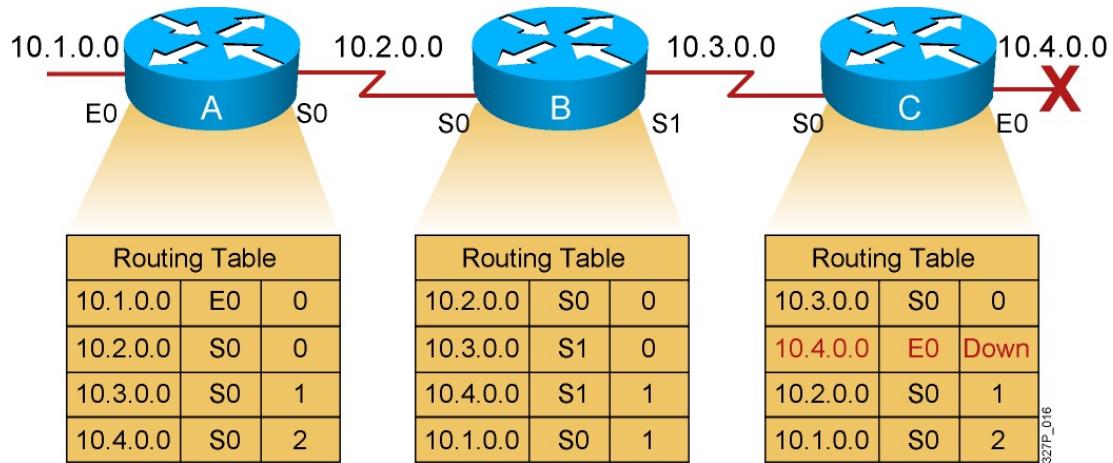
13.4.1 Algoritmus Bellman-Ford

- Pro formální vysvětlení a pseudokód viz otázku „Hledání nejkratších cest ze zdrojového uzlu do všech ostatních uzlů grafu“.
- Jde o algoritmus hledání nejkratších cest ze zdrojového uzlu do všech ostatních uzlů. To znamená, že ho provádí každý uzel (směrovač).
- Využívá jej např. směrovací protokol RIP (*Routing Protocol*).
 - Komunikuje přes UDP na portu 520.
 - Metrika: *hop count*.
 - Proti *Infinity Count* (viz dále) se brání vylepšením *Split Horizon*: Nikdy se nepošle informace o cestě zpátky rozhraním, ze kterého přišla.



Obrázek 13.4: Příklad výpočtu nejkratší cesty pomocí algoritmu Bellman-Ford.

Infinity Count V případě, že vypadne některá ze sítí, může nastat problém počítání ceny cesty do nekonečna. Jelikož si uzly vyměňují pravidelně směrovací informace. Viz následující obrázek.

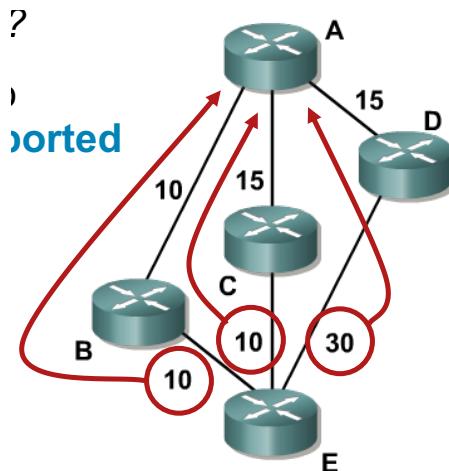


Obrázek 13.5: Příklad *Infinity Count*. Síť 10.4.0.0 vypadne a pro uzel C je najednou nedosažitelná. Uzel B se cestu do sítě 10.4.0.0 naučil dříve od C. Uzel C čerpá informaci o cestě do sítě z uzel B. Tímto způsobem se cena do sítě zacyklí a jde k ∞ .

13.4.2 Diffusing update algorithm (DUAL)

- Algoritmus DUAL (diffusing update algorithm) je algoritmus společnosti Cisco, který zajišťuje, že daná trasa je přepočítána globálně, kdykoli by mohla způsobit směrovací smyčku.

- Klíčová je tzv. **Feasibility Condition**, která zajišťuje, že jsou vždy vybírány pouze trasy bez smyček. Pokud podmínka platí, nemohou vzniknout žádné smyčky, ale podmínka může za určitých okolností odmítnout všechny trasy k cíli, přestože některé jsou bez smyček.
 - Vzdálenost sousedů k cílové síti nazveme *RC* (*reported distance*).
 - Nejlepší vzdálenost k danému uzlu nazveme *FD* (*feasible distance*).
 - Pokud platí $RD < FD$, tak cesta neobsahuje smyčku.
- Využívá jej směrovací protokol EIGRP (*Enhanced Interior Gateway Routing Protocol*) od společnosti Cisco.
 - Pro hledání nejkratší cesty používá algoritmus Bellman-Ford + *Feasibility Condition*.
 - Komunikace: vlastní protokol zabalený do IP paketu (*Reliable Transport Protocol*).
 - Kompozitní metrika – jedno číslo spočítaný na základě několika parametrů (šířka pásma, zpoždění, rychlosť, ...).



Obrázek 13.6: Příklad výpočtu nejlepší cesty a uplatnění feasibility condition.

Cesta z A do E:

$$\begin{array}{lll} RD(B) = 10 & RD(C) = 10 & RD(D) = 30 \\ via(B) = 20 & via(C) = 25 & via(D) = 45 \end{array}$$

Potom z pohledu A:

- Cesta $via(B)$ je *FD* a platí podmínka $FD > RD(B)$, tudíž cesta přes B neobsahuje smyčku.
- Pro cestu $via(C)$ platí podmínka $FD > RD(C)$, tudíž cesta přes C neobsahuje smyčku.
- Pro cestu $via(D)$ neplatí podmínka $FD > RD(D)$, tudíž cesta přes D může obsahovat smyčku a není brána v potaz.

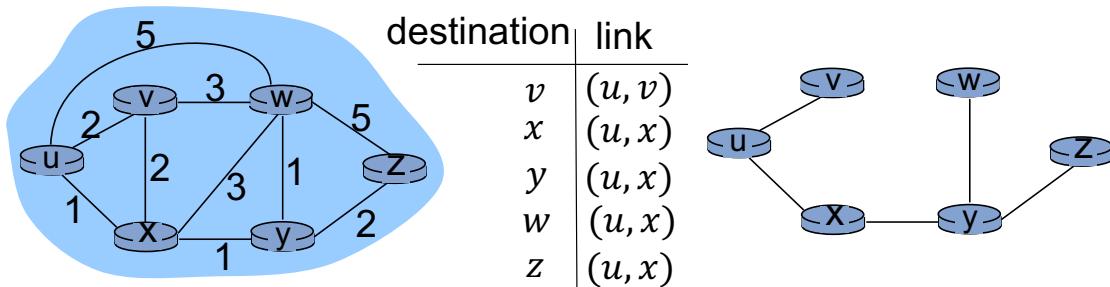
13.5 Link State přístup

- Globální směrovací informace – Každý uzel zná celou topologii sítě. Na začátku si uzly vymění informace o topologii sítě.
- Používá se pro směrování uvnitř autonomních systémů.
- *Single-metric.*

Dijkstrův algoritmus

- Pro formální vysvětlení a pseudokód viz otázku „Hledání nejkratších cest ze zdrojového uzlu do všech ostatních uzlů grafu“.
- Jde o algoritmus hledání nejkratších cest ze zdrojového uzlu do všech ostatních uzlů. To znamená, že ho provádí každý uzel (směrovač).
- Používá ho např. směrovací protokol OSPF (*Open Shortest Path First*).
 - Komunikace: vlastní protokol nad IP.
 - Metrika: odvozena od rychlosti linky.

| Step | N' | $D(v), p(v)$ | $D(w), p(w)$ | $D(x), p(x)$ | $D(y), p(y)$ | $D(z), p(z)$ |
|------|----------|--------------|--------------|--------------|--------------|--------------|
| 0 | u | $2, u$ | $5, u$ | $1, u$ | ∞ | ∞ |
| 1 | ux | $2, u$ | $4, x$ | | $2, x$ | ∞ |
| 2 | uxy | $2, u$ | $3, y$ | | | $4, y$ |
| 3 | $uxyv$ | | $3, y$ | | | $4, y$ |
| 4 | $uxyw$ | | | | | $4, y$ |
| 5 | $uxyvwz$ | | | | | |



Obrázek 13.7: Příklad výpočtu nejkratší cesty pomocí algoritmu Dijkstra. Step značí iteraci, N' je množina již prozkoumaných uzlů, D je pole vzdáleností do uzlu, p je pole předchůdzů uzlu.

13.6 Path Vector přístup

- Globální směrovací informace.
- Na síť je pohlíženo jako na množinu autonomních systémů.
- Používá se pro směrování mezi autonomními systémy.
- *Multi-metric.*

Path Vector algoritmus

- Posílá celou cestu (posloupnost uzlů – autonomních systémů), ne jenom vzdálenost.
- Slouží také pro detekci smyček v rámci cesty.
- Cíl je projít přes co nejméně autonomních systémů.
- Umožňuje tzv. *flexible policies*, můžu se rozhodnout jakou cestou budu provoz směrovat. Část provozu můžu přes AS_1 , jinou část pres AS_2 , apod.
- Známe celou cestu směrování, víme přes co pakety půjdou, můžeme některý uzel vyřadit.
- Směrovací protokol Border Gateway Protocol (BGP).

Kapitola 14

PDS – Řízení toku dat (flow-control) a prevence zahlcení (congestion-control) na transportní vrstvě (MP-TCP, QUIC, SCTP, DCCP).

14.1 Metadata

- Předmět: Přenos dat, počítačové sítě a protokoly (PDS)
- Přednáška:
 - 02-transportni-protokoly.pdf
- Záznam:
 - 2021-02-19

14.2 Úvod a kontext

Viz. „Úvod a kontext“ v předchozích otázkách z tohoto předmětu.

Transportní vrstva Transportní vrstva je název čtvrté vrstvy modelu vrstvové síťové architektury (ISO/OSI model). Leží mezi vrstvou síťovou (L3) a aplikační (L7).

- Činnost na straně odesílatele: obdrží data z aplikační vrstvy, nasegmentuje je a ke každému segmentu/datagramu přidá L4 hlavičku.
- Činnost na straně příjemce: obdrží segmenty/datagramy, které uspořádá do správného pořadí (*out-of-order* doručení), předá je aplikační vrstvě.
- Komunikace mezi vrstvami L4 a L7 probíhá pomocí socketů.
- Zodpovědnost za *end-to-end* spojení.
- Adresuje aplikace pomocí portových čísel.
- *Quality of Service* – zotavení se z chyb, spolehlivost, řízení toku, řízení zahlcení.

Multi homing – Využití více bodů připojení zároveň (pro jednu komunikaci). Pokud je podporováno a změní se IP adresa, nemusí se navázat nové spojení.

Connection-oriented – Typ spojení, kdy je nejprve nutné, navázat komunikaci. Typicky tři fáze: navázání spojení, přenos dat, ukončení spojení.

Connection-less – Typ spojení, které nerozlišuje jestli spojení existuje, nebo nikoliv. Pouze fáze přenosu dat.

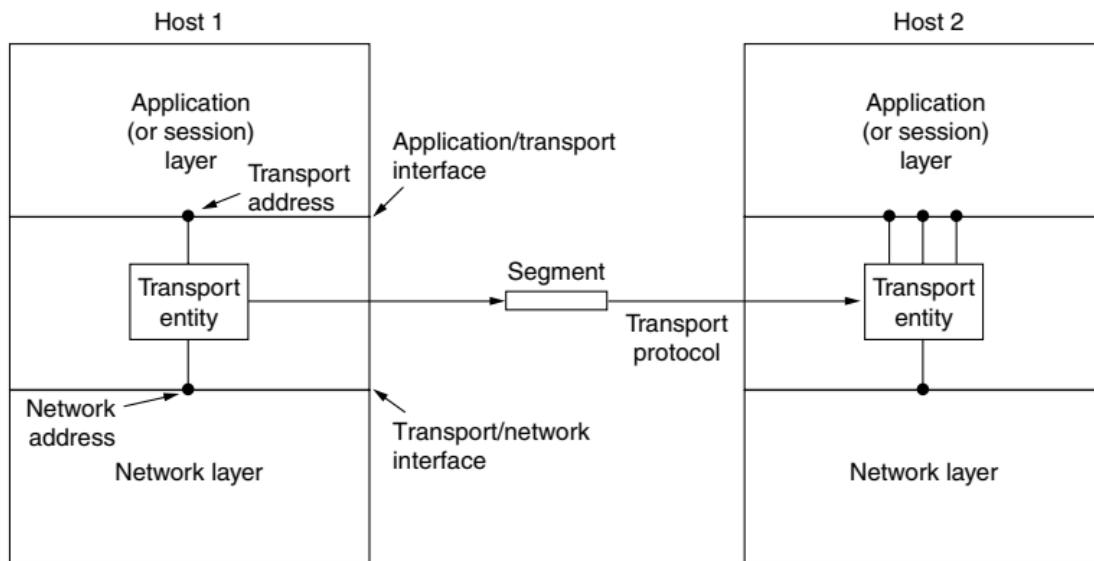
Round trip time (RTT) RTT je doba, která uplyne od vyslání paketu z jednoho komunikujícího uzlu k druhému po návrat potvrzení zpátky na první uzel.

bandwidth, throughput, goodput

- *Bandwidth* – šířka pásma (maximální přenos dat danou cestou v jeden moment)
- *Throughput* – propustnost komunikace na úrovni protokolu
- *Goodput* – propustnost komunikace na úrovni aplikace (throughput – režie)

Bitové chyby

- „Přeskočení bitu“ (*bit swapping*), jde o chyby na úrovni zpracování informace uzlem v síti (typicky spojené s přenosovým médiem)
- Zabýváme se jimi na úrovni L2
- Řešením je obvykle redundancy (Hammingův kód, CRC, Viterbiho algoritmus, ...)



Obrázek 14.1: Transportní vrstva (L4).

14.3 Paketové chyby

- Jde o chyby na úrovni přenosu informace mezi dvěma uzly.
- Zabýváme se jimi na transportní vrstvě.
- Řešením je obvykle znovuzaslání paketu.
- Můžeme měřit:

- PER (*packet error rate*)

$$PER = \frac{\text{počet chybných paketu}}{\text{počet prenesených paketu}}$$

- BER (*bit error rate*)

$$BER = \frac{\text{počet chybných bitu}}{\text{počet prenesených bitu}}$$

- Dále jsou představeny druhy paketových chyb.

Ztráta (zpoždění) paketu Jak může dojít ke ztrátě (zpoždění) paketu?

- Neopravitelná bitová chyba (rámcem je zahozen na L2)
- Zahlcení linky (směrovače jsou přetíženy a některé pakety zahazují)
- Špatné směrovací tabulky (paket jede špatnou cestou, nebo cyklí a vyprší mu TTL)

Ztráta fragmentovaných dat Fragment se může ztratit ze stejného důvodu jako standardní nefragmentovaný paket.

Duplicita paketu Příjemce obdržuje stále paket se stejným sekvenčním číslem. Odesílatel si myslí, že se paket k příjemci nikdy nedostal. Proč?

- Ztratí se potvrzení o přijetí paketu.

Vložení paketu Příjemce obdrží paket, který do spojení nepatří. Jak k tomu může dojít?

- Přijetí zpožděného paketu, který dorazil až po skončení jednoho a začátku dalšího separátního datového toku.
- Podvrhávání paketu útočníkem snažícím se narušit integritu sítě.
- Paket je chybou „zvláštně zmrzačen“ tak, že chybu nelze rozpoznat (např. je změněm bit v cílové IP adrese) a je směrován jinému příjemci.

Změna pořadí Vychází z designu počítačových sítí, kdy paketu mohou proudit různými cestami a s různým zpožděním. Data jsou *out-of-order* a příjemce je musí přeuspěřádat.

14.4 Řízení toku

Řízení toku (*flow control*) je řízení komunikaci mezi dvěma uzly na straně koncového systému. Řízení zahlcení (*congestion control*) je řízení komunikace v rámci sítě. Alternativně, se za řízení toku považuje koncept detekce a korekce paketových chyb. Často jsou tyto termíny vykládány různě.

14.4.1 Detekce paketové chyby

Paketové chyby detekujeme pomocí sekvenčních čísel. Sekvenční číslo je jedinečný identifikátor paketu v rámci datového toku, který identifikuje jeho pořadí. Tímto poznáme ztrátu, duplicitu i změnu pořadí.

Jak rozpoznat ztrátu od zpoždění?

- **Timeout** – Může být fixní, ale v lepším případě se odvozuje od RTT.
- **Negativní potvrzování** – Paket jsem dostal (ACK) vs. paket jsem nedostal (NACK). Pokud je příčina zpoždění zahlcení, tak tento přístup linku ještě více zahltí.

14.4.2 Korekce paketových chyb

Pokud paket nedorazil, odesílatel ho pošle znova. Tzv. *Automatic Repeat/Request* (ARQ) – Čeká se na ACK, když nepřijde do timeoutu, dojde k znovuzaslání. Princip tzv. klouzavého okna (*sliding window*). Existují 3 strategie.

Stop and wait

- Klouzavé okno o velikosti 1.
- Odesílatel pošle paket a čeká na potvrzení. Po přijetí potvrzení pošle další paket.
- Špatná efektivita využití pásma.

Go back n

- Buffer na straně odesílatele.
- Příjemce potvrzuje naposledy přijatým paketem (např. příjemce pošle ACK2 znamená, že dostal pakety 0, 1, 2).
- Plýtvání při znovuzaslání (např. odesílatel pošle 5 paketů, ztratí se první, musí poslat znova všech 5).

Selective repeat

- Buffery jsou na obou stranách.
- Efektivní využití šířky pásma, ale složitější implementace.
- Příjemce potvrzuje:
 - Fast Retransmit – Příjemce posílá potvrzení s naposledy přijatým paketem.
 - Bitová maska – V ACK je zavedeno další políčko (bitová maska), které obsahuje informace o tom, co příjemce přijal a co ne.
 - NACK – Příjemce posílá ACK, pokud balík dostal a NACK, pokud nikoliv.

14.5 Řízení zahlcení

Zahlcení je detekováno, co se s tím dá dělat?

- Zvýšit kapacitu sítě,
- Snížit množství provozu.

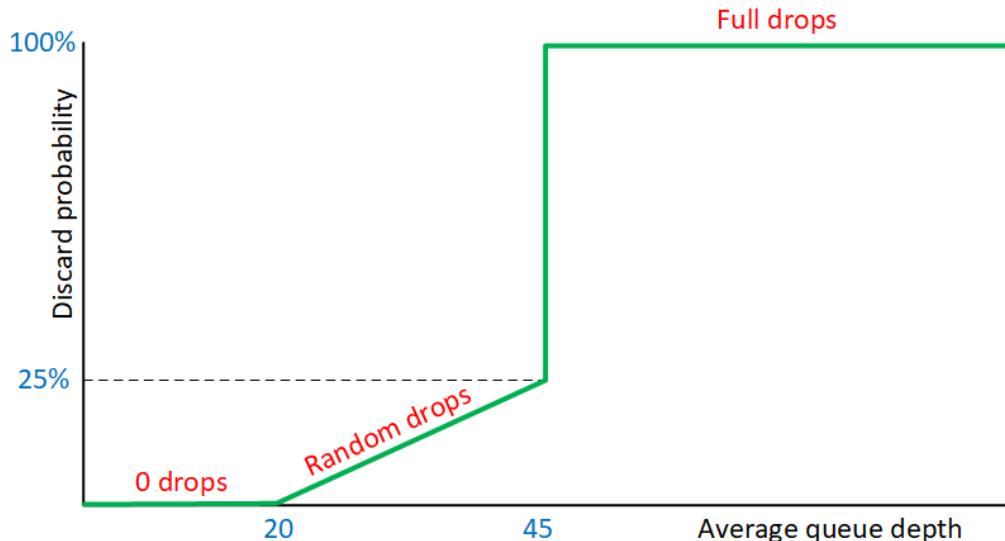
Jak se to dá dělat?

- Opravit zahlcený provoz (*repair*) – Je detekováno zahlcení a řešíme, co budeme dělat.
- Předcházet problému zahlcení (*avoidance*) – Snažíme se zahlcení předcházet (*prevention*).

14.5.1 Zahazování paketů (opravování provozu)

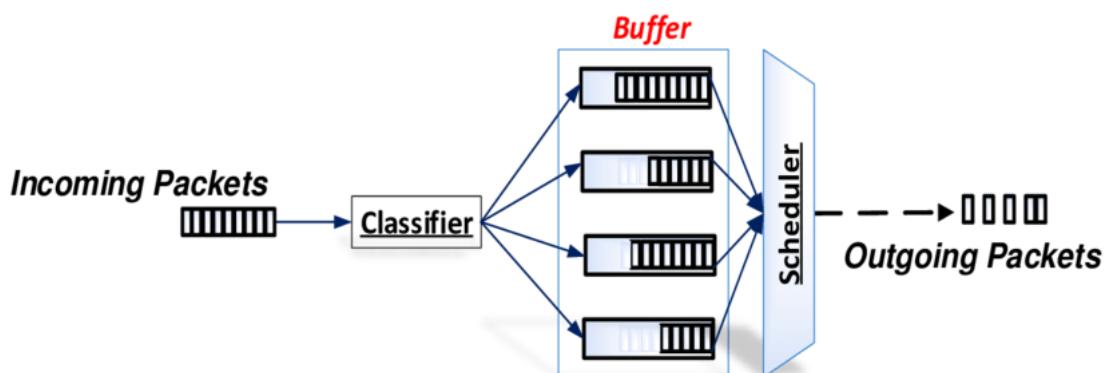
Každý uzel (nejčastěji směrovač) v síti kontroluje délku své fronty, pokud příliš narůstá, tak začne preventivně některé pakety zahazovat (to však vede k jejich znovuzaslání).

RED (Random Early Detection) Algoritmus RED sleduje velikost fronty a zahazuje pakety na základě pravděpodobnosti. Pokud je fronta téměř prázdná, pravděpodobnost zahození je 0. Jak fronta roste, roste i pravděpodobnost zahození příchozího paketu. Když je fronta plná, pravděpodobnost dosáhne hodnoty 1 a všechny příchozí pakety jsou zahazovány.



Obrázek 14.2: Vizualizace RED.

WRED (Weighted Random Early Detection) Algoritmus WRED je rozšířením RED o přidání váh pro jednotlivé třídy paketů (máme frontu pro každou třídu). WRED nezahazuje všechny pakety se stejnou pravděpodobností, ale rozlišuje jejich důležitost (pomocí hodnot IP precedence nebo DSCP). Některé pakety chceme prioritizovat, např. VoIP (*Voice over Internet Protocol*), RTP (*Real-time Transport Protocol*).



Obrázek 14.3: Vizualizace WRED.

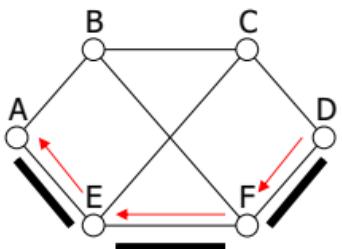
14.5.2 Choke pakety (opravování provozu)

Příjemce v síti má problém zpracovávat provoz (nestíhá), tak dá vědět odesílateli (odešle choke paket), aby zpomalil.

Plain choke packets Viz obrázek 14.4. Uzel D detekuje zahlcení způsobené provozem od A. Uzel D pošle choke paket A. Proti zahlcení „zabuje“ pouze odesílatel (A).

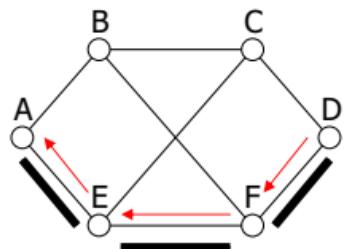
Hop-by-hop choke packets Viz obrázek 14.4. Uzel D detekuje zahlcení způsobené provozem od A. Uzel D pošle choke paket A. Každý uzel po cestě „bojuje“ proti zahlcení (F, E, A).

Plain Choke packets



A heavy flow is established
Congestion is noticed at D
A Choke packet is sent to A
The flow is reduced at A
The flow is reduced at D

Hop-by-hop Choke packets



A heavy flow is established
Congestion is noticed at D
A Choke packet is sent to A
The flow is reduced at F
The flow is reduced at D

Obrázek 14.4: Plain choke packets vs. hop-by-hop choke packets.

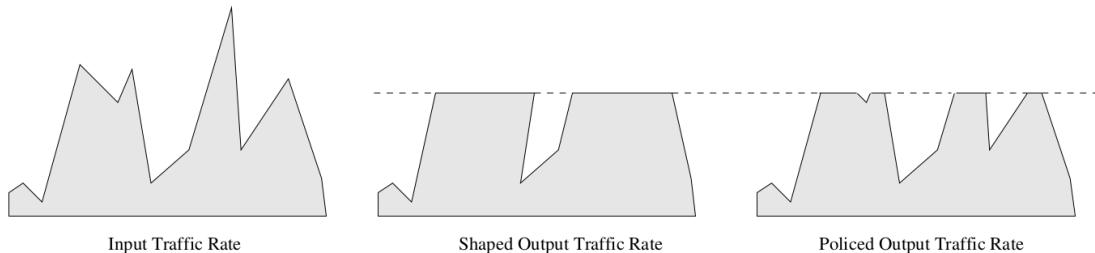
14.5.3 Ořezání a rozložení provozu (předcházení zahlcení)

Jedná se o předcházení zahlcení ze strany odesílatele. Odesílatel je „ohleduplný“ a myslí na to, aby příjemce nebyl zahlcen. Odesílatel má pevně stanovenou hranici kolik může maximálně posílat. Implementace pomocí tzv. *token bucket*.

Ořezání provozu (policing) Cokoliv nad hranici je „oříznuto“ a zahozeno.

Rozložení provozu (shaping) Cokoliv nad hranici je uloženo do bufferu a odesláno později, až je menší provoz.

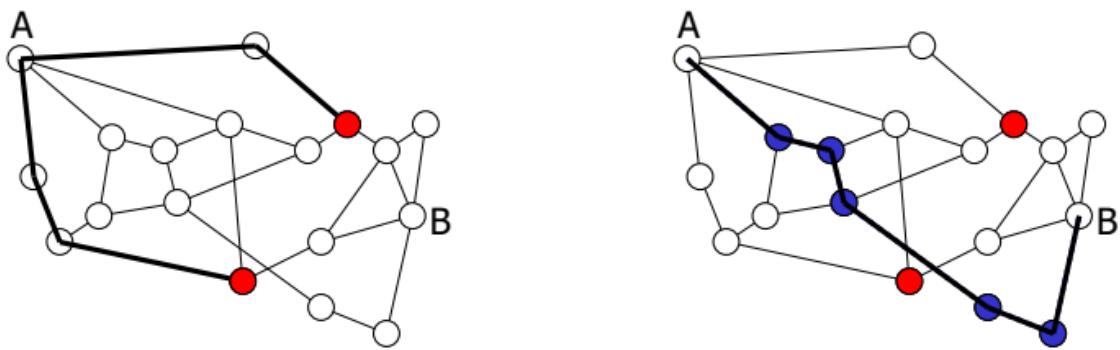
Čím se liší rozložení provozu (policing) a omezení provozu (shaping)?



Obrázek 14.5: Ořezání provozu vs rozložení provozu.

14.5.4 Rezervace (předcházení zahlcení)

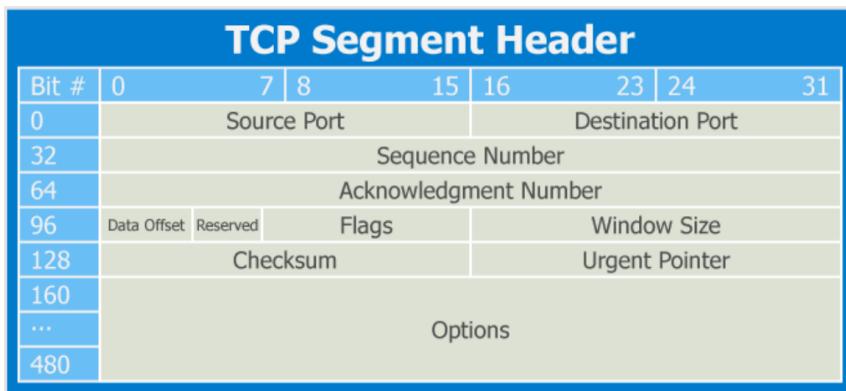
Zařízení v síti (směrovače) jsou schopni se domluvit a vytvořit si tzv. „virtuální okruh“ pro konkrétní komunikaci. Každé zařízení na okruhu ví o provozu a rezervuje potřebné pásmo (předem se deklarují kapacity). Je nutné více signaliace (režie).



Obrázek 14.6: Příklad rezervace.

14.6 TCP (*Transmission Control Protocol*)

- Garantuje spolehlivé doručení.
- Garantuje doručení v pořadí.
- Řízení toku a zahlcení.
- *Connection-oriented* – Navázání spojení pomocí *three-way handshake*.
- Pracuje s *byte stream* – nezohledňuje hranice aplikačních dat.
- Sekvenční číslo závisí na tom, kolik bajtů se posílá. Např. $seq = 92, data = 8 B, seq = 100, data = 20 B, seq = 120, data = 13 B, \dots$

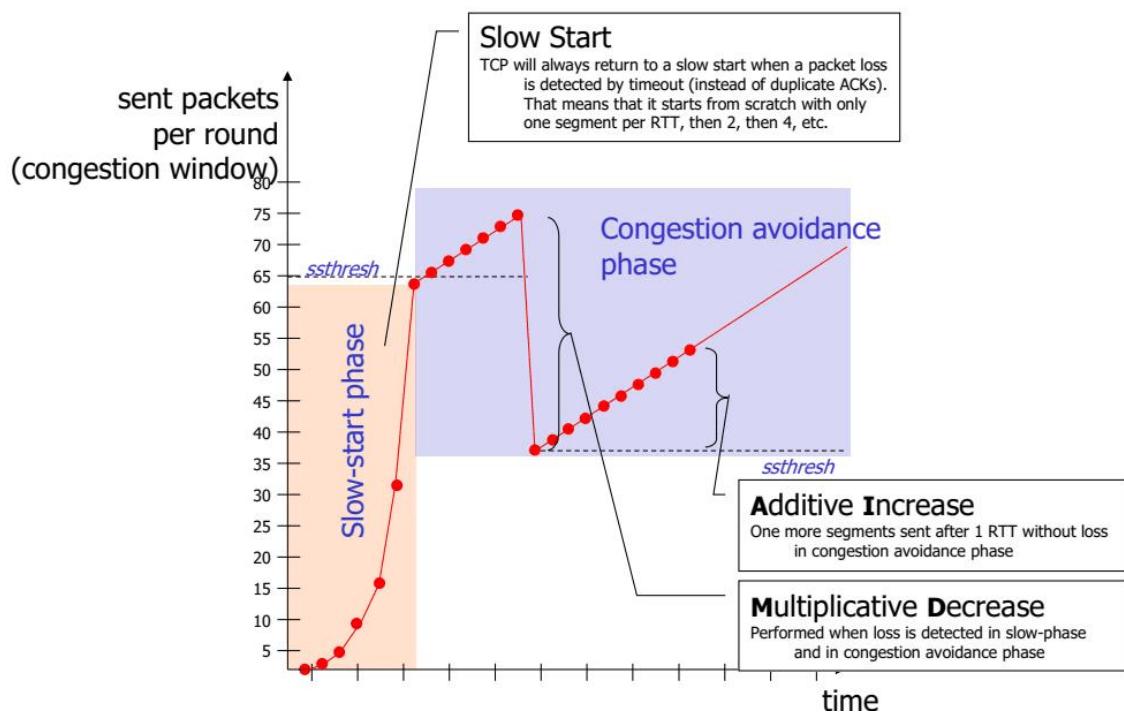


Obrázek 14.7: Hlavička TCP.

Řízení toku a zahlcení V TCP jsou rozlišovány 3 fáze přenosu dat.

- *Slow start* – Exponenciální zvyšování rychlosti odesílání až po dosažení nějaké hranice.
- *Congestion avoidance (additive increase)* – Lineární zvyšování počtu odeslaných paketů.
- *Congestion avoidance (multiplicative decrease)* – Skokové snížení rychlosti odesílání. Nastane, pokud je detekována ztráta paketu.

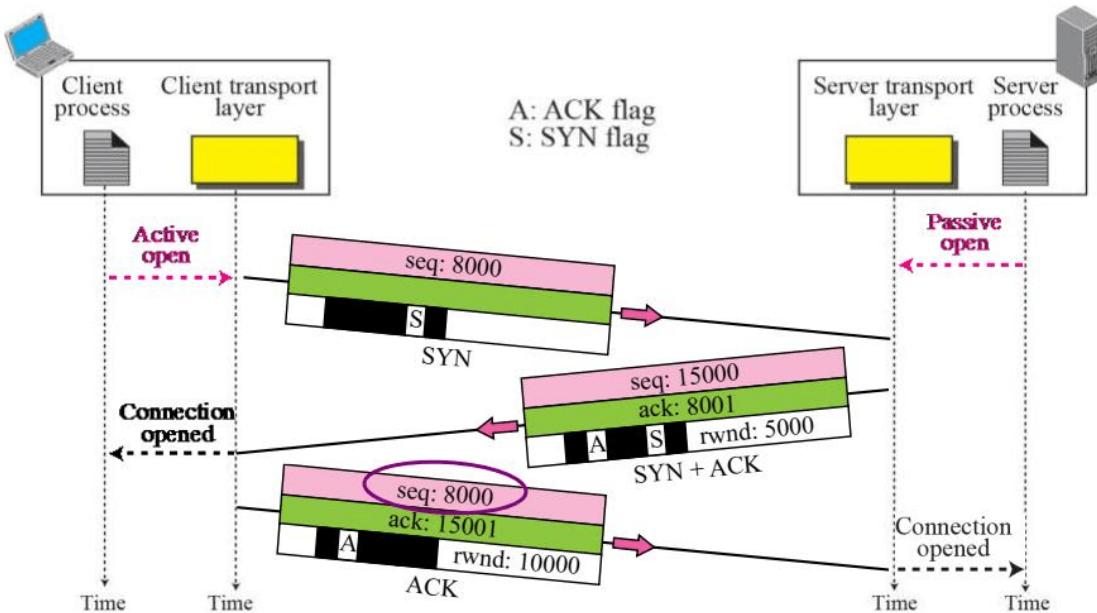
V TCP existuje spoustu algoritmů pro řízení toku a zahlcení, které implementují tento princip (Tahoe, Reno, New Reno, Vegas, CUBIC, Westwood, ...).



Obrázek 14.8: Jednotlivé fáze řízení zahlcení v TCP.

Nevýhody a problémy

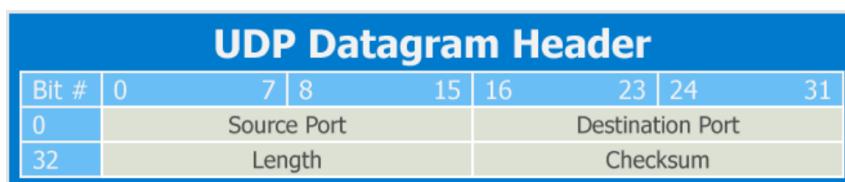
- Nepodporuje *multi homing*.
- *Head of line blocking* – Pokud dojde ke ztrátě paketu, tak vše za ním je pozdrženo
- Velká režie potvrzování (až 35 % veškerého provozu je režie)



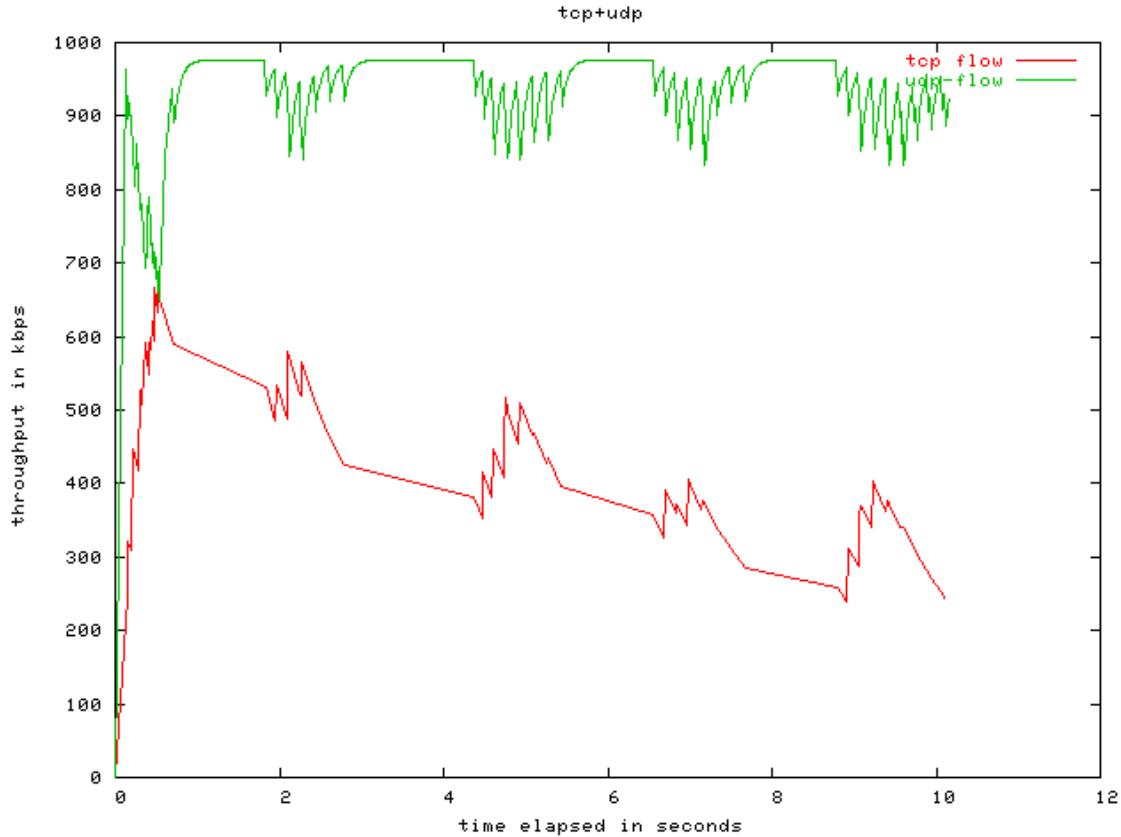
Obrázek 14.9: Zahájení spojení pomocí tzv. *Three-Way Handshake*.

14.7 UDP (*User Datagram Protocol*)

- Negarantuje spolehlivé doručení (*best effort*).
- Negarantuje doručení v pořadí.
- *Connection-less*.
- Pracuje s *byte stream* – nezohledňuje hranice aplikačních dat.
- Má minimální režii.



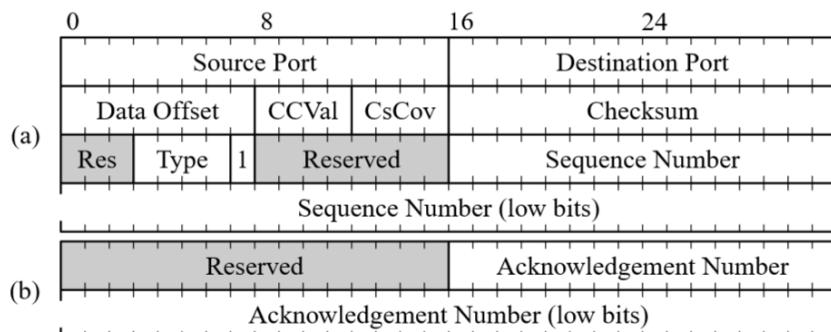
Obrázek 14.10: Hlavička UDP.



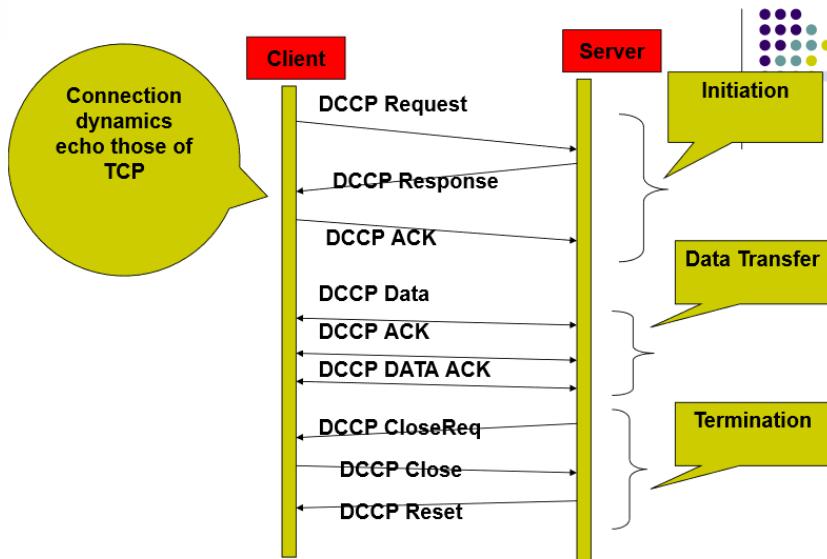
Obrázek 14.11: Problém UDP, kdy si pro sebe zabere celé pásmo, protože TCP se chová „zodpovědně“ a snižuje svůj provoz.

14.8 DCCP (Datagram Congestion Control Protocol)

- Cílem je zajistit podporu řízení zahlcení pro UDP.
- Negarantuje spolehlivé doručení (*best effort*).
- Negarantuje doručení v pořadí.
- *Connection-oriented* – Navázání spojení pomocí *three-way handshake*.
- Využití pro audio/video konference.



Obrázek 14.12: Hlavice DCCP. Sekvenční a potvrzovací čísla souvisí pouze s účtováním provozu kvůli řízení zahlcení, nikoliv pro zajištění spolehlivého přenosu.

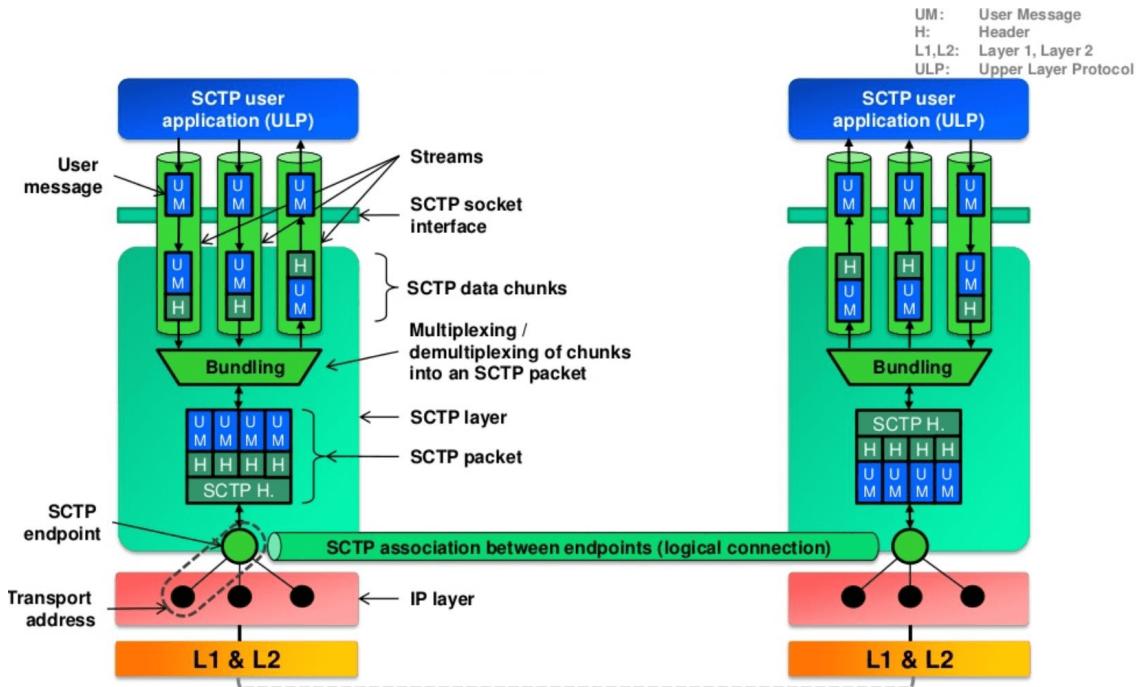


Obrázek 14.13: Komunikace pomocí DCCP. Obsahuje několik typů zpráv: Request, Response, ACK, Data, Data ACK, CloseReq, Close, Reset.

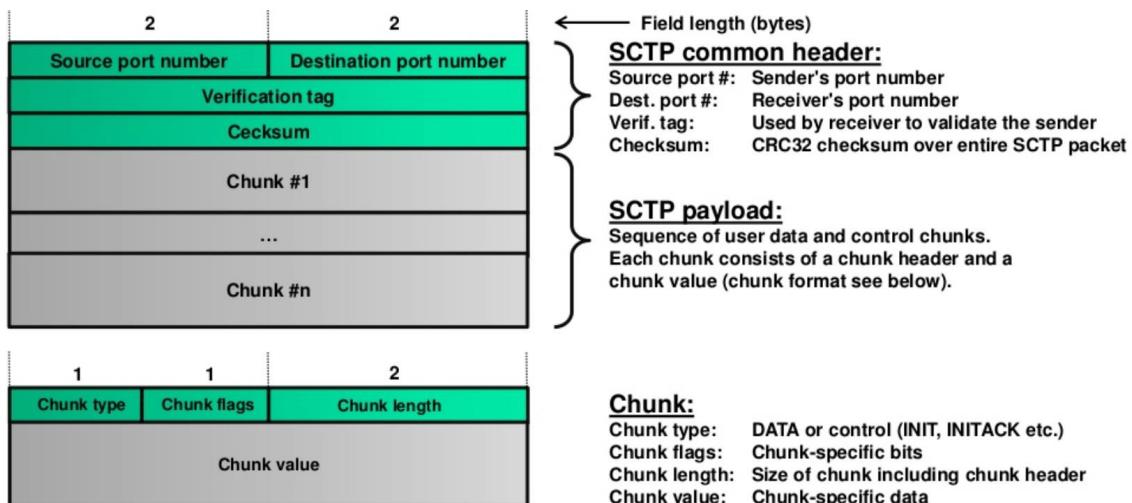
Jak je zahlcení realizováno V potvrzeních se nastavují tzv. ECN (*Explicit Congestion Notification*) bity. Pomocí nich dává příjemce odesílateli vědět, jak moc je zahlcen. V průběhu navazování spojení se pomocí CCID (*Congestion Control ID*) dohodne, jaký typ řízení zahlcení se bude používat (algoritmy jako pro TCP, nebo vlastní).

14.9 SCTP (*Stream Control Transmission Protocol*)

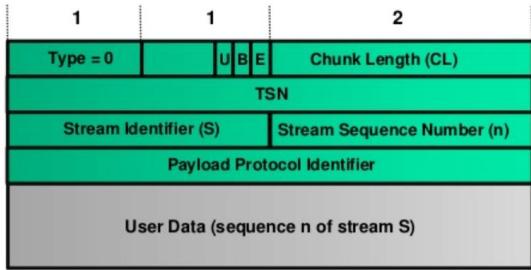
- Garantuje spolehlivé doručení.
- Garantuje doručení v pořadí.
- Pracuje s *message stream* – Respektuje hranice aplikačních dat.
- *Connection-oriented* – Navázání spojení pomocí *four-way handshake*.
- *Path MTU discovery* – Obrana proti fragmentaci (od zdroje k cíli jsou prozkoumány MTU všech uzlů).
- Netrpí problémem *head of line blocking*.
- Podporuje *multi homing*.
- Řízení zahlcení.
- Aplikace je třeba přeprogramovat, aby místo TCP socketů využívali SCTP sockety.



Obrázek 14.14: Komunikace pomocí SCTP. Komunikace s aplikací probíhá přes sockety. Každé aplikační zprávě je přidán základní hlavičkou – vzniká *data chunk*. Několik data chunku je zabaleno dohromady a opatřeno SCTP hlavičkou – vzniká SCTP paket. Ten je poté možno posílat přes více IP spojení.



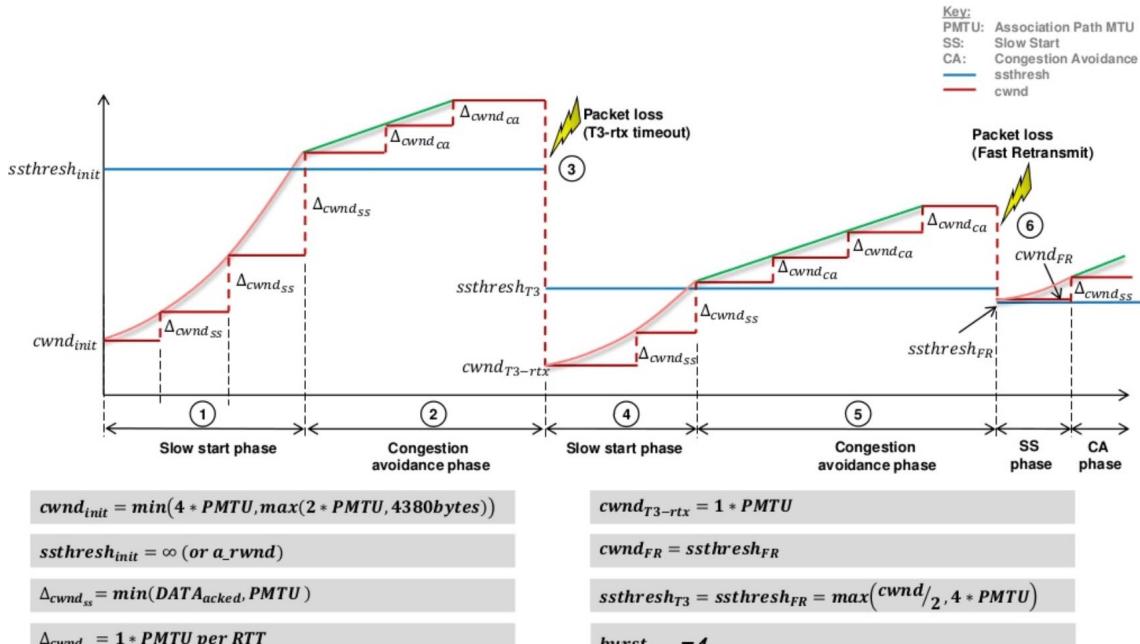
Obrázek 14.15: Hlavička SCTP. SCTP paket a *data chunk*.



Data chunk:

- U bit: If set to 1, indicates that this is an unordered chunk. Unordered chunks are transmitted and sent to the receiving application as is without re-ordering based on the sequence number.
- B bit: Begin of fragment bit. If set to 1 this is the first fragment of a larger fragmented user data message.
- E bit: End of fragment bit. If set to 1 this is the last fragment of a larger fragmented user data message.
- Stream identifier: Identifies the stream to which this chunk belongs.
- Stream seq. no.: Sequence number of user data within this stream. In case of fragmentation this number is identical for all fragments.
- Payload proto. id.: Identifies the upper (application) layer protocol (e.g. HTTP).
- User data: Application user data (e.g. HTTP header and payload).

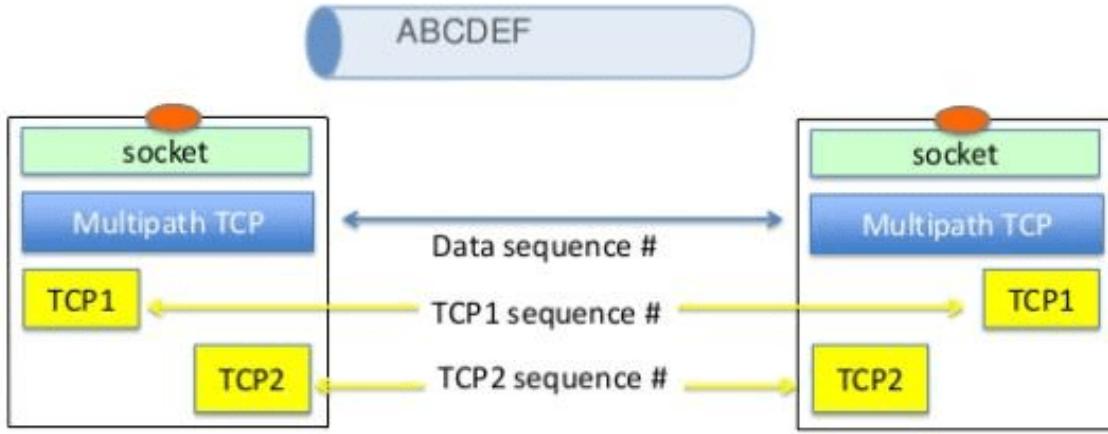
Obrázek 14.16: Detailní hlavička *data chunk*. Identifikátor streamu a identifikátor paketu v rámci streamu.



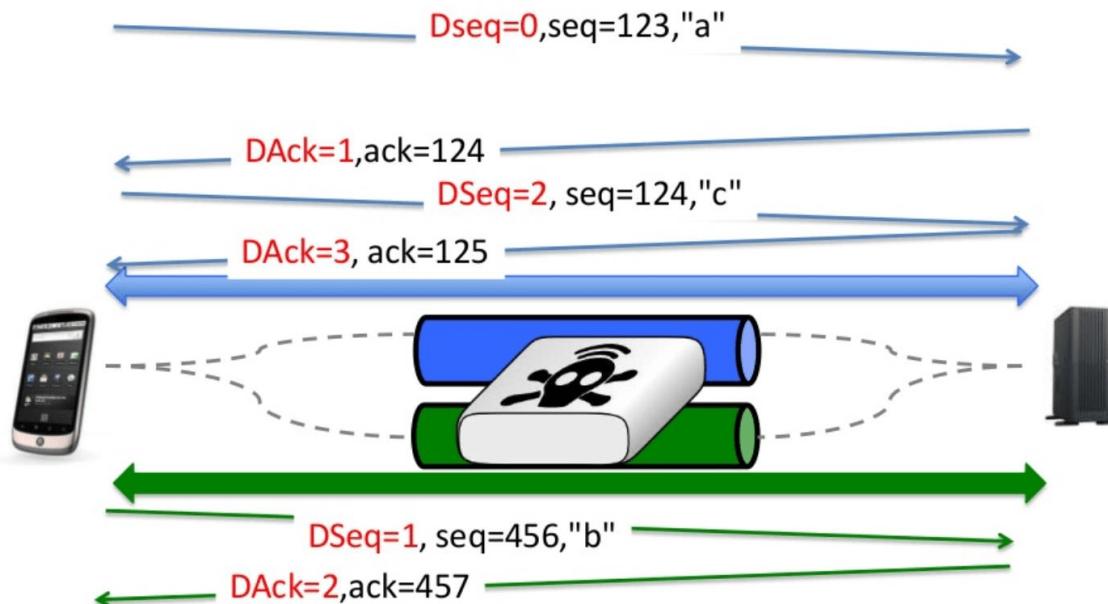
Obrázek 14.17: Řízení zahlcení (*slow start*, *congestion avoidance (additive increase)*, *congestion avoidance (multiplicative decrease)*) je podobné jako TCP Westwood.

14.10 MP-TCP (*Multipath TCP*)

- Cílem je vzít co nejvíce vlastností z SCTP a poskytnout je v TCP bez toho, aniž by se musely přeprogramovávat aplikace (má obdobné sockety jako TCP).
- Má menší signalizační náročnost oproti SCTP.
- Chová se jako TCP pouze, podporuje *multi homing* a má vlastní řízení zahlcení.
- Rozšíření TCP hlavičky o tzv. *options*.



Obrázek 14.18: MPTCP komunikace, rozdělení aplikačních dat na paralelní TCP streamy. Rozlišujeme datová sekvenční čísla a sekvenční čísla TCP.

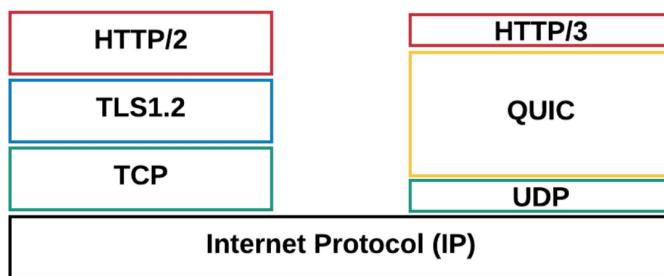


Obrázek 14.19: MPTCP řízení toku. Dvě TCP komunikace a posílání řetězce „abc“. Navázání spojení podobné jako v TCP. Může vypadnout i celý komunikační kanál.

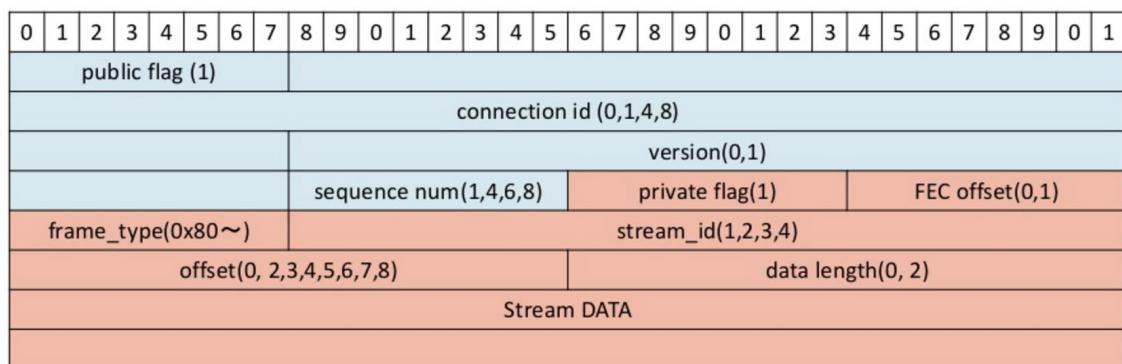
14.11 QUIC (*Quick UDP Internet Connections*)

- Motivace pro QUIC: dnes většina provozu probíhá šifrované přes HTTPS. Navázání bezpečného HTTPS spojení má velkou režii (HTTPS = TCP + SSL + HTTP). Až $3 \times \text{RTT}$ přes započetím komunikace!
 - TCP: syn, syn ack, ack
 - TLS: client hello, server hello, change cipher spec, encrypted handshake message
- Staví nad UDP (aby všechny zařízení v internetu nemuseli implementovat nový transportní protokol).

- Paralelní datové streamy, aby bylo zabráněno *head of line blocking*.
- Má zabudovanou *Forward Error Correction* (FEC), přidává bitovou redundanci, tím umožňuje velké množství ztracených paketů dopočítat. Pokud to není možný, tak „chytrý“ selective repeat.



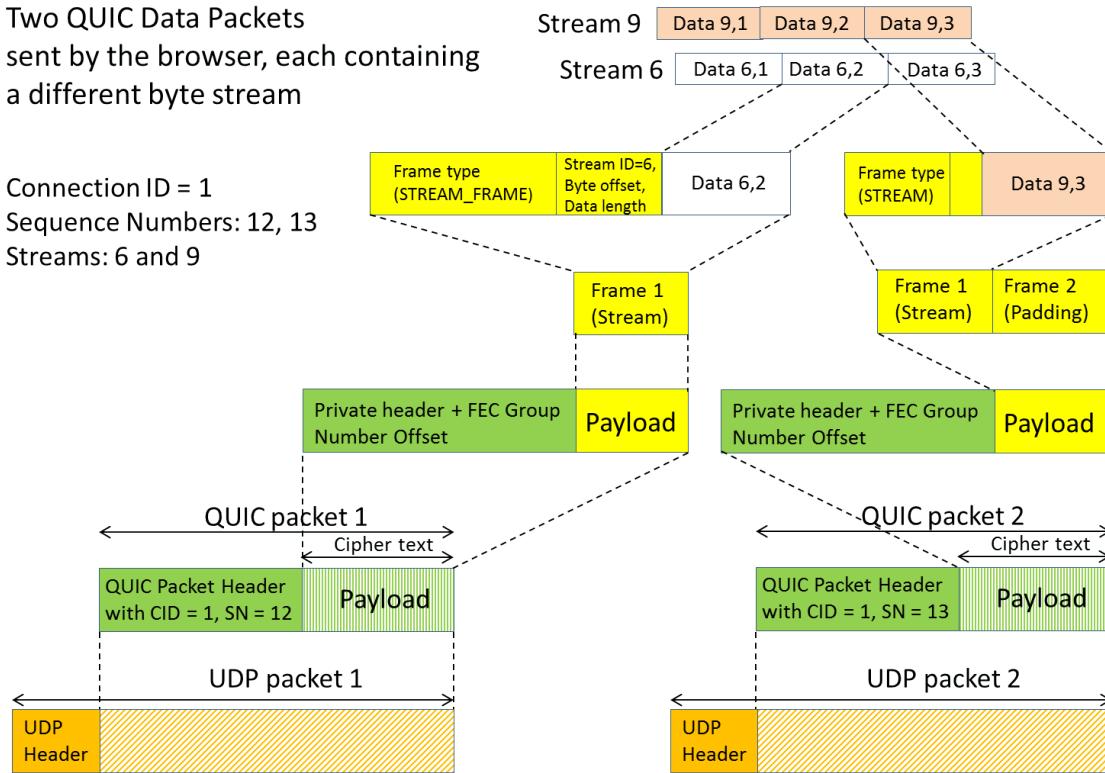
Obrázek 14.20: Standardní TCP stack vs. QUIC.



Obrázek 14.21: Hlavička QUIC. Spousta políček má proměnlivou velikost.

Two QUIC Data Packets
sent by the browser, each containing
a different byte stream

Connection ID = 1
Sequence Numbers: 12, 13
Streams: 6 and 9



Obrázek 14.22: Složení QUIC zprávy.