

Dana Pardubská

Katedra informatiky, FMFI UK, Bratislava

+421 2 602 95 158

[pardubska@fmph.uniba.sk](mailto:pardubska@fmph.uniba.sk)

<http://www.dcs.fmph.uniba.sk/~pardubska/brno>

## Čo je distribuovaný systém?

**Prepojenie** sústavy **autonómnych** počítačov, procesov, procesorov

**Autonómny**- vlastné riadenie

**Prepojenie**- komunikácia/ výmena informácie

## Prečo distribuované systémy?

**pôvodne** veľa procesov rozdistribuovaných vo veľkej geografickej ploche

**neskôr** lokálne siete, multiprocessorové systémy so spoločne zdieľanou pamäťou

- Výmena informácií (WAN)
- Zdieľanie zdrojov
- Zvýšenie spoľahlivosti replikáciou
- Zvýšenie výkonu využitím paralelizmu
- Zjednodušenie návrhu využitím špecializácie

## rozdiely DA

### spôsob medziprocesorovej komunikácie

shared memory / point-to-point posielanie správ / broadcasting / remote procedure calls

### časovanie

- **synchrónne systémy** (spoločný čas)
- **asynchrónne systémy** (rôzne rýchlosti, rôzne poradie krokov)
- **čiasťočne synchronizované** (čiasťočné info o časovaní udalostí)

### model chýb

- spoľahlivý vs. nespoľahlivý hardware
- tolerancia nejakého množstva chybného správania
  - $\left\{ \begin{array}{ll} \textit{chyby procesora} & \text{stop / crash / byzantínske} \\ \textit{chyby kanálov} & \text{strata / duplikácia / modifikácia správ} \end{array} \right.$

**DA** - sústreďujeme sa na algoritmy s vysokým **stupňom neurčitosti** a veľkou **nezávislosťou** jednotlivých **akcií**

- neznámy počet procesorov
- neznáma topológia
- nezávislé vstupy na rôznych miestach
- viacero programov (rôzny čas začiatku, rôzne rýchlosti, ..)
- nedeterminizmus procesorov
- neistý čas doručenia správ
- neznáme poradie doručovania správ
- chyby procesorov a komunikácie/liniek

## **Distribuované vs. Centralizované algoritmy**

znalosť globálneho stavu

pojmem globálneho času

nedeterminizmus/determinizmus

Kvôli dôkazom korektnosti/ zložitosti / neriešiteľnosti používame **model**

## Príklad - spoľahlivá výmena informácií cez nespoľahlivé médium

Procesy a,b; procedúry kontroly siete NCP -A,B; správa m

Naviazanie a ukončenie komunikácie

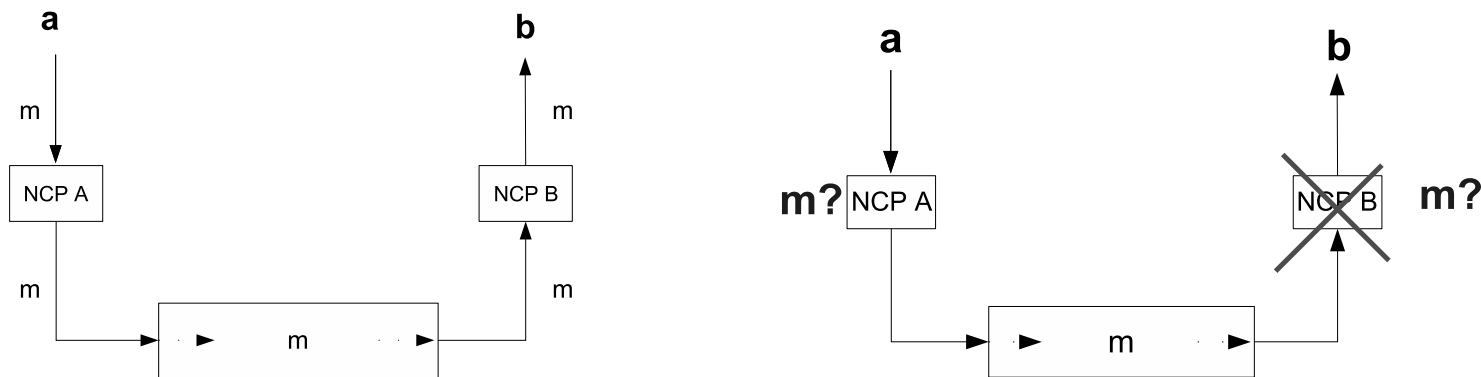
Informačná jednotka sa môže stratiť alebo duplikovať, NCP môže zlyhať (reštart v uzavretom stave)

### Spoľahlivá výmena nie je dosiahnuteľná

→ Inicializácia komunikácie procesom a

→ NCP A a NCP B začnú komunikáciu, počas ktorej B doručí m do b

→ NCP B sa zrúti a je reštartovaný v uzavretom stave



V tejto situácii ani A ani B nevedia, či m bolo doručené

## Konverzácia jednou správou- $\langle \text{data}, m \rangle$

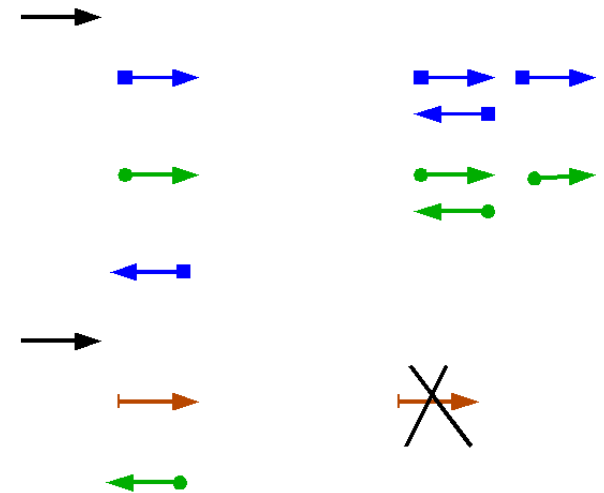
1. NCP A **send**  $\langle \text{data}, m \rangle$ , notify, close
  2. NCP B receive  $\langle \text{data}, m \rangle$ , deliver m, close
- strata, keď sa správa nedoručí  $\Rightarrow$  zavádza sa **potvrdenie príjmu** - nedochádza k duplikácii

## Konverzácia dvomi správami- $\langle \text{data}, m \rangle$ , - $\langle \text{ack} \rangle$

1. NCP A **send**  $\langle \text{data}, m \rangle$
  2. NCP B receive  $\langle \text{data}, m \rangle$ , deliver m, **send**  $\langle \text{ack} \rangle$ , close
  3. NCP A receive  $\langle \text{ack} \rangle$ , notify, close
- zavádza sa **vypršanie času a opakované poslanie správy**
1. NCP A **send**  $\langle \text{data}, m \rangle$
  2. NCP B receive  $\langle \text{data}, m \rangle$ , deliver m, **send**  $\langle \text{ack} \rangle$ , close
  3. DN  $\langle \text{ack} \rangle$  sa stratí
  4. NCP A **timeout**, **send**  $\langle \text{data}, m \rangle$
  5. NCP B receive  $\langle \text{data}, m \rangle$  deliver m, **send**  $\langle \text{ack} \rangle$ , close
  6. NCP A receive  $\langle \text{ack} \rangle$ , notify, close
- opätovné poslanie správy  $\rightarrow$  možnosť duplikácie

## Pri zavedení potvrdenia sa môžu strážené správy stratiť

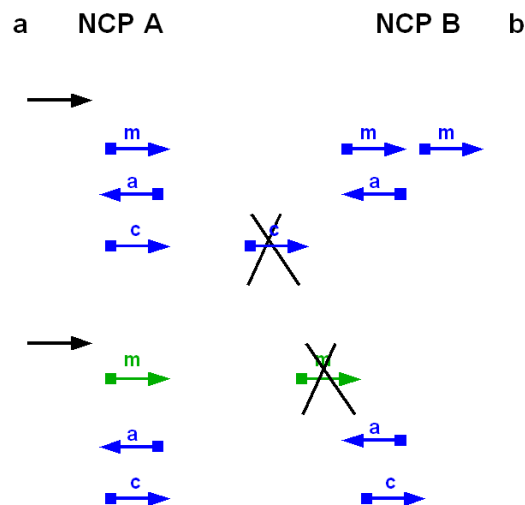
1. NCP A **send** < data, m1 >
2. NCP B receive <data, m1> deliver m1,  
**send** < ack >, close
3. NCP A timeout, **send** < data, m1 >
4. NCP B receive <data, m1> deliver m1,  
**send** < ack >, close
5. NCP A receive <ack>, notify, close
6. NCP A **send** < data, m2 >
7. DN <data, m2> lost
8. NCP A receive <ack>, notify, close



Dvojnásobné potvrdenie m1 spôsobí, že sa nedozvieme o strate m2

## Konverzácia tromi správami- $\langle \text{data}, m \rangle$ , $\langle \text{ack} \rangle$ , $\langle \text{close} \rangle$

1. NCP A **send**  $\langle \text{data}, m \rangle$
2. NCP B receive  $\langle \text{data}, m \rangle$ , deliver  $m$ , **send**  $\langle \text{ack} \rangle$
3. NCP A receive  $\langle \text{ack} \rangle$ , notify, **send**  $\langle \text{close} \rangle$ , close
4. NCP B receive  $\langle \text{close} \rangle$ , close



1. NCP A **send**  $\langle \text{data}, m1 \rangle$
2. NCP B receive  $\langle \text{data}, m1 \rangle$ , deliver  $m1$ , **send**  $\langle \text{ack} \rangle$
3. NCP A receive  $\langle \text{ack} \rangle$ , notify, **send**  $\langle \text{close} \rangle$ , close
4. DN  $\langle \text{close} \rangle$  sa stratí
5. NCP A **send**  $\langle \text{data}, m2 \rangle$
6. DN  $\langle \text{data}, m2 \rangle$  sa stratí
7. NCP B retransmit  $\langle \text{ack} \rangle$  (krok 2)
8. NCP A receive  $\langle \text{ack} \rangle$ , notify, **send**  $\langle \text{close} \rangle$ , close
9. NCP B receive  $\langle \text{close} \rangle$ , close

strata  $\langle \text{ack} \rangle$  nevedie k duplikovaniu; strata  $\langle \text{close} \rangle$  môže spôsobiť re-send  $\langle \text{ack} \rangle$



## Konverzácia tromi číslovanými správami- $\langle \text{data}, m, x \rangle$ , $\langle \text{ack}, x, y \rangle$ , $\langle \text{close}, x, y \rangle$

1. NCP A **send**  $\langle \text{data}, m, x \rangle$
2. NCP B receive  $\langle \text{data}, m, x \rangle$ , deliver m, **send**  $\langle \text{ack}, x, y \rangle$
3. NCP A receive  $\langle \text{ack}, x, y \rangle$ , notify, **send**  $\langle \text{close}, x, y \rangle$ , close
4. NCP B receive  $\langle \text{close}, x, y \rangle$ , close

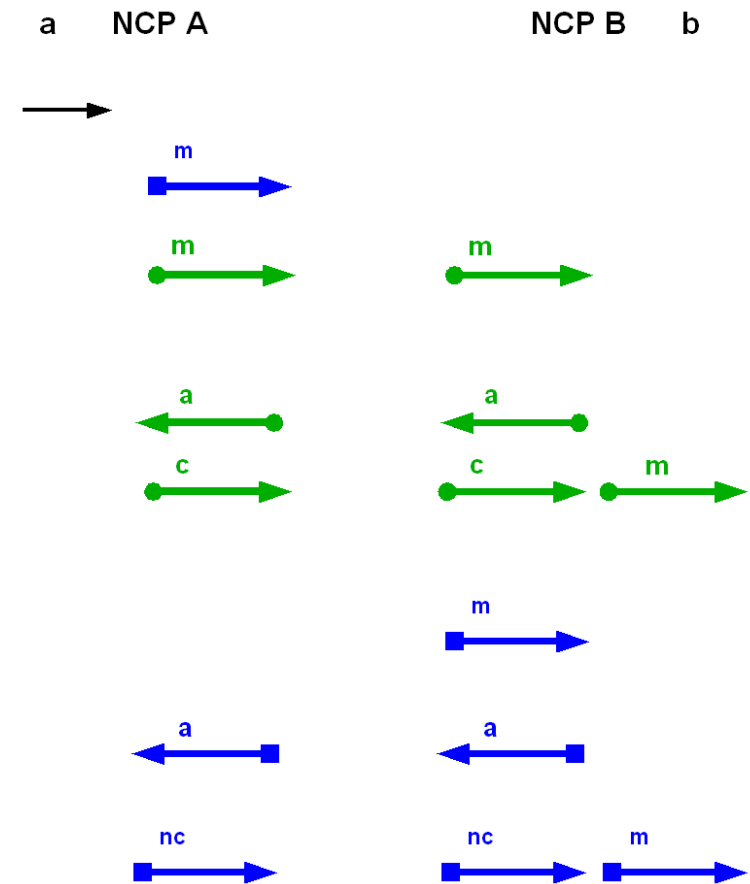
### Overíme platnosť $\langle \text{data}, m, x \rangle$

1. NCP A **send**  $\langle \text{data}, m, x \rangle$
2. NCP B receive  $\langle \text{data}, m, x \rangle$  **send**  $\langle \text{ack}, x, y \rangle$
3. NCP A receive  $\langle \text{ack}, x, y \rangle$ , notify, **send**  $\langle \text{close}, x, y \rangle$ , close
4. DN  $\langle \text{close}, x, y \rangle$  sa stratí
5. NCP B timeout, retransmit  $\langle \text{ack}, x, y \rangle$
6. NCP A receive  $\langle \text{ack}, x, y \rangle$ , reply  $\langle \text{nocon}, x, y \rangle$
7. NCP B receive  $\langle \text{nocon}, x, y \rangle$ , deliver m, close

**B musí doručiť dáta aj keď A nepotvrdí spojenie s x,y - môže nastať duplikácia**

Môže nastať **duplikácia**

1. NCP A **send**  $\langle \text{data}, m, x \rangle$
2. NCP A timeout,  
**retransmit**  $\langle \text{data}, m, x \rangle$
3. NCP B receive  $\langle \text{data}, m, x \rangle$  (step 2)  
**send**  $\langle \text{ack}, x, y1 \rangle$
4. NCP A receive  $\langle \text{ack}, x, y1 \rangle$ , notify,  
**send**  $\langle \text{close}, x, y1 \rangle$ , close
5. NCP B receive  $\langle \text{close}, x, y1 \rangle$ ,  
deliver m, close
6. NCP B receive  $\langle \text{data}, m, x \rangle$  (step 1),  
**send**  $\langle \text{ack}, x, y2 \rangle$
7. NCP A receive  $\langle \text{ack}, x, y2 \rangle$ ,  
**reply**  $\langle \text{nocon}, x, y2 \rangle$
8. NCP B receive  $\langle \text{nocon}, x, y2 \rangle$ ,  
deliver m, close



## Konverzácia štyrmi správami

**send** < open, x, y >, **send** < data, x, y >, **send** < agree, x, y >, **send** < ack, x, y >

Vzájomná dohoda na id. číslach konverzácie pred doručením dát

1. NCP A **send** < data, m, x >
2. NCP B receive <data, m,x>, **send** < open, x, y >
3. NCP A receive <open,x,y>, **send** < agree, x, y >
4. NCP B receive <agree,x,y>, deliver m, **send** < ack, x, y >, close
5. NCP A receive <ack, x,y>, notify, close

Stále dochádza k duplikácii, keď NCP padne. Dá sa to zmodifikovať tak, aby A zazna-  
menalo a ukončilo po prijatí <noncon,x,y>; tým zabránime duplicitu, ale môžu nastať  
straty

... □

# Model distribuovaných výpočtov

## Distribuovaný systém

sieť –  $G=(V, E)$ ,  $|V|=n$  procesorov,  $|E|=m$  kanálov

## Processor

lokálna pamäť

komunikácia so susedmi výmenou správ

správanie popísané stavmi, riadené správami

## Formálny popis

Burns(1980), Lynch, Fischer(1981), Angluin(1980)

# Predpoklady modelu distribuovaných výpočtov

## ⇒ Asynchrónnosť

procesory

komunikácia

poradie správ

## ⇒ Nie je centrálné riadený

## ⇒ Procesory majú jednoznačné identifikačné čísla

(anonymné siete – bez IČ)

## ⇒ Komunikácia procesorov výmenou správ

## ⇒ Čas lokálneho výpočtu zanedbateľný vzhľadom k času potrebnému na prenos

## ⇒ Znalosť o topológii

striktne lokálna (zoznam susedov)

lokálna + štruktúra siete (kruh, strom, ...)

## ⇒ Plne distribuovaný systém - každý má rovnaký algoritmus

**Prechodový systém**  $S = (C, \rightarrow, I)$ , kde

$C$  je množina konfigurácií

$\rightarrow$  je prechodová relácia,  $\gamma \rightarrow \delta$

$I, I \subseteq C$ , je množina počiatočných konfigurácií

**Vykonanie** (execution)  $S$  je maximálna postupnosť  $E = (\gamma_0, \gamma_1, \gamma_2, \dots)$ , kde

$\gamma_0 \in I, \forall i \gamma_i \rightarrow \gamma_{i+1}$

**Terminálna konfigurácia**  $\gamma$ - neexistuje  $\delta : \gamma \rightarrow \delta$

$\delta$  je **dosiahnuteľná z  $\gamma$**  ( $\gamma \rightsquigarrow \delta; \gamma \rightarrow^* \delta$ )

$$\exists \gamma = \gamma_0, \gamma_1, \gamma_2, \dots, \gamma_k = \delta, \gamma_i \rightarrow \gamma_{i+1} \forall 0 \leq i < k$$

$\delta$  je **dosiahnuteľná**, ak  $\gamma \rightsquigarrow \delta$  pre  $\gamma \in I$

- DS= množina procesov + komunikačný podsystem
- Prechodový systém-proces; konfigurácia-**stav**; prechod-**udalosť**
- Udalosť – interná, send, receive

### Lokálny algoritmus procesu

$(Z, I, \mapsto^i, \mapsto^s, \mapsto^r)$ , kde

$Z$ -množina stavov

$I$ -počiatočné stavy ( $I \subseteq Z$ )

$\mapsto^i \subseteq Z \times Z$  **interná udalosť**

$\mapsto^s \subseteq Z \times M \times Z$  **send**;  $M$  je množina možných správ

$\mapsto^r \subseteq Z \times M \times Z$  **receive**

$$c \mapsto d \Leftrightarrow (c, d) \in \mapsto^i \text{ alebo } \exists m \in M ((c, m, d) \in \mapsto^s \cup \mapsto^r)$$

**Distribovaný algoritmus** pre množinu procesov  $\mathbb{P} = \{ p_1, \dots, p_N \}$  je množina lokálnych algoritmov, jeden pre každý procesor z  $\mathbb{P}$

Prechodový systém indukovaný distribuovaným algoritmom pre  $\mathbb{P} = \{ p_1, \dots, p_N \}$  pod **asynchrónnou komunikáciou** je  $\mathbf{S} = (\mathbf{C}, \rightarrow, \mathbf{I})$ , kde

$$1. \mathbf{C} = \{(c_{p1}, \dots, c_{pN}, M) : (\forall p \in \mathbb{P} : c_p \in Z_p), M \in \mathcal{M}(M)\}$$

$$2. \rightarrow = (\cup_{p \in \mathbb{P}} \rightarrow_p), \text{ kde } \rightarrow_p \text{ je množina dvojíc}$$

$$(2a) (c_{pi}, c'_{pi}) \in \mapsto_{pi}^i, M1 = M2$$

$$(2b) \exists m \in M (c_{pi}, m, c'_{pi}) \in \mapsto_{pi}^s, M2 = M1 \cup m$$

$$(2c) \exists m \in M (c_{pi}, m, c'_{pi}) \in \mapsto_{pi}^r, M1 = M2 \cup m$$

$$3. \mathbf{I} = \{(c_{p1}, \dots, c_{pN}, M) : (\forall p \in \mathbb{P} c_p \in I_p) \wedge M = \emptyset\}$$



Prechodový systém indukovaný distribuovaným algoritmom pre  $\mathbb{P} = \{ p_1, \dots, p_N \}$  pod **synchronnou komunikáciou** je  $\mathbf{S} = (\mathbf{C}, \rightarrow, \mathbf{I})$ , kde

$$1. \mathbf{C} = \{(c_{p1}, \dots, c_{pN}) : (\forall p \in \mathbb{P} : c_p \in Z_p)\}$$

$$2. \rightarrow = (\cup_{p \in \mathbb{P}} \rightarrow_p) \cup (\cup_{p, q \in \mathbb{P} : p \neq q} \rightarrow_{pq}), \text{ kde}$$

$\rightarrow_p$  je množina dvojíc

$$(c_{p1}, \dots, c_{pi}, \dots, c_{pN}), (c_{p1}, \dots, c'_{pi}, \dots, c_{pN}) : (c_{pi}, c'_{pi}) \in \mapsto_{pi}^i$$

$\rightarrow_{pq}$  je množina dvojíc

$$(\dots, c_{pi}, \dots, c_{pj}, \dots), (\dots, c'_{pi}, \dots, c'_{pj}, \dots) : \\ \exists m \in M(c_{pi}, m, c'_{pi}) \in \mapsto_{pi}^s \ \& \ (c_{pj}, m, c'_{pj}) \in \mapsto_{pj}^r$$

$$3. \mathbf{I} = \{(c_{p1}, \dots, c_{pN}) : (\forall p \in \mathbb{P} c_p \in I_p)\}$$

## Overovanie vlastností prechodových systémov

Bezpečnosť a životnosť sú predikátmi na množine konfigurácií

**Požiadavka bezpečnosti**/fairness (tvrdenie  $P$  je pravdivé v **každej** konfigurácii každého vykonania algoritmu)

**Požiadavka životnosti**/liveness (tvrdenie  $P$  je pravdivé v **nejakej** konfigurácii každého vykonania algoritmu)

## Invarianty

$S = (C, \rightarrow, I)$ ,  $P, Q$  sú predikáty. Píšeme

$$\{P\} \rightarrow \{Q\} \Leftrightarrow (\forall \gamma \rightarrow \delta \text{ ak } P(\gamma) \text{ tak } Q(\delta))$$

**Definícia 1** *Tvrdenie  $P$  je **invariantom** systému  $S$ , ak*

$$- \forall \gamma \in I, P(\gamma)$$

$$- \{P\} \rightarrow \{P\}$$

**Fakt 1** *Ak  $P$  je invariantom  $S$ , tak  $P$  platí v každej konfigurácii každého vykonania  $S$*

**Fakt 2** *Nech  $Q$  je invariant  $S$ ,  $Q \Rightarrow P(\forall \gamma \in C)$ . Potom  $P$  platí v každej konfigurácii každého vykonania*

**Definícia 2** Nech  $S$  je prechodový systém,  $P, Q$  tvrdenia. Hovoríme, že  $P$  je ***Q-derivát***, ak

1.  $\forall \gamma \in I, Q(\gamma) \Rightarrow P(\gamma)$
2.  $\{Q \wedge P\} \rightarrow \{Q \Rightarrow P\}$

**Fakt 3** Ak  $Q$  je invariant a  $P$  je  $Q$ -derivát, tak  $Q \wedge P$  je invariant.

$Q$  je invariant,  $Q(\gamma)$  platí  $\forall \gamma \in I$  & (def)  $\forall \gamma \in I, Q(\gamma) \Rightarrow P(\gamma)$

↓

$\{P(\gamma)\}$ , preto  $\{Q(\gamma) \wedge P\gamma\}$

Nech  $\gamma \rightarrow \delta$ ,  $Q(\gamma) \wedge P(\gamma)$ ;  $Q$  je invariant (teda  $Q(\delta)$ )

Keďže  $\{Q \wedge P\} \rightarrow \{Q \Rightarrow P\}$ , tak  $Q(\delta) \Rightarrow P(\delta)$

Teda  $Q(\delta) \wedge P(\delta)$

**Životnosť** - tvrdenie je pravdivé v **nejakej** konfigurácii každého vykonania

Nech  $S$  je prechodový systém,  $P$  (cieľ) predikát.

**term**-predikát, ktorý je (ne)pravdivý v (ne)terminálnych konfiguráciách

**Definícia 3** *Systém  $S$  **končí korektne** ak predikát  $(term \Rightarrow P)$  v  $S$  vždy platí.*

**Uviaznutie** - terminálna konfigurácia, cieľ  $P$  sa nedosiahol

**Dobre založená množina** - množina s usporiadaním, pre kt. neexistuje nekonečná klesajúca postupnosť

Čiastočné usporiadanie  $(W, <)$  je dobre založené, ak neexistuje nekonečná klesajúca postupnosť  $w_1 > w_2 > \dots$

Nech  $S$  je prechodový systém,  $P$  tvrdenie. Funkcia  $f$  z  $C$  do dobre založenej  $W$  sa nazýva **norma** (vzhľadom na  $P$ ) ak  $\forall \gamma \rightarrow \delta : (f(\gamma) > f(\delta) \text{ alebo } \mathbf{P}(\delta))$

**Fakt 4** *Nech  $S = (C, \rightarrow, I)$ ,  $P$  (cieľ) predikát. Ak  $S$  skončí korektne a norma  $f$  existuje tak  $P$  platí v niektorej konfigurácii každého vykonania.*

Dôkaz: Nech  $E = (\gamma_0, \gamma_1, \gamma_2, \dots)$  je vykonanie.

1.  $E$  je **konečná**,  $E = (\gamma_0, \gamma_1, \gamma_2, \dots, \gamma_k)$ , potom  $P(\gamma_k)$  platí

2.  $E$  je **nekonečná**.

Uvažujme najdlhší prefix  $E' = (\gamma_0, \gamma_1, \gamma_2, \dots)$  taký, že  $P$  neplatí.

Nech  $F = (f(\gamma_0), f(\gamma_1), f(\gamma_2), \dots)$ .

**f je norma**  $\Rightarrow f(\gamma_0) > f(\gamma_1) > f(\gamma_2) > \dots$

Je to klesajúca postupnosť, preto je konečná;

$F = (f(\gamma_0), f(\gamma_1), f(\gamma_2), \dots, f(\gamma_k))$  a  $P(\gamma_{k+1})$

## Poradie(causality order)

Nech  $E = (\gamma_0, \gamma_1, \gamma_2, \dots)$  je vykonanie. Relácia  $\prec$ , ktorej hovoríme **poradie udalostí** v  $E$ , je najmenšia relácia, ktorá spĺňa

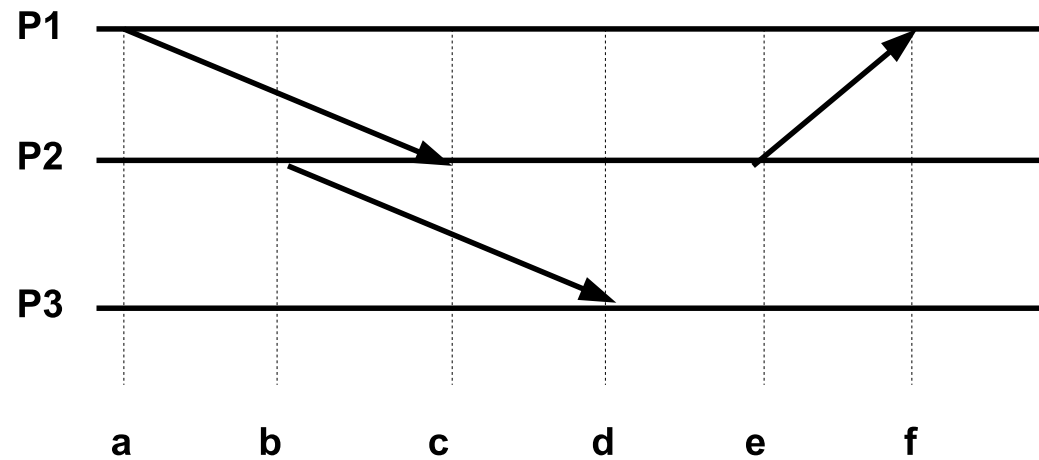
- Ak  $e, f$  sú rôzne udalosti toho istého procesora,  $e$  sa vykonalo pred  $f$ , tak  $e \prec f$
- Ak  $s$  je send,  $r$  odpovedajúca receive udalosť, tak  $s \prec r$
- $\prec$  je tranzitívna

Čiastočné usporiadanie  **$\mathbf{a} \leq \mathbf{b}$**  :  $a \prec b \vee a = b$

Súbežne bežiacie procesy ( **$\mathbf{a} || \mathbf{b}$** ) ani  $a \leq b$ , ani  $b \leq a$

## Time-space diagram

Ku každému vykonaniu  $E = (\gamma_0, \gamma_1, \gamma_2, \dots)$  existuje odpovedajúca postupnosť udalostí  $F = (e_0, e_1, e_2, \dots)$  takých, že  $e_i$  spôsobilo prechod  $\gamma_i \rightarrow \gamma_{i+1}$



$$F = (a, b, c, d, e, f)$$

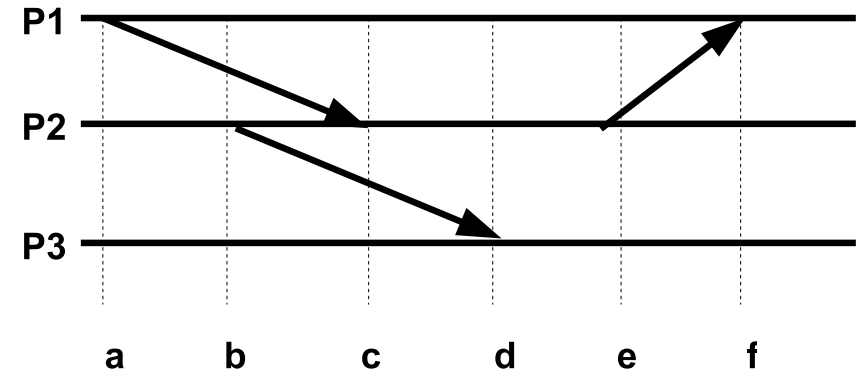
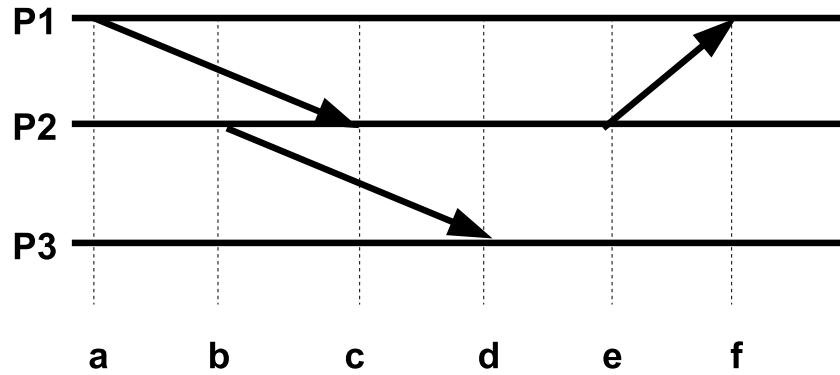


## Závislé a nezávislé udalosti

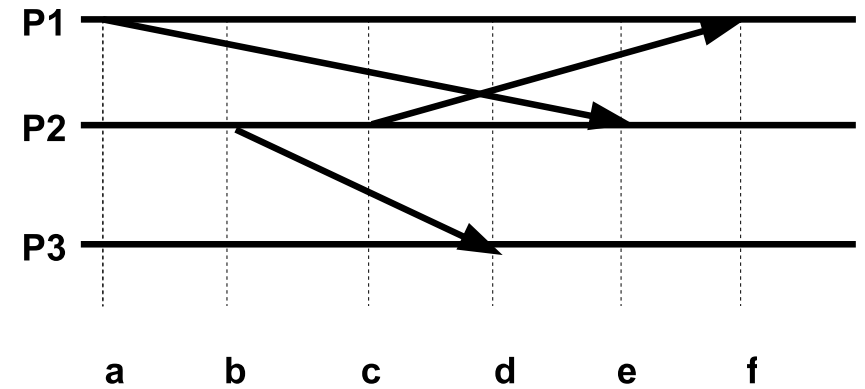
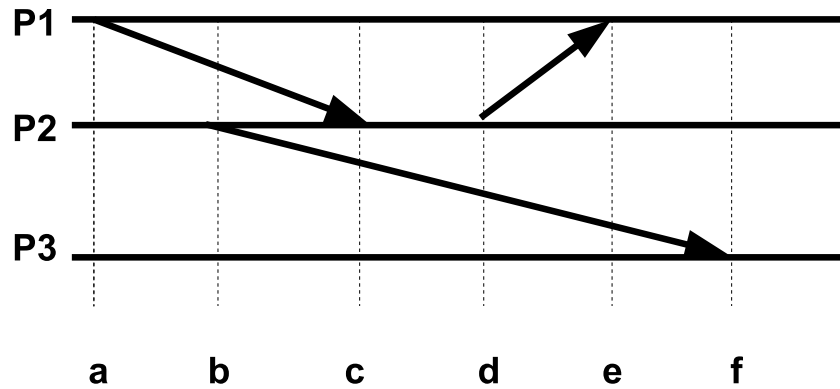
**Tvrdenie 1** *Nech  $\gamma$  je konfigurácia,  $e_p$  a  $e_q$  sú udalosti rôznych procesov  $p, q$ , aplikovateľné v  $\gamma$ . Potom  $e_p$  je aplikovateľné v  $e_q(\gamma)$ ,  $e_q$  je aplikovateľné v  $e_p(\gamma)$  a  $e_p(e_q(\gamma)) = e_q(e_p(\gamma))$ .*

- Správy sú jednoznačne spojené s konkrétnym procesorom
- $e_p, e_q$  nemôžu byť odpovedajúce send/receive udalosti

# Rôzne vykonania



$$a < c, b < d, e < f, c < e$$



## Ekvivalencia vykonaní

$E = (\gamma_0, \gamma_1, \gamma_2, \dots)$  je vykonanie,

$\mathcal{E} = (e_0, e_1, e_2, \dots)$  odpovedajúca postupnosť udalostí,

$\mathcal{F} = (f_0, f_1, f_2, \dots)$  konzistentná permutácia  $\mathbf{f}_i \leq \mathbf{f}_j \Rightarrow \mathbf{i} \leq \mathbf{j}$

**F vykonanie s postupnosťou udalostí  $\mathcal{F}$ .**

**Tvrdenie 2**  $\mathcal{F}$  definuje jednoznačné vykonanie  $F$ , ktoré začína v počiatočnej konfigurácii  $E$ , má rovnaký počet udalostí ako  $E$ , posledná konfigurácia  $F$  je poslednou konfiguráciou  $E$ .

Ekvivalentné vykonania definujú triedu ekvivalencie  $\approx$

**Výpočet distribuovaného algoritmu** je trieda ekvivalencie vykonaní toho algoritmu

**Logické hodiny** Môžeme definovať hodiny, ktoré budú vyjadrovať poradie/causality

Nech  $\Theta$  je funkcia z množiny udalostí do **usporiadanej** množiny.  $\Theta$  sú **logické hodiny** ak  **$a \prec b$  implikuje  $\Theta(a) < \Theta(b)$**

- **Poradie**- ak je vykonanie definované postupnosťou udalostí  $\mathcal{E} = (e_0, e_1, e_2, \dots)$ , definujeme  $\Theta(e_i) = i$
- **Real-time hodiny** – procesor má hardwareové hodiny (vlastne to kazí definíciu DS)
- **Vektorové hodiny** –  $a \prec b \Leftrightarrow \Theta(a) < \Theta(b)$  (súbežne bežiacie procesy majú neporovnateľné hodnoty)

$$\Theta_v(a) = (a_1, a_2, \dots, a_n), \text{ kde}$$

$a_i$  je maximálne číslo udalosti  $e$  v procesore  $p_i$ , pre ktorú  $e \leq a$

## Lamportove hodiny

$\Theta(a)$  definujeme ako dĺžku najdlhšej postupnosti  $(e_0, e_1, e_2, \dots, e_k)$  takej, že  
 $e_0 \prec e_1 \prec e_2 \prec \dots \prec e_k = a$

- Ak  $b$  je **interná alebo send**,  $a$  predchádzajúca udalosť toho istého procesora, tak  $\Theta_L(\mathbf{b}) = \Theta_L(\mathbf{a}) + 1$
- Ak  $r$  je **receive** udalosť,  $i$  predchádzajúca udalosť toho istého procesora,  $s$  odpovedajúca send udalosť, tak  $\Theta_L(\mathbf{r}) = \max\{\Theta_L(\mathbf{i}), \Theta_L(\mathbf{s})\} + 1$
- Ak  $a$  je **prvá udalosť**, tak  $\Theta_L(a) = 0$

Realizácia

$\Theta_p$  číslo poslednej udalosti v procesore  $p$

$\Theta_m$  číslo udalosti-posielania správy  $m$

## Topológia siete

Sieť = množina procesorov + komunikačný subsystém (graf)

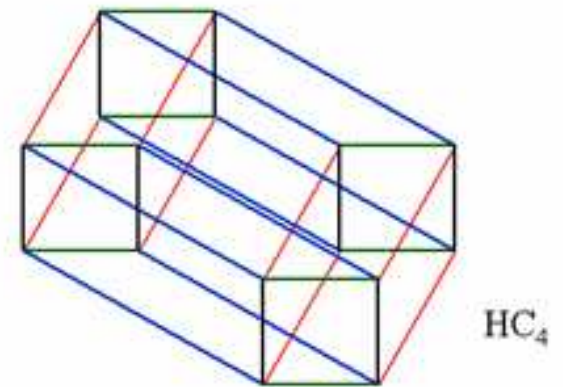
- kruh (tokren ring),
- strom (rýchle výpočty; kostra - prenositeľné)
- hviezda (master-slave; zraniteľnosť centra)
- úplný graf
- hyperkocka

$$HC_n = (V, E), N = 2^n, V = \{0, 1\}^n$$

$$((a_0, a_1, \dots, a_n), (b_0, b_1, \dots, b_n)) \in V$$



$$\exists! i \ a_i \neq b_i$$



## Vlastnosti kanálov

- **spoľahivosť** (správy sa nestrácajú/neduplikujú/nevznikajú-predpokladáme)
- **FIFO** (zachováva poradie, v akom doň vstupujú-nepredpokladáme)
- **kapacita** (def. predpokladá neobmedzenú kapacitu)

## Vedomosti procesu/processora

### Topologická informácia

- Čísla procesorov, polomer, topológia
- Zmysel pre orientáciu(konzistentné označenie hrán so smermi)

### Identita procesu (pomenovaná sieť, anonymná sieť)

### Susedné identity (nepredpokladáme; priame vs. nepriame adresovanie)

## Zložitosť DA

**Počet správ** (total)

**Bitová**(väčšinou správy dĺžky logn)

## Časová

idealizované meranie času

- Čas spracovania je nulový
- Prenosový čas najviac jednotka času

Časová zložitosť je čas spotrebovaný vzhľadom na tieto predpoklady

**Priestorová**(vzhľ. na jeden proces)



# Predpoklady

Ak nepovieme inak

- Sieť je silne súvislá
- Komunikácia je asynchrónna
- Kanály sú FIFO
- Kanály nestrácajú správy
- Správy sú doručené v konečnom čase

**Routing** – smerovanie packetov do cieľa

**routovacia tabuľka**

$RT_u(v) = w \Leftrightarrow$  packet smerujúci do vrchola  $v$  je z vrchola  $u$  posunutý do  $w$

**routovanie**

1. výpočet routovacích tabuliek
2. samotné posúvania packetov podľa nich

**Kritériá dobrého routovania**

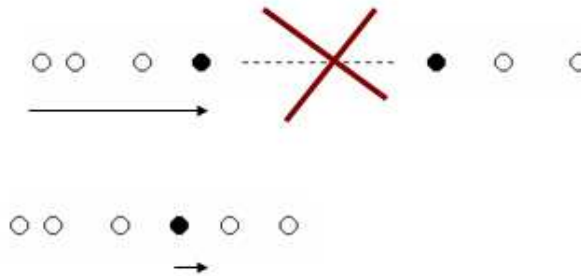
- **Korektnosť** (každý packet je v konečnom čase doručený)
- **Zložitosť** počítania tabuliek (čas, pamäť, správy)
- **Efektívnosť** (používanie optimálnych ciest)
- **Robustnosť** (prepočítavanie tabuliek pri zmene v topológii)
- **Adaptabilita** tabuliek, ktorá eliminuje/znižuje preťaženosť kanálov.
- Predpokladáme, že vrcholy sú **obsluhované rovnomerne**

## Smerovanie packetov podľa cieľa

optimálny algoritmus existuje, ak

- cena posielania je nezávislá od momentálneho používania
- cena zreťazenia ciest je súčet ciest jednotlivých ciest
- graf neobsahuje cyklus zápornej dĺžky

**Fakt 5** *Ak existuje cesta z  $u$  do  $v$ , tak existuje jednoduchá optimálna cesta.*



## sink tree

**Fakt 6** *Nech  $G = (V, E)$  je súvislý graf. Pre každý vrchol/cieľ  $d \in V$  existuje strom  $T_d = (V, E_d)$ ,  $E_d \subseteq E$  tak, že pre každý vrchol  $v \in V$  je  $v - d$  cesta v strome  $T_d$  najkratšou  $v - d$  cestou v grafe  $G$ .*

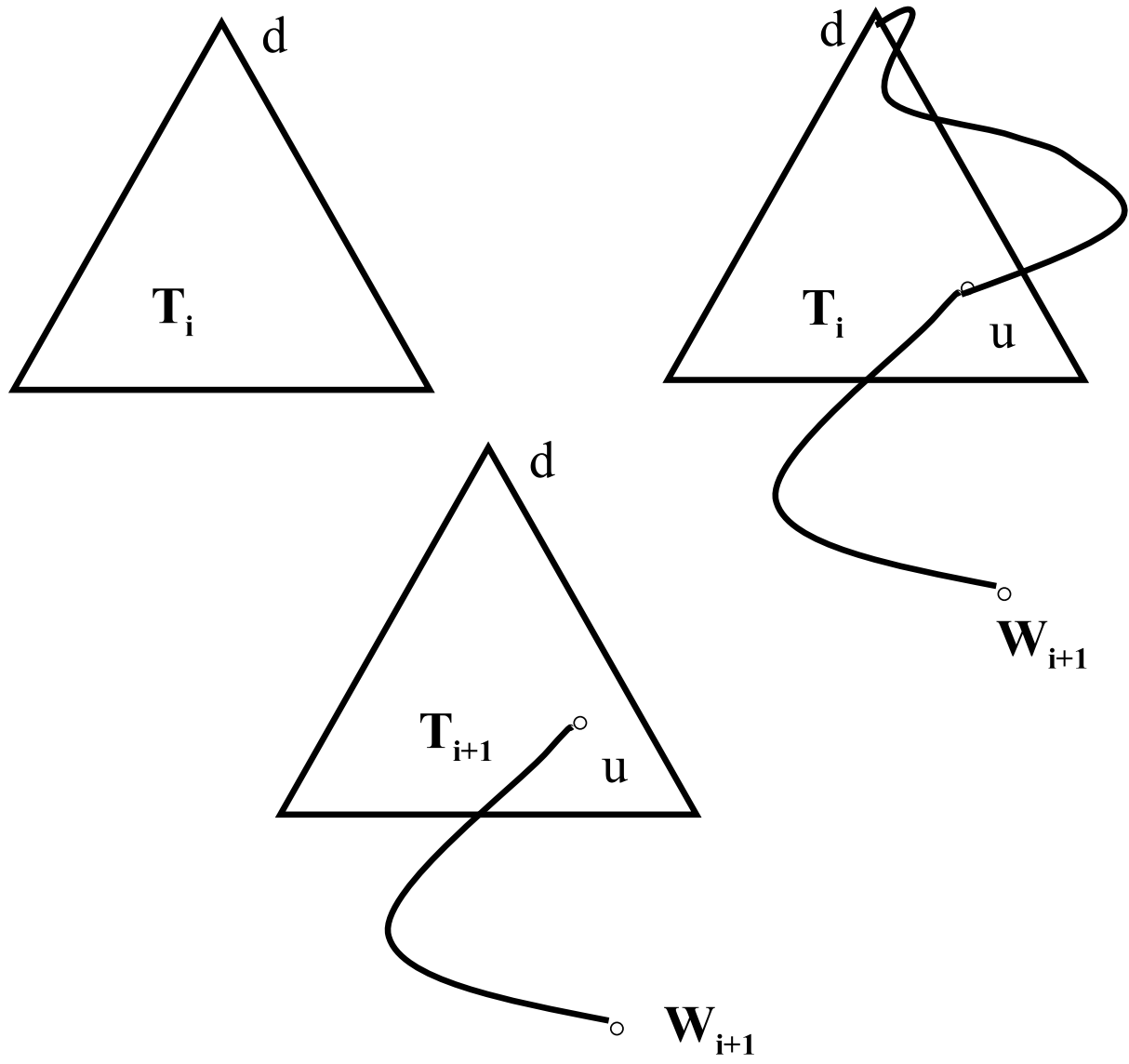
**Konštrukcia** – iteratívne  $T_0, T_1, \dots, T_m$

- $T_0 = (\{d\}, \emptyset)$
- $T_i = (V_i, E_i)$ , kde  $T_i$  je podgraf grafu  $G$ ,  $T_i$  je podgraf grafu  $T_{i+1}$
- $\forall w \in V_i$  je  $w - d$  cesta v  $T_i$  optimálnou  $w - d$  cestou v  $G$ .

Nech  $w_{i+1} \notin V_i$  a nech  $\mathbf{w}_{i+1} - \mathbf{d}$  je optimálna cesta v  $G$ . Nech  $u$  je prvý vrchol na tejto ceste, ktorý *patrí* do  $V_i$ .

$T_{i+1}$  vznikne z  $T_i$  **pripojením** cesty  $\mathbf{w}_{i+1} - \mathbf{u}$  k vrcholu  $\mathbf{u}$ .

- proces skončí vytvorením stromu  $T_m$  so všetkými vrcholmi. Vtedy  $T_d := T_m$



## Destination based routing

- Vypočítaj sink-tree  $T_d \forall d \in V$
- $T_u[d]$  je otec vrchola  $u$  v sink-tree s koreňom  $d$

**Algoritmus**(packet prijatý/vygenerovaný vrcholom  $u$  smeruje do  $d$ )

**if**  $d = u$  **then** *deliver* lokálne  
**else** send packet do  $T_u[d]$

## všetky najkratšie cesty

$G = (V, E)$  je ohodnotený orientovaný graf bez záporných cyklov; váhy  $\omega(\mathbf{u}, \mathbf{v})$  môžu byť aj záporné

Vychádzame zo **sekvenčného Floyd-Warshall alg.**

$d^S(\mathbf{u}, \mathbf{v})$  dĺžka najkratšej  $u - v$  cesty s *vnútornými* vrcholmi z množiny  $S$

**inicializácia:**

$$d^S(u, u) \leftarrow \emptyset$$

$$d^\emptyset(u, v) \leftarrow \omega(u, v) \text{ pre } (u, v) \in E, u \neq v$$

$$d^\emptyset(u, v) \leftarrow \infty \text{ pre } (u, v) \notin E, u \neq v$$

**iterácia:**

$$d^{S \cup \{w\}}(u, v) = \min\{d^S(u, v), d^S(u, w) + d^S(w, v)\}$$

## Floyd-Warshall

$S \leftarrow \emptyset$

for all  $u, v \in V$  do

  if  $u=v$  then  $D_U[v] \leftarrow \emptyset$

  else if  $(u, v) \in E$  then  $D_u[v] \leftarrow \omega(u, v)$

  else  $D_u[v] \leftarrow \infty$

▷ inicializácia

while  $S \neq V$  do

  pick  $w \in V \setminus S$ ;

▷ w-pivot

  for all  $u, v \in V$  do

$D_u[v] \leftarrow \min\{D_u[v], D_u[w] + D_w[v]\}$

▷ nie je kratšie ísť cez  $w$  ?

$S \leftarrow S \cup \{w\}$

▷ zväčšenie množiny  $S$

**sekvenčná zložitosť**

$$O(n^2) + n \cdot O(n^2) = \mathbf{O}(n^3)$$



## Paralelná verzia - **TOUEG**ov algoritmus

### Predpoklady

- cykly sú **pozitívne**
- na začiatku výpočtu všetky **vrcholy poznajú  $V$**
- vrcholy **poznajú** svojich susedov ( **$\text{Neigh}_u$** )
- výber pivota **podľa dopredu zvolenej** stratégie

**$S_u$**  množina (vnútorných) vrcholov

**$D_u$**  pole váh/vzdialeností

**$Nb_u$**  pole vrcholov(susedov na najkratších cestách)

## Jednoduchý Toueg - pre vrchol $u$ ;

$S_u := \emptyset$ ;

for all  $v \in V$  do

*inicializácia*;

while  $S \neq V$  do

*voľba pivota  $w$* ;

*distribúcia  $D_w[v]$  po sieti*;

for all  $v \in V$  do

*aktualizácia  $D_u[v], Nb_u[v]$* ;

$S_u := S_u \cup \{w\}$

## algoritmus jednoduchýTOUEG

### inicializácia

```
 $S_u \leftarrow \emptyset;$   
for all  $v \in V$  do  
  if  $u = v$  then  
     $D_u[v] \leftarrow 0;$   
     $Nb_u[v] \leftarrow \perp$   
  else if  $(u, v) \in E$  then  
     $D_u[v] \leftarrow \omega(u, v);$   
     $Nb_u[v] \leftarrow v$   
  else  
     $D_u[v] \leftarrow \infty;$   
     $Nb_u[v] \leftarrow \perp$ 
```

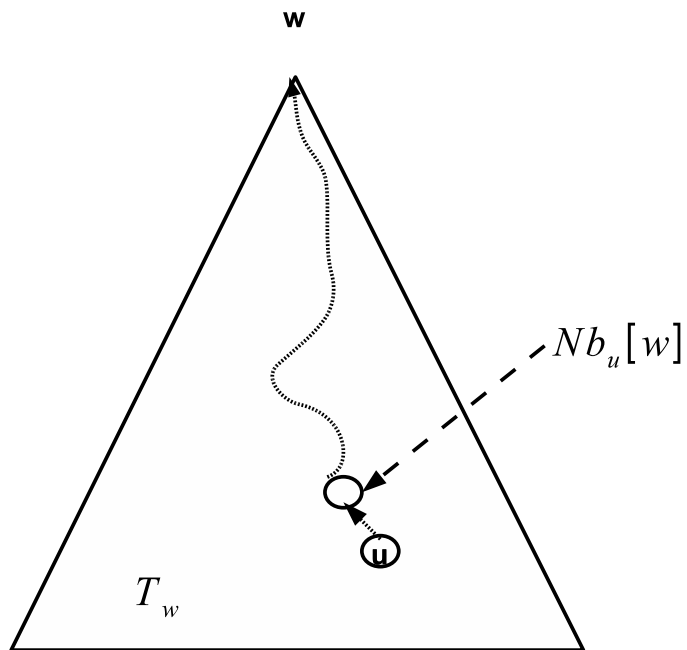
### iterácia

```
while  $S \neq V$  do  
  pick  $w \in V \setminus S$   
  
  if  $u = w$  then broadcast  $D_w$   
  else receive  $D_w$   
  
  for  $v \in V$  do  
    if  $D_u[w] + D_w[v] < D_u[v]$  then  
       $D_u[v] \leftarrow D_u[w] + D_w[v]$   
       $Nb_u[v] \leftarrow Nb_u[w]$   
  
   $S_u \leftarrow S_u \cup \{w\}$ 
```

**Fakt 7** Algoritmus jednoduchýToueg skončí po  $N$  opakovaníach hlavného cyklu. Po jeho skončení je v  $D_u[v]$  dĺžka najkratšej  $u$ - $v$  cesty,  $Nb_u[v]$  obsahuje identifikáciu prvého vrcholu na nej (ak cesta existuje) alebo má nedefinovanú hodnotu.

## Vylepšenie:

- ak  $D_u[w] = \infty$ ,  $u$  nemaní svoje tabuľky vo fáze s pivotom  $w$   
 $\Rightarrow$  broadcast  $D_w[v]$  iba po strome  $T_w$
- ak  $D_u[w] < \infty$ ,  $u$  pozná otca, ale nie synov v  $T_w$ , preto  
 $v$  oznámi svojmu susedovi  $u$ , či je jeho synom v  $T_w$  alebo nie



$u$  oznámi  $Nb_u[w]$ , že je jeho syn vo fáze s pivotom  $w$

## algoritmus vylepšený Toueg (náčrt pre vrchol $u$ )

**Inicializácia**  $D_u, Nb_u$

**Začni fázu (pre pivota  $w$ )**

**for**  $x \in Neigh_u$  **do**

**if**  $Nb_u[w] = x$  **then** send  $\langle ys, w \rangle$  do  $x$

**else** send  $\langle nys, w \rangle$  to  $x$

$numRec_u \leftarrow 0$

**while**  $numRec_u < |Neigh_u|$  **do**

    receive  $\langle ys, w \rangle$  alebo  $\langle nys, w \rangle$

$numRec_u \leftarrow numRec_u + 1$

**if**  $D_u[w] < \infty$  **then**

**if**  $u \neq w$  **then** receive  $\langle dtab, w, d \rangle$  od  $Nb_u[w]$

**for**  $x \in Neigh_u$  **do**

**if**  $\langle ys, w \rangle$  bol prijatý od  $x$  **then**

            send  $\langle dtab, w, d \rangle$  do  $x$  ;

$S_u \leftarrow S_u \cup \{w\}$

▷ Vytvor v  $T_w$  spojenie otec-syn

▷  $u$  musí prijať  $|Neigh_u|$  správ

▷ Zúčastni sa fázy (s pivotom  $w$ )

▷ lokálny update

$\langle ys, w \rangle \setminus \langle nys, w \rangle$  – **som** \ **nie som** tvoj syn vo fáze s pivotom  $w$

$\langle dtab, w, D \rangle$  – broadcast riadka  $D_w$  vo fáze s pivotom  $w$

**Veta 1** *Algoritmus vylepšený Toueg skončí. Po jeho skončení je v  $D_u[v]$  dĺžka najkratšej  $u - v$  cesty,  $Nb_u[v]$  obsahuje identifikáciu prvého vrchola na nej (ak cesta existuje) alebo má nedefinovanú hodnotu.*

*Nech  $\mathbf{W}$  je počet bitov postačujúci na kódovanie mien aj váh. Potom počas vykonania algoritmu sa prekomunikuje*

- $\mathbf{O}(\mathbf{N})$  správ - kanál
- $\mathbf{O}(\mathbf{N}|E|)$  správ - celkovo
- $\mathbf{O}(\mathbf{N}^2\mathbf{W})$  bitov - kanál
- $\mathbf{O}(\mathbf{N}^3\mathbf{W})$  bitov - celkovo

$$N^2 \cdot NW + 2N|E| = O(N^3W)$$

*Pamäť potrebná v jednotlivých procesoroch je  $\mathbf{O}(\mathbf{NW})$  bitov.*

## Od sekvenčných algoritmov k distribuovaným

- premenné sekvenčného algoritmu sa distribuujú medzi vrcholy; výpočty sú lokálne
- ak treba vzdialenú premennú, vyžaduje sa komunikácia
- znižovanie komunikácie využitím vlastností sekvenčného algoritmu

Ešte stále nám ostali **zlé vlastnosti Touegovho** algoritmu

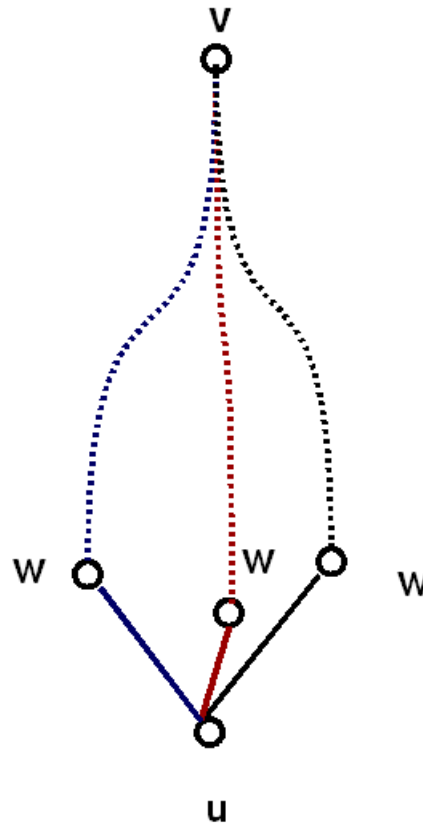
- kvôli pivotovi vyžadujeme znalosť všetkých vrcholov
- vyžadujeme informáciu, ktorá nie je dostupná ani vo vrchole ani v jeho susedoch  
?  $d(u, w) + d(w, v) < d(u, v)$ ?



Chandy-Misra

$$d(u, v) = \begin{cases} 0, & \text{ak } u=v; \\ \min_{w \in \text{Neigh}_u} \{ \omega_{uw} + d(w, v) \}, & \text{inak.} \end{cases}$$

- Min hop
- Lokálna komunikácia
- Výpočty nezávislé
- Vzhľadom na lokálnosť jednoduchšia adaptabilita





## sekvenčný Dijkstra – z každého do $s$

**inicializácia:**

$T \leftarrow \{s\}; S := V \setminus s; d_s \leftarrow 0; Nb_s \leftarrow s$

$d_u \leftarrow \omega(u, s); Nb_u \leftarrow s$  pre  $(u, s) \in E$

$d_u \leftarrow \infty; Nb_u \leftarrow \perp$  pre  $(u, s) \notin E$

**výpočet**

**while**  $S \neq \emptyset$  **do**

    nech  $v$  je taký, že  $d_v = \min\{d_w, w \in S\}$ ;

$T \leftarrow T \cup \{v\}; S \leftarrow S - \{v\}$ ;

**for**  $w \in S$  **do**

**if**  $d_w < \omega(w, v) + d_v$  **then**

$d_w \leftarrow \omega(w, v) + d_v$

$Nb_w \leftarrow v$

budujeme vlastne sink-tree s koreňom  $s$ :

- vo vrchole  $u$  bude informácia  $d_u$  a  $Nb_u$
- $S, T$  sú v grafe uložené implicitne; vrchol, ktorý "dostáva trvalú značku" oznámi všetkým ostatným, aká je jeho vzdialenosť ku koreňu
- každé vylepšenie oznamuje vrchol všetkým ostatným

**Inicializácia** -  $D_u[w] := \infty; Nb_u := \perp$

**Výpočet pre vrchol  $v$**  (koreň  $v$  kostry sa "zobudil")

$D_v[v] := 0;$

$\forall w \in Neigh_v$  send  $\langle mydist, v, 0 \rangle$  do  $w$

**while** nie je koniec **do**

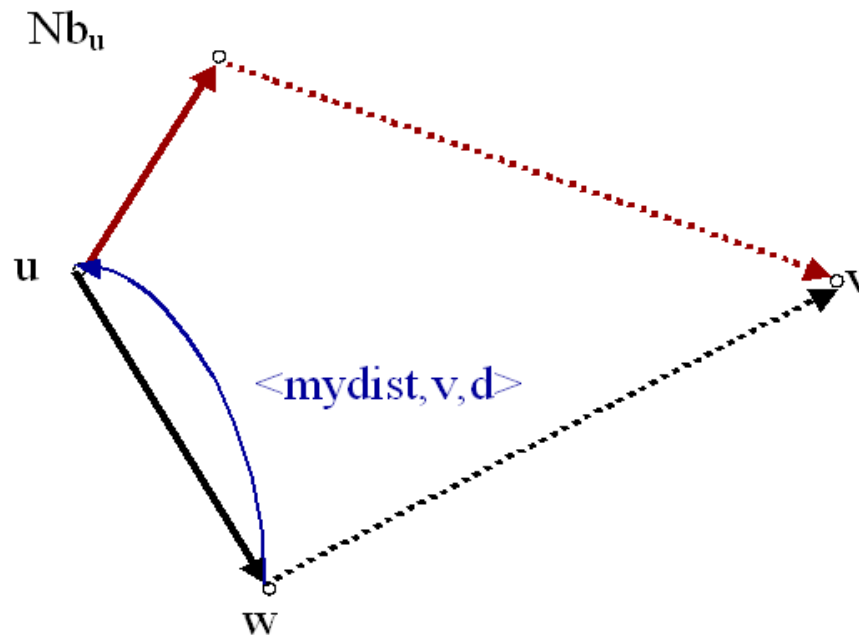
**for**  $w \in V$  **do**

        prijmi správu;

        spracuj prijatú správu

## Spracovanie $\langle \text{mydist}, v, d \rangle$

```
receive  $\langle \text{mydist}, v, d \rangle$  od  $w$ ;  
if  $d + \omega(uv) < D_u[v]$  then begin  
   $D_u[v] := d + \omega(uv)$ ;  $Nb_u[v] := w$ ;  
  forall  $x \in \text{Neigh}_u$  do send  $\langle \text{mydist}, v, D_u[v] \rangle$  do  $x$   
end;
```

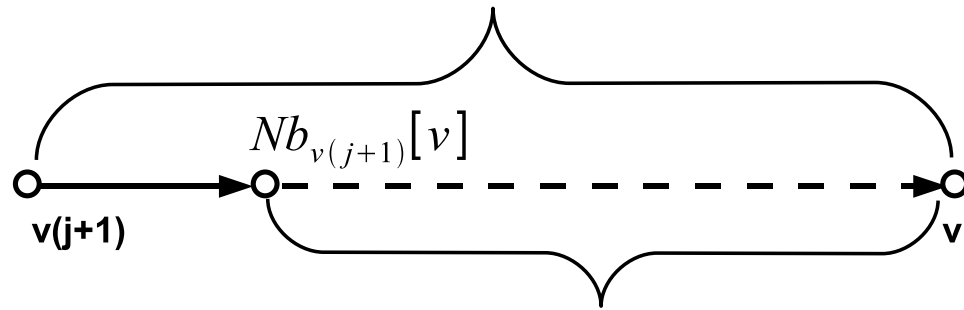


**Fakt 8** V každej realizácii algoritmu Chandy-Misra (CHM) sa dosiahne konfigurácia, v ktorej  $\forall u \in V D_u[v] \leq d(u, v)$ .

Kvôli dôkazu tohto tvrdenia si očísľujeme vrcholy v  $T_v$  tak, že otec má menšie číslo ako syn. Potom tvrdenie dokážeme matematickou indukciou.

**IP**  $\forall j \leq N - 1$  dosiahneme konfiguráciu, v ktorej

$$D_{v(j+1)}[v] = \omega(v(j+1), v_i) + d(v_i, v) = d(v(j+1), v)$$



$$Nb_{v(j+1)}[v] = v_i, 1 \leq j \wedge d_{v_i}[v] = d(v_i, v)$$

## zložitosť algoritmu Chandy-Misra

Každý vrchol posiela správy svojim susedom vždy, keď sa dozvie o vylepšení. Keďže každý vrchol si môže vylepšiť info max  $\mathbf{N} - \mathbf{1}$  krát, dostávame pre **počet správ** vymenených pri výpočte **najkratších ciest do v**

$$O(N \sum_{u \in V} |N_{eig_u}|) = \mathbf{O}(\mathbf{N}|\mathbf{E}|)$$

Ak počítame najkratšie cesty ku každému vrcholu (teda  $\mathbf{N}$  krát opakujeme vysvetlený postup), dostávame

- $O(N^2|E|)$  správ /  $O(N^2|E|W)$  bitov pre výpočet **všetkých najkratších ciest**
- $O(N^2)$  správ /  $O(N^2W)$  bitov **na kanál**

## Netchange - najkratšie cesty v prítomnosti chýb kanálov

### predpoklady

- Vrcholy poznajú veľkosť siete
- Kanály sú FIFO
- Vrcholy sa dozvedia o poruche/oprave susedných kanálov
- Cena cesty je počet kanálov na nej

### Od algoritmu vyžadujeme

1. Ak po konečnom počte zmien ostane topológia siete nezmenená, algoritmus v konečnom čase skončí.
2. V okamihu, keď algoritmus skončí, pre vrchol  $u$  platí:
  - $Nb_u[v] = \text{local}$ , ak  $u = v$
  - $Nb_u[v] = w$ , kde  $w$  je prvý vrchol na najkratšej  $u - v$  ceste
  - $Nb_u[v] = \perp$  ak  $u - v$  cesta neexistuje

## premenné

$\mathbf{Neigh}_u$  - aktuálna množina susedov vrchola  $u$

$\mathbf{D}_u[\mathbf{v}]$  - odhad  $d(u, v)$

$\mathbf{Nb}_u[\mathbf{v}]$  - prvý vrchol na ceste s  $D_u[v]$

$\mathbf{ndis}_u[\mathbf{w}, \mathbf{v}]$  - odhad  $d(w, v)$

## Inicializácia v $u$ :

for all  $w \in \mathit{Neigh}_u, v \in V$  do  $\mathit{ndis}_u[w, v] := N$

for all  $v \in V$  do  $D_u[v] \leftarrow N; \mathit{Nb}_u[v] \leftarrow \perp$

$D_u[u] \leftarrow 0; \mathit{Nb}_u[u] \leftarrow \text{local};$

for all  $w \in \mathit{Neigh}_u$  do send  $\langle \mathit{mydist}, u, 0 \rangle$  do  $w$ ;

- $\langle \mathbf{mydist}, \mathbf{v}, \mathbf{d} \rangle$  – vrchol, od ktorého prišla, oznamuje svoju vzdialenosť  $d$  od vrchola  $v$
- $\langle \mathbf{fail}, \mathbf{v} \rangle$  – správa prijatá kanálom od vrchola  $w$  oznamuje poruchu kanálu  $(v, w)$
- $\langle \mathbf{repair}, \mathbf{v} \rangle$  – kanál  $(v, w)$  bol obnovený

$\langle \mathbf{repair}, \mathbf{v} \rangle$

```

receive  $\langle \mathit{repair}, v \rangle$ 
 $Neigh_u \leftarrow Neigh_u \cup \{w\}$ 
for all  $v \in V$  do
   $ndist_u[w, v] \leftarrow N$ 
  send  $\langle \mathit{mydist}, v, D_u[v] \rangle$  do  $w$ 

```

$\langle \mathbf{fail}, \mathbf{v} \rangle$

```

receive  $\langle \mathit{fail}, w \rangle$ 
 $Neigh_u \leftarrow Neigh_u - \{w\}$ 
for all  $v \in V$  do
  Recompute( $\mathbf{v}$ )

```

$\langle \mathbf{mydist}, \mathbf{v}, \mathbf{d} \rangle$

```

receive  $\langle \mathit{mydist}, v, d \rangle$ 
 $ndist_u[w, v] \leftarrow d$ 
Recompute( $\mathbf{v}$ )

```



## Recompute( $v$ )

**if**  $v = u$  **then**

$D_u[v] \leftarrow 0$

$Nb_u[v] \leftarrow local$

**else**

$D \leftarrow 1 + \min\{ndist_u[w, v] \mid w \in Neigh_u\}$

▷ nech  $w_D$  je ten vrchol, pre ktorý sa dosahuje minimum

**if**  $D < N$  **then**

$D_u[v] \leftarrow D$

$Nb_u[v] \leftarrow w_D$

**else**

$D_u[v] \leftarrow N$

$Nb_u[v] \leftarrow \perp$

**if**  $D_u[v]$  sa zmenilo **then**

**for all**  $x \in Neigh_u$  **do** send  $\langle mydist, D \rangle$  do  $x$

## Netchange – korektnosť

hovoríme, že  $up(u, w)$  platí  $\Leftrightarrow$  obojsmerný kanál  $(u, w)$  pracuje korektne

$Q_{uw}, Q_{wu}$  – rady, ktoré modelujú kanál  $(u, w)$

Korektnosť sa ukáže pomocou invariantov.

$\mathbf{P}(\mathbf{u}, \mathbf{v}, \mathbf{w}) \equiv$

1.  $up(u, w) \Leftrightarrow w \in Neigh_u$
2.  $up(u, w) \wedge Q_{uw}$  obsahuje  $\langle mydist, v, d \rangle$  potom v poslednej  $d = D_w[v]$
3.  $up(u, w) \wedge Q_{uw}$  neobsahuje  $\langle mydist, v, d \rangle$  potom  $ndist_u[w, v] = D_w[v]$

$\mathbf{L}(\mathbf{u}, \mathbf{v}) \equiv$

4.  $u = v \Rightarrow (D_u[v] = 0 \wedge Nb_u[v] = local)$
5.  $(u \neq v \wedge \exists w \in Neigh_u ndist_u[w, v] < N - 1) \Rightarrow$   
 $(D_u[v] = 1 + \min\{ndist_u[w, v]\} = 1 + ndist_u[Nb_u[v], v])$
6.  $(u \neq v \wedge \forall w \in Neigh_u ndist_u[w, v] \geq N - 1) \Rightarrow (D_u[v] = N \wedge [Nb_u[v] = \perp])$

## Netchange – stabilizácia

predikát **stable**  $\equiv \forall u, w (up(u, w) \Rightarrow Q_{uw})$  neobsahuje  $\langle mydist, v, d \rangle$

Hovoríme, že **konfigurácia**  $\gamma$  je **stabilná**, ak  $stable(\gamma)$ .

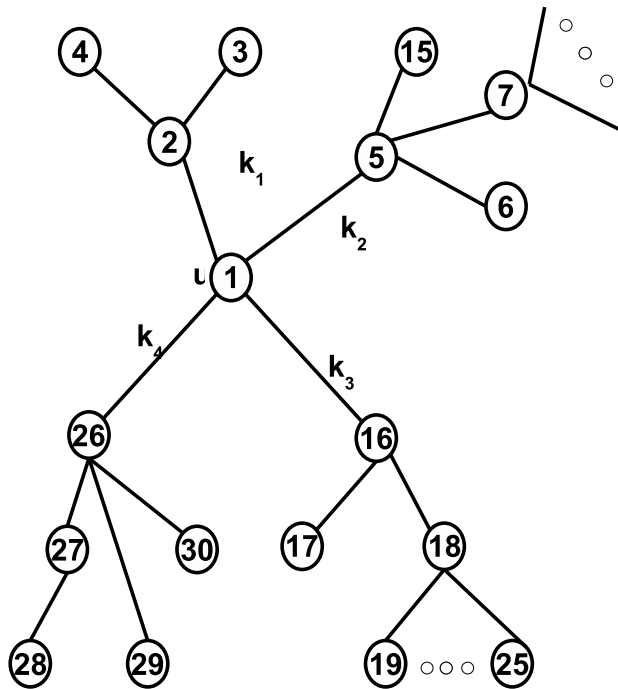
Nech  $t_i(\gamma) = \# \langle mydist, \dots, i \rangle v \gamma + \#(u, v)$ , pre ktoré v  $\gamma$   $D_u[v] = \dots$

**norma**  $f(\gamma) = (t_0(\gamma), t_1(\gamma), \dots, t_N(\gamma))$

**Fakt 9** *Spracovanie  $\langle mydist, v, i \rangle$  znižuje hodnotu  $f$ .*

**Fakt 10** *Ak po konečnom počte topologických zmien ostane topológia siete nemenná, algoritmus dosiahne stabilnú konfiguráciu.*

## Routing s kompaktnými tabuľkami



- také adresovanie, ktoré minimalizuje veľkosť routovacej tabuľky
- Tree-labeling - kanálom routujeme interval adres
- topológia je strom

vo vrchole 1

adresa	2	3	...	20	...	30
kanál	$k_1$	$k_1$		$k_3$		$k_4$

alebo kompaktnejšie

kanál	$k_1$	$k_2$	$k_3$	$k_4$
adresy	2,4 (2...4)	5, 15 (5...15)	16,25 (16...25)	26,30(26...30)

## Tree labeling — Santoro/Khatib

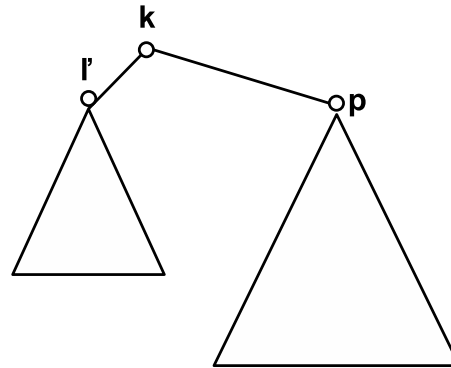
vrcholy  $0, 1, \dots, N - 1$

$$\text{cyklický interval } [a, b) = \begin{cases} \{a, a + 1, \dots, b - 1\} & \text{ak } a < b; \\ \{a, a + 1, \dots, N, 0, 1, \dots, b - 1\} & \text{ak } b < a. \end{cases}$$

Ak  $[a_1, b_1), \dots, [a_k, b_k)$  sú cyklické intervaly na odchádzajúcich kanáloch vrchola  $u$ , tak na routovaciu tabuľku v  $u$  stačí pamäť  $\mathbf{k} \cdot \log \mathbf{N}$

**Fakt 11** *Vrcholy stromu  $T$  možno očíslovať tak, že množina adries routovaných cez jeden výstupný kanál tvorí cyklický interval.*

napr. preorder – koreň, ľavý, pravý



## Tree labeling – Správa $(m, d)$ pre $d$ (prijatá/generovaná) v $u$

**if**  $d = u$  **then** spracuj lokálne  
**else**

- v stromoch je optimálne
- v súvislých grafoch použijeme kostru.

nech  $w$  je také, že  $d \in [a_w, b_w)$   
send  $(m, d)$  po kanáli  $v$

$G$  je graf,  $T$  jeho kostra,  $d_G(u, v)$ ,  $d_T(u, v)$   
cena najkratšej cesty medzi  $u, v$  v grafe  $G$   
resp. kostre  $T$ .

**Fakt 12** : *Neexistuje ohraničenie pomeru  $\frac{d_T(u,v)}{d_G(u,v)}$ . A to ani v prípade min-hop.*

### *kruh*

Majme súvislý grafu so symetrickým ohodnotením hrán ( $d(u, v) = d(v, u)$ ). Označme  $D_G$  polomer grafu  $G$ . Potom

**Fakt 13** *Existuje taká kostra  $T$  grafu  $G$ , že  $\forall u, v \in V$   $d_T(u, v) \leq 2 \cdot D_G$ .*

Nech  $T$  je optimálny sink-tree pre vrchol  $w$  z centra. O vzdialenosti  $d_T(u, v)$  platí

$$\begin{aligned} d_T(u, v) &\leq d_T(u, w) + d_T(w, v) \\ &= d_G(u, w) + d_G(w, v) \leq 2 \cdot D_G \end{aligned}$$

## Nevýhody routovania po strome

- nepoužívame hrany mimo stromu (mrhanie prostriedkami)
- používame len stromové hrany (zahrtenie)
- výpadok jednej hrany spôsobí kolaps systému

**ILS**(interval labeling scheme) je značkovanie vrcholov a hrán siete značkami zo  $Z_N$  tak že

- vrcholy majú rôzne značky
- v uzle majú všetky kanály rôzne značky

Hovoríme, že ILS je **platná**, ak routovaním podľa nej každá správa v konečnom čase dosiahne svoj cieľ.

**Veta 2** *Ku každému súvislému grafu existuje platná ILS.*

$l_w$  značka vrchola  $w$

$\alpha_{uv}$  značka hrany  $(u, v)$

$k_w = l_w + |T(w)|$  označuje podstrom s koreňom  $w$

- vrcholy značkujeme preorderom( koreň  $w$ , podstrom  $l_w, l_{w+1}, \dots, l_{w+|T(w)|-1}$ )
- značka hrana  $(u, w)$ 
  - $(u, w)$  je frond (nestromová hrana), tak  $\alpha_{uw} = l_w$
  - $w$  je  $T$ -syn vrchola  $u$ , tak  $\alpha_{uw} = l_w$
  - $w$  je  $T$ -otec vrchola  $u$ ,  $k_u \neq N$  alebo  $u$  nemá frond do koreňa, tak  $\alpha_{uw} = k_u$
  - $w$  je  $T$ -otec vrchola  $u$ ,  $k_u = N$  a  $u$  má frond do koreňa, tak  $\alpha_{uw} = l_w$ .



Nech  $\mathbf{lca}(u,v)$  je najbližší spoločný predchodca vrcholov  $u, v$ . Potom správe, ktorá sa nachádza vo vrchole  $u$  a smeruje do vrchola  $v$ , priradíme hodnotu

$$f_v(u) = (-lca(u, v), l_u)$$

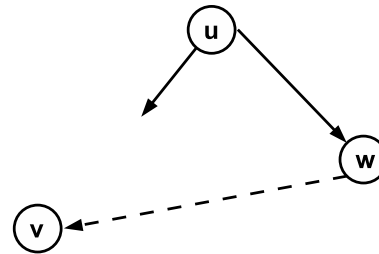
**vrchol  $u$  posúva správu pre vrchol  $v$  do  $w$**

**Lema 1** *V situácii, keď podľa DFS ILS posúva vrchol  $u$  správu pre vrchol  $v$  do vrchola  $w$  platí*

$$l_u > l_v \quad \Rightarrow \quad l_w < l_u$$

$$l_u < l_v \quad \Rightarrow \quad l_w \leq l_v$$

$$l_u < l_v \quad \Rightarrow \quad f_v(w) < f_v(u)$$



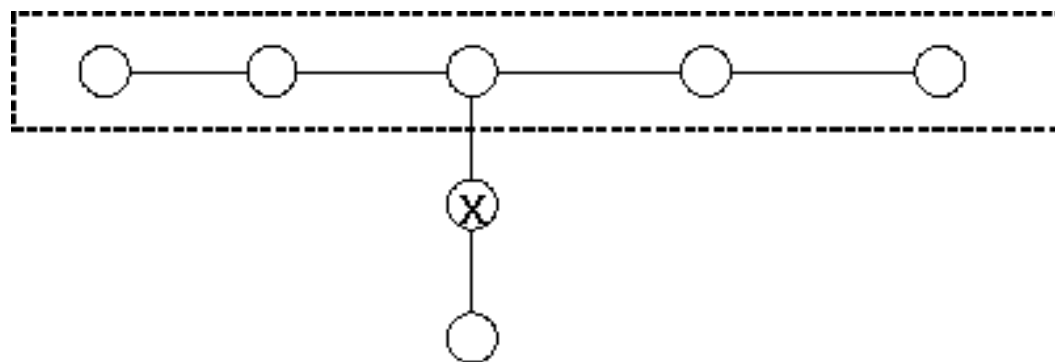
**Fakt 14** *Existuje sieť  $G$  taká, že pre každú platnú ILS existujú v  $G$  dva vrcholy  $u, v$  také, že na doručenie packetu z  $u$  do  $v$  treba aspoň  $\frac{3}{2}D_G$  hopov.*

**Fakt 15** *Existuje min-hop ILS pre kruh.*

**Fakt 16** *Existuje min-hop ILS pre  $n \times n$  mriežku.*

lineárne ILS schémy – intervaly priradené jednotlivým kanálom sú lineárne

**Fakt 17** *Existuje sieť, pre ktorú neexistuje platná lineárna ILS*



**Fakt 18** *Existuje min-hop ILS pre hyperkocku*

**PLS** (prefix labeling scheme) nad abecedou  $\Sigma$  je značkovanie vrcholov a hrán siete reťazcami nad  $\Sigma^*$  tak, že

- vrcholy majú rôzne značky
- v uzle majú všetky kanály rôzne značky

PLS je *platná*, ak routovaním podľa nej každá správa v konečnom čase dosiahne svoj cieľ  
**packet pre  $d$  (prijatý/generovaný) v  $u$**

**if  $d = l_u$  then** doruč lokálne  
**else**

nech  $\alpha_i$  je najdlhšie číslo kanála také, že  $\alpha_i < d$ ;  
send packet pre  $d$  po kanáli  $\alpha_i$

**Veta 3** *Ku každému súvislému grafu  $G$  existuje platná PLS.*

**Konštrukcia tree-PLS** Nech  $T$  je kostra grafu  $G$ .

- koreň je označený  $\epsilon$
- ak  $w$  je syn  $u$  tak  $l_w = l_u.a_i$   
ak  $u_1, \dots, u_k$  sú rôzni synovia vrchola  $u$ , tak  $a_1 \neq \dots \neq a_k$  a  $l_{u_i} = l_u.a_i$
- ak  $uw$  je frond, tak  $\alpha_{uw} = l_w$
- ak  $w$  je syn  $u$ , tak  $\alpha_{uw} = l_w$
- ak  $w$  je otec  $u$ , tak  $\alpha_{uw} = \epsilon$  okrem prípadu, keď  $u$  má frond do koreňa; vtedy  $\alpha_{uw} = l_u$

**Fakt 19** *Pre popísanú tree-PLS platí*

1. ak  $u \neq v$  tak v  $u$  existuje kanál označený prefixom  $l_v$
2. v situácii, keď  $u$  posiela správu do  $v$  cez  $w$ 
  - (a) ak  $u \in T[v]$ , tak  $w$  je predchodca  $u$ , resp.  $v$  je predchodca  $w$
  - (b) ak  $u$  je predchodca  $v$  tak  $w$  je predchodca  $v$  bližšie k  $v$  (ako  $u$ )
  - (c) ak  $u \notin T[v]$ , tak  $w$  je predchodca  $v$  alebo  $d_T(w, v) < d_T(u, v)$

**Dôsledok 1** *Pre každý graf  $G$  s priemerom  $D_G$  existuje PLS, ktorá každý packet doručí  $\leq 2D_G$  hopov.*