

2-INF-132 Úvod do distribuovaných algoritmov

Rasto Královič

Katedra informatiky, FMFI UK Bratislava
kralovic@dcs.fmph.uniba.sk



- ▶ **Gerard Tel:** *Introduction to Distributed Algorithms*, Cambridge University Press, 2000, ISBN 0521794838
- ▶ **Nancy Lynch:** *Distributed Algorithms*, Morgan Kaufmann Publishers, 1996, ISBN 1558603484
- ▶ **Frank Thomson Leighton:** *Introduction to Parallel Algorithms and Architectures: Arrays, Trees, Hypercubes*, Morgan Kaufmann Publishers, 1991, ISBN 1558601171

sekvenčné výpočty

- ▶ algoritmus (počítanie funkcie)
- ▶ čas
 - problém: presný, ale neprenositelný
- ▶ riešenie: počet inštrukcií v modeli
 - problém: aký model?
- ▶

jednoduchý "assembler"

- registre $r_0, r_1, r_2, r_3, \dots$
- vstup $a_0, a_1, a_2, a_3, \dots$
- program = postupnosť inštrukcií
- pozícia v programe

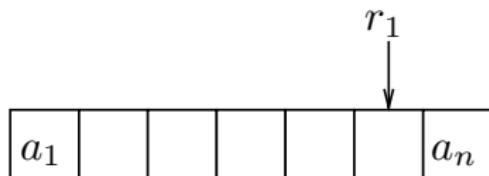
inštrukcie

priradenie	$r_X := r_Y$	X, Y je konštanta alebo register
	$r_X := a_Y$	
	$r_X := c$	c je konštanta

výpočet	$r_X := r_Y \square r_Z$	\square je $+, -$
skok	goto i	
podmienka	if $r_X \leq r_Y$ then goto i	

príklad: čo robí tento program?

- 1: vstup: $a_0 = n$ a_1, a_2, \dots, a_n , $a_i \geq 0$
- 2: $r_0 := -1$
- 3: $r_1 := a_0$
- 4: $r_1 := r_1 + 1$
- 5: $r_1 := r_1 - 1$
- 6: **if** $r_1 = 0$ **then goto** 10
- 7: **if** $r_0 \geq a_{r_1}$ **then goto** 5
- 8: $r_0 := a_{r_1}$
- 9: **goto** 5
- 10: výstup: r_0



r_0

maximum

sekvenčné výpočty

- ▶ algoritmus (počítanie funkcie)
- ▶ čas
 - problém: presný, ale neprenositelný
- ▶ riešenie: počet inštrukcií v modeli
 - problém: aký model?
- ▶ riešenie: zaujíma nás asymptotický rast

$$f \in O(g) \Leftrightarrow \exists c, n_0 \ \forall n \geq n_0 : f(n) \leq c \cdot g(n)$$

$$f \in \Omega(g) \Leftrightarrow \exists c, n_0 \ \forall n \leq n_0 : f(n) \leq c \cdot g(n)$$

distribuované algoritmy

dôvody

- ▶ lepší návrh
OS, browser, ...
- ▶ rýchlejšie riešenie problémov
scheduling, FEM, ...
- ▶ veľké dátá
- ▶ využitie hardvéru
GPU, seti@home, ...
- ▶ fyzická vzdialenosť
ATM, social networking, ...

výzvy

- ▶ rôzne ciele
- ▶ rôzne úlohy
- ▶ rôzny hardvér
- ▶ narastanie zložitosti

taxonómia

- ▶ Flynn (SISD, SIMD, MIMD)
- ▶ tightly/loosely coupled

nie je jeden model

bez zdieľanej pamäte – procesy

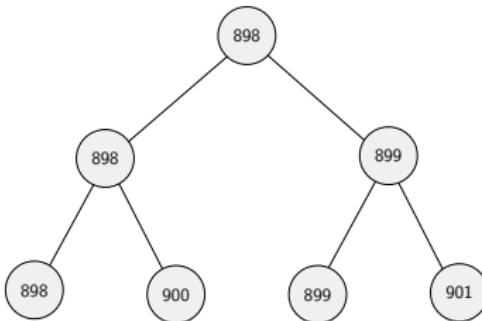
```
void do_child(int data_pipe[]) {                                int main() {  
    int c,rc;                                         int data_pipe[2];  
    close(data_pipe[1]);                                int pid,rc;  
    while ((rc = read(data_pipe[0], &c, 1)) > 0)          rc = pipe(data_pipe);  
        putchar(c);                                     pid = fork();  
    exit(0);                                         switch (pid) {  
}                                                               case 0:  
                                                               do_child(data_pipe);  
void do_parent(int data_pipe[]) {                               default:  
    int c,rc;                                         do_parent(data_pipe);  
    close(data_pipe[0]);                                }  
    while ((c = getchar()) > 0)                         }  
        rc = write(data_pipe[1], &c, 1);  
    close(data_pipe[1]);  
    exit(0);
```

```
int main(void) {
    int value1, value2;
    int a=42;

    printf("%d started (a=%d)\n",
           getpid(),a);
    value1 = fork();
    printf("%d: value1=%d (a=%d)\n",
           getpid(),value1,a);

    value2 = fork();
    printf("%d: value2=%d (a=%d)\n",
           getpid(),value2,a);
    return 0;
}
```

```
898 started (a=42)
898: value1=899 (a=42)
899: value1=0 (a=42)
898: value2=900 (a=42)
899: value2=901 (a=42)
900: value2=0 (a=42)
901: value2=0 (a=42)
```



bez zdieľanej pamäte – procesy

```
void proc(int p,int fd[2]) {           int main(void) {
    int i=0;
    close(fd[0]);
    while(1) {
        sleep(rand()%4);
        sprintf(line,"h%d from %d\n",
                ++i,p);
        write(fd[1], line, strlen(line));
    }
}

int fds[NPROC][2];
int n,c,i;
fd_set r;
int maxf = 0;
for (i=0;i<NPROC;i++) {
    pipe(fds[i]);
    if (fork()==0) proc(i,fds[i]);
    else {
        close(fds[i][1]);
        if (fds[i][0]>maxf) maxf=fds[i][0];
    }
}
while(1) {
    FD_ZERO(&r);
    for (i=0;i<NPROC;i++)
        FD_SET(fds[i][0],&r);
    select(maxf+1,&r,NULL,NULL,NULL);
    for (i=0;i<NPROC;i++)
        if (FD_ISSET(fds[i][0],&r)) {
            n = read(fds[i][0],line,MAXLINE);
            printf("%s",line);
        }
}}
```

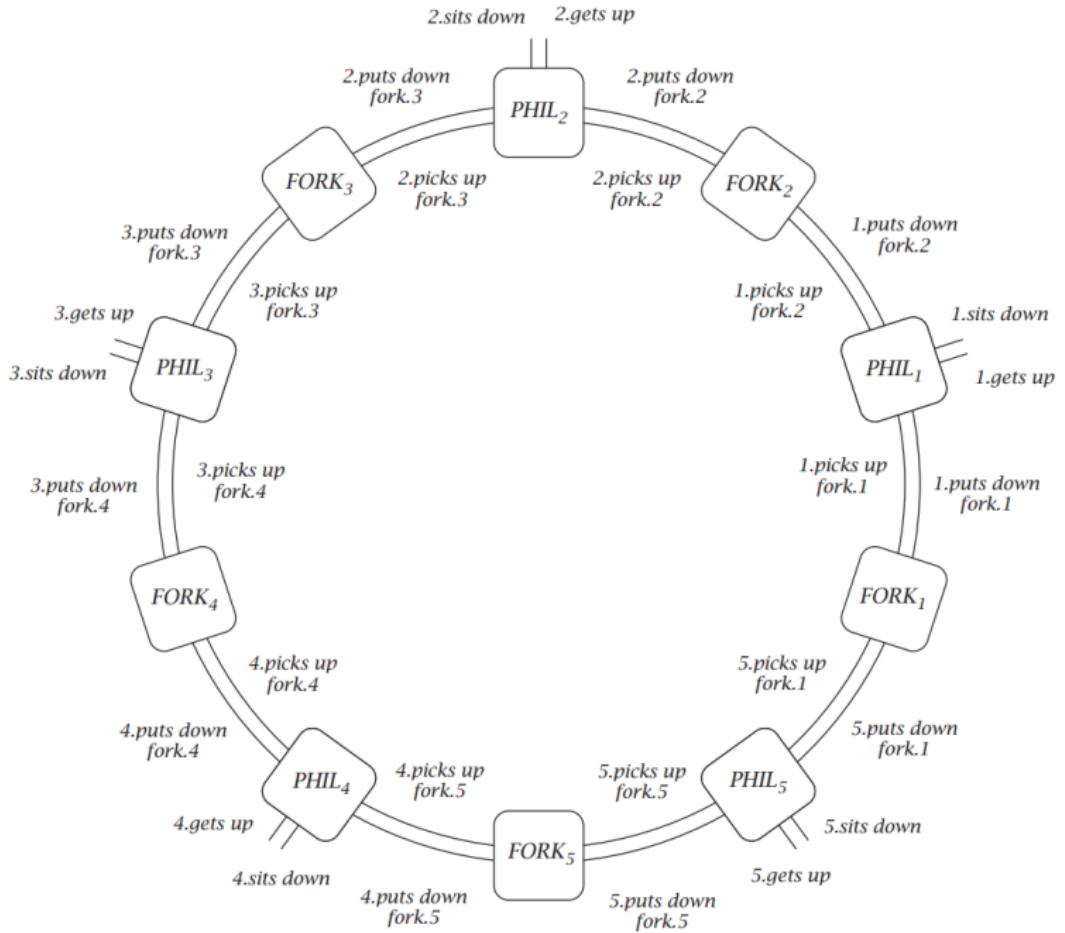
CSP (Communicating Sequential Processes)

- ▶ sémantika
- ▶ pozorovanie procesu: postupnosť udalostí
- ▶ proces: množina pozorovaní (**trace**)
- ▶ vytváranie procesov:
 - ▶ **prefix:** $(x \mapsto P)$
 - ▶ **rekurzia:** $\text{CLOCK} = (\text{tick} \mapsto \text{CLOCK})$
 - ▶ $F \equiv \mu X. F(X)$
 - ▶ **deterministický výber:** $(x \mapsto P \mid y \mapsto Q)$
 - ▶ **paraleлизmus:** $(P \parallel Q)$ všetky "premiešania" pozorovaní
- ▶ **komunikácia:** zdieľaním udalostí (rendezvous)

$$K = \mu X. \text{minca} \mapsto (\text{káva} \mapsto X \mid \text{čaj} \mapsto X)$$

$$Z = \mu X. (\text{čaj} \mapsto X \mid \text{káva} \mapsto X \mid \text{minca} \mapsto \text{čaj} \mapsto X)$$

$$(Z \parallel K) = \mu X. (\text{minca} \mapsto \text{čaj} \mapsto X)$$



$$\begin{aligned}FORK_i = & (i.\text{picks up fork}.i \rightarrow i.\text{puts down fork}.i \rightarrow FORK_i \\& | (i \oplus 1).\text{picks up fork}.i \rightarrow (i \oplus 1).\text{puts down fork}.i \rightarrow FORK_{i+1})\end{aligned}$$
$$\begin{aligned}LEFT_i = & (i.\text{sits down} \rightarrow i.\text{picks up fork}.i \rightarrow \\& i.\text{puts down fork}.i \rightarrow i.\text{gets up} \rightarrow LEFT_i)\end{aligned}$$
$$\begin{aligned}RIGHT_i = & (i.\text{sits down} \rightarrow i.\text{picks up fork}.(i \oplus 1) \rightarrow \\& i.\text{puts down fork}.(i \oplus 1) \rightarrow i.\text{gets up} \rightarrow RIGHT_i)\end{aligned}$$
$$PHIL_i = LEFT_i \parallel RIGHT_i$$
$$PHILOS = (PHIL_0 \parallel PHIL_1 \parallel PHIL_2 \parallel PHIL_3 \parallel PHIL_4)$$
$$FORKS = (FORK_0 \parallel FORK_1 \parallel FORK_2 \parallel FORK_3 \parallel FORK_4)$$
$$COLLEGE = PHILOS \parallel FORKS$$

$$\begin{aligned}FORK_i = & (i.picks\ up\ fork.i \rightarrow i.puts\ down\ fork.i \rightarrow FORK_i \\& | (i \ominus 1).picks\ up\ fork.i \rightarrow (i \ominus 1).puts\ down\ fork.i \rightarrow FORK_i)\end{aligned}$$
$$\begin{aligned}LEFT_i = & (i.sits\ down \rightarrow i.picks\ up\ fork.i \rightarrow \\& i.puts\ down\ fork.i \rightarrow i.gets\ up \rightarrow LEFT_i)\end{aligned}$$
$$\begin{aligned}RIGHT_i = & (i.sits\ down \rightarrow i.picks\ up\ fork.(i \oplus 1) \rightarrow \\& i.puts\ down\ fork.(i \oplus 1) \rightarrow i.gets\ up \rightarrow RIGHT_i)\end{aligned}$$
$$PHIL_i = LEFT_i \parallel RIGHT_i$$
$$PHILOS = (PHIL_0 \parallel PHIL_1 \parallel PHIL_2 \parallel PHIL_3 \parallel PHIL_4)$$
$$FORKS = (FORK_0 \parallel FORK_1 \parallel FORK_2 \parallel FORK_3 \parallel FORK_4)$$
$$COLLEGE = PHILOS \parallel FORKS$$

$$U = \bigcup_{i=0}^4 \{i.gets\ up\} \quad D = \bigcup_{i=0}^4 \{i.sits\ down\}$$

$$FOOT_0 = (x : D \rightarrow FOOT_1)$$
$$FOOT_j = (x : D \rightarrow FOOT_{j+1} \mid y : U \rightarrow FOOT_{j-1}) \quad \text{for } j \in \{1, 2, 3\}$$
$$FOOT_4 = (y : U \rightarrow FOOT_3)$$
$$NEWCOLLEGE = (COLLEGE \parallel FOOT_0)$$

nenastane deadlock

$$U = \bigcup_{i=0}^4 \{i.\text{gets up}\} \quad D = \bigcup_{i=0}^4 \{i.\text{sits down}\}$$

$(newcollege/s) \neq STOP$ for all $s \in \text{traces(NEWCOLLEGE)}$

- ▶ $\text{seated}(s) := \#(s \upharpoonright D) - \#(s \upharpoonright U)$
- ▶ $\text{seated} \leq 4$, lebo $FOOT_0$
- ▶ $\text{seated} = \leq 3$, môže sa posadiť
- ▶ nech $\text{seated} = 4$
 - ▶ niekto je jediaci, môže prestaviť
 - ▶ max 3 zdvihnuté vidličky \Rightarrow aspoň 1 môže zdvihnuť ľavú
 - ▶ 4 zdvihnuté \Rightarrow jeden môže zobrať pravú
 - ▶ 5 zdvihnutých \Rightarrow jeden je jediaci

bez zdieľanej pamäte – procesy

```
void do_child(int data_pipe[]) {                                int main() {  
    int c,rc;                                         int data_pipe[2];  
    close(data_pipe[1]);                                int pid,rc;  
    while ((rc = read(data_pipe[0], &c, 1)) > 0)          rc = pipe(data_pipe);  
        putchar(c);                                     pid = fork();  
    exit(0);                                         switch (pid) {  
}                                                               case 0:  
                                                               do_child(data_pipe);  
void do_parent(int data_pipe[]) {                               default:  
    int c,rc;                                         do_parent(data_pipe);  
    close(data_pipe[0]);                                }  
    while ((c = getchar()) > 0)                         }  
        rc = write(data_pipe[1], &c, 1);  
    close(data_pipe[1]);  
    exit(0);
```

v sieti: posielanie správ (sokety)

server

```
int main(int argc, char *argv[]) {  
    int sockfd, newsockfd, portno;  
    socklen_t clilen;  
    char buffer[256];  
    struct sockaddr_in serv_addr, cli_addr;  
    int n;  
  
    sockfd = socket(AF_INET, SOCK_STREAM, 0);  
    bzero((char *) &serv_addr, sizeof(serv_addr));  
    portno = atoi(argv[1]);  
  
    serv_addr.sin_family = AF_INET;  
    serv_addr.sin_addr.s_addr = INADDR_ANY;  
    serv_addr.sin_port = htons(portno);  
  
    bind(sockfd, (struct sockaddr *) &serv_addr,  
          sizeof(serv_addr));  
    listen(sockfd,5);  
    clilen = sizeof(cli_addr);  
    newsockfd = accept(sockfd,  
                      (struct sockaddr *) &cli_addr, &clilen);  
    bzero(buffer,256);  
  
    n = read(newsockfd,buffer,255);  
    printf("Here is the message: %s\n",buffer);  
  
    n = write(newsockfd,"I got your message",18);  
    close(newsockfd);  
    close(sockfd);  
}
```

klient

```
int main(int argc, char *argv[]) {  
    int sockfd, portno, n;  
    struct sockaddr_in serv_addr;  
    struct hostent *server;  
    char buffer[256];  
  
    portno = atoi(argv[2]);  
    sockfd = socket(AF_INET, SOCK_STREAM, 0);  
    server = gethostbyname(argv[1]);  
  
    bzero((char *) &serv_addr, sizeof(serv_addr));  
    serv_addr.sin_family = AF_INET;  
    bcopy((char *)server->h_addr,  
          (char *)&serv_addr.sin_addr.s_addr,  
          server->h_length);  
    serv_addr.sin_port = htons(portno);  
  
    connect(sockfd,(struct sockaddr *) &serv_addr,  
            sizeof(serv_addr));  
  
    printf("Please enter the message: ");  
    bzero(buffer,256);  
    fgets(buffer,255,stdin);  
    n = write(sockfd,buffer,strlen(buffer));  
    bzero(buffer,256);  
    n = read(sockfd,buffer,255);  
    printf("%s\n",buffer);  
    close(sockfd);  
    return 0;  
}
```

v sieti: posielanie správ (MPI)

```
int main(int argc, char *argv[]) {
    const int tag = 47;
    ...
    MPI_Init(&argc, &argv);
    MPI_Comm_size(MPI_COMM_WORLD, &ntasks);
    MPI_Comm_rank(MPI_COMM_WORLD, &id);

    if (id == 0) {
        MPI_Get_processor_name(msg, &length);
        printf("Hello World from process %d running on %s\n", id, msg);
        for (i=1; i<ntasks; i++) {
            MPI_Recv(msg, 80, MPI_CHAR, MPI_ANY_SOURCE,
                     tag, MPI_COMM_WORLD, &status);
            source_id = status.MPI_SOURCE;
            printf("Hello World from process %d running on %s\n", source_id, msg);
        }
    }
    else {
        MPI_Get_processor_name(msg, &length);
        MPI_Send(msg, length, MPI_CHAR, 0, tag, MPI_COMM_WORLD);
    }
    MPI_Finalize();
    if (id==0) printf("Ready\n");
}
```

sietový model

- ▶ nezávislé zariadenia (procesy, procesory, uzly)
- ▶ posielanie správ
- ▶ lokálny pohľad (číslovanie portov)
- ▶ asynchronné, spoľahlivé správy
- ▶ topológia siete
- ▶ wakeup/terminácia

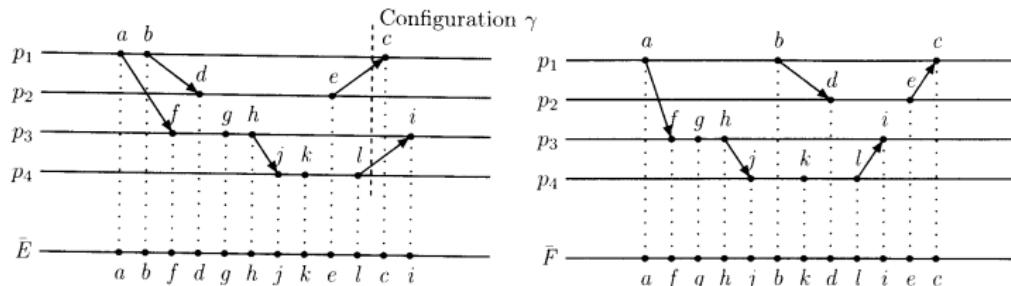
zložitosť: v závislosti od počtu procesorov

- ▶ počet správ / bitov
- ▶ čas (beh normovaný na max. dĺžku 1)

sietový model

formálny model: prechodové systémy

- systém je trojica $(\mathcal{C}, \mapsto, \mathcal{I})$
- beh (execution) je postupnosť $\gamma_0 \mapsto \gamma_1 \mapsto \gamma_2 \mapsto \dots$, kde $\gamma_0 \in \mathcal{I}$
- $\mathcal{C} = \mathcal{C}_L^n \times \mathcal{M}$: stavy systému = stavy procesorov + správy na linkách
- \mapsto : krok jedného procesora (výpočet, posanie správy, prijatie správy)
- fair scheduler?



Sieťový model

formálny model: prechodové systémy

- ▶ (ne)závislosť udalostí:
 - 1) poslanie pred prijatím
 - 2) časovanie v rámci procesora
 - 3) tranzitivita
- ▶ výpočet (computation): trieda ekvivalencie behov

logické hodiny (Lamport)

```
var  $\theta_p$  : integer      init 0 ;  
  
(* An internal event *)           (* A receive event *)  
   $\theta_p := \theta_p + 1$  ;          receive ⟨messg,  $\theta$ ⟩ ;  $\theta_p := \max(\theta_p, \theta) + 1$  ;  
  Change state                   Change state  
  
(* A send event *)  
   $\theta_p := \theta_p + 1$  ;  
  send ⟨messg,  $\theta_p$ ⟩ ; Change state
```

sietový model

horný odhad pre problém

Existuje algoritmus, ktorý **pre všetky** topológie, vstupy, časovania,
... pracuje správne a vymení najviac $f(n)$ správ

dolný odhad pre problém

Pre každý algoritmus, **existuje kombinácia** topológie, vstupu, časovania,
....
... že bud' nepracuje správne alebo vymení aspoň $f(n)$ správ

sietový model: broadcast a convergecast (maximum)

```
const: deg : integer
      ID : integer
      Neigh : [1...deg] link
      parent : link
var: msg : text
Init:  
  
if parent = NULL  
    send <dispatch, msg> to self
```

Code:

loop forever

On receipt $\langle \text{dispatch}, \text{new_msg} \rangle$ from parent or self :

```
msg := new_msg  
for all  $l \in \text{Neigh} - \{\text{parent}\}$  do  
    send <dispatch, msg> to  $l$   
skonči algoritmus
```

sietový model: broadcast a convergecast (maximum)

const: *deg* : integer
 ID : integer
 Neigh : [1...*deg*] link
 parent : link
 msg : integer
var: *count* : integer
 max : integer

Init:

count = 0
max = *msg*
if *deg* = 1
 send ⟨my_max, *max*⟩ **to** parent

Code:

loop forever

On receipt ⟨my_max, *x*⟩ from *Neigh*[*i*]:

max = maximum{*max*, *x*}
 count ++
 if *count* = *deg* - 1
 send ⟨my_max, *max*⟩ **to** parent
 skonči algoritmus

cvičenia

- ▶ voľba šéfa na (anonymných) stromoch
- ▶ voľba šéfa na (anonymnej) mriežke

prehľadávane

- ▶ na začiatku je jeden iniciátor
- ▶ treba informovať všetkých

shout-and-echo

- ▶ iniciátor pošle **shout**
- ▶ ak príde **shout** do nového vrchola
 - ▶ označí hranu
 - ▶ pošle po neoznačených **shout**
 - ▶ čaká na všetky **echo**
 - ▶ pošle **echo**
- ▶ ak príde **shout** do navštíveného –
okamžite **echo**

prehľadávane

- ▶ na začiatku je jeden iniciátor
- ▶ treba informovať všetkých

shout-and-echo

- ▶ iniciátor pošle **shout**
- ▶ ak príde **shout** do nového vrchola
 - ▶ označí hranu
 - ▶ pošle po neoznačených **shout**
 - ▶ čaká na všetky **echo**
 - ▶ pošle **echo**
- ▶ ak príde **shout** do navštíveného – okamžite **echo**

zložitosť

- ▶ $4m$ správ
- ▶ $2diam(G)$ čas

traverzovanie

- ▶ na začiatku je jeden iniciátor
- ▶ treba informovať všetkých
- ▶ v každom okamihu je len jeden *token*

$f(x)$ -traverzovanie

Na objavenie x vrcholov treba $\max\{f(x), n\}$ správ.

prehľadávanie do hĺbky

```
var usedp[q] : boolean    init false for each  $q \in \text{Neigh}_p$  ;  
                           (* Indicates whether  $p$  has already sent to  $q$  *)  
fatherp : process      init undef ;
```

For the initiator only, execute once:

```
begin fatherp := p ; choose  $q \in \text{Neigh}_p$  ;  
     usedp[q] := true ; send ⟨tok⟩ to q  
end
```

For each process, upon receipt of ⟨tok⟩ from q_0 :

```
begin if fatherp = undef then fatherp := q0 ;  
      if  $\forall q \in \text{Neigh}_p : \text{used}_p[q]$   
          then decide  
      else if  $\exists q \in \text{Neigh}_p : (q \neq \text{father}_p \wedge \neg \text{used}_p[q])$   
          then begin if fatherp  $\neq q_0 \wedge \neg \text{used}_p[q_0]$   
                  then  $q := q_0$   
                  else choose  $q \in \text{Neigh}_p \setminus \{\text{father}_p\}$   
                        with  $\neg \text{used}_p[q]$  ;  
                  usedp[q] := true ; send ⟨tok⟩ to q  
          end  
      else begin usedp[fatherp] := true ;  
              send ⟨tok⟩ to fatherp  
          end  
end
```

- zložitosť $2m$ (čas aj správy)

- aké súčetné $x - ?$

Awerbuch: $4m$ správ, $4n - 2$ čas

```
var  $used_p[q]$  : boolean      init false for each  $q \in Neigh_p$  ;  
                               (* Indicates whether  $p$  has already sent to  $q$  *)  
 $father_p$  : process        init  $undef$  ;
```

For the initiator only, execute once:

```
begin  $father_p := p$  ; choose  $q \in Neigh_p$  ;  
      forall  $r \in Neigh_p$  do send  $\langle vis \rangle$  to  $r$  ;  
      forall  $r \in Neigh_p$  do receive  $\langle ack \rangle$  from  $r$  ;  
       $used_p[q] := true$  ; send  $\langle tok \rangle$  to  $q$   
end
```

For each process, upon receipt of $\langle tok \rangle$ from q_0 :

```
begin if  $father_p = undef$  then  
      begin  $father_p := q_0$  ;  
            forall  $r \in Neigh_p \setminus \{father_p\}$  do send  $\langle vis \rangle$  to  $r$  ;  
            forall  $r \in Neigh_p \setminus \{father_p\}$  do receive  $\langle ack \rangle$  from  $r$   
      end ;  
      if  $p$  is the initiator and  $\forall q \in Neigh_p : used_p[q]$   
      then decide  
      else if  $\exists q \in Neigh_p : (q \neq father_p \wedge \neg used_p[q])$   
      then begin if  $father_p \neq q_0 \wedge \neg used_p[q_0]$   
              then  $q := q_0$   
              else choose  $q \in Neigh_p \setminus \{father_p\}$   
                    with  $\neg used_p[q]$  ;  
               $used_p[q] := true$  ; send  $\langle tok \rangle$  to  $q$   
      end  
      else begin  $used_p[father_p] := true$  ;  
                send  $\langle tok \rangle$  to  $father_p$   
      end  
end
```

For each process, upon receipt of $\langle vis \rangle$ from q_0 :

```
begin  $used_p[q_0] := true$  ; send  $\langle ack \rangle$  to  $q_0$  end
```

► otázka: treba čakať na ack ?

prehľadávanie do šírky

Cheung'83: kubická komunikácia, lineárny čas

- ▶ každá správa obsahuje počítadlo hopov
- ▶ na začiatku iniciátor: všetkým susedom pošle 1
- ▶ každý vrchol p : lokálnu vzdialenosť $dist_p := \infty$
- ▶ on recv i : $dist_p := \min\{i, dist_p\}$, ak sa zmenila, pošli všetkým

Cheung,Zhu'87: kvadratická komunikácia aj čas

- ▶ buduje kostru po vrstvách

kombinácia?

voľba šéfa: úplné grafy

const: N : integer
 ID : integer
 $Neigh$: $[1\dots N-1]$ link

var: $leader$: boolean
 $count$: integer
 i : integer

Init:

$count := 0$
 $leader := \text{false}$

Code:

```
for  $i = 1$  to  $N - 1$  do
    send ⟨elect,  $ID$ ⟩ to  $Neigh[i]$ 
while  $count < N - 1$  wait
for  $i = 1$  to  $N - 1$  do
    send ⟨leader,  $ID$ ⟩ to  $Neigh[i]$ 
 $leader := \text{true}$ 
```

On receipt ⟨elect, id_i ⟩ from $Neigh[i]$:

if $id_i > ID$ **send** ⟨accept⟩ **to** $Neigh[i]$

On receipt ⟨accept⟩ from $Neigh[i]$:

$count ++$

On receipt ⟨leader, id_i ⟩ from $Neigh[i]$:

Skonči algoritmus

vol'ba šéfa: úplné grafy II

const: N : integer
 ID : integer
 $Neigh$: [1...N-1] link

var: $leader$: boolean
 $state$: {active, captured, killed}
 $level$: integer
 $parent$: link
 msg : {victory, defeat}
 i : integer

Init:

$state := \text{active}$
 $level := 0$
 $leader := \text{false}$

Code:

```
for  $i = 1$  to  $N - 1$  do
    send ⟨capture, [level, ID]⟩ to Neigh[i]
    receive ⟨accept⟩ from Neigh[i]
    level ++
leader := true
for  $i = 1$  to  $N - 1$  do
    send ⟨leader, ID⟩ to Neigh[i]
```

Dead:

loop forever

On receipt ⟨capture, [level_i, id_i]⟩ from Neigh[i]:

```
if state ∈ {active, killed} and [leveli, idi] > [level, ID]
    state := captured
    parent := Neigh[i]
    send ⟨accept⟩ to parent
    goto Dead
else if state = captured
    send ⟨help, [leveli, idi]⟩ to parent
    receive msg from parent
    if msg = defeat
        send ⟨accept⟩ to Neigh[i]
        parent := Neigh[i]
```

On receipt ⟨help, [level_i, id_i]⟩ from Neigh[i]:

```
if [leveli, idi] < [level, ID]
    send ⟨victory⟩ to Neigh[i]
else
    send ⟨defeat⟩ to Neigh[i]
    if state = active
        state := killed
    goto Dead
```

On receipt ⟨leader, id_i⟩ from Neigh[i]:

Skončí algoritmus

voľba šéfa: úplné grafy II - analýza

Lema 1

V ľubovoľnom výpočte existuje pre každý level $l = 0, \dots, N - 1$ aspoň jeden proces, ktorý bol počas výpočtu na leveli l .

Lema 2

Nech v je aktívny proces (*state = active*) s levelom l . Potom existuje l zjazatých procesov ktoré patria v (t.j. ich premenná *parent* ukazuje na v).

⇒ práve jeden proces je šéf

Lemma 3

Ľubovoľnom výpočte je najviac $N/(l + 1)$ procesov, ktoré niekedy dosiahli level l .

⇒ maximálne $\sum_{l=1}^{N-1} \frac{N}{l+1} = N(\mathbf{H}_N - 1) \approx N \log N$ postupov o level (=správ)

dolný odhad

Treba $\Omega(n \log n)$ správ

- ▶ “globálny” algoritmus
- ▶ graf indukovaný novými správami
- ▶ udržujem výpočet:
 - ▶ jeden súvislý komponent, ostatné vrcholy izolované
 - ▶ $e(k)$ – koľko viem vynútiť na komponente k
 - ▶ $e(2k + 1) = 2e(k) + k + 1$

jednosmerné kruhy

Chang, Roberts

const: ID : integer
 l_{in}, l_{out} : link
var: $leader$: integer

Init:

$leader := NULL$

Code:

send $\langle ID \rangle$
wait until $leader <> NULL$

On receipt $\langle i \rangle$:

if $i < ID$ **then send** $\langle i \rangle$
if $i = ID$ **then**
 $leader := ID$
 send $\langle leader, ID \rangle$

On receipt $\langle leader, x \rangle$:

$leader := x$
send $\langle leader, ID \rangle$

- ▶ $\Omega(n^2)$ správ v najhoršom prípade
- ▶ $O(n \log n)$ v priemernom

Chang Roberts - analýza

Koľkokrát sa pohla správa i ?

$$P(i, k) = \frac{\binom{n-i}{k-1} \cdot (i-1)}{\binom{n-1}{k-1} \cdot (n-k)}$$

$$E[X_i] = \sum_{k=1}^{n-i+1} k \cdot P(i, k)$$

$$E[X] = n + \sum_{i=2}^n E[X_i] = n + \sum_{i=2}^n \sum_{k=1}^{n-i+1} k \cdot P(i, k)$$

Chang Roberts - analýza

lema

$$\sum_{j=k}^{n-1} \frac{j!}{(j-k)!} = k! \binom{n}{k+1}$$

dôkaz

$$\sum_{j=k}^{n-1} \frac{j!}{(j-k)!} = \sum_{j=k}^{n-1} \frac{k!j!}{k!(j-k)!} = k! \sum_{j=k}^{n-1} \binom{j}{k} = k! \binom{n}{k+1}$$

$$\begin{aligned}
n + \sum_{i=2}^n \sum_{k=1}^{n-i+1} k \cdot \frac{\binom{n-i}{k-1} \cdot (i-1)}{\binom{n-1}{k-1} \cdot (n-k)} &= n + \sum_{j=1}^{n-1} \sum_{k=1}^j k \cdot \frac{\binom{j-1}{k-1} \cdot (n-j)}{\binom{n-1}{k-1} \cdot (n-k)} = \\
&= n + \sum_{j=1}^{n-1} \sum_{k=1}^j \frac{k(j-1)!(n-j)(k-1)!(n-k)!}{(k-1)!(j-k)!(n-1)!(n-k)} = \\
&= n + \sum_{k=1}^{n-1} \left[\frac{k(n-k-1)!}{(n-1)!} \sum_{j=k}^{n-1} \frac{(j-1)!(n-j)}{(j-k)!} \right] = \\
&= n + \sum_{k=1}^{n-1} \left[\frac{k(n-k-1)!}{(n-1)!} \left(\sum_{j=k}^{n-1} \frac{n(j-1)!}{(j-k)!} - \sum_{j=k}^{n-1} \frac{j!}{(j-k)!} \right) \right] = \\
&\quad \sum_{j=k}^{n-1} \frac{(j-1)!}{(j-k)!} = (k-1)! \binom{n-1}{k} \text{ a } \sum_{j=k}^{n-1} \frac{j!}{(j-k)!} = k! \binom{n}{k+1} \\
&= n + \sum_{k=1}^{n-1} \frac{k(n-k-1)!}{(n-1)!} \left[n \cdot (k-1)! \binom{n-1}{k} - k! \binom{n}{k+1} \right] = \\
&\quad = n + \sum_{k=1}^{n-1} \frac{n}{k+1}
\end{aligned}$$

dvojsmerné kruhy

Hirschberg-Sinclair

- ▶ level l : dobývať územie 2^l
- ▶ $\log n$ levelov
- ▶ $n/2^l$ vrcholov na leveli
- ▶ každý vrchol pošle 2^l správ

dvojsmerné kruhy

Hirschberg-Sinclair

- ▶ level l : dobývať územie 2^l
- ▶ $\log n$ levelov
- ▶ $n/2^l$ vrcholov na leveli
- ▶ každý vrchol pošle 2^l správ

Franklin

- ▶ level l : poraziť susedov (na rovnakom leveli; "synchronizácia")
- ▶ $\log n$ levelov
- ▶ n správ na level

dvojsmerné kruhy

Hirschberg-Sinclair

- ▶ level l : dobývať územie 2^l
- ▶ $\log n$ levelov
- ▶ $n/2^l$ vrcholov na leveli
- ▶ každý vrchol pošle 2^l správ

Franklin

- ▶ level l : poraziť susedov (na rovnakom leveli; "synchronizácia")
- ▶ $\log n$ levelov
- ▶ n správ na level

Dolev – simulácia na 1-smernom kruhu

- ▶ idea: presunúť identitu

dolný odhad

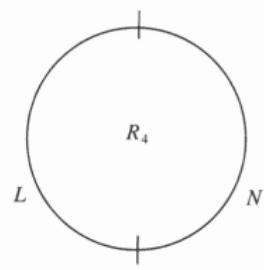
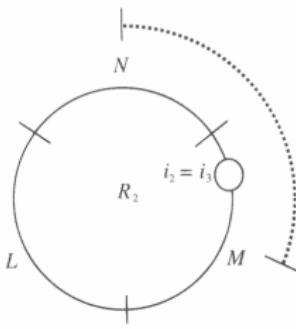
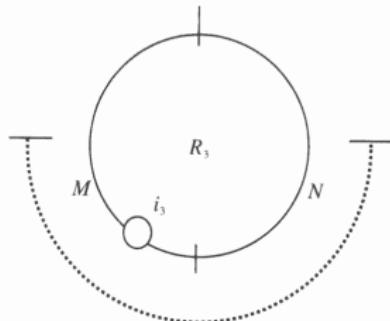
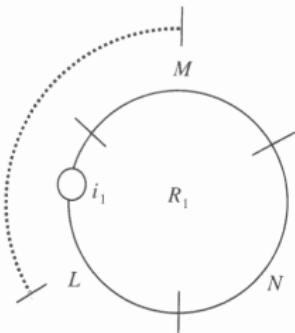
zapojiť do čiary L , vymení sa $C(L)$ správ

lema

Pre každé r existuje nekonečne veľa čiar dĺžky 2^r , kde $C(L) \geq r2^{r-2}$

- ▶ indukcia
- ▶ dve vrecia: vyberám trojice, chcem dve spojiť item pri spojení: dĺžka 2^{r+1} , počet správ $\geq r2^{r-1}$
- ▶ ešte treba 2^{r-1} správ; sporom

dolný odhad



GHS

- ▶ ľubovoľná topológia
- ▶ buduje sa kostra
- ▶ "segmenty"
- ▶ spájanie po najlacnejšej odchádzajúcej hrane:
 - A menší sa pripojí k väčšiemu
 - B rovnakí sa spoja
 - C väčší čaká
- ▶ veľkosť ⇒ level

GHS

```
var statep : (sleep, find, found) ;
    stachp[q] : (basic, branch, reject) for each q ∈ Neighp ;
    namep, bestwtp : real ;
    levelp : integer ;
    testchp, bestchp, fatherp : Neighp ;
    recp : integer ;
```

- (1) As the first action of each process, the algorithm must be initialized:

```
begin let pq be the channel of p with smallest weight ;
    stachp[q] := branch ; levelp := 0 ;
    statep := found ; recp := 0 ;
    send ⟨connect, 0⟩ to q
end
```

- (2) Upon receipt of ⟨connect, L⟩ from q:

```
begin if L < levelp then (* Combine with Rule A *)
    begin stachp[q] := branch ;
        send ⟨initiate, levelp, namep, statep⟩ to q
    end
    else if stachp[q] = basic
        then (* Rule C *) process the message later
        else (* Rule B *) send ⟨initiate, levelp + 1, ω(pq), find⟩ to q
end
```

- (3) Upon receipt of ⟨initiate, L, F, S⟩ from q:

```
begin levelp := L ; namep := F ; statep := S ; fatherp := q ;
    bestchp := undef ; bestwtp := ∞ ;
    forall r ∈ Neighp : stachp[r] = branch ∧ r ≠ q do
        send ⟨initiate, L, F, S⟩ to r ;
    if statep = find then begin recp := 0 ; test end
end
```

GHS

- (4) **procedure** *test*:
- ```
begin if $\exists q \in \text{Neigh}_p : \text{stack}_p[q] = \text{basic}$ then
 begin $\text{testch}_p := q$ with $\text{stack}_p[q] = \text{basic}$ and $\omega(pq)$ minimal ;
 send $\langle \text{test}, \text{level}_p, \text{name}_p \rangle$ to testch_p
 end
 else begin $\text{testch}_p := \text{udef}$; report end
end
```
- (5) Upon receipt of  $\langle \text{test}, L, F \rangle$  from *q*:
- ```
begin if  $L > \text{level}_p$  then (* Answer must wait! *)
    process the message later
    else if  $F = \text{name}_p$  then (* internal edge *)
        begin if  $\text{stack}_p[q] = \text{basic}$  then  $\text{stack}_p[q] := \text{reject}$  ;
            if  $q \neq \text{testch}_p$ 
                then send  $\langle \text{reject} \rangle$  to q
            else test
        end
        else send  $\langle \text{accept} \rangle$  to q
    end
end
```
- (6) Upon receipt of $\langle \text{accept} \rangle$ from *q*:
- ```
begin $\text{testch}_p := \text{udef}$;
 if $\omega(pq) < \text{bestwt}_p$
 then begin $\text{bestwt}_p := \omega(pq)$; $\text{bestch}_p := q$ end ;
 report
end
```
- (7) Upon receipt of  $\langle \text{reject} \rangle$  from *q*:
- ```
begin if  $\text{stack}_p[q] = \text{basic}$  then  $\text{stack}_p[q] := \text{reject}$  ;
    test
end
```

GHS

- (8) **procedure** *report*:
- ```
begin if $rec_p = \#\{q : stach_p[q] = branch \wedge q \neq father_p\}$
 and $testch_p = undef$ then
 begin $state_p := found$; send $\langle report, bestwt_p \rangle$ to $father_p$ end
 end
```
- (9) Upon receipt of  $\langle report, \omega \rangle$  from  $q$ :
- ```
begin if  $q \neq father_p$ 
      then (* reply for initiate message *)
          begin if  $\omega < bestwt_p$  then
              begin  $bestwt_p := \omega$  ;  $bestch_p := q$  end ;
                   $rec_p := rec_p + 1$  ; report
              end
          else (*  $pq$  is the core edge *)
              if  $state_p = find$ 
                  then process this message later
              else if  $\omega > bestwt_p$ 
                  then changeroot
                  else if  $\omega = bestwt_p = \infty$  then stop
          end
      end
```
- (10) **procedure** *changeroot*:
- ```
begin if $stach_p[bestch_p] = branch$
 then send $\langle changeroot \rangle$ to $bestch_p$
 else begin send $\langle connect, level_p \rangle$ to $bestch_p$;
 $stach_p[bestch_p] := branch$
 end
 end
```
- (11) Upon receipt of  $\langle changeroot \rangle$ :
- ```
begin changeroot end
```

analýza

správnosť

ukázať, že sa zvolí práve jeden šéf: nenastane deadlock

počet správ

- ▶ testovacie správy: jeden test po každej hrane
- ▶ kostrové správy: fragment s n_i vrcholmi pri postupe o level $O(n_i)$ správ
- ▶ postupy na level l : dizjunktné vrcholy

LE - summary

- ▶ pre kruhy, úplné grafy $\Theta(n \log n)$ správ
- ▶ bez promisu: $\Theta(n \log n + |E|)$ správ

ako vplývajú zmeny modelu na zložitosť?

- ▶ ("komprimovaná") znalosť mapy?
- ▶ synchrónnosť?

KKM

- ▶ $f(x)$ -traverzovanie
- ▶ tokeny traverzujú/označujú územia
- ▶ leveley: keď sa stretnú dva, vznikne nový
- ▶ naháňanie

KKM

```
var  $lev_p$  : integer      init -1 ;
     $cat_p, wait_p$  :  $\mathcal{P}$       init undef ;
     $last_p$         :  $Neigh_p$     init undef ;

begin if  $p$  is initiator then
    begin  $lev_p := lev_p + 1$  ;  $last_p := trav(p, lev_p)$  ;
            $cat_p := p$  ; send  $\langle annex, p, lev_p \rangle$  to  $last_p$ 
    end ;
    while ... (* Termination condition, see text *) do
        begin receive token  $(q, l)$  ;
            if token is annexing then  $t := A$  else  $t := C$  ;
            if  $l > lev_p$  then (* Case I *)
                begin  $lev_p := l$  ;  $cat_p := q$  ;
                        $wait_p := undef$  ;  $last_p := trav(q, l)$  ;
                       send  $\langle annex, q, l \rangle$  to  $last_p$ 
                end
            else if  $l = lev_p$  and  $wait_p \neq undef$  then (* Case II *)
                begin  $wait_p := undef$  ;  $lev_p := lev_p + 1$  ;
                        $last_p := trav(p, lev_p)$  ;  $cat_p := p$  ;
                       send  $\langle annex, p, lev_p \rangle$  to  $last_p$ 
                end
            else if  $l = lev_p$  and  $last_p = undef$  then (* Case III *)
                 $wait_p := q$ 
            else if  $l = lev_p$  and  $t = A$  and  $q = cat_p$  then (* Case IV *)
                begin  $last_p := trav(q, l)$  ;
                   if  $last_p = \text{decide}$ 
                     then  $p$  announces itself leader
                   else send  $\langle annex, q, l \rangle$  to  $last_p$ 
                end
            else if  $l = lev_p$  and  $((t = A \text{ and } q > cat_p) \text{ or } t = C)$  then (* Case V *)
                begin send  $\langle chase, q, l \rangle$  to  $last_p$  ;  $last_p := undef$  end
            else if  $l = lev_p$  then (* Case VI *)
                 $wait_p := q$ 
```

KKM

počet správ

- ▶ naháňacie: 1 na vrchol a level, spolu n za level
- ▶ objavovacie: $\sum_i f(n_i)$
- ▶ ak f je konvexná, t.j. $f(a) + f(b) \leq f(a + b)$, tak je $O(\log n(n + f(n)))$ správ

K_n : vplyv orientácie

zmysel pre orientáciu

Porty sú číslované podľa vzdialenosť vo fixnej Hamiltonovej kružnici.

- ▶ algoritmus: zajať fixný pattern $\{i[1..k], i[2k], i[3k], \dots, i[n - k]\}$
- ▶ prvá fáza $i[1..k]$
 - ▶ $level$, ID
 - ▶ $level$ je počet zajatých
 - ▶ pripája sa celé územie
- ▶ druhá fáza $i[2k], i[3k], \dots, i[n - k]$
 - ▶ nastav $owner$ vrcholom $i[1..k]$, ack
 - ▶ pošli $elect$ do $i[2k], i[3k], \dots, i[n - k]$
 - ▶ $elect$: ak je v prvej fáze, alebo je slabší $\Rightarrow accept$

K_n : vplyv orientácie

počet správ

- ▶ prvá fáza: $O(n)$: odpoveď zajatému 1x, dizjunknosť území
- ▶ druhá fáza: max $O(n/k)$ kandidátov, každý pošle n/k správ $\Rightarrow O(n^2/k^2)$ správ

čas

- ▶ po zobudení (najsilnejšieho) $O(k)$, v najhoršom $O(n)$
- ▶ pridať wakeup fázu (poslať od $i[1], i[k]$) $\Rightarrow O(k + n/k)$
- ▶ $k := \sqrt{n}$

vplyv synchrónnosti

algoritmy založené na porovnaniach

- ▶ ekvivalentné okolia
- ▶ c -symetrický reťazec: pre každé $\sqrt{n} \leq l \leq n$ a pre každý segment S je $\left\lfloor \frac{cn}{l} \right\rfloor$ ekvivalentných
- ▶ bit-reversal je $1/2$ -symetrický
- ▶ c -symetrický reťazec, alg. nemôže skončiť po k kolách pre $\left\lfloor \frac{cn}{2k+1} \right\rfloor \geq 2$
- ▶ počet správ: $k = \left\lfloor \frac{cn-2}{4} \right\rfloor$, je aspoň $k + 1$ kôl
- ▶ aspoň $\left\lfloor \frac{cn}{2r-1} \right\rfloor$ aktívnych v r -tom kole pošle správu

vplyv synchrónnosti

algoritmy využívajúce integer ID

- ▶ rôzne rýchlosťi
- ▶ v i -tom kroku test

cvičenie

V toruse rozmerov $n \times n$ (t.j. zacyklenej mriežke) sú na začiatku zobudené dva vrcholy. Napíšte algoritmus, pomocou ktorého sa každý z nich dozvie identifikátor druhého s použitím $O(n)$ správ.

Najdite asymptoticky optimálny (čo do počtu správ) algoritmus na voľbu šéfa v úplnom bipartitnom grafe $K_{n,n}$. Dokážte jeho zložitosť a optimalitu.

★

broadcasting a voľba šéfa na (ne?)orientovanej hyperkocke s lineárnym počtom správ

odolnosť voči chybám – strata správ

Problém dohody

- ▶ synchrónny systém
- ▶ známe identifikátory
- ▶ každý má na vstupe 0/1
- ▶ správy sa môžu strácať
- ▶ každý proces sa musí rozhodnúť
- ▶ treba zaručiť
 - ▶ **Dohoda:** všetky procesy sa rozhodnú na tú istú hodnotu
 - ▶ **Terminácia:** každý proces sa rozhodne v konečnom čase
 - ▶ **Netrivialita:**
 1. Ak všetci začnú s hodnotou 0, musia sa dohodnúť na 0.
 2. Ak všetci začnú s hodnotou 1 a správy sa nestrácajú, musia sa dohodnúť na 1.

neexistuje deterministické riešenie

- ▶ 2 vrcholy, 1 linka
- ▶ sporom, nech existuje a trvá r kôl
- ▶ výpočet, kde začnú obaja s hodnotou 1 a nestrácajú sa správy
- ▶ dohodnú sa na 1
- ▶ stratí sa posledná správa, jeden z nich to nezistí
- ▶ výpočet, kde neprejde ani jedna správa a dohodnú sa na 1
- ▶ jeden z nich dostane na vstup 0
- ▶ aj druhý

randomizované riešenie (úplný graf)

komunikačný pattern

zoznam trojíc (i, j, t) : v čase t sa nestratí správa z $i \mapsto j$

(fixný) adversary = vstup a komunikačný pattern

Dohoda: $\Pr[\text{nejaké dva procesy sa rozhodnú na rôznu hodnotu}] \leq \varepsilon$

algoritmus s $\varepsilon = 1/r$

daný adversary γ : dvojice (i, t) , kde i -procesor, t -čas majme usporiadanie:

1. $(i, t) \leq_\gamma (i, t')$, kde $t \leq t'$
2. ak $(i, j, t) \in \gamma$, tak $(i, t - 1) \leq_\gamma (j, t)$
3. tranzitivita

úroveň informovanosti

1. $\text{level}_\gamma(i, 0) = 0$
2. ak $t > 0$ a existuje $j \neq i$ také, že $(j, 0) \not\leq_\gamma (i, t)$, tak $\text{level}_\gamma(i, t) = 0$
3. nech I_j je $\max\{\text{level}_\gamma(j, t') \mid (j, t') \leq_\gamma (i, t)\}$
potom $\text{level}_\gamma(i, t) = 1 + \min\{I_j \mid j \neq i\}$

algoritmus

- ▶ prvý proces vygeneruje náhodný kľúč
- ▶ procesy si počítajú level
- ▶ rozhodnutie 1, ak všetci majú 1 a môj level je aspoň kľúč

rounds := rounds + 1

let (L_j, V_j, k_j) be the message from j , for each j from which a message arrives
if some $k_j \neq \text{undefined}$ then $key := k_j$

for all $j \neq i$ do

 if some $V_{i'}(j) \neq \text{undefined}$ then $val(j) := V_{i'}(j)$

 if some $L_{i'}(j) > level(j)$ then $level(j) := \max \{L_{i'}(j)\}$

$level(i) := 1 + \min \{level(j) : j \neq i\}$

if $rounds = r$ then

 if $key \neq \text{undefined}$ and $level(i) \geq key$ and $val(j) = 1$ for all j then

$decision := 1$

 else $decision := 0$

dôkaz

- ▶ **Dohoda:** $\Pr[\text{nejaké dva procesy sa rozhodnú na rôznu hodnotu}] \leq \varepsilon$
- ▶ **Terminácia:** každý proces sa rozhodne v konečnom čase
- ▶ **Netrivialita:**
 1. Ak všetci začnú s hodnotou 0, musia sa dohodnúť na 0.
 2. Ak všetci začnú s hodnotou 1 a správy sa nestrácajú, musia sa dohodnúť na 1.

terminácia a netrivialita sú zrejmé
pre fixný pattern, aká je pravdepodobnosť nezhody?
leveley sa líšia max o 1, preto jediný problém je ak $key = \max\{l_i\}$

dolný odhad

ľubovoľný r -kolový algoritmus má pravdepodobnosť nezhody aspoň $\frac{1}{r+1}$

orez

pre adversary B s patternom γ a proces i , $B' = \text{prune}(B, i)$

1. ak $(j, 0) \leq_{\gamma} (i, r)$ tak sa vstup j zachová, inak znuluje
2. trojica (j, j', t) je v kom. patterne B' , akk je v γ a $(j', t) \leq_{\gamma} (i, r)$

$$P^B[i \text{ sa rozhodne } 1] = P^{\text{prune}(B, i)}[i \text{ sa rozhodne } 1]$$

lema

Ak majú na vstupe všetci 1, $P[i \text{ sa rozhodne } 1] \leq \varepsilon(\text{level}(i, r) + 1)$

lema

Ak majú na vstupe všetci 1, $P[i \text{ sa rozhodne } 1] \leq \varepsilon(\text{level}(i, r) + 1)$

- ▶ indukcia na $\text{level}(i, r)$: nech $\text{level}(i, r) = 0$:
- ▶ $B' = \text{prune}(B, i) = \text{prune}(B', i)$
- ▶ $P^B[i \text{ sa rozhodne } 1] = P^{B'}[i \text{ sa rozhodne } 1]$
- ▶ od j -čka neprišla správa, $B'' = \text{prune}(B', j) = \text{prune}(B'', j)$ je triviálny adversary
- ▶ $P^{B'}[j \text{ sa rozhodne } 1] = P^{B''}[j \text{ sa rozhodne } 1]$
- ▶ lenže $P^{B''}[j \text{ sa rozhodne } 1] = 0$, takže $P^{B'}[j \text{ sa rozhodne } 1] = 0$
- ▶ pôť nezhody je ε , $\Rightarrow |P^{B'}[i \text{ sa rozhodne } 1] - P^{B'}[j \text{ sa rozhodne } 1]| \leq \varepsilon$
- ▶ preto $P^{B'}[i \text{ sa rozhodne } 1] \leq \varepsilon$ a $P^B[i \text{ sa rozhodne } 1] \leq \varepsilon$

lema

Ak majú na vstupe všetci 1, $P[i \text{ sa rozhodne } 1] \leq \varepsilon(\text{level}(i, r) + 1)$

- ▶ indukcia na $\text{level}(i, r)$: nech $\text{level}(i, r) > 0$
- ▶ $B' = \text{prune}(B, i) = \text{prune}(B', i)$
- ▶ existuje j , že $\text{level}_{B'}(j, r) \leq l - 1$
- ▶ podľa i.p. $P^{B'}[j \text{ sa rozhodne } 1] \leq \varepsilon(\text{level}(j, r) + 1) \leq \varepsilon l$
- ▶ pôsť nezhody je ε , $\Rightarrow |P^{B'}[i \text{ sa rozhodne } 1] - P^{B'}[j \text{ sa rozhodne } 1]| \leq \varepsilon$
- ▶ preto $P^{B'}[i \text{ sa rozhodne } 1] \leq \varepsilon(l + 1)$ a $P^B[i \text{ sa rozhodne } 1] \leq \varepsilon(l + 1)$

chyby procesov – stop chyby

Problém dohody

- ▶ synchrónny systém
- ▶ známe identifikátory
- ▶ každý má na vstupe 0/1
- ▶ proces môže havarovať (uprostred posielania správ)
- ▶ maximálne f havarovaných procesov
- ▶ každý proces sa musí rozhodnúť
- ▶ treba zaručiť
 - ▶ **Dohoda:** všetky procesy (ktoré nehavarovali) sa rozhodnú na tú istú hodnotu
 - ▶ **Terminácia:** každý proces (ktorý nehavaroval) sa rozhodne v konečnom čase
 - ▶ **Netrivialita:** ak všetci začnú s rovnakou hodnotou i , musia sa rozhodnúť na i .

chyby procesov – stop chyby

algoritmus

flood počas $f + 1$ kôl; ak je iba jedna hodnota, rozhodni sa, inak (default) 0

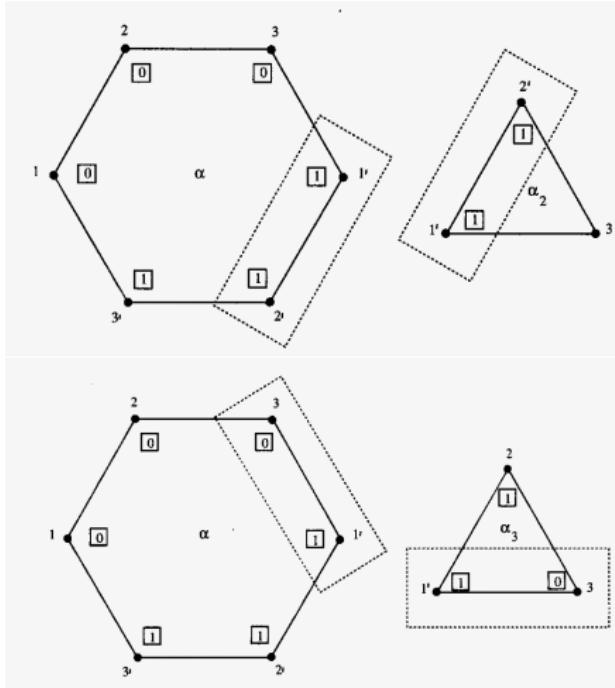
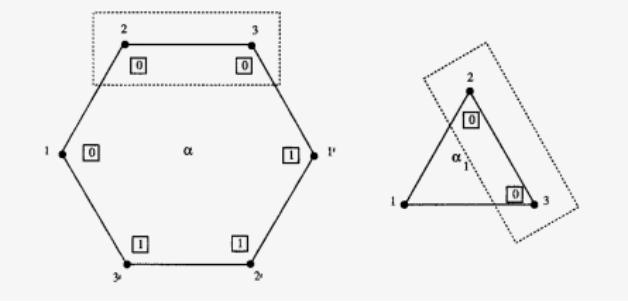
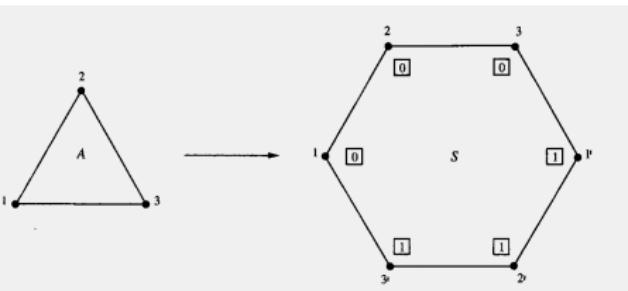
- ▶ existuje kolo, v ktorom nikto nehavaruje; potom sa udržiavajú rovnaké hodnoty
- ▶ $(f + 1)n^2$ správ
- ▶ zlepšenie: posielat' iba keď sa zmení hodnota $\Rightarrow O(n^2)$ správ

chyby procesov – byzantínske chyby

Problém dohody

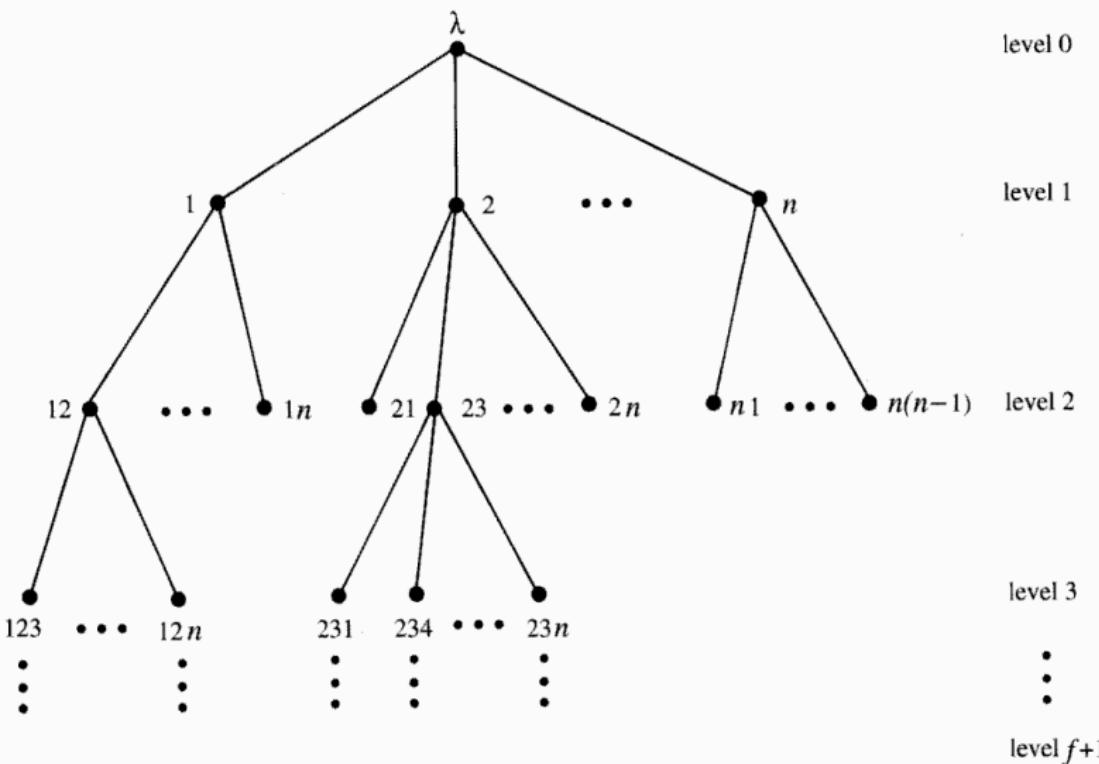
- ▶ synchrónny systém
- ▶ známe identifikátory
- ▶ každý má na vstupe 0/1
- ▶ niektoré procesy sú zlé
- ▶ maximálne f zlých procesov
- ▶ každý proces sa musí rozhodnúť
- ▶ treba zaručiť
 - ▶ **Dohoda:** všetky dobré procesy sa rozhodnú na tú istú hodnotu
 - ▶ **Terminácia:** každý dobrý proces sa rozhodne v konečnom čase
 - ▶ **Netrivialita:** ak všetci dobrí začnú s rovnakou hodnotou i , všetci dobrí sa musia dohodnúť na i .

dolný odhad na počet zlých: jeden zlý spomedzi troch



pre viac: simulácia

EIG algoritmus



$newval(x)$: väčšina z $newval(xj)$

dôkaz

lema

Po $f + 1$ kolách platí: nech $j, j, k, i \neq j$ sú tri dobré procesy. Potom $\text{val}(xk)_i = \text{val}(xk)_j$ pre všetky x .

lema

Po $f + 1$ kolách platí: nech k je dobrý proces. Potom existuje v , že $\text{val}(xk)_i = \text{newval}(xk)_i = v$ pre všetky dobré procesy i

lema

ked' všetci začnú s rovnakou hodnotou, musia sa na nej dohodnúť

dôkaz

vrchol x je spoľahlivý, ak všetky dobré procesy i majú po $f + 1$ koláčach
 $\text{newval}(x)_i = v$ pre nejaké v

lema

Po $f + 1$ koláčach je na každej ceste z koreňa do listu spoľahlivý vrchol

lema

Po $f + 1$ koláčach: ak existuje pokrytie podstromu vo vrchole x dobrími vrcholmi, potom x je dobrý.

polynomiálny počet správ

konzistentný broadcast

- ▶ ak dobrý proces i poslal (m, i, r) v kroku r , dobrí ju akceptujú najneskôr v $r + 1$
- ▶ ak dobrý proces i neposlal (m, i, r) v kroku r , nikto dobrý ju neakceptuje
- ▶ ak je správa (m, i, r) akceptovaná dobrým j v r' , najneskôr v $r' + 1$ ju akc. všetci dobrí

polynomiálny počet správ

algoritmus

- ▶ i pošle $(init, m, i, r)$ v kole r
- ▶ ak dobrý dostane $(init, m, i, r)$ v kole r , pošle $(echo, m, i, r)$ všetkým dobrým v kole $r + 1$
- ▶ ak pred kolom $r' \geq r + 2$ dostane dobrý od $f + 1$ echo, pošle $(echo, m, i, r)$ v r'
- ▶ ak dostal echo od $n - f$, akceptuje

polynomiálny počet správ

dohoda

- ▶ dvojkrokové fázy
- ▶ v prvom kole bcastujú všetci s 1
- ▶ v kole $2s - 1$ pošlú tí, čo akceptovali od $f + s - 1$ a ešte nebcastovali
- ▶ ak po $2(f + 1)$ kolách i akceptoval od $2f + 1$ procesov, tak 1, inak 0

routovanie

cieľ: doručovať srávy medzi ľubovoľnou dvojicou

modely

- ▶ destination based
- ▶ splittable
- ▶ connections (wormhole)
- ▶ buffering
- ▶ selfish

ciele

- ▶ statické váhy (najkratšie cesty)
- ▶ dynamické váhy (hot potato)
- ▶ deadlock

najkratšie cesty

```
begin (* Initialize  $S$  to  $\emptyset$  and  $D$  to  $\emptyset$ -distance *)
     $S := \emptyset$  ;
    forall  $u, v$  do
        if  $u = v$  then  $D[u, v] := 0$ 
        else if  $uv \in E$  then  $D[u, v] := \omega_{uv}$ 
        else  $D[u, v] := \infty$  ;
    (* Expand  $S$  by pivoting *)
    while  $S \neq V$  do
        (* Loop invariant:  $\forall u, v : D[u, v] = d^S(u, v)$  *)
        begin pick  $w$  from  $V \setminus S$  ;
            (* Execute a global  $w$ -pivot *)
            forall  $u \in V$  do
                (* Execute a local  $w$ -pivot at  $u$  *)
                forall  $v \in V$  do
                     $D[u, v] := \min(D[u, v], D[u, w] + D[w, v])$  ;
                 $S := S \cup \{w\}$ 
            end (*  $\forall u, v : D[u, v] = d(u, v)$  *)
        end
    end
```

najkratšie cesty

```
var  $S_u$  : set of nodes ;
       $D_u$  : array of weights ;
       $Nb_u$  : array of nodes ;

begin  $S_u := \emptyset$  ;
       forall  $v \in V$  do
         if  $v = u$ 
           then begin  $D_u[v] := 0$  ;  $Nb_u[v] := undef$  end
         else if  $v \in Neigh_u$ 
           then begin  $D_u[v] := \omega_{uv}$  ;  $Nb_u[v] := v$  end
           else begin  $D_u[v] := \infty$  ;  $Nb_u[v] := undef$  end ;
       while  $S_u \neq V$  do
         begin pick  $w$  from  $V \setminus S_u$  ;
               (* All nodes must pick the same node  $w$  here *)
               if  $u = w$ 
                 then “broadcast the table  $D_w$ ”
                 else “receive the table  $D_w$ ”
               forall  $v \in V$  do
                 if  $D_u[w] + D_w[v] < D_u[v]$  then
                   begin  $D_u[v] := D_u[w] + D_w[v]$  ;
                          $Nb_u[v] := Nb_u[w]$ 
                   end ;
                  $S_u := S_u \cup \{w\}$ 
             end
         end
end
```

najkratšie cesty

```
var  $S_u$  : set of nodes ;
 $D_u$  : array of weights ;
 $Nb_u$  : array of nodes ;

begin  $S_u := \emptyset$  ;
    forall  $v \in V$  do
        if  $v = u$ 
            then begin  $D_u[v] := 0$  ;  $Nb_u[v] := undef$  end
        else if  $v \in Neigh_u$ 
            then begin  $D_u[v] := \omega_{uv}$  ;  $Nb_u[v] := v$  end
        else begin  $D_u[v] := \infty$  ;  $Nb_u[v] := undef$  end ;
    while  $S_u \neq V$  do
        begin pick  $w$  from  $V \setminus S_u$  ;
            (* Construct the tree  $T_w$  *)
            forall  $x \in Neigh_u$  do
                if  $Nb_u[w] = x$  then send  $\langle ys, w \rangle$  to  $x$ 
                    else send  $\langle nys, w \rangle$  to  $x$  ;
                 $num\_rec_u := 0$  ; (*  $u$  must receive  $|Neigh_u|$  messages *)
            while  $num\_rec_u < |Neigh_u|$  do
                begin receive  $\langle ys, w \rangle$  or  $\langle nys, w \rangle$  message ;
                     $num\_rec_u := num\_rec_u + 1$ 
                end ;
            if  $D_u[w] < \infty$  then (* participate in pivot round *)
                begin if  $u \neq w$ 
                    then receive  $\langle dtab, w, D \rangle$  from this  $Nb_u[w]$  ;
                    forall  $x \in Neigh_u$  do
                        if  $\langle ys, w \rangle$  was received from  $x$ 
                            then send  $\langle dtab, w, D \rangle$  to  $x$  ;
                    forall  $v \in V$  do (* local  $w$ -pivot *)
                        if  $D_u[w] + D[v] < D_u[v]$  then
                            begin  $D_u[v] := D_u[w] + D[v]$  ;
                                 $Nb_u[v] := Nb_u[w]$ 
                            end
                end ;
             $S_u := S_u \cup \{w\}$ 
        end
    end
```

Netchange

```

var Neighu : set of nodes ; (* The neighbors of u *)
Du : array of 0.. N ; (* Du[v] estimates d(u, v) *)
Nbu : array of nodes ; (* Nbu[v] is preferred neighbor for v *)
ndisu : array of 0.. N ; (* ndisu[w, v] estimates d(w, v) *)

```

Initialization:

```

begin forall w ∈ Neighu, v ∈ V do ndisu[w, v] := N ;
  forall v ∈ V do
    begin Du[v] := N ; Nbu[v] := undef end ;
    Du[u] := 0 ; Nbu[u] := local ;
    forall w ∈ Neighu do send (mydist, u, 0) to w
  end

```

Procedure *Recompute* (*v*):

```

begin if v = u
  then begin Du[v] := 0 ; Nbu[v] := local end
  else begin (* Estimate distance to v *)
    d := 1 + min{ndisu[w, v] : w ∈ Neighu} ;
    if d < N then
      begin Du[v] := d ;
        Nbu[v] := w with 1 + ndisu[w, v] = d
      end
    else begin Du[v] := N ; Nbu[v] := undef end
  end ;
  if Du[v] has changed then
    forall x ∈ Neighu do send (mydist, v, Du[v]) to x
end

```

Processing a (**mydist**, *v*, *d*) message from neighbor *w*:

```

{ A (mydist, v, d) is at the head of Qwv }
begin receive (mydist, v, d) from w ;
  ndisu[w, v] := d ; Recompute (v)
end

```

Upon failure of channel *uw*:

```

begin receive (fail, w) ; Neighu := Neighu \ {w} ;
  forall v ∈ V do Recompute (v)
end

```

Upon repair of channel *uw*:

```

begin receive (repair, w) ; Neighu := Neighu ∪ {w} ;
  forall v ∈ V do
    begin ndisu[w, v] := N ;
      send (mydist, v, Du[v]) to w
    end
end

```

Netchange

```
var Neighu : set of nodes ; (* The neighbors of u *)
Du : array of 0.. N ; (* Du[v] estimates d(u, v) *)
Nbu : array of nodes ; (* Nbu[v] is preferred neighbor for v *)
ndisu : array of 0.. N ; (* ndisu[w, v] estimates d(w, v) *)
```

Initialization:

```
begin forall w ∈ Neighu, v ∈ V do ndisu[w, v] := N ;
  forall v ∈ V do
    begin Du[v] := N ; Nbu[v] := undef end ;
    Du[u] := 0 ; Nbu[u] := local ;
    forall w ∈ Neighu do send ⟨mydist, u, 0⟩ to w
  end
```

Procedure *Recompute* (v):

```
begin if v = u
  then begin Du[v] := 0 ; Nbu[v] := local end
  else begin (* Estimate distance to v *)
    d := 1 + min{ndisu[w, v] : w ∈ Neighu} ;
    if d < N then
      begin Du[v] := d ;
        Nbu[v] := w with 1 + ndisu[w, v] = d
      end
    else begin Du[v] := N ; Nbu[v] := undef end
  end;
  if Du[v] has changed then
    forall x ∈ Neighu do send ⟨mydist, v, Du[v]⟩ to x
end
```

Processing a ⟨mydist, v, d⟩ message from neighbor w:

```
{ A ⟨mydist, v, d⟩ is at the head of Qwv }
begin receive ⟨mydist, v, d⟩ from w ;
  ndisu[w, v] := d ; Recompute (v)
end
```

Upon failure of channel uw:

```
begin receive ⟨fail, w⟩ ; Neighu := Neighu \ {w} ;
  forall v ∈ V do Recompute (v)
end
```

Upon repair of channel uw:

```
begin receive ⟨repair, w⟩ ; Neighu := Neighu ∪ {w} ;
  forall v ∈ V do
    begin ndisu[w, v] := N ;
      send ⟨mydist, v, Du[v]⟩ to w
    end
end
```

korektnost'

lexikograficky klesá hodnota $[t_0, t_1, \dots, t_N]$

kde t_i je počet správ $\langle \text{mydist}, i \rangle$ + počet dvojic u, v kde $D_u[v] = i$

hierarchické routovanie

cieľ: minimalizovať počet rozhodnutí

veta

Pre siet' s N vrcholmi stačí $O(\sqrt{N})$ rozhodnutí pri použití 3 farieb.

s -klastre:

- ▶ každý je súvislý, pokrývajú všetky vrcholy
 - ▶ každý obsahuje aspoň s vrcholov a má polomer najviac $2s$
- kostra spájajúca centrá klastrov: m listov $\Rightarrow m - 2$ vetvení

hierarchické routovanie

veta

Pre siet s N a pre $f \leq \log N$ stačí $O(f \cdot N^{1/f})$ rozhodnutí a $2f + 1$ farieb

po i klastrovaniach s parametrom s : m_i listov, max. $m_i - 2$ vetvení \Rightarrow

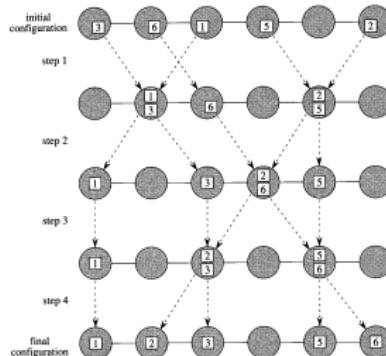
$$m_{i+1} = m_i(2/s)$$

$$m_f + fs \text{ rozhodnutí}$$

$$s \approx 2N^{1/f}$$

packet routing

- ▶ synchrónny režim
- ▶ vrcholy majú pakety (uložené v bufferoch)
- ▶ v jednom kroku po jednej linke ide max. jeden paket
- ▶ algoritmus = odchádzajúce linky + priorita bufferov
- ▶ celkový čas



packet routing na mriežke $\sqrt{N} \times \sqrt{N}$

Každý vrchol má 1 paket, do každého smeruje 1 paket (permutation routing)

Najprv riadok, potom stĺpec. Prednosť má ten s najdlhšou cestou.

analýza: stačí $2\sqrt{N} - 2$ krokov

- ▶ po $\sqrt{N} - 1$ krokoch je každý v správnom stĺpci (nebrzdia sa)
- ▶ routovanie v stĺpci ide v $\sqrt{N} - 1$ krokoch
 - ▶ pre každé i platí: po $N - 1$ krokoch sú koncové pakety na koncových miestach
 - ▶ dôvod: zdržujú sa iba navzájom

veľkosť buffra v najhoršom prípade: $2/3\sqrt{N} - 3$

veľkosť buffra: priemerný prípad I

setting

Každý vrchol má jeden paket s **náhodným cieľom**

max. veľkosť buffra \approx počet zahnutí vo vrchole

psť, že aspoň r zahne $\leq \binom{\sqrt{N}}{r} \left(\frac{1}{\sqrt{N}}\right)^r < \left(\frac{e}{r}\right)^r$

pre $r = \frac{e \log N}{\log \log N}$ je psť $o(N^{-2})$

veľkosť buffra: priemerný prípad II

wide-channel: nepredbiehajú sa

lema

pst', že vo wch prejde aspoň $\alpha\Delta/2$ paketov cez hranu e počas $t+1, t+2, \dots, t+\Delta$ je najviac $e^{(\alpha-1-\alpha \ln \alpha)\Delta/2}$

očakávaný počet paketov na hrane $(i,j) \mapsto (i+1,j)$ je

$$\frac{2i(\sqrt{N} - i)\Delta}{N} \leq \frac{\Delta}{2}$$

chceme ukázať, že s veľkou pstoou ich neprejde príliš viac

Černovov odhad

lema

Majme n nezávislých Bernoulihho náh. prem. X_1, \dots, X_n , pričom
 $\Pr[X_k = 1] \leq P_k$. Potom

$$\Pr[X \geq \beta P] \leq e^{(1 - \frac{1}{\beta} - \ln \beta) \beta P}$$

kde $X = \sum X_i$, $P = \sum P_i$

$$E[e^{\lambda X_k}] \leq 1 + P_k(e^\lambda - 1) \leq e^{P_k(e^\lambda - 1)}$$

$$E[e^{\lambda X}] \leq e^{P(e^\lambda - 1)}$$

$$\Pr[e^{\lambda X} \geq e^{\lambda \beta P}] \leq \frac{E[e^{\lambda X}]}{e^{\lambda \beta P}} \leq e^{P(e^\lambda - 1) - \lambda \beta P}$$

veľkosť buffra: priemerný prípad II

lema

Majme n nezávislých Bernoulihho náh. prem. X_1, \dots, X_n , pričom $Pr[X_k = 1] \leq P_k$. Potom $Pr[X \geq \beta P] \leq e^{(1 - \frac{1}{\beta} - \ln \beta)\beta P}$ kde $X = \sum X_i$, $P = \sum P_i$

lema

pst', že vo wch prejde aspoň $\alpha\Delta/2$ paketov cez hranu e počas $t+1, t+2, \dots, t+\Delta$ je najviac $e^{(\alpha-1-\alpha \ln \alpha)\Delta/2}$

očakávaný počet paketov na hrane $(i, j) \mapsto (i+1, j)$ je

$$\frac{2i(\sqrt{N}-i)\Delta}{N} \leq \frac{\Delta}{2}$$

chceme ukázať, že s veľkou pstoou ich neprejde príliš viac
 $n = 2i\Delta$, $P_k = \frac{\sqrt{N}-i}{N}$, $P = \frac{2i(\sqrt{N}-i)\Delta}{N}$, $\beta = \frac{\alpha N}{4i(\sqrt{N}-i)}$

veľkosť buffra: priemerný prípad II

lema

ak je paket vo vzd. d od hrany e v čase T , a p prejde cez e v čase $T + d + \delta$, potom v každom kroku $[T + d, T + d + \delta]$ prejde paket cez e

lema

ak počas $[T + 1, T + \Delta]$ prejde cez e x paketov v št., tak pre nejaké t prejde $x + t$ paketov cez e v čase $[T + 1 - t, T + \Delta]$ vo wch.

lema

psť, že cez e prejde viac ako $\alpha\Delta/2$ paketov počas konkrétneho okna Δ krokov je najviac $O(e^{(\alpha-1-\alpha \ln \alpha)\Delta/2})$

dosledok

s pravdopođenosťou $1 - O(\frac{1}{N})$ neprejde po e viac ako $c \log N$ paketov v posebe idúcich krokoch, kde $c = \frac{5 \ln 2}{2 \ln 2 - 1} < 9$.

dynamické routovanie

model

V každom kroku sa v každom vrchole s pravdepodobnosťou λ narodí paket s náhodným cieľom.

stabilita

Pre $\lambda \geq 4/\sqrt{N}$ je systém nestabilný

veta

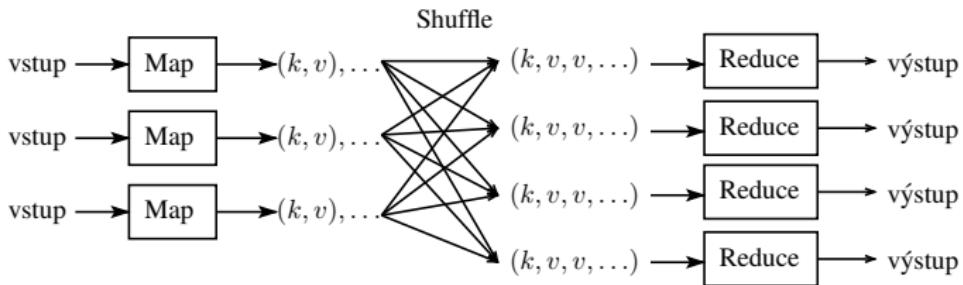
Ak je $\lambda \leq 0.99\frac{4}{\sqrt{N}}$, tak pravdepodobenosť zdržania konkrétneho paketu o Δ krokov je $e^{-O(\Delta)}$.

W.h.p. stačí buffer $O(1 + \frac{\log T}{\log N})$.

Apache Hadoop (Google MapReduce)

- ▶ odolnosť voči chybám, synchronizácia, scheduling
- ▶ menej flexibility: komunikácia Map → Shuffle → Reduce

- ▶ Map: vstup $\rightarrow (k_1, v_1), (k_2, v_2), \dots$
- ▶ Shuffle: $\rightarrow (k_1, v, v, v, \dots), (k_2, v, v, v, \dots)$
- ▶ Reduce: $(k, v_1, v_2, \dots) \rightarrow$ výstup



Príklad

	Map	Shuffle	Reduce
škrtia sa krty	\rightarrow (škrtia, 1), (sa, 1), (krty, 1)	... (škrtia, 1) ... (sa, 1)	\rightarrow škrtia: 1 \rightarrow sa: 1
krt krta škrtí	\rightarrow (krt, 1), (krta, 1), (škrtí, 1)	... (krt, 1, 1) ... (krta, 1, 1)	\rightarrow krt: 1 \rightarrow krta: 1 \rightarrow škrtí: 1
ked' krt krta zaškrtí	\rightarrow (ked', 1), (krt, 1), (krta, 1), (zaškrtí, 1)	... (škrtí, 1)	\rightarrow krt: 2 \rightarrow krta: 2 \rightarrow zaškrtí: 1
do rána ho zmaškrtí	\rightarrow (do, 1), (rána, 1), (ho, 1), (zmaškrtí, 1)	... (ked', 1) ... (zaškrtí, 1) ... (do, 1) ... (rána, 1) ... (ho, 1) ... (zmaškrtí, 1)	\rightarrow ked': 1 \rightarrow zaškrtí: 1 \rightarrow do: 1 \rightarrow rána: 1 \rightarrow ho: 1 \rightarrow zmaškrtí: 1

```
Map(String input_line):
```

```
    for each word w in input_line: Emit(w, 1);
```

```
Reduce(String key, Iterator values):
```

```
    int result = 0;  
    for each v in values: result += v;  
    Emit(key + ": " + result);
```

Ked' treba viacero map-shuffle-reduce fáz

Zoznam slov s > 1000 výskytmi

- ▶ 1. MR: slovo → počet výskytov
- ▶ 2. MR: filter

- ▶ Správa pipeline: skript / iný nástroj
- ▶ Ošetrenie chýb (reštart)
- ▶ Komplikovaná optimalizácia
- ▶ ⇒ MR je príliš nízkoúrovňová

Flume: Abstrakcia nad MR

- ▶ Knižnica: abstrakcia pre paralelizovanú kolekciu dát
- ▶ **PCollection<typ>**:
veľmi veľká, distribuovaná kolekcia
obmedzený spôsob manipulácie
- ▶ **PTable<typ_kľúča, typ_hodnoty>** =
PCollection< KV<typ_kľúča, typ_hodnoty> >

Premenná typu PCollection:

- ▶ nedá sa meniť (immutable)
- ▶ nedá sa k nej priamo pristupovať (not materialized)
- ▶ malá interná reprezentácia

Operácie nad kolekciami

ParDo

- ▶ PCollection $<t_1>$ vstup
- ▶ PCollection $<t_2>$ výstup =
vstup.ParDo(new Funkcia())

```
class Funkcia extends DoFn< $t_1, t_2>$  {  
    public void Do( $t_1$  vstup, EmitFn< $t_2>$  emit) {  
         $t_2$  výstup = ...  
        emit.Emit(výstup)  
    }  
}
```

- ▶ Funkcia beží paralelne \Rightarrow nesmie komunikovať stavom
- ▶ Funkcia nesmie mať vedľajšie efekty

Operácie nad kolekciami

GroupByKey

- ▶ PTable< t_1, t_2 > vstup
- ▶ PTable< t_1 , list< t_2 >> výstup =
vstup.GroupByKey()

Flatten

- ▶ PCollection< t > vstup1, ..., vstupK
- ▶ PCollection< t > výstup =
Flatten(vstup1, ..., vstupK)

Operácie nad kolekciami

CombineValues

- ▶ PTable $<t_1, t_2>$ vstup
- ▶ PTable $<t_1, t_2>$ výstup =
vstup.CombineValues(new Funkcia())

```
class Funkcia extends CombineFn $<t_1, t_2>$  {  
    public  $t_2$  Combine( $t_1$  klúč, list $<t_2>$  hodnoty) {  
         $t_2$  výstup = ...  
        return výstup  
    }  
}
```

- ▶ Ako GroupByKey + ParDo
- ▶ Asociativita; nemusí vidieť všetky hodnoty naraz

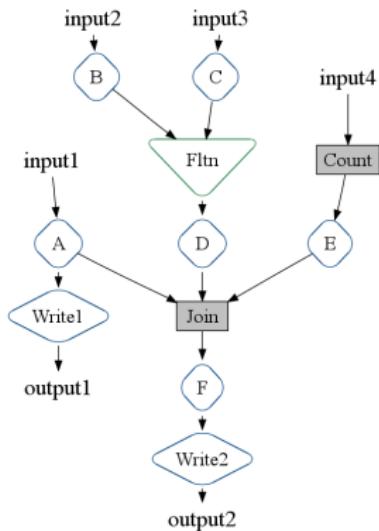
Ako prebieha výpočet

- ▶ `ReadFromFile()`, `ParDo()`,
`GroupByKey()`, ...,
`WriteToFile()`
- ▶ `Flume::Run()`

Ako prebieha výpočet

- ▶ `ReadFromFile()`, `ParDo()`,
`GroupByKey()`, ... ,
`WriteToFile()`
- ▶ `Flume::Run()`

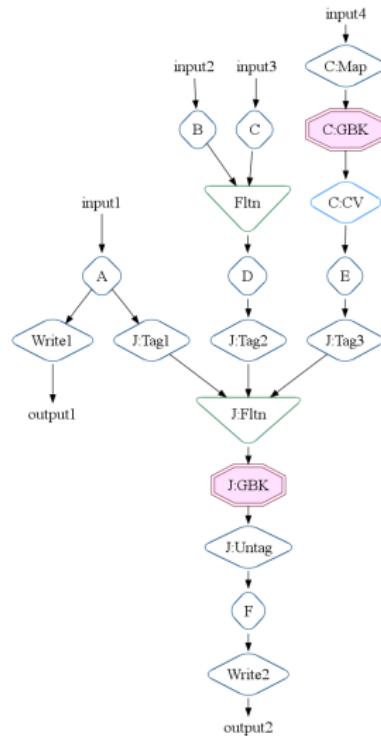
- ▶ Strom operácií
- ▶ Lazy evaluation



Ako prebieha výpočet

- ▶ `ReadFromFile()`, `ParDo()`,
`GroupByKey()`, ... ,
`WriteToFile()`
- ▶ `Flume::Run()`

- ▶ Strom operácií
- ▶ Lazy evaluation



\mathcal{MRC} teória

\mathcal{MRC}^i

- ▶ algoritmus je postupnosť $\langle \mu_1, \rho_1, \dots, \mu_R, \rho_R \rangle$
- ▶ μ_r je (randomizovaný) RAM s $O(n^{1-\varepsilon})$ pamäťou a poly časom
- ▶ ρ_r rovnako
- ▶ pamäť $\sum_{\langle k; v \rangle} (|k| + |v|)$ z každého reducera je $O(n^{2-2\varepsilon})$
- ▶ počet kôl je $O(\log^i n)$

prefix sums

- ▶ [1, 4, 0, 0, 3, 0, 2]:
in: $\{(1, 1), (2, 4), (5, 3), (7, 2)\}$,
out: $\{(1, 1, 1), (2, 4, 5), (5, 3, 8), (7, 2, 10)\}$
- ▶ blok m^ε