

# EXTERNÍ TŘÍDĚNÍ

Třídění na páskách (50. a 60. léta):

sekvenční přístup k datům

algoritmy na principu Mergesortu (většinou)

Třídění na discích:

přímý přístup k blokům dat

Míra časové složitosti:

vstupní a výstupní operace (čtení a zápis dat)

operace v interní paměti jsou zanedbatelné

# TRÍDĚNÍ NA PAŠKÁCH - MERGESORT

Přímý MERGESORT - dvoufázové mergesort

používá 3 pásky T1, T2, T3

nenáročný na vnitřní paměť (stačí 2 prvky)

Příklad:

0) T1: 44, 55, 12, 42, 94, 18, 6, 67 vstup

T2:

T3:

1a) T1: 44, 55, 12, 42 vstup se rozdělí

T2: 94, 18, 6, 67

T3:

1b) T1:

T2: *mergují se dvojice*

T3: 44, 94 | 18, 55 | 6, 12 | 42, 67

2a) T1: 44, 94 | 18, 55

T2: 6, 12 | 42, 67

*rozdělí se*

T3:

2b) T1:

T2: *mergují se čtveřice*

T3: 6, 12, 44, 94 | 18, 42, 55, 67

3a) T1: 6, 12, 44, 94

T2: 18, 42, 55, 67

T3:

rozděli se

3b) T1:

T2:

T3: 6, 12, 18, 42, 44, 55, 67, 94

Vypočet má 3 stupně (obecně  $\log n$ )

V každém stupni 2 fáze : a) rozdělení (kopírování)

b) mergesort

V každé fázi pro každý prvek 1 čtení

1 zápis

Složitost:  $O(n \log n)$

Fáze a) je z hlediska třídění neproduktivní

(nepřeuspořádává prvky)

spotřebává polovinu operací

## Jednotlivé' mergování'

používá 4 pásky

výsledek mergování' se hned rozděluje na 2 pásky

Vývožné' mergování' - na každé' pásce stejný' počet běhů'

Příklad:

0a) T1: 44, 55, 12, 42, 94, 18, 6, 67

T2:

T3:

T4:

0b): T1: 44, 55, 12, 42

T2: 94, 18, 6, 67

rozdělení'

T3:

T4:

1) T1:

T2:

mergování' + rozdělení'

T3: 44, 94 | 6, 12

T4: 18, 55 | 42, 67

2) T1: 18, 44, 55, 94

T2: 6, 12, 42, 67

T3:

T4:

3) T1:

T2:

T3: 6, 12, 18, 42, 44, 55, 67, 94

T4:

### Přirozený MERGESORT

pro soubory, které obsahují delší uspořádané úseky  
(běhy)

mergují se dvojice běhů

rozdělování se provádí podle počtu běhů, ne podle  
počtu prvků

#### Příklad:

0a) T1: 44, 55, 12, 42, 94, 18, 6, 67

T2:

T3:

T4:

0b) T1: 44, 55 | 18

T2: 12, 42, 94 | 6, 67

T3:

T4:

rozdělení

1) T1:

T2:

mergování + rozdělení

T3: 12, 42, 44, 55, 94

T4: 6, 18, 67

2) T1: 6, 12, 18, 42, 44, 55, 67, 94

T2:

T3:

T4:

Složitost:  $O(n \log r)$ kde  $r$  je počet běhů ve vstupuMergesort s využitím interní paměti: $m$  = kapacita interní paměti1) vytvoření běhů délky  $m$ :ze vstupní pařky se načte  $m$  prvků, setřídí se,

vypíší se na pařku T1

načte se dalších  $m$  prvků, setřídí se, vypíší se na T2

— " —

— " —

— " — T1

2) mergování běhů

Složitost: krok 1)  $O(n)$

krok 2) -  $\log \frac{n}{m}$  průchodů

celkem  $O(n \log \frac{n}{m})$

Technické vylepšení:

kyž se dá pásek ošit oběma směry, střídá se  
vzestupně a sestupně mergesování

ušetří se převijení pásek na začátek po každém  
průchodu

Příklad:

ma'm T1: 44, 94 | 6, 12

T2: 18, 55 | 42, 67

T3:

T4:

jsem na konci pásek T1, T2

jdou odzadu

mergují sestupně

na T3 a T4 zapisují odpředu

T1:

T2:

T3: 67, 42, 12, 6

T4: 94, 55, 44, 18

jsem na konci T3, T4

jdou odzadu

mengují vzestupně

na T1 a T2 zapisují odpředu

T1: 6, 12, 18, 42, 44, 55, 67, 94

T2:

T3:

T4:

Poznámka - může se stát, že výsledek bude  
setříděný obráceně

(pak se ještě jednou přehopíruje)



Zobecnění - vícecestné mergesort:

maíme  $2N$  pásek,  $r = \frac{n}{m}$  běhů

běhy se rozdělí na  $N$  pásek

mergeje se současně  $N$  běhů ze vstupních pásek

$x$  mergesortované běhy se rozdělují na  $N$  výstupních pásek

vyvážené mergesort - na každé pásece stejný počet  
běhů

Složitost:  $O(n \log_N (\frac{n}{m}))$

nevýhoda - velké množství pásek

Vicetázové mergesort

maíme  $N$  vstupních pásek, 1 výstupní

běhy se na pásy rozdělí nevyváženě

Příklad:

1)  $T_1$ : 13 běhů

$T_2$ : 8 běhů

$x$  mergeji se  $T_1, T_2$  na  $T_3$

$T_3$ : 0

2) T1: 5 bēhu°

T2: 0

T3: 8 bēhu°

zmerguji' se T1, T3 na T2

3) T1: 0

T2: 5 bēhu°

T3: 3 bēhy

zmerguji' se T2, T3 na T1

4) T1: 3 bēhy

T2: 2 bēhy

T3: 0

atd.

5) T1: 1 bēh

T2: 0

T3: 2 bēhy

6) T1: 0

T2: 1 bēh

T3: 1 bēh

7) T1: 1 bēh

T2: 0

T3: 0

konec

Příklad:

T1	T2	T3	T4	T5	T6	pošky
16	15	14	12	8	0	počty běhů
8	7	6	4	0	8	
4	3	2	0	4	4	
2	1	0	2	2	2	
1	0	1	1	1	1	
0	1	0	0	0	0	konec

Podstatné je počáteční rozdělení běhů na pošky

1. příklad - jsou to dvě po sobě jdoucí Fibonacciho čísla

obecně - Fibonacciho čísla odpovídají jeho řádky :

$$F_{i+1} = F_i + F_{i-1} + \dots + F_{i-k}$$

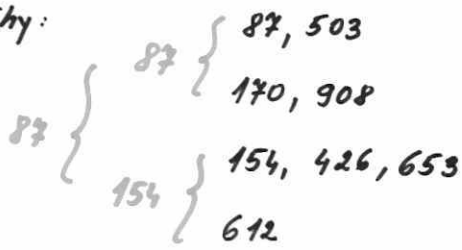
Když se počáteční počet běhů nedá takto rozdělit, doplň se hypotetickými prázdnými běhy.

Současné mergování více běhů

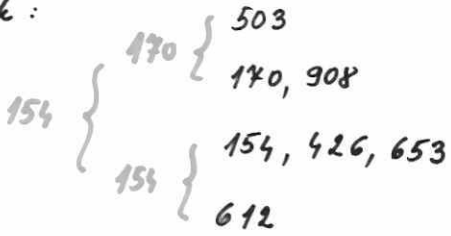
výběr minima probíhá v interní paměti;  
použije turnajový strom nebo haldy

Turnajový strom:

ma'm běhy:

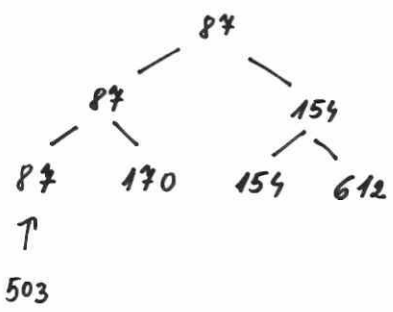


dalsi' krok:



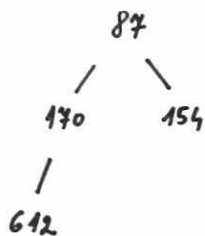
atd.

neboli:



atd.

Min - halda :



provedu DELETEMIN (87)  
INSERT (503)  
DELETEMIN  
atd.

### Počáteční vytvoření běhu

- 1) načíst do interní paměti maximální množství prvků  
interně je seřadit  
vznikají běhy stejné délky  $m$
- 2) metoda Replacement selection  
mohou vznikat běhy delší než  $m$ , v průměru  $2m$   
princip: jakmile zapíšu první prvek z interní  
paměti na pásku, uvolní se tím místo  
a může se hned načíst další

## Postup:

načtu  $m$  prvků do paměti, vytvořím min-haldu  
provedu DELETĚMIN

minimum zapíšu na výstupní pásku

načtu další prvek ze vstupní pásky -

je-li větší než odebrané minimum, přijde do  
těže haldy (a těžože bůhu)

je-li menší, dáím ho stranou

provedu další DELETĚMIN

atd.

až vyprázdním celou haldu, postavím novou

s prvky, které jsem odložila

(můžu ji stavět průběžně)

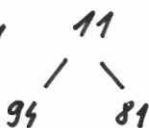
## Příklad:

tridím 11, 94, 81, 96, 12, 35, 17, 99, 28, 58, 41, 75

$m = 3$

1) načtu 11, 94, 81

vytvořím haldu



2) vypišu 11, načtu 96

vytroším haldu 81  
/ \  
94 96

3) vypišu 81, načtu 12

halda: 94 12 tam nepatří'  
/  
96

4) vypišu 94, načtu 35

1. halda: 96 2. halda: 12  
/  
35

5) vypišu 96, načtu 17

hotový běh: 11, 81, 94, 96

nová halda: 12  
/ \  
35 17

6) vypišu 12 - začátek nového běhu  
načtu 99

halda: 17  
/ \  
35 99

atd.

výsledek - běhy : 11, 81, 94, 96

12, 17, 28, 35, 41, 58, 99

75



# PŘIHRÁDKOVÉ TRĚDĚNÍ NA PÁSKAČN

Příklad:

trídím 0, 1, 2, 3, 7, 6, 5, 4

tj. v binárním zápise:

000, 001, 010, 011, 111, 110, 101, 100

použiji se 4 pásky - 2 pro vstup

2 pro výstup

0) T1: 000, 001, 010, 011, 111, 110, 101, 100 vstup

T2:

T3:

T4:

1) T1:

T2:

*rozdělení podle poslední cifry*

T3: 000, 010, 110, 100

T4: 001, 011, 111, 101

2) T1: 000, 100, 001, 101

T2: 010, 110, 011, 111

T3: *rozdělení podle prostřední cifry*

T4:

*v pořadí T3, T4*

- 3) T1: rozdělení podle 1. cifry  
 T2:  
 T3: 000, 001, 010, 011  
 T4: 100, 101, 110, 111

- 4) T1: 000, 001, 010, 011, 100, 101, 110, 111  
 T2:  
 T3: spojení T3, T4  
 T4:

Složitost:  $k$  průchodů pro čísla v rozsahu  
 $0, \dots, 2^k - 1$

V každém průchodu se každé číslo 1x kopíruje  
 celkem:  $O(k \cdot n)$

## QUICKSORT NA 2 PÁŠKÁCH

mám pásky 0 a 1 - fungují jako zásobníky  
jsou obousměrné

vstupní soubor je na páse 0

(vejde-li se do vnitřní paměti, setřídí se interně  
a vrátí na pásku)

pivot = poslední prvek souboru =  $k$

pásku 0 se přeodradu -

prvky  $> k$  se kopírují na pásku 1

pásku 0 se přeodpředu -

prvky  $= k$  se kopírují na pásku 1

pásku 0 se přeodradu -

prvky  $< k$  se kopírují na pásku 1

výsledek:  $T_1$ : prvky  $> k$ , prvky  $= k$ , prvky  $< k$

$T_0$ : prázdná

rekurze:

zpracují se prvky  $< k$  a setříděné se zapíší na  $T_0$

prvky  $= k$  se překopírují na  $T_0$

zpracují se prvky  $> k$  a setříděné se zapíší na  $T_0$

Poznámka: Při rozdělování se střídají posky, ale  
také nerovnosti -

tj. když rozdělají soubor (část) z T1 na T0,  
zapisují ho na T0 v pořadí:

prvky  $< k$ , prvky  $= k$ , prvky  $> k$

Složitost: nejhorší  $O(n^2)$

očekávána  $O(n \log n)$

Příklad:

T0: 11, 94, 81, 96, 12, 35, 17, 99, 28, 58, 41, 75

T1:

rozdělím podle 75

T0:

T1: 99, 96, 81, 94, 75, 41, 58, 28, 17, 35, 12, 11

$\underbrace{\hspace{10em}}_{> 75}$                        $\underbrace{\hspace{10em}}_{< 75}$

zpracuji prvky  $< 75$ , rozdělím podle 11

T0: 11, 12, 35, 17, 28, 58, 41

T1: 99, 96, 81, 94, 75

zpracuji prvky  $> 11$ , rozdělím podle 41

T0: 11

T1:  $\underbrace{99, 96, 81, 94, 45, 58, 41, 28, 17, 35, 12}_{>41}$

prvky  $< 41$  rozdělím podle 12

T0: 11, 12,  $\underbrace{35, 17, 28}_{>12}$

T1:  $\underbrace{99, 96, 81, 94, 45, 58, 41}$

prvky  $> 12$  rozdělím podle 28

T0: 11, 12

T1:  $\underbrace{99, 96, 81, 94, 45, 58, 41, 35, 28, 17}_{>28}$

zkopiruji 12 na T1

zpracuji prvky  $< 12$  (žádné nejsou)

zkopiruji 11 na T1

zpracuji prvky  $< 11$  (žádné nejsou)

T0:

T1:  $\underbrace{99, 96, 81, 94, 45, 58, 41, 35, 28, 17, 12, 11}_{>45}$

kopiruji na T0: prvky  $< 11$

$$= 11$$

$$< 12$$

$$= 12$$

$$< 28$$

$$= 28$$

$$< 41$$

$$= 41$$

$$< 75$$

$$= 75$$

T0: 11, 12, 17, 28, 35, 41, 58, 75

T1: 99, 96, 81, 94

dotřídím prvky  $> 75$ , rozdělím podle 94

T0: 11, ..., 75, 81, 94, 96, 99  
                                    $\underbrace{\hspace{2em}}$                    $\underbrace{\hspace{2em}}$   
                                    $< 94$                                    $> 94$

T1:

zpracuji prvky  $> 94$ , rozdělím podle 99

T0: 11, ..., 75, 81, 94

T1: 99, 96  
                    $\underbrace{\hspace{1em}}$   
                    $< 99$

zkopiruji 94 na T1

zpracuji prvky  $< 94$  (tj. 81)

T0: 11, ..., 75

T1: 99, 96, 94, 81

kopiruji na T0: prvky  $< 94$

= 94

$< 99$

= 99

$> 99$  (žádné nejsou)

výsledek:

T0: 11, 12, 17, 28, 35, 41, 58, 75, 81, 94, 96, 99

T1:

# TŘÍDĚNÍ NA DISKÁCH

na páskách: prvky se přenášejí po jednom

na disku: -11- po blocích

značení:  $N$  počet prvků v souboru

$M$  počet bloků, které se vejdou do interní  
paměti

$B$  počet prvků v jednom bloku

$P$  počet bloků, které se dají přenášet  
současně (pro nás  $P=1$ )

platí:  $1 \leq B \leq M < N$ ,  $1 \leq P \leq \lfloor \frac{M}{B} \rfloor$

složitost (počet I/O operací):

$$\Omega\left(\frac{N}{PB} \cdot \frac{\log\left(1 + \frac{N}{B}\right)}{\log\left(1 + \frac{M}{B}\right)}\right)$$

v nejhorším i průměrném případě

Aggarwal, Vitter (1988)



# MERCESORT

$\frac{N}{B}$  bloků se postupně načte do interní paměti

(vejde se jich  $\frac{M}{B}$ )

obsah interní paměti se setřídí

vytvorí se  $\frac{N}{M}$  běhů

běhy se po blocích uloží zpět na disk

složítost:  $\frac{N}{B}$  čtení,  $\frac{N}{B}$  zápisů

další fáze: vždy  $\frac{M}{B}$  běhů se mergesuje do jednoho

(z každého běhu se načte 1 blok, když se při mergesování vyčerpá, načte se další blok z téhož běhu)

počet fází:  $\log_{\frac{M}{B}} \frac{N}{B}$

v každé se přečte a zapíše  $\frac{N}{B}$  bloků

celkem:  $2 \frac{N}{B} + 2 \frac{N}{B} \log_{\frac{M}{B}} \frac{N}{B} = O\left(\frac{N}{B} \cdot \frac{\log \frac{N}{B}}{\log \frac{M}{B}}\right)$

vylepsení použitím Replacement selection

nevýhoda - netřídí na místě

# EXQUISIT (Six, Wegner 1984)

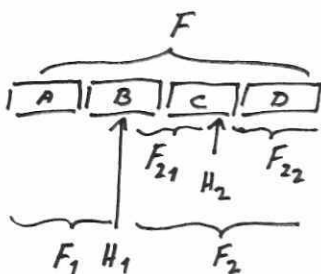
externí varianta Quicksortu  
stačí  $M=2$

Algoritmus:

- 1) načte se první a poslední blok vstupního souboru  $F$   
z jejich prvků se vypočte pivot  $H$
- 2)  $F$  se rozdělí na část  $F_1$  (prvky menší než  $H$ ) a  $F_2$  (prvky větší než  $H$ )  
pointer  $L$  jde od 1 do  $B$  v levém bloku  
-||-  $R$  -||-  $B$  do 1 v pravém -||-  
prvky se porovnávají s  $H$  a případně se prohází  
když pointer dosáhne hranice bloku, blok se  
zapiše na disk a načte se sousední
- 3) když se pointery setkají, třídí se rekurzivně  
 $F_1$  a  $F_2$

Príklad:

disk:



interní paměť:

[A] [D]

začátek

[B] [D]

L přešel hranici mezi A a B

[B] [C]

R — " — D a C

[B] [A]

L a R se setkali v B

dotřídňuje se  $F_1$  (dokončí se A)

[B] [D]

začíná se třídít  $F_2$

[C] [D]

L přešel z B do C

[C] [B]

L a R se setkali v C

dotřídňuje se  $F_{21}$  (dokončí se B)

[C] [D]

dotřídňuje se  $F_{22}$  (dokončí se C, D)

## Složitost

průměrný případ - řeší se rekurze

$$T(N) = \lfloor \frac{N}{B} \rfloor + 1 + \frac{1}{N-1} \sum_{k=1}^{N-1} T(k) + T(N-k), \quad N > B+1$$

$$T(N) = 0 \text{ pro } N \leq B+1$$

řešení:

$$T(N) \leq 2 \frac{N}{B} \log\left(\frac{N}{B+1}\right) + \frac{3N}{B} + 1 - \frac{N}{B(B+1)}$$

nejhorší případ:

$$\text{není } O\left(\left(\frac{N}{B}\right)^2\right), \text{ ale } O\left(\frac{N^2}{B}\right)$$

tj. když se příkrodeím rozdělení oddělí 1 prvek

volba pivotu

je ještě významnější než v interním Quicksortu  
možnosti je méně (počítá se jen ze 2 bloků)

doporučení:  $H = \lfloor \frac{F(L) + F(R)}{2} \rfloor$  fiktivní pivot

výhoda - není třeba ho po rozdělení umístit  
na správnou pozici

## Poznámky:

- 1) Blok se neukládá na disk, dokud není nutné uvolnit místo v interní paměti
- 2) Blok se nepřepisuje (na disk), pokud v něm neproběhla žádná výměna  
(tj. např. u setříděného souboru není žádný zápis na disk, pouze čtení)
- 3) Pořadí volání rekurze - třídi se dřív levá část
- 4) Rekurze se zastaví, když se příslušná část souboru vejde do interní paměti
- 5) Je-li  $M > 2$ , může se pracovat s logickými bloky místo fyzických
- 6) Jiný způsob rozdělování v interní paměti - pomocí dvoukoncevých haldy  
(operace MIN a MAX v konstantním čase, ostatní v logaritmickeém)

# HEAPSORT

Wiedermann (1983)

Wiedermann, Šimo, Neruda (1997)

poslívá se  $k$ -regulární haldá blokú

a) každý uzel  $v$  obsahuje  $M$  záznamú

(logický blok  $S_v$  velikosti  $M$ )

b)  $k \approx \frac{M}{B}$  (tj. počet fyzických blokú, které se  
vejdou do interní paměti)

c) haldová podmínka :

když  $v$  je synem  $u$  v haldě  $T$ , pak pro všechna

$x \in S_v, y \in S_u$  platí  $x \leq y$

záznamy v uzlech jsou vzestupně setříděné

Da se reprezentovat polem.

Budování haldy

soubor se rozdělí na logické bloky velikosti  $M$   
každý se načte do interní paměti

tam se setřídí

zapiše zpátky na disk

úprava haldy - podle Floydů (shora dolů):

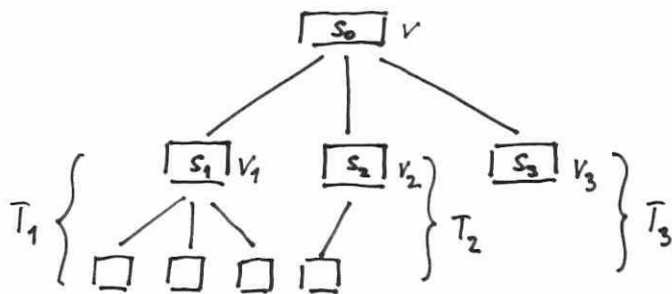
$q$  = počet vrcholů v haldě

$$i := \frac{q-1}{k}$$

while  $i \geq 1$  do restore( $i, k$ );  $i := i-1$  enddo

restore( $i, k$ ) - obnoví haldovou podmínku v haldě  
s kořenem  $i$

Příklad: máme  $T$



$T_1, T_2, T_3$  už jsou správně vytvořené 3-regularní haldy

Není-li splněna haldova podmínka pro koře  $v$ , pak:

1) pro  $i=0, 1, \dots, k$  označím  $s_i = \min \{S_i\}$

postupnost  $s_0, s_1, \dots, s_k$  setřídím vzestupně

tj. najdu permutaci  $\pi$ , že  $s_{\pi(0)} \leq s_{\pi(1)} \leq \dots \leq s_{\pi(k)}$

(tomu odpovídá pořadí hald  $T_{\pi(0)}, T_{\pi(1)}, \dots, T_{\pi(k)}$ )

2) vytvořím  $U = S_0 \cup S_1 \cup \dots \cup S_k$

setřídím  $U$  vzestupně pomocí  $(k+1)$ -cestného mergesortu - výsledek bude  $\Sigma$

3)  $\Sigma$  rozdělím na úseky  $G_0, G_1, \dots, G_k$  délky  $M$



- 4) protože  $s_{\pi(i)} \leq t_i = \min \{G_i\}$  pro všechna  $i$ , můžeme přiřadit:  $G_i$  kořeni  $T_{\pi(i)}$  pro  $i=1, \dots, k$   
 $G_0$  kořeni  $T$

haldova' podmínka bude splněna v  $T_1, \dots, T_k$

- 5) není-li haldova' podmínka splněna v kořeni  $T$ ,  
 vyměníme  $G_0$  a  $G_k$  (je v kořeni  $T_{\pi(k)}$ )  
 rekurzivně obnovíme haldovou podmínku v  $T_{\pi(k)}$

Složitost:

při vhodné implementaci  $O\left(\frac{N}{B}\right)$  pro každé  $k \geq 2$

Heapsort

- vybudovat haldu
- opakovat DELETMAX

Složitost Heapsortu:

$$O\left(k \cdot \frac{N}{B} \cdot \frac{\log \frac{N}{B}}{\log k}\right) \text{ pro každé } k \geq 2$$

Pro  $k \geq 2$  se fáze b) provede jinak a tím zmizí multiplikační konstanta  $k$ .

Příklad: třídím postupně

21, 1, 2, 38, 53, 50, 57, 59, 55, 52, 11, 17, 8, 10, 63, 66, 7, 37,  
13, 49, 61, 62, 70, 68

$N = 24$

necht  $B = 1, M = 4$ , použijí binární haldy

- 1) setřídím úseky délky 4
- 2) naházím je do haldy



3) upravím haldy

a) prohodím uzly  $n_2$  a  $n_5$

b) uzel  $n_1$  nelze prohodit s žádným z jeho synů

najdu  $s_1 = \min \{n_1\} = 50$

$s_2 = \min \{n_3\} = 8$

$s_3 = \min \{n_4\} = 7$

platí  $s_3 < s_2 < s_1$

vytvorím  $U = \{n_1\} \cup \{n_3\} \cup \{n_4\}$

seřídím  $U$  pomocí 3-cestného mergeování:

$\underbrace{7, 8, 10, 13}_{G_0}, \underbrace{37, 49, 50, 53}_{G_1}, \underbrace{57, 59, 63, 66}_{G_2}$

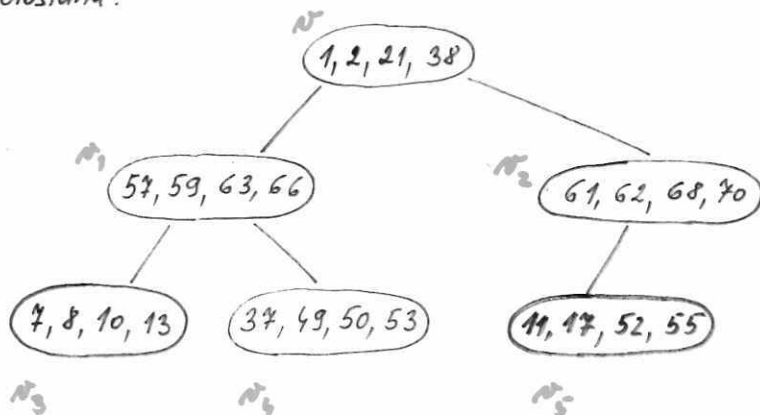
$G_0$  přijde do  $n_1$

$G_1$  —||—  $n_3$

$G_2$  —||—  $n_3$

prohodím uzly  $n_1$  a  $n_3$

dostanu:



c)  $n$  nelze prohodit s  $n_1$  ani s  $n_3$

$$\text{najdu } S_1 = \min \{n\} = 1$$

$$S_2 = \min \{n_1\} = 57$$

$$S_3 = \min \{n_2\} = 61$$

plati'  $s_1 < s_2 < s_3$

vytvorim  $U = \{n\} \cup \{n_1\} \cup \{n_2\}$ , setridim:

$$1, 2, 21, 38, 57, 59, 61, 62, 63, 66, 68, 70$$

$\underbrace{\hspace{10em}}_{G_0} \quad \underbrace{\hspace{10em}}_{G_1} \quad \underbrace{\hspace{10em}}_{G_2}$

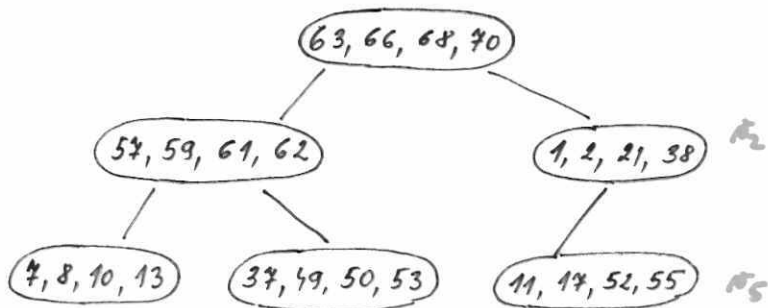
$G_0$  prijde do  $n$

$G_1$  -"-  $n_1$

$G_2$  -"-  $n_2$

prohodim  $n$  a  $n_2$

cestanu:



d) upraví se halda s kořenem  $n_2$

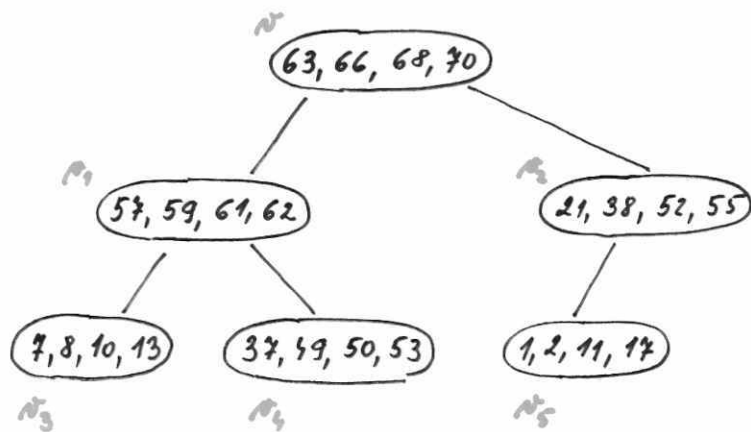
setřídím  $U = \{n_2\} \cup \{n_5\}$

rozdělím na 2 části

menší prvky přijdou do  $n_2$ , větší do  $n_5$

prohodi se  $n_2$  a  $n_5$

Výsledná haldá:



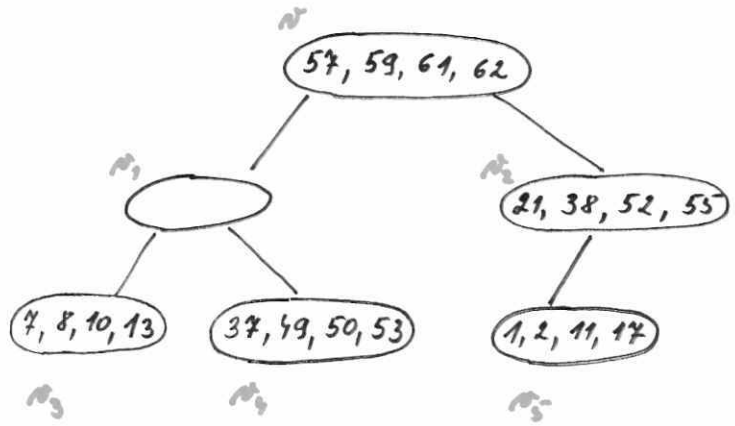
Třídění:

a) odeberu  $n_1$ , nahradím ho  $n_5$ , upravím haldu atd.

hodi se pro binární haldy, ne pro více regulární

b) odeberu  $n$

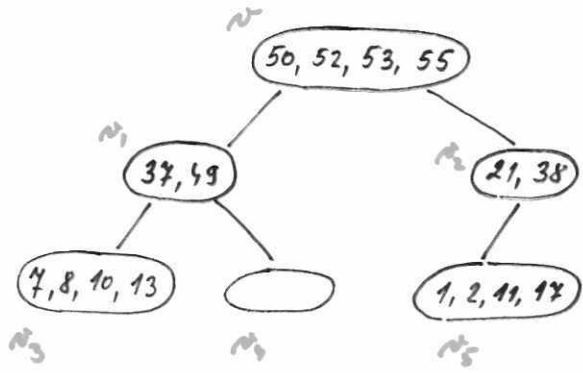
do kořene namerguji 4 největší prvky z  $n_1$  a  $n_2$



do  $n_1$  namerguji 4 největší prvky z  $n_3$  a  $n_4$  (j.  $n_3$ )

odeberu kořen

namerguji do něj 4 největší prvky z  $n_3$  a  $n_2$



doplňím uzly  $n_1$  a  $n_2$   
odeberu kořeny  
atd.

Pozor! Když se vyprázdní nějaký uzel a chce ho doplnit z jeho synů, musím nejdřív vyčistit interní paměť (mohlo v ní něco zůstat z předchozího mergeování - to se musí vrátit na disk).

### Prostorová složitost

Algoritmus se dá implementovat tak, aby trídil na místě, ale je to technicky složitě.

### Poznámky:

- 1) Existují externí verze i dalších interních třídících algoritmů
- 2) Prvky externího třídění se uplatňují i v interních algoritmech, pokud vnitřní paměť je hierarchická!

Míra složitosti: počet operací (porovnání, výměny, ...)  
+ počet výpadků z cache