

## Milí řešitelé a řešitelky!

Léto se přehouplo v podzim, venku se každou chvíli vájí mlha, a když se roztrhá, ukáže vám, jak se v uličkách valí spadané listy. Přichází podzimní deště, sychravé večery i první kola lejačkových soutěží. KSPčko má první sérii už za sebou, ale i do druhé se může zapojit kdokoliv, kdo má zájem lémat si hlavu nad úložkami, třeba v teple krbu s hněčkem horkého čaje. Krb i čaj si musíte zajistit sami, ale úložky vám s radostí dodáme. Račte se začíst.

Také bychom vás rádi pozvali v podčečr 22. listopadu na **Kalíšek** (setkání v čajovně) v Praze a hned následující den na **Den otevřených dveří na Matfyzu**. Více informací se brzy objeví na našich stránkách i na Facebooku.

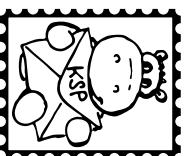
A pokud stále váháte, připomínáme, že každému řešiteli, který v tomto ročníku z každé série dostane alespoň 5 bodů, darujeme KSP propiskku, blok, pláček a možná i další překvapení.

Za řešení KSP je také možné být přijat na MFF UK bez přijímacích zkoušek. Na získání osvědčení úspěšného řešitele je letos třeba získat v hlavní kategorii alespoň 150 bodů. Maturnanti, kteří by chtěli osvědčení využít letos, jej musí získat nejpozději za čtvrtou sérii.

**Termín série: Pondělí 5. prosince 2016 v 8:00**

**Odevzdávání:** Přes web na adrese <https://ksp.mff.cuni.cz/submit/>

**Odměna série: Sladkou odměnu pošleme každému, kdo získá alespoň polovinu bodů z pěti odevzdaných úloh.**



## Druhá série dvacátého devátého ročníku KSP

*„To ten den skutečně začíná,“ prošlo Erice hlavou, když chuvině přesouvala svůj notebook z kolejniho stolu do baru. Holdny na záti nanzodory ujhřezájným pohledem ukazovaly patnáct minut do začátku výuky, která ovšem probíhala přes půl hodiny odtud.*

*Budk ji měl vzbudit už o půl osmné, ale když pak na chvilku zavřela oči, musela usnout, protože najednou bylo skoro devět. Koněně měla vše potřebné a mohla vyběhnout. Napadlo ji, kolik vlastně existuje stejně rychlých cest do školy, ale raději si zakázala experimentovat.*

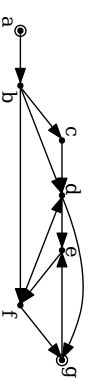
**29-2-1 Cesty do školy 10 bodů**

Studentka se chce dostat do školy za právě  $K$  minut – ne více, ale ani ne méně (to by musela zbytečně čekat na chodbě). Zajímalo by ji, kolika různými způsoby to jde udělat. Protože venku už začíná být docela zima, nechce se ani po cestě nikde zastavovat. Raději celých  $K$  minut stráví chůzí, i kdyby to znamenalo trochu si zajít, nebo třeba jít kus tam a zpátky.

$K$  dispoziční má mapu, ve které je zakresleno  $N$  významných míst – kromě startovního a cílového bodu (koleje a školy) také různé křižovatký a další místa, přes která je možné procházet. Dále mapa obsahuje seznam povolených přesunů: každá jeho položka říká, že z místa  $i$  jde dojit na jiné místo  $j$ , a to za přesně jednu minutu. Jinudy než podle povolených přesunů se pohybovat nelze. Pozor, u přesunů záleží na směru. Může se stát, že je povoleno jít z  $i$  do  $j$ , ale ne z  $j$  do  $i$  (třeba protože v opaktem směru je to moc do kopce).

Formálníji: je dán orientovaný graf. Jeho vrcholy odpovídají místům a hrany povoleným přesunům. Zajímá nás, kolik v něm existuje různých sledů mezi danými dvěma vrcholy  $s$  a  $t$  dlouhých přesně  $K$  hran. Sled je něco jako cesta, jen se na něm mohou opakovat vrcholy a hrany.

Uvažujme například následující mapu a  $K = 6$ ,  $s = a$ ,  $t = g$ :



V ni existuje právě pět sledů délky 6 vedoucích z  $a$  do  $g$ . Jsou to *abcdefg*, *abdfefg*, *abdefgf*, *abdfeffg*, *abfdfeffg*.

Obvyklý způsob dopravy dnes zafungoval. Erice se povedlo chytit autobus tramvají (přes hlavní lakt viděla přiskakující minutu) i další tramvají a krátce před půl už přebíhala Mlostranské náměstí, až skoro smetla nějakou podobně starou dítku. Schodky vzala po více najednou, přebíhla celou chodbu a zkusila se dobyt do učebny — klient se ovšem tlačila tiše, a hlavně zamčeně.

Erice hned vyřídila mobil a na Hangoutu napsala svému společníkovi: *Ahoj, prosím Te, analýza se nám pávesanaula?* Vzápětí získala neověřené koukat na čas 8:28 vedle své zprávy. Odpověď přišla záhy. *Ne ale začíná přece v 9. A pak přišla další zpráva: Vas ze je zámní cas?*

Jen malým zárukem nedošlo k našim na neviněm telefontu. Msto toho se Erice vydala skomrat nástěnký na chodbách. Na jednu někdo vylepřl papír se zašifrovaným textem a výzvou k rozluštění. Spíš než opravdové řešení teď ale Ericka toužila najít něco jako „pomsta vgnulcezi letního času“.

**29-2-2 Hledání pomsty 13 bodů**

Předpokládáme, že text na nástěnce je zašifrovaný jednoduchou substitucí šifrou. Ta funguje tak, že pro každé písmeno abecedy nahradíme všechny jeho výskytý v původním textu nějakým jiným písmenem (např. každé  $A$  zneháme na  $X$ , každé  $B$  na  $U$  atd.).

Návíc platí, že dvě různá písmena nikdy nenahradíme stejným, protože pak by se text nedal jednoznačně dešifrovat.

Předpřisn, které písmenko nahrazujeme kterým, se říká *klic*. Konkrétními zplosoby, jak zašifrovat slovo PAPIR, jsou například UWTXI nebo PEZNL, ale nikoli SPFGH (každé P jsme nahradili za něco jiného) nebo naopak CLKXC (P i R jsme nahradili stejným písmenem).

Dostanete zašifrovaný text (neznaným klíčem) a hledaný řečivec (nezašifrovaný). Vaším úkolem je najít všechny pozice, na kterých se v původním textu mohli nacházet řečivec vyskytovat. Různé výskyty mohou předpokládat různé klíče.

Například slovo POTOPA v textu ZAGHAGZGHLLQWOW můžeme najít na dvou místech:

POTOPA  
ZAGHAGZGHLLQWOW  
POTOPA

V prvním případě klic přečkáá P→G, O→H, T→A, A→Z. Ve druhém je správné přičtení P→H, O→G, T→Z, A→L. Stačilo si rozmyslet, že jinde už se toto slovo vyskytovat nemůže.

Po přednášce, která proběhla překvapivě v klidu (jen jeden nejasný dotaz a jen dvě rypnutí od spoleháka, kteří změnu času zaznamenali), čekala Erika ještě programová cí cvičení. Obvykle ho měla ráda, ale dnes se jí ponocilo jedním střechákem navíc vyrobil nekonečný cyklus. Přitom na puvod chybly přišla až po hodině, takže se dnes z výložně těšila.

Začali když si na dveřce naplínovovali svaz s kamarádkou ze střední! Potkal se na zastávce Karlova náměstí znilo jako skvělý nápad, dokud Erika nezjistila, že těch zastánek je více kus od sebe. Zauvážil kamarádce znilo jako skvělý nápad, dokud telefon suse neoznámil, že „volání ústánek není dostupný“.

A tak Erika nekonkrétně přebhala mezi jednotlivými zastávkami. Přitom si usmíla, že tu stále visí nejvzácnější nolební reklama.

### 29-2-3 Billboardová většina 13 bodů

Erika probhla ulicemi, kde visí spousta volební reklamy: billboardy, plakáty... Víme, v jakém pořadí okolo reklamních materiálů probhla a které strany propagovaly.

Rádi bychom uměli pro lhovohlou část její cesty zjistit, jestli v tomto úseku měla nějaká strana nadpoloviční většinu reklamních materiálů (a tedy by přesvědčila lidi, kteří projdou jen tuto část cesty).

Na vstupu dostanete nejdřív posloupnost  $N$  přirozených čísel  $(a_0, \dots, a_{N-1})$ . Ta udává, které strany propagují jednotlivé plakáty, v pořadí, v jakém je Erika mješla (tedy  $a_i$  je číslo strany propagované  $i$ -tou reklamou). Čísla stran mohou být lhvohlavě velká.

Poté hude nastédovat  $Q$  dotazů. Každý dotaz je tvořen dvojicí čísel  $(k, l)$ . Úkolem vašeho algoritmu je pro každý dotaz určit, zda v úseku  $a_k, a_{k+1}, \dots, a_{l-1}$  posloupnosti má nějaké číslo strany nadpoloviční zastoupení.

Například pro posloupnost

$a_0$	$a_1$	$a_2$	$a_3$	$a_4$	$a_5$	$a_6$	$a_7$	$a_8$	$a_9$
1	4	4	4	1	4	1	1	7	

a dotazy  $(1, 6)$ ,  $(3, 6)$ ,  $(2, 5)$ ,  $(5, 10)$ ,  $(0, 10)$  jsou správnými mi odpovědná 4, -4, 1, - (kde - značí, že v daném úseku nemá většinu míklo). První úsek  $(1, 6)$  je výše znázorněn podtržením.

Vyhodnocovat každý dotaz zvlášť by bylo pomalé. Zkusíte si na začátku pro posloupnost něco předpčítat, abyste pak zvládli dotazy vřizovat rychleji. Předpokládejte, že počet dotazů bude řádově stovnatelny s  $N$ .

Při třetím náratu na původní zastávku se oškerá zadání a obě dísky se konečně poškaly. V blízké kavárně pak naštené propovídky dít hodiny jako nic. Při loučení se si sblížily, že se zase brzy uvidí, a pak už každá vyrazila za svým dalším programem.

V Eričině případě to znamenalo vrátit se na koleji a skatit si uše, co by mohla potřebovat na hodiné uvození. K jejímu provozování se nechala přesvědčit už na začátku se mestru Hankou z koleje, která nechtěla chodit sama. Když Erika dorazila na svaz s Hankou, bylo už dost pozdě. Přeslo se Hanka nejvíce ze všeho toužila zmateně.

„Máš určitě úšlechto?“ zeptala se nejistě. „No jasně,“ měla Erika rýkou. „ ve které, jak si právě uedámila, měla sportovní tašku. „Eh, ne, počkej, hned jsem zpátky!“ Hanka měla z Eriky náramnou legraci a dohrála si jí i po letech, kdy se Erika chystala vydat za dalšími kamarádky. „Nemám Te raději doprovodit na místo?“ ptala se. „Huš,“ kontrolovala Erika na mobila jízdní řády, „za chvílika mi jede autobus a přestoupit na metro zvládnu.“

Dívky se rozloučily a Erika za chvilí skatečně nastoupila do autobusu. Když ale ani po pěti zastávkách nebyla v dočasném cíli, znejsvěkla a na další zastávce raději vystoupila. Zaujalo ji náměstěčko prokhané chodníčky. Mezi nimi se nacházely květinové záhony podmlouhdi neprumdelných tvarů. Sedmáhléinkový záhon, kdo to kdy viděl!

### 29-2-4 Nejsložitější záhon 9 bodů

Na náměstěčku tvaru  $N$ -dlehňka jsou různé chodníčky a kyty, které jsou ale vedené tak, že se navzájem nekřičí a vycházejí jen z pomyslných rohků, nikoliv ze samoných stran.

Oblasti mezi chodníčky tvoří květinové záhony. Nás by zajímalo, který záhon má nejsložitější tvar, tedy je tvořen mnohotělníkem s nejvíce stranami.

To je praktická open-data úloha. V odvezdkávcím systému si necháte vygenerovat výstupy a odevzdláte přišlešné výstupy. Záleží jen na vás, jak výstupy vytvoříte.

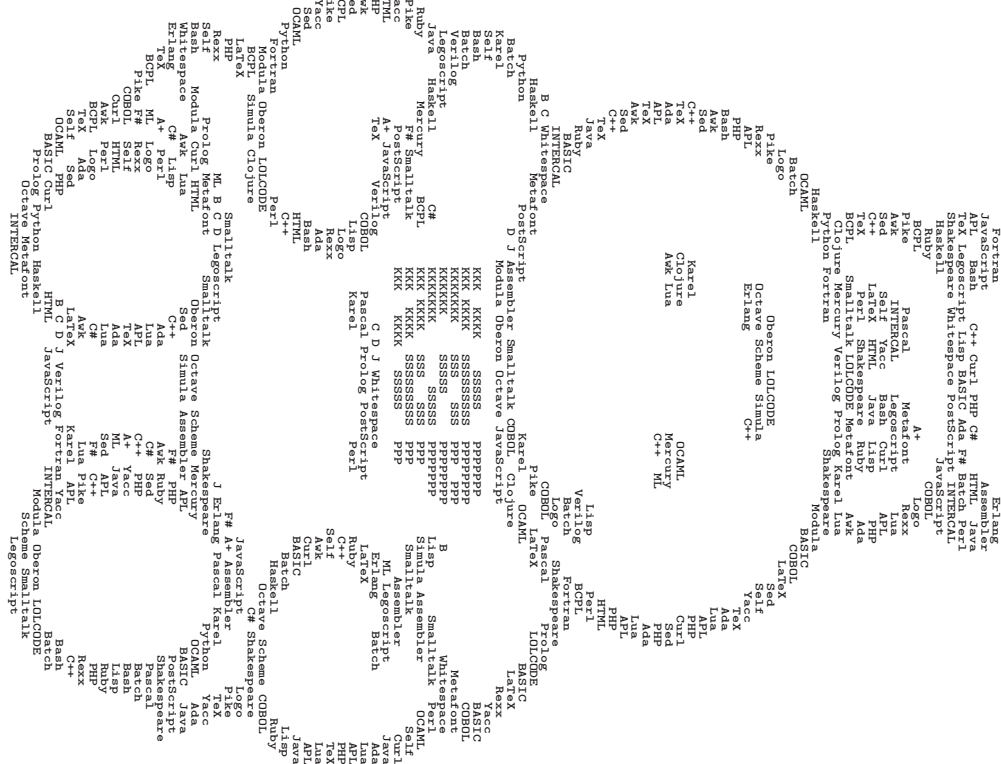
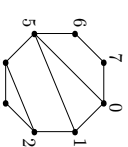
Formát výstupů: Na prvním řádku vstupu dostanete čísla  $N$  a  $K$  udávající počet vrcholů  $N$ -dlehňka a počet chodníčků. Na dalších  $K$  řádcích následuje popis chodníčků – každý chodníček je určen dvěma čísly udávajícími, mezi kterými dvěma vrcholy náměstěčka vede. Vrcholy číslyjme od 0 do  $N-1$  v pořadí na obvodu.

Formát výstupů: Na výstupu na prvním řádku počet vrcholů na obvodu největší souvislé oblasti (záhonu) a na druhém řádku nezserou oddělená čísla zajímavých vrcholů ohraničujících tuto oblast – to jsou takové vrcholy, kde předcházející z jednoho chodníčku na druhý. Čísla vyjste v rostoucím pořadí, v jakém se vyskytují na obvodu této oblasti. Pokud existuje více takových oblastí, vyberte lhvohlou.

Ukázkový vstup: Ukázkový výstup:

8 3	4
5 1	0 5
0 5	
2 4	

Nejsložitější záhon má tvar čtyřtělníku – můžeme si vybrat  $(1, 2, 4, 5)$  nebo  $(0, 5, 6, 7)$ . Na prvním jsou významné všechny body, na druhém jen body 0 a 5.







podle definice rekurzí, kde si již spočítaná  $\check{S}(a, b)$  zapamätujeme do tabulky, abychom je nepočítali znovu.

Můžeme postupovat i jiným způsobem, a to spočítat všechna  $\check{S}(a, a+k)$  iterativně po „vrstvách“, tedy nejprve všechna  $\check{S}(a, a)$ , potom všechna  $\check{S}(a, a+1), \dots$ , až nakonec získáme hledané  $\check{S}(0, n-1)$ . Při takovém způsobu počítání nám stačí si pamatovat předchozí vrstvu, tabulka tedy není nutná.

Jestliže bychom navíc chtěli znát postup, jak jsme šli rozšiřovali, budeme si kromě  $\check{S}(a, b)$  pamatovat i  $D(a, b)$  říkájící, přidáním kterého prvku (nebo ze kterého směru) byl nejlepší šit mezi  $a$  a  $b$  rozšíření. Projítím  $D(0, n-1)$  potom získáme postup rozšiřování šitů v opacém pořadí.

Nyní stačí spočítat časovou a paměťovou složitost tohoto algoritmu. Všech dvojic  $a, b$  je  $\mathcal{O}(n^2)$ , každou spočítáme nejvýše jednou. Na vypočítání každého  $\check{S}(a, b)$  spotřebujeme konstantní čas. Celková časová složitost je tedy  $\mathcal{O}(n^2)$ . Jestliže nám stačí znát pouze nejlepší součet zachráněných hraček, vystačíme si s  $\mathcal{O}(n)$  pamětí při použití iterativní metody. Jinak je paměťová složitost shodná s časovou složitostí  $\mathcal{O}(n^2)$ .

Program (C):  
<http://ksp.mff.cuni.cz/viz/29-1-6.c>

Vladislav Konečný

## 29-1-7 Stromy kolem nás

### Úkol 1: Příklady stromů

Děkujeme za pěkné příklady stromů: rozdomek (já a všichni mi biologičtí předkové), řeka se svými přítoky, adresářová struktura na disku, nebo třeba všechny možnosti, jak se může vyvíjet partit deskové hry. Přidáme k nim strukturu webových stránek (do sebe zaměřené elementy HTML), programů v programovacích jazycích a nádrvkem ještě vztahy mezi větvými členy či větvami v souvětí. Stačí? :)

### Úkol 2: Strom z postřívováno výpisu

Úkol byl formulován trochu neoprádně: pokud o vrcholech stromu vůbec nic nevíme, tvar stromu se z jeho postřívováno výpisu určit nedá. Pokud nám ale někdo řekne, kolik má mít který vrchol synů, už to možné je. U stromů vyráží je tato podmínka triviálně splněna: čísla jsou listy, operace mají právě dva syny. Pro jednodušlost budeme nadále předpokládat, že pracujeme se stromem vyrazím.

Postřívový zápis budeme procházet zleva doprava a postupně ho překládat na stromy: pokud jsme z 1234+\*\* přičetli 1234+, máme zatím hotové jednovrcholový stromy 1 a 2 a třívrcholový strom s kořenem +, pokud nímž jsou listy 3 a 4. Tyto stromy postupně slepíme do jednoho velkého, ale z toho, co jsme zatím přičetli, dosud není jasné jak.

Budeme si tedy pamatovat nějakou posloupnost rozpracovaných stromů, říkáme jí *aleje*. Kdyžkoliv ze vstupu přičteme číslo, vytvoříme jednovrcholový strom s tímto číslem a přidáme ho na konec aleje. Přečteme-li napaek nějakou operaci, odpojíme z konce aleje dva stromy, spojíme je pod nové založený kořen s danou operací a výsledok opět zapíšeme na konec aleje. Obecně objví-li se na vstupu zápis vrcholu stromu  $s_i$ , spojujeme posledních  $s$  stromů z aleje.

Časová složitost tohoto algoritmu je zřejmá lineární. Když bychom chtěli počítat dokázat korektnost, pak třeba takto: Uvažme průběh prohlédávání stromů do hloubky. Právě stejně v nějakém vrcholu  $v$  a clystáme se ho opírat, tedy do

postřívového zápisu toto  $v$  vypsát. Vrcholy na cestě z kořene do  $v$  jsme už objvíli, ale na vypsání teprve čekají. Vše vlevo od této cesty jsme už vypsali, vše napravo napaek ani neobjvíli.

O dekodování algoritmu pak platí následující invariant: při zpracování  $v$  se v aleji nachází posloupnost podstromů všich vlevo od cesty do  $v$  (v pořadí shora dolů) následovaných podstromy ležícími pod  $v$  (zleva doprava). Jakmile algoritmus uvídn  $v$ , podstromy správně poslepuje a invariant bude nadále platit.

Za odhmem vám říkáme ještě jedno, mírně magické řešení: Postřívový zápis očíme, čímž se z něj stane přehledový (až na opacní pořadí synů). Ten můžeme poskládat do stromu jednodušším rekurzivním algoritmem: pokud narazíme na číslo, vytvoříme z něj jednovrcholový strom a skončíme. Pokud na operaci, dvakrát se rekurzivně zavoláme a získáme dva stromy spojíme pod společný kořen. Hotovo,  $\mathcal{O}(n)$ . Formální důkaz složitosti by se veel stejně jako u minulého algoritmu.

### Úkol 3: Nejednoznamenitý inkřový zápis

Prosím, milý Watsone: stačí uvážít zápis 1+2+3. Ten můžeme uvažokovat dvěma způsoby: (1+2)+3 a 1+(2+3). Pokudé vyjde jiný strom. Přehledové tyto stromy zapíšeme +123 a +1+23, postřívové 12+3+ a 123++ – vše bez závorek, z předchozho úkolu už přeci víme, že nejšon třeba.

### Úkol 4: Průměr stromu

Jak zadání naznačuje, k měření délky nejdelší cesty se skutečně bude hodit něco počítat během DFS. Kdyžkoliv se budeme vracet z podstromem s kořenem  $v$ , spočítáme dvě čísla:  $h(v)$  udávající hloubku podstromu a  $l(v)$  – maximální délku cesty v podstromu.

Vypočet  $h(v)$  už známe ze zadání: v listu platí  $h(v) = 0$ , ve vnitřním vrcholu se syny  $s_1, \dots, s_k$  musí být  $h(v) = \max(h(s_1), \dots, h(s_k)) + 1$ .

Jak tedy s  $l(v)$ ? Rozmysleme si možné polohy vrcholu  $v$  vzhledem k této cestě:

1.  $v$  je list – tehdy je celý podstrom jednovrcholový, pročež  $l(v) = 0$
2.  $v$  vřbec na cestě *neleží* – tehdy jsme cestu již započítali do některé z dětek  $l(s_1)$  až  $l(s_k)$ .
3.  $v$  je koncovým vrcholem cesty – cesta vede z  $v$  dolů do nejvzdálenějšího listu, takže měří  $h(v)$ .
4.  $v$  je „zlomením“ cesty – cesta přichází zespoda z některého  $s_i$  a zase odchází dolů do nějakého jiného  $s_j$ . První část určité vede z listu, takže měří  $h(s_i)$ , pak následuje hraha  $s_i v$ , hraha  $v s_j$  a závěrečná část do listu, délky  $h(s_j)$ . Vrcholy  $s_i$  a  $s_j$  samozřejmě volíme tak, abychom použili největší a druhé největší  $h(s_i)$ .

Pro vrcholy s jedním synem může nastat druhá a třetí možnost, tedy  $l(v) = \max(l(s_1), h(v) + 1)$ . Je-li synů více, třetí možnost je vždy horší než čtvrtá, takže spočítáme  $l(v) = \max(l(s_1), \dots, l(s_k), m_1 + m_2 + 2)$ , kde  $m_1$  a  $m_2$  je první a druhé největší  $h(s_i)$ .

Nejdelší cestu v celém stromu pak najdeme v  $l(v)$  kořene. Našími výpočty strávíme v každém vrcholu lineární čas s počtem synů, což odpovídá složitosti samotného DFS. Celý algoritmus tedy poběží v lineárním čase.

### Průměr stromu podruhé

Předevedme ještě jeden způsob, jak změřit nejdelší cestu. Strom zakotíme v libovolném vrcholu  $a$ . Pak najdeme

Úkol 2 [4b]: Vytvořte datovou strukturu pro strom s obarvenými vrcholy. Na počátku jsou všechny vrcholy zelené. Chceme umět prováčet tyto operace:  $SetColor(v, c)$  – nastaví barvu vrcholu  $v$  na červenou, zelenou nebo modrou;  $CountColor(v, c)$  – zjistí, jaká barva je v podstromu pod vrcholem  $v$  nejčastější (je-li to nerozhodné, odpoví, že nerozhodné).

Úkol 3 [6b]: Jestliže jedna datová struktura. Na začátku dostaneme strom s vrcholy obdomecenými přirozenými čísly, Chceme umět operaci  $Touch(v)$ , která v podstromu pod  $v$  provede následující: nalozme minimum z nenulových čísel ve vrcholech, toto minimum od všech nenulových čísel odečte a nakonec ohlási, jestli už se povedlo všechna čísla v podstromu vymazat.

### Dvojprozměrné intervalové stromy

Další zajímavé triky můžeme prováčet s dvojprozměrnými intervalovými stromy. Do těch vkládáme dvojice čísel, které si můžeme představovat jako body v rovině. Roli intervalů pak hrají libovolné obdélníky.

Pro naše účely postacím velmi jednoduchá podoba 2D intervalových stromů, která umí takovéto operace:

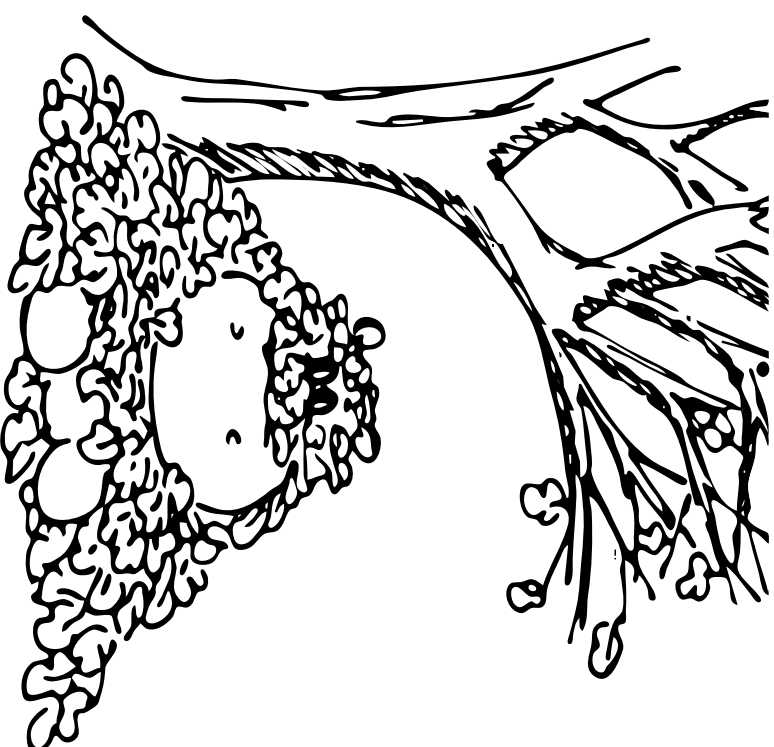
- $Init(x_1, y_1), \dots, (x_n, y_n)$  – *inicializace* –  $\mathcal{O}(n \log n)$  vytvoří nový strom s body  $o$  zadaných souřadnicích
- $RangeCount(x_{min}, x_{max}, y_{min}, y_{max})$  – *obdélníkový počet* –  $\mathcal{O}(\log^2 n)$  spočítá, kolik ze zadaných bodů leží v daném obdélníku, tedy pro kolik různých  $t$  je  $x_{min} \leq x_t \leq x_{max}$  a současně  $y_{min} \leq y_t \leq y_{max}$ .

Jak takový 2D strom sestavit, najdete například ve vzorovém řešení úlohy 24-4-7-3 dokonce včetně mazaného triku, který zrychlil intervalový dotaz na  $\mathcal{O}(\log n)$ .

V následujícím úkolu nás ale konkrétní implementace nemusí zajímat, stačí použít 2D strom jako černou skříňku.

Úkol 4 [3b]: Dostaneme strom s nestromovými hranami. Chceme si něco předpocítat tak, abychom uměli rychle odpovídat na dotazy typu „existuje nestromová hraha mezi podstromem pod  $u$  a podstromem pod  $v$ ?“.

Martin „Měkký“ Mareš



Řetězec je v podstatě jakákoliv posloupnost symbolů zapísaná za sebou a s nimi budeme v této kapitole pracovat. Každého napadne „vyhledávání v textu“ nebo „hledání jmen v telefonním seznamu“, ale řetězec najdeme i na nižších úrovních informatiky. Například celé číslo zakódované v binární soustavě, které dostaneme na vstupu programu, je také jen řetězec nul a jedniček.

Jiný příklad použití řetězce (a algoritmu s nimi pracujícího) najdeme v biologii. Například DNA není o mnoho více, než čtyřtý uložení posloupnosti čtyř znaků/nukleových bazí – dleme-li hledat vzory anebo konkrétní podposloupnosti, bude se nám hodit znalost základních algoritmů pro práci s řetězci.

Nemáme bohužel šanci vysvětlit *všechny* algoritmy s řetězci, protože je příliš mnoho možných věcí, co s řetězci dělat. Převážněm řetězci na čísla (hesování) jsme se věnovali v jiné kapitole, v této se budeme soustředit na algoritmy, které se objevují spíše v práci s textem.

Kromě úvodu do práce s řetězci popíšeme dva *starší* *kačňáky* textových algoritmů, což bude datová struktura pro slovníky – *trie* a vyhledání v textu s předprogramovaným hledaným slovem a jeho rozšíření pro více slov. S jejich znalostí pak bude mnohem snazší vymýšlet řešení složitějších, reálnějších problémů.

### Jak řetězec chápát

Když programátor dělá první křížky, často moc netuší, co s těmi řetězci vlastně může a nesmí dělat. V programování jazyce to je jasné – něco mu jazyk dovolí a na něco nejsoň prostředky. Ale jak to je na úrovni ryze teoretické?

Jak jsme si řekli na začátku, řetězec bude posloupnost nějakých symbolů, kterým říkáme *znaky*. Tyto znaky jsou z nějaké množiny, které říkáme *abeceda*. Abeceda může být jen  $\{0, 1\}$  pro čísla v binárním zápisu, klasické  $\{A-Z, a-z\}$  pro anglickou abecedu anebo plyš rozsaš univerzální znakové sadu Unicode, která má až 2<sup>31</sup> znaků. Nezapomínejme, že nejenom písmena a číslice, ale i mezery a interpunkce jsou znaky!

Vidíme, že zanedbat velikost abecedy při odhadu složitosti by bylo příliš troufalé, a tak budeme velikost abecedy označovat  $|Σ|$ . Abeceda sama se v textech o řetězci často značí řeckým  $Σ$ .

O značká samotyřch předpokládáme, že jsou dostatečně malé, abychom s nimi mohli pracovat v konstantním čase, podobně jako s celými čísly v ostatních kapitolách.

Nyní hlavní otázka – máme chápát řetězec jako pole znaků, nebo jako spojový seznam? Salanumská odporově: můžeme s ním pracovat tak i tak. Když budeme potřebovat převést řetězec na spojový seznam (protože se nám hodí rychle připojování řetězci), tak si její převádeme. Tento převod nás samozřejmě bude stát čas lineární závislý na  *délce* řetězce. Budeme ji dále značit  $L$ ; časová složitost převodu bude  $O(L)$ .

Standardně se ale počítá s tím, že řetězec je uložěn v poli někde v paměti (již od začátku algoritmu), takže ke každému znaku můžeme přistupovat v konstantním čase.

Jelikož jsme řetězec definovali jako posloupnosti, nesmíme zapomínat ani na *prázdný řetězec*  $ε$ . A když už máme řetě-

zec, určitě máme i *podřetězec* – souvislou podposloupnost znaků jiného řetězce. Například **BAR**, **RET**,  $ε$  i **KABARET** jsou podřetězce slova (řetězce) **KABARET**; **KAT** však podřetězcem není.

Často nás budou zajímat dva zvláštní druhy podřetězci. Pokud ze slova odstraníme nějaký souvislý úsek na konci, vznikne *podřetězec*, kterým říkáme *prefix* (řeky předpona), a pokud odstraníme nějaký souvislý úsek ze začátku, dostaneme *suffix* neboli příponu. **RET** je suffix slova **KABARET**, **KABA** je zase jeho prefixem.

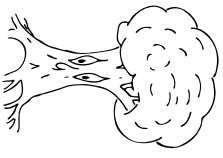
Terminologie tvorby zepředu i zezadu odstranění prázdný řetězec – to znamená, že slovo je samo sobě prefixem i suffixem. Pokud chceme mluvit o prefixech, suffixech nebo obecně podřetězci, kde jsme museli alespoň jeden znak odstranit, označíme takové podřetězce jako *vlástin*.

Pro některá použití řetězci je důležité, abychom je mohli porovnávat – když máme řetězec  $R$  a  $S$ , chceme umět rozhodnout, který je menší a který je větší. Jak přesně toto uspořádání bude, závisí na naší aplikaci, ale mnohdy se používá tzv. *lexikografické uspořádání*.

Pro lexikografické uspořádání potřebujeme nejprve zadané vlástin uspořádat na znacích. Tím se myslí takové, kde jsou všechny prvky navzájem porovnatelné, a v podstatě to znamená, že jsou uspořádaný v jedné řadě za sebou (kromě binárního  $0 < 1$  se často používá „telefonní“  $A = a < B = b < \dots < Z = z$ , které je ovšem lineární až na velikost znaků).

Když máme zadané uspořádání na znacích, na všechny řetězce je rozšíříme následovně: Nejkratší je prázdný řetězec. Ostatní řetězce řídíme podle jejich prvního znaku. Jestliže se první znak shoduje, tak podle druhého znaku, atd. Pokud přitom znaky jednoho z řetězci dojdou dříve, prohlásíme tento řetězec za menší.

Plati tedy třeba,  $ε < A < AUTO < AUTOBUS < AUTOGRAM < AUTOR < BAMBETKA < BARNABAS < Z$ .



### Adresář pomoci trie

Typický „textový“ problém je udržování množiny řetězci – můžete si představit třeba slovník. Slova ve slovníku si dleme šikovně předpracovanou, abychom pak mohli efektivně odpovídat na otázky typu: „Je slovo  $S$  obsaženo ve slovníku?“ Můžeme také po předpracování chtít přidávat nové položky, nebo i odebrat staré.

Pokud bychom neměli odebrat slova, můžeme použít hesování, které je rychle a účinné. Více o něm najdete v hesovací kapitole.<sup>4</sup> Má však tu nevýhodu, že při velkém zaplnění se začne chovat pomaleji a mírně nepřehledně.

Pro každého lučšičníka  $k$  teď chceme najít co největší nekřížící se skupinu  $A_k$ , pro níž platí, že  $k$  do ní patří (tj. bude střeřit) a jinak se skládá jen ze střeřící naloze od  $k$ . Všimněte si, že v této skupině míří na nejpravější cíl právě lučšičník  $k$ . Projdeme teď lučšičníky zleva a pro každého budeme chtít spočítat číslo  $d_k$ , což je velikost skupiny  $A_k$ .

Jistě  $A_1 = \{1\}$  a tedy  $d_1 = 1$ . Pro  $k > 1$  vytvoříme skupinu  $A_k$  tak, že prodloužíme nějakou předšou skupinu (každá větší skupina je nutně rozšířením nějaké předšle skupiny: stačí si uvědomit, že odebráním střeřce nejvíce napravo do stápneme kratší skupinu). Jistě chceme vzít takovou s co největším počtem střeřící: stále ale musí platit, že se dráhy nekříží. Pokud dáme dohromady všechna pravidla, vyjde nám, že  $d_k = \max d_i + 1$ , kde  $i < k$  a  $c_k < c_i$ . Právě druhé pravidlo zajišťi, že dráha  $k$ -tého střeřce není v konfliktu s ostatními dráhami.

Po dobehnutí cyklu pak nejvyšší  $d_k$  označuje maximální počet lučšičníků, kteří mohou střeřit nalethon. Pokud dleme znát konkrétní čísla střeřci, stačí algoritmus jednoduše rozšířit: vždy po vypočtení nového  $d_k$  si zapamatujeme číslo střeřce, rozšířením jehož skupiny vznikla  $A_k$ . Na konci cyklu pak najdeme nejvyšší  $d_k$  a od něj tyto zpětné odkazy projdeme a vypíšeme.

Zbyvá určit složitost. U  $k$ -tého lučšičníka musíme projít všech  $k - 1$  předchůdci, dohromady nám to tedy zabere  $O(1 + 2 + \dots + (N - 1)) = O(\frac{N(N+1)}{2})$  času. Paměťová složitost je  $O(N)$ , u každého střeřce si ukládáme pouze konstantní počet hodnot.

Existuje však lepší řešení. Nejprve si pro jakoukoliv nekřížící se skupinu lučšičníků definujeme *první obrov*, což je číslo cile, kam nyní lučšičník nejvíce napravo. To je zároveň nejvyšší  $c_k$  skupiny.

Může se stát, že máme více stejně velkých skupin, a chceme je rozšířit o lučšičníka, který stojí napravo od ostatních střeřci v skupinách. Pak upřičňostňujeme skupinu s nejnížším prvním okrajem ( $c_k$  nového lučšičníka musí být větší, než  $c_k$  ostatních).

Zavedme posloupnost  $m_i, i \in \{0, \dots, N\}$ . Pokud si vezmeme všechny nekřížící se skupiny velikosti  $i, m_i$  bude nejmenší z jejich pravy okrajů. V případě, že skupina velikosti  $i$  neexistuje, pak  $m_i = \infty$ . Platí, že  $m_{i-1} \leq m_i$ ; rozmyslete si, že ze skupiny  $s$  i střeřci lze vytvořit skupinu  $s + 1$  střeřci takovou, že pravy okraj zůstane stejný. Pokud pro zadané lučšičníky a jejich terče dokážeme vypočítat posloupnost  $m_i$ , je maximálním počtem lučšičníků nejvyšší takové  $i$ , že  $m_i < \infty$ .

A jak tedy  $m_i$  získat? Opět projdeme střeřce zleva doprava; na začátku zvolíme  $m_0 = -1$ , pro  $k > 0$  bude  $m_k = \infty$ . Pro  $k$ -tého střeřce potom vyepíšeme posloupnost takto: protože posloupnost  $m_i$  je celou dobu neklesající, dokážeme najít takové  $i$ , že  $m_i < c_k \leq m_{i+1}$ .

Všimněte si, že přidáním  $k$ -tého lučšičníka nevyepíšeme  $m_i$ , kde  $i \leq i$ , pro skupiny velikosti  $i$  jsme už našli menší pravy okraje. Můžeme však zlepšit  $m_{i+1}$ : představte si, že ze skupiny o  $l + 1$  lučšičnicích s prvním okrajem  $m_{l+1}$  odstraníme střeřce nejvíce napravo a přidáme  $k$ -tého lučšičníka – tím  $m_{l+1}$  snížíme, nebo necháme nezmenšit. Pro  $i > l + 1$  toto udělat nemůžeme, museli bychom odstranit dva a více lučšičníků a velikost skupiny by se změnila.

Akoliv i zde uděláme  $N$  kroků, najít správné  $l$  lze pomocí binárního vyhledávání, následně už jen upravíme  $m_{l+1}$ . Celkový čas je tedy  $O(N \log N)$ .

Na závěr důležité pozorování: pokud  $c_0, c_1, \dots, c_N$  zapíšeme jako posloupnost, odporově nalezneme řešení nejdelší rostoucí podposloupnosti v  $c_k$ . Nemí těžké převést tyto úlohy mezi sebou: ostatně většina řešitelů s plymým počtem bodů si této spojitosti všimla.

Kloba Maroušek

Program s řešením v  $O(N^2)$  (Python 3):  
<http://ksp.mff.cuni.cz/viz/29-1-5-pomale.py>  
 Program s řešením v  $O(N \log N)$  (Python 3):  
<http://ksp.mff.cuni.cz/viz/29-1-5-rychle.py>

### 29-1-6 Štítové konzo

Nabízí se triviální algoritmus, který vyzkouší všechny možnosti. Jenomže toto nebude fungovat vůbec rychle. Pojďme se podívat proč.

Jednotlivé kroky rozšřřování štítu můžeme zapsat do posloupnosti  $\{L, P\}^{n-1}$ , kde  $n$  je počet všech hradeb. Nechtěť pak v této posloupnosti  $L$  reprezentují rozšřřování štítu vlevo a  $P$  rozšřřování štítu vpravo. V této definici můžeme jednoduše získat počáteční úsek hradeb spočítáním všech  $L$  (můžeme začít dostatečně vpravo, aby rozšřřování vlevo mělo smysl), takto každá posloupnost jednoznačně určuje postup rozšřřování štítu.

Počet všech takových posloupností čini  $2^{n-1}$ . Vyzkoušení všech možností by nám tedy trvalo vždy  $\Omega(2^n)$  času. Jak se z toho vyhrádit?

Hledové řešení nefunguje. Vybrat vždy nejvyšší úsek hradeb je již vyvárcem příkladem v zadání. Stejně tak nás zradí rozšřřování směrem, kde je větší součet stojících hradeb. Ukážeme si to na jednoduchém příkladu. Májíme následující hradeby:

10 10 0 0 0 0 0 100

Zde se hledový algoritmus rozhodne rozšřřovat štít vpravo. Jenže než se dostane k nejpravějšímu úseku, oba úseky nalevo ztratí příliš mnoho hodnoty. Nakonec zadržání jen  $(100 - 5) + (10 - 6) + (10 - 7) = 102$  hradeb. Kdybychom nejprve ochránili část hradeb vlevo a až potom se vydali vpravo, zachráníme jich  $10 + 9 + (100 - 7) = 112$ , tedy mnohem více.

Hlavní potíž při zkoumání všech možností spočívá ve skutečnosti, že takto spouští společných postupů mnohočet zbytečné znovu počítáme. Třeba tyto dvě varianty aktivních štítů  $A B C D E$  a  $A B C D E$  vyžadují ze štítu  $A B C D E$  každý při zkoumání všech možností celý přepočítáme dvakrát. Takový problém řeší princip *dynamického programování*.<sup>8</sup>

Označme  $\tilde{S}(a, b)$  nejlepší možný součet výšek hradeb v úseku od  $a$  do  $b$ , který chrání náš štít v kroce  $k = b - a$ . Dále si označme  $V(k, i)$  výšku  $k$ -tého nechráněného úseku hradeb v  $k$ -tém kroce, jejíž hodnota se bude rovnat  $\max(0, h[i] - k)$ .

Všimněte si, že  $\tilde{S}(a, b)$  se rovná buďto  $\tilde{S}(a, b - 1) + V(k, b)$ , anebo  $\tilde{S}(a + 1, b) + V(k, a)$ , podle té z možností poskytnující vyšší součet chráněných hradeb. Tedy vzorec pro  $\tilde{S}(a, b)$  je  $\max \{ \tilde{S}(a, b - 1) + V(k, b), \tilde{S}(a + 1, b) + V(k, a) \}$ .

Nyní můžeme spočítat  $\tilde{S}(0, n - 1)$ , jenž odporově optimálnímu rozšřřování štítu nad celou hradbou. Můžeme postupovat

<sup>4</sup> <http://ksp.mff.cuni.cz/viz/kucharka/hesovani>

<sup>8</sup> <http://ksp.mff.cuni.cz/viz/kucharka/dynamika>



hustolesa, musela by být vlozena též v SZKH, ale pak bychom se našim obdelínkem vůbec nemohli zabývat, protože mezi  $u_0$  a  $u_1$  by ještě existovalo nějaké  $u_i$ , takže by  $u_0$  a  $u_1$  nebyly línou, po sobě jdoucí polohy zametací přímky, tedy máme spory, kterým jsme dokázali, že tento obdelník je horních a dolních hran hustolesi zcela prost.

Pokud tedy nějaký hustoles v obdelníku leží, vždy obsáhne celou jeho výšku, a tedy se vždy musí protnout s úsečkou jdoucí po jeho šířkovém. Stejně tak obráceně, pokud nějaký hustoles v obdelníku neleží, tak se rozhodně mezi souřadnicemi  $u_0$  a  $u_1$  nemůže protnout se zkomponovanou úsečkou.

Kromě hustolesi tedy budeme mít na naší zametací přímce ještě bod, který bude označovat průsečík zkomponované úsečky s aktuální přímkou, a při každém posunutí přímky s ním poměrně na správné místo a zkontrolujeme, jestli jsme tím náhodou netrhli hustoles.

Co když ale úsečka vede shora dolů, takže musíme projít všechny polohy zametací přímky? Jednou by to nevedlo, ale mohly by takto vést třeba všechny úsečky a pak bychom si se složitostí vůbec nepamohli. Štále totiž každý posun zametací přímky zabere  $O(\log H)$  času, takže bychom si zhoršili složitost na  $O(U \cdot H \log H)$ . To nám je platné jak nrvěmu zinnku.

Všimneme si tedy, že při posouvání zametací přímky sem-tam počítáme pořád dokola toleže – její stav.

Pořídíme si tedy datovou strukturu, která dokáže nějak rozumně uložit všechny možné polohy zametací přímky, a zároveň nebude přehnaně velká. Pokud jsme ochotni obětovat až  $O(H^2)$  paměti, stačí si postupně uložit všechny daty zametací přímky. Jenomže na vygenerování  $O(H^2)$  dat potřebujeme také  $O(H^2)$  času, takže jsme zase nahraní.

Ne tak docela. Ukládáním všech stavů zametací přímky bychom zase bedacos ukládali redundantně, takže si pořídíme nějaký vhodný druh komprese. Každý hustoles budeme reprezentovat dvěma seznamy ukazatelů: seznamem levých sousedů a pravých sousedů. Seznam levých sousedů říká pro každý interval souřadnice  $y$  nejbližší sousední hustoles směrem vlevo. Analogicky vypadá seznam pravých sousedů. Seznamy postavíme tak, že projdeme zametací přímkou přes všechny události v SZKH. Na začátku není v oblasti žádný hustoles, ale vytvoříme si dva virtuální se souřadnicemi  $x_A = -\infty$ ,  $x_B = \infty$ .

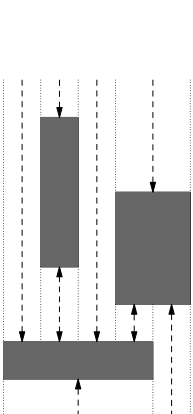
Když potkáme událost „začátek hustolesa“, podíváme se, mezi které dva jiné hustolesy jej máme přidat. (Vždycky tam budou alespoň ty dva virtuální.) Do seznamu levých sousedů pravého hustolesa a do seznamu pravých sousedů levého hustolesa tedy přidáme odkaz na právě přidávané hustoles, do seznamu levých sousedů právě přidávaného hustolesa vložíme odkaz na levý hustoles a do seznamu pravých sousedů právě přidávaného hustolesa vložíme odkaz na pravý hustoles.

Když potkáme událost „konec hustolesa“, podíváme se také, mezi kterými dvěma jinými hustolesy byl. Seznamy levých a pravých sousedů odebráme hustolesa uzavřeme, do seznamu levých sousedů pravého hustolesa přidáme levý hustoles a do seznamu pravých sousedů levého hustolesa přidáme pravý hustoles.

Jak dlouho to trvá? Každý posun zametací přímky potřebuje  $O(\log H)$  času kvůli aktualizaci jejího stavu. K tomu vytvoříme 4 nebo 2 nové odkazy, což je konstanta, která logarithmem nelme. Celkem nás bude vytvoření datové struk-

tury stát  $O(H \log H)$  času, což je stále stejně jako počítání potřeba seřadit.

Jak velká bude zkonstruovaná datová struktura? To se dá spočítat z algoritmu, kterým ji vytváříme. Každý hustoles vygeneruje v SZKH jednu událost začátku a jednu událost konce hustolesa. Zpracováním těchto dvou událostí vznikne v datové struktuře právě 6 nových odkazů. Datová struktura tedy spotřebuje celkem  $O(H)$  paměti. To vypadá nádherně.



Na vstupech, které generovalo naše odevzdávátko, už tato úprava běžela dostatečně rychle na získání plného počtu bodů, protože dříve většina zadáních úseček byla v porovnání s velikostí oblasti velmi krátká. Při zpracování každé úsečky totiž stačilo projít několik málo odkazů – prázdných oblastí, přes které zrovna procházela, i když byla třeba poměrně dlouhá.

Stále se nám síce může stát, že budeme při zpracování každé úsečky projít přes  $O(H)$  odkazů, takže jsme si pomohli zpátky na  $O(U \cdot H)$  v nejlhorším případě.

Bylo by fér spočítat, že v průměrném případě na tom budeme lín, nicméně počítat statistiku nad úsečkami a obdelníky v rovině je poměrně ošemetné, jak naznačuje kupř. Bertrandův paradox, tak si to raději odpustíme.

Program s řešením v čase  $O(N^2)$  (C):  
<http://ksp.mff.cuni.cz/viz/29-1-4-slow.c>

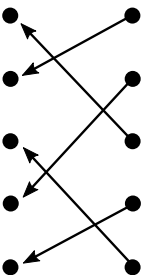
Program s rychlejším řešením v čase až  $O(N^2 \log N)$  (C):  
<http://ksp.mff.cuni.cz/viz/29-1-4-fast.c>

Jan „Mlosky“ Matějka

### 29.1-5 Lucičkiní

Na vstupu máme  $N$  lucičkinů. Zleva doprava si je očísujeme od 1 do  $N$ , stejně jako cíle, na které se míří. Platí, že cíl s číslem  $k$  se nachází naproti střelci  $k$ . Jako  $c_k$  si označíme cíle lucičkinů  $k$ .

Mnoho z vás se snažilo najít nejlepší řešení tímto způsobem: pro každého lucičkina našlo počet dráh, s nimiž se ta jeho kříží, a následně postupně odstranovalo lucičkinů s největším možným počtem křížení. Takový postup ale obecně nefungoval. Například v zadání na obrázku byste mohli odstranit i druhého lucičkina, který do správného řešení ovlivňuje páří, ačkoli se jeho dráha kříží se dvěma dalšími dráhami.

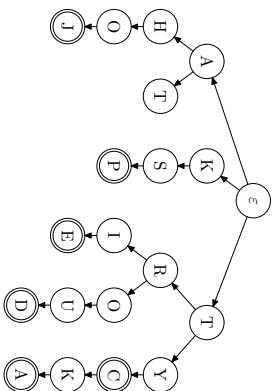


Uvedeme si řešení božic v čase  $O(N^2)$  a to následně zlepšujeme. Zavedeme pojem *nekrýžící se skupina* pro takovou podmnožinu střelců, že jejich dráhy se navzájem nekrýží.

Ukažeme si jiné řešení, které je také asymptoticky rychlé a není ani příliš náročné na paměť. Využívá stromové struktury a říká se mu *trie* (vyslovujeme česky „trýpě“ a anglicky jako část slova „retrieval“, z něhož slovo *trie* vzniklo). V českém se občas používá také označení „písmenkový strom“.

Trie bude zakoreněný strom. V prvním patře se bude větvit podle prvního písmene slova, ve druhém podle dalšího, a tak dále.

Obrázek vyvá za tisíc definic, pojďme se podívat, jak vypadá trie pro slova AHOJ, AT, KSP, TRIE, TROUD, TYG, TVČKA. Pro přehlednost písmena místo na hrany kreslíme do následujících vrcholů:



Všimněte si, že vrcholy v hloubce  $h$  (tedy v  $h$ -tém patře trie) odpovídají prefixům délky  $h$  zadáních slov. Například prefixy délky 2 jsou AH, AT, KS, TR a TY. Hrana mezi prefixy vede právě tělady, lze-li jeden z druhého získat připsáním písmene na konec.

Jak bychom takovou trii postavili algoritmem? Přesně, jak jsme ji definovali: každé slovo ze slovníku budeme procházet znak po znaku, a bude-li nějaká hrana chybět, tak ji vytvoříme a pokračujeme dále podle slova.

Z takto popsané trie bohnžel nepoznáme, kde končí slovo ze slovníku a kde končí jen jeho prefix. Standardní způsob, jak to vyřešit, jsou dva: buď si do každého vrcholu přidáme informaci o tom, je-li koncem celého slova, nebo si (jak je to naznačeno dvojitými kroužky v obrázku), nebo si rozšíříme abecedou o speciální znak, který se v ní předtím nevyškrtyval – třeba \$ – a pak všem slovům přilepíme tento \$ na konec.

Budeme-li se později ptát, bylo-li slovo ve slovníku, do příkladu trii zkontrolujeme ještě, jestli z konečného vrcholu vede hrana odpovídající znaku \$.

Jestli jsme si nerozmysleli, jak budeme v jednotlivých vrcholech trie reprezentovat hrany do dalších prefixů. Abychom mohli vyhledávat skutečně lineárně, potřebovali bychom umět v konstantním čase odpovědět na otázku „má vrchol  $P$  potomka přes hranu se znakem  $c$ ?“.

Abychom zajistili konstantní čas odpovědi, museli bychom mít v každém vrcholu pole indexované znaky abecedy. To ovšem znamená, že takové pole budeme muset vytvořit, a tedy alokovat  $|Σ|$  políček v každém znaku.

To zvyšší paměťovou náročnost trie (a časovou náročnost její stavby) na  $O(D \cdot |Σ|)$ , kde  $D$  znací velikost vstupu, čili součet délek všech slov ve slovníku. To je naprosto přijatelné pro malé abecedy, ale už pro {A-Z, a-z} je tento faktor roven 52 a pro Unicode je taková alokace nemyslitelná.

Jak z toho ven? Můžeme ožádat konstantní rychlost dotazu

<http://mj.ucw.cz/vyuka/ga/>

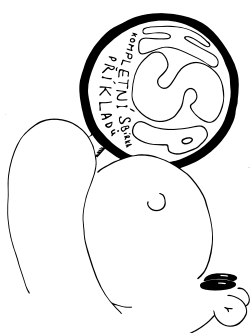
a použít namísto pole třeba binární vyhledávací strom nebo hešovací tabulku všech znaků, kterými aktuální prefix může pokračovat. Nebo můžeme každý znak velké abecedy zapsat pomocí několika znaků menší abecedy. Tou menší abecedou může být třeba {0, 1}. Těldy nahradíme každé znaky pívodní abecedy [log<sub>2</sub> |Σ|] novými (jeho zápisem ve dvojkové soustavě). Tím se časová složitost konstrukce zlepši na  $O(D \cdot \log |Σ|)$  a časová složitost dotazu na slovo délky  $L$  zhoší na  $O(L \cdot \log |Σ|)$ .

### Poznámky

- Chcete-li algoritmus konstruace trie vidět napřes v Pascahu, podívejte se do knihy *Algorithms a programming techniques*.
- Třím se také říká *prefixové stromy*, což popisuje, že každý vrchol odpovídá prefixu některého slova ve slovníku.
- Kdybychom chtěli, mohli bychom pomocí trie vyhledávat v textu v lineárním čase. Můžeme přeci postavit slovník ze všech slov v daném textu, a pak procházet příslušnou trii. Má to ale pár háčků: jednak je často hledány řetězce krátké, ale text se nevejde do paměti, jednak pokud bychom použili jako oddělovací mezeru, mohli bychom hledat jen jednotlivá slova, a mnohí jejich konce nebo delší kusy věty.
- Asi se po posledním poznance ptáte – existuje nějaká modifikace trie, která umí hledat libovolnou část textu? Ano, jmenuje se *sufixový strom* a jdu s ní dělat spousty krásných konstat. Říká se, že každou řetězovou lhotu lze řešit v lineárním čase pomocí sufixových stromů. Ví se o nich dočtete třeba v knize *Krajinnou grafových algoritmů*.<sup>5</sup>

### Čvčtění

- Řekneme, že chceme slovník na vstupu seřadit v lexikografickém pořadí (definovaném v sekci „Jak řetězce chápát“). Problém pro klasické třídící algoritmy je to, že porovnání dvou řetězců není bohnžel konstantně rychlé. Vymyslete způsob, jak seřadit takový slovník rychle pomocí trie.
- Komprese trie. Co kdybychom chtěli odstranit přebytečné vrcholy trie, tedy ty, v nichž se slova nerýví? Rozmyslele si, jestli by něčemu vadilo místo takovýto cest mít jen jednotlivé hrany. Zkusíte se konstruace nebo vyhledávání? Mnohočetem, je celkem jasné, že takováto komprimovaná trie přinese jen konstantní zrychlení dotazů i prostoru, a tak na seřazích apod. stačí použít základní variantu.



Začátek situace je asi zřejmý – máme na vstupu zadán dlouhý text a krátké slovo. Slovo si můžeme nějak předpracovat, načež projdeme co nejrychleji text a nahlásneme jeden nebo více shodných slov. Zaujímají nás při tom i vyskytnutí, které se nazývají překryvaty: v textu MAMAMA se slovo MAMA vyskytuje dvakrát. Často se hovoří o „hledání jehly v kuce se sena“, proč se textu předává *sema* a hledáme textu slovo *jehla*. Délku jehly označíme  $J$  a délku textu  $S$ .

Představme si nejdříve hledané slovo jako spojový seznam, třeba slovo INSTINKT:



Mohl bychom text začít procházet znak po znaku a kontrolovat, zda se text shoduje s našim slovem/spojovým seznamem. Pokud by si znaky odpovídaly, skočíme na další znak z textu a i na další znak v seznamu. Co když se ale neshodují? Pak nemůžeme jen skočit na další znak textu – co kdybychom v textu narazili na slovo INSTINKTINKT?

Musíme se tedy vrátit nejen na začátek spojového seznamu, ale i zpátky v textu na druhý znak, který jsme označili jako odpovídající, a zkoušet porovnávat s jehlou znova od začátku. To už naznačuje, že takto získaný algoritmus nebude lineární, protože se musí vracet zpět v textu o délku jehly.

Sice je předchozí popis skutečně v nejlhorším případě složitý  $O(S \cdot J)$ , avšak stačí malá úprava a složitost přijde na lineární  $O(S+J)$ . Je skutečností algoritmus nepomaha-vatlo v určení se – za špatnou složitost mohl fakt, že jsme se vraceli *přihřívá zpátky*.

Třeba v našem příkladu s textem INSTINKT se nemáme vracet ve spojovém seznamu na začátek, jakmile narčíme INSTINK. Měli jsme se vrátit jen na druhý znak, tedy do prvního N, a pak kontrolovat, jaký znak pokrčuje dál. Když následuje S jako v našem případě, můžeme pokračovat dále v čtení a nevracíme se v textu. Kdyžby text byl jiný, třeba INSTINB, vrátili bychom se po načtení B na začátek spojového seznamu a v textu bychom pokračovali dále bez zastavení.

Pro každý znak ve spojovém seznamu si tedy určíme políčko spojového seznamu, na které skóčíme, pokud se následující znak v textu liší od toho očekávaného. Pořadové číslo tohoto políčka nám poradí tzv. *zpětná funkce F*, což bude funkce definovaná pomocí, kde  $F[i]$  bude pořadové číslo políčka, na které se má skočit z políčka číslo  $i$ . Porovnávat pak budeme s následujícím znakem. Pokud  $F[i] = 0$ , znamená to, že máme začít porovnávat úplně od prvního znaku jehly.

Pokud máme rádi grafovou terminologii, můžete se na náš spojový seznam dívat jako na graf a hovořit o *zpětných hranách*.

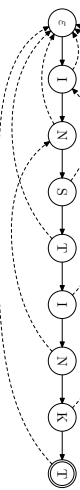
Zatím jsme ale přisane nepopsali, na které políčko přesně bude zpětná funkce ukazovat. Někdy chceme určit zpětné políčko pro druhé N ve slově INSTINKT. Pracujeme teď s přefixem INSTIN. Selský řečeno, chceme najít „konec slova INSTIN takový, že je stejný jako začátek slova INSTIN“.

Abychom náš požadavek upřesnili, zamysleme se nad zpětným políčkem pro jiné slovo. Co kdyby jehlou bylo slovo ABABABC a my určovali zpětné políčko pro ABABAB? Když bychom ukázali na první písmenko B, nebylo by to správné,

protože pak bychom pro text ABABABABC nezaháslili výskyt jehly, což je jasná chyba. Musíme se vrátit už na ABAB! Zajímá nás tedy ne libovolný sufix, který je stejný jako začátek, ale nejdelší takový konec/sufix. A ještě navíc nejen ten nejdelší, ale nejdelší „neutriviální“ – slovo INSTIN by samo sobě přechkem a suffixem, ale zpětná funkce pro N by se neměla cyklit, měla by vesti zpátky.

Reknuvme to tedy znova, zcela formálně: pokud bychom právě určovali hodnotu zpětné funkce pro znak číslo  $i$ , kterému odpovídá prefix  $P$ , pak její hodnota bude *délka nejdelšího netriviálního suffixu slova P*, pro který ještě platí, že je zároveň prefixem  $P$ .

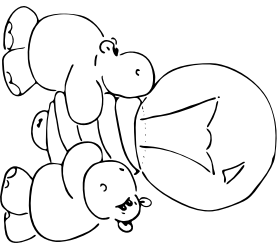
Pro slovo INSTINKT vypadá spojový seznam se zpětnou funkcí (zakroužčenou pomocí ukazatelů) takto:



Nyní vysvětlíme dvě otázky: Jakou to má celé časovou složitost? A jak spočítat zpětnou funkci?

Poporne se nejdříve s tou první. Pro každý znak vstupního textu mohou nastat dva případy: Buď znak rozšiřuje aktuální prefix, nebo musíme použít zpětnou funkci. První případ má jasně konstantní složitost, druhý je horší, neboť zpětná funkce může být pro jeden znak volána až  $J$ -krát. Při každém volání však klesne pořadové číslo aktuálního stavu (políčka) alespoň o jedna, zatímco kdykoliv stav prodlužujeme, roste jen o jeden znak. Proto všechny zkrácení dohromady může být nejvýše tolik, kolik bylo všech prodloužení, čili kolik jsme přechli znaky textu. Celkem je tedy počet kroků automatu lineární v délce textu, tj.  $O(S)$ .

Konstruici zpětné funkce provedeme malým trikem. Všimněme si, že  $F[i]$  je přesně číslo stavu, do něž se dostaneme při spuštění našeho vyhledávacího algoritmu na textězec, který tvoří prefix délky  $i$  z jehly bez prvního znaku.



Proč to tak je? Zpětná funkce říká, jaký je nejdelší vlastní suffix daného stavu, který je také stavem, zatímco políčko, ve kterém po  $i$  krocích skončíme, označuje nejdelší sufix textu, který je stavem. Tyto dvě věci se předtím liší jen v tom, že ta druhá připisuje i nevlasní sufixy, a právě tomu zabránilme odstraněním prvního znaku.

Takže  $F$  získáme tak, že spustíme vyhledávání na část samotné jehly. Jenže z hledávání zase potřebujeme funkci  $F$ . Jak z toho ven? Budeme zpětnou funkci vytvářet postupně od nejkratších prefixů. Zřejmě  $F[1] = 0$ . Pokud již

Dlouho podotknout, že celkový počet možností samozřejmě nezávisí na tom, v jakém pořadí si gardisty vybírají kopí. Jen když si představujeme, že si vybírají v pořadí od nejmenšího, je o dost snazší možnosti spočítat.

Teď v podstatě stačí dosadit do vzorce  $K$ . Tomu bychom ale nejdřív potřebovali spočítat jednotlivá  $m_i$ . To je celkem jednoduché. Na začátku si kromě gardistů i kopí seřadíme vzestupně podle délky. Teď budeme postupně procházet jednotlivé gardisty a přiřůžně si udržovat index největšího kopí menšího než aktuální gardista (tento index odpovídá právě  $m_i$ ).

Když přjdeme na následujícího gardistu, tento index zvyšujeme, dokud nenarazíme na kopí větší než on. Takto spočítáme všechna  $m_i$  jedním průchodem přes obě pole v lineárním čase. Podrobněji ve vzorovém programu. Jde o podobný princip jako při séřádění seřizdřených posloupností.

Celková časová složitost je  $O(N \log(N))$ , protože tolik času strávíme tříděním a zbyřek stiháme v lineárním čase. Paměťová složitost bude lineární.

Program (Python 3):

`http://ksp.mff.cuni.cz/viz/29-1-3.py`

Věšina z vás na základní myšlenku přišla. Co jste často opomjeli, je zřihvodit, proč řešení funguje. Z výše uvedeného vzorečku není bez komentáře vůbec vidět, že by měl platit. Po přechtení spousty řešení nebylo ani jasné, proč vlastně vstup třídí, když se poe nikdy nezmenily, k čemu to využívá, či dokonce že bez třídění by vzoreček neplatil. Zřihvodnění jsme tentokrát porovžovali za poměrně dlouhým souřezím tlavy (už proto, že algoritmus samotný je poměrně jednoduchý) a patřičně sřihávali body.

Také jste se malokdy zamysleli, co se stane, když mexistuje záhada možnost, jak gardistům kopí rozdat (například některá kopí jsou vyšší než všichni gardisté). V takovém případě se může stát, že některé ze závorek  $m_i - i$  vyjdou nulové či dokonce záporné.

V tomhle případě platí, že pokud některá z mř vyjde záporná, bude mezi nimi i nějaká nulová (rozmyslete si), tedy desazvaním do vzorečku dostaneme správný výsledek: nulu. Ale není to vůbec zřejmé a slušelo by se nad tím v řešení zamyslet.

Jesté se zmíníme o jednom detailu, který je asi přijatelné na úrovni KSP přihř neřšit: výsledný počet možnosti může být velmi velké číslo. Počítáme součin  $N$  čísel velkých až  $N$  (v nejlhorším případě jsou všechna kopí kratší než všichni gardisté), tedy výsledek může být velký řádově až  $N^N$ , což se kromě nejmenších vstupů do běžných celočíselových typů nevejde.

S tím se musíme nějak vypořádat v zřivřlosti na tom, čeho přesně chceme dosáhnout. Pokud by nám výsledek například stačil přibližně či řádově, nejštedrošším řešením je použít nějaký typ s plovcou částkou (float, double), který umí uchovávat velmi velké čísla, být s omezenou přesností, a zacházet s nimi v konstantním čase a prostoru.

Druhou možností je počítat s velkými čísly poचितě přesně. Ale to už obecně nezřidchne v konstantním čase: číslo velké řádově  $K$  zabere v paměti místo  $L = O(\log K)$  a vynásobením dvou takovýchto čísel zřidchne poměrně jednoduchým algoritmem v čase zhruba  $O(L^{1.5})$  (jde to i rychleji).

`http://ksp.mff.cuni.cz/viz/kuchacky/geometry`

V našem případě  $L = O(\log N^N) = O(N \log N)$ , tedy jedno násobení stiháme v čase  $O(N^{1.5} \log^2 N)$ , počítáme  $N$  součinů, takže celková časová složitost vzroste na  $O(N^{2.5} \log^2 N)$ , paměťová na  $O(N \log N)$ .

*Přip. Stěškovský*

**29-1-4 Zběsilý útek**

Uloha byla velmi jednoduchá, až triviální. Prosit pro každou úsečku projít všechny hustloty; zjistím, se kterým se protíná, spočítám dráhu, kterou urazím hustlotou a kterou řídíkolesem, obojí vyřídím přísušnou vrcholostí a výsledné časy sečtu.

Takové řešení má časovou složitost  $O(U \cdot H)$ , kde  $U$  je počet úseček a  $H$  je počet hustlot. Tyto hodnoty mají být řádově stejné, aneb  $U = O(N)$ ,  $H = O(N)$ , celková složitost je tedy  $O(N^2)$ . Není to moc? Je to moc. Aha. Uloha nebyla tak jednoduchá a triviální, jak by se na první pohled mohla bodit. Přesto i na triviálním řešení šlo nasbírat nemálo bodů.

Připomeňte si techniku zametání roviny přímkou? Projdeme celou oblast například zdola nahoru a poznamenejme si, kde začíná a kde končí který hustlote. Na to si musíme seznam hustlot seřidit, což zabere  $O(H \log H)$ . Seznam začátků a konců hustlot seřidit si označíme jako SZKH.

Tahle metoda se nám hodí při procházení lesem. Na začátku stojím v bodě  $y_0$ , spočítám si, které hustloty procházejí přímkou o souřadnici  $y = y_0$ , a budu tedy znát jejich seznam. Bude se hodit, aby byl seřidřený; zabere nám to jen  $O(H \log H)$ .

Pak posunu přímku do  $y_1$ . Abych ji nemusel počítat celou odznova, podřívám se do SZKH a zpracuju všechny udlosti, které jsou mezi  $y_0$  a  $y_1$ . A celou dobu se přitom dívám, jestli se na přímce nahodou neobjeví nějaký hustlote se souřadnicí mezi  $y_0$  a  $y_1$ , a pokud ano, tak jej hned podřívám zkoumáním, jestli nahodou nebyl protát právě zpracovávanou úsečkou.

Pak posunu stejným způsobem přímku do  $y_2$  (a zpracovávám druhou úsečku), pak do  $y_3$  ... Postupně tak zpracuju všechny úsečky na vstupu.

Implementace se pro účely ukládání stavu na přímce bude patrně hodit použít nějaký vyvažovaný vyhledávací strom, třeba červenocerný. Složitost zpracování každé úsečky pak bude  $O(K \log H)$ , kde  $K$  je počet kroků, které je potřeba udeřat.

Tím jsme si ale vůbec nepomohli, protože úsečka zrovna mohla vstát z levého okraje území na pravý, takže stejně musíme vzřkousšet všechny hustloty, i když víme, že prostout má nejvýše jeden.

Neuvědomili jsme si ale, že vůbec nemusíme zkoušet všechny hustloty! Změní se  $x_0$  a  $x_1$ , ale často jen mály úsek. Někdy aktuální vodovorná přímka má souřadnici  $y_0$ , a nejbližší následující má souřadnici  $y_1$ . Zkoumaná úsečka mezi body  $(x_0, y_0)$  a  $(x_1, y_1)$  protne tyto dvě přímky na souřadnicích  $x_0 + \frac{(x_1 - x_0)(y_0 - y_0)}{y_1 - y_0}$  a  $x_1 + \frac{(x_1 - x_0)(y_1 - y_0)}{y_1 - y_0}$ .

Co se stane v obědřím, který jsme si vyřidili souřadnicemi  $(x_0, y_0, x_1, y_1)$ ? Především vřitit něj nemůžeme ležet. Záhdá vodovorná hrana žádchého hustlotosa. Kdyby v našem obědřím snad chtěla ležet horní nebo dohí hrana nějakého



## 29-1-2 Kupecké počty

Každý kupce v $i$ , kolik zaplatil, tedy se domluví, napíšou všechny částky na papír a spočítají z nich průměr. To je částka, kterou každý z nich chce nakonec zaplatit. Některé jí zaplatili víc, jiní méně. V seznamu kupců tedy od výše každé investice odečteme průměr.

Tím dostaneme pro každého kupce buďto kladné číslo – o tolik zaplatil do projektu  $v$ , a tedy tolik má od ostatních dohranady dostat – nebo záporné číslo – o tolik zaplatil méně, a tedy tyto peníze nějak rozdělí svým kolegům.

Mžou nastat výteř i nulový kupce zaplatit právě tolik, kolik zaplatil měl, a nemusí se tedy s nikým vyrovnávat. Vyřadíme jej ze seznamu, neboť takových kupců by klidně mohlo být docela dost a zbytečně by kazili časovou i paměťovou složitost zbytekn algoritmu.

Kupci se mezi sebou bndou vyrovnávat tak, že se vždy potká nějaký dluzník (to je ten, co má platit), s nějakým věřitelem (to je ten, co má peníze dostat) a zaplatí si nějakou částku.

Kolik si mžou zaplatit? Označme-li celkovou výši dlhu jako  $D$  a celkovou výši pohledávky jako  $P$ , mžou se vyrovnat nejvýše o min( $D, P$ ). Kdyžby si dřeli snad předat víc peněz, buď dojde k předlužení, nebo k přepřacení, což zadání zakazuje.

Jaké je nejmenší množstvř transakcí, které by teoreticky mohlo proběhnout? Jedna transakce má vždy dvě strany, ovřivň tedy dva kupce. Jednou transakcí tedy dojde k vyrovnování nejvýše dvou kupců. Je-li tedy kupců celkem  $N$ , je nejmenší množstvř transakcí k vyrovnání  $\lceil \frac{N}{2} \rceil$ .

Zadání nám tedy dovoluje provést  $N$  transakcí. To je ale velmi mále, protože každou transakci dokážeme alespoň jednou kupce vyrovnat. Stačí když převedeme nejvyšší množno částku: pokud min( $D, P$ ) =  $D$ , právě vyrovnová dluzníka; pokud min( $D, P$ ) =  $P$ , vyrovnová jsme věřitele.

Rýsuj se nám tedy jednoduchý algoritmus:

Nejprve si spočítáme, kolik má každý zaplatit, a rozdělíme vstup na dlhu a pohledávky, to vše v čase  $O(N)$ . Poté, dokud budou nějaké pohledávky a dlhuzy zbyvat, vybereme libovolný dlhu a libovolnou pohledávku (třeba první ze seznamu) a převedeme maximální možnou částku. Když už žádné pohledávky a dlhuzy nejsou, máme hotovo a mžžeme se jít klouzat.

Proč to funguje? Každým převodem vymlnujeme alespoň jednoho kupce, tedy s konkrétnř množstvřm převodů budou vyrovnováni všichni. Tež z toho plyne, že převodů provedeme nejvýše  $N - 1$ , takže jsme splnili zadání.

Jak je to rychlé? Předpracování trvá  $O(N)$ . Každý krok pak trvá konstantně čas: výběr dlhu, výběr pohledávky i samotné vyrovnání.

Paněťové složitost zahrnuje uložení několika proměnných, vstupu velikosti  $N$  a dvou seznamů o celkové délce takéž  $N$ , celkem tedy  $O(N)$ .

Některě řešitelé tvrdili, že je potřeba vstup seřadit, nedokázali ale přijít s žádným rozumným argumentem, proč by to mělo být potřeba. Je to lákavé, ale mžžeme neřdit a díky tomu si nezavřect do složitosti logaritmus  $O(N \log N)$  místo  $O(N)$ , a to za to stojí, ne?

[https://en.wikipedia.org/wiki/Subset\\_sum\\_problem](https://en.wikipedia.org/wiki/Subset_sum_problem)

Jiní se mě zase snažili přesvědčit, že po každé transakci musíme proběhnout celé pole a vyřadit z něj ty, kdo už platili. To však není vřbec potřeba. Vždyť jediné dva, kterých se to týká, jsme právě měli v ruce. Takovým postupem se dalo dosáhnout až (pro tuto úlohu jistě impozantní) složitosti  $O(N^2)$ .

Tež se objevilo několik řešení, kde řešitelé pouštěvali podmínky dané v zadání, například všechno přepočítávali, s dokonec porfilišli banku, která všechny zprostředkovala. Někteří se tohoto prohřšku dopustili skrytě, jiní to dokonce drze deklarovali. Pro odvážn postavě se organizátorům a hrně řešř jinou úlohu. Jest však neměl valného pochopení a udeloval pouze body třetky.

Úlohu jsme zadali úmyslně s požadavkem na maximálně dvě dvojnásobný počet převodů. Dřvod je prostř. Optimální řešení totiž získáme tak, že bychom našli rozdělení množiny kupců do co nejvíce podmnožin tak, aby součet v každé z podmnožin byl nulový. Nicméně už jenom otázka, jestli vřbec existuje podmnožina, jejíž součet je nulový, je NP-ŕpný problém (anglicky Subset Sum Problem),<sup>6</sup> jistě pochopřte, že po vás nemůžeme chtřt vymyslet takové dlhuzy v polyonomálním čase.

*Jan „Moskylor“ Matějka*

### 29-1-3 Strřřání zbrání

Na začátku si seřadíme gardisty podle výšky od nejnižšího po nejvyššího. Představme si, že si v tomto pořadí chodí vybrat kopř. Označme si  $m_i$  počet kopř kratšř než  $i$ -ř gardista (počítáme od nuly). První gardista si mžže vybrat ze všech kopř, které je schopen používat, má tedy  $m_0$  možností.

Druhý gardista je schopen používat  $m_1$  různých kopř. Ale vybrat si mžže pouze z  $m_1 - 1$ , protože jedno z nich už si vzal první gardista (první gardista je menšř než druhý, tedy kopř, které si vzal, by mohl používat i druhý, ten tak o jednu možnost přišel).

Analogický třetř gardista dokáže používat  $m_2$  kopř, ale dvě z nich už si zabrali první dva gardisté, má tedy na výběr z  $m_2 - 2$  možností. Obecně  $i$ -ř gardista si mžže vybrat z  $m_i - i$  kopř.

Počet možností, které má  $i$ -ř gardista, nezavřřs na tom, jaká kopř si vybrali předchozí. Tady je vidět, proč je dlhšířté gardisty brát v pořadí od nejmenšřho. Kdyžby byl první gardista výšř než druhý, počet kopř, které zbudou na druhého gardistu, by závisel na tom, jak vysoké kopř si vzal první.

Takhle víme, že pro každou z  $m_0$  možností výběru prvního gardisty existuje právě  $m_1 - 1$  možností volby druhého. Celkem tedy je  $m_0 \cdot (m_1 - 1)$  možností, jak si mžže vybrat první dvojici.

Obecně když provádřme několik výběrů po sobě a počet možností, ze kterých v každém kroku vybřřáme, nezavřřs na volběř v předchozřm krociř, je celkový počet možností prostě součinem počtů možností v jednotlivých krociř. Tomuto principu se řřka *pravidlo kombinatorická součř-mu*.

Dostáme tedy celkem

$$m_0 \cdot (m_1 - 1) \cdot (m_2 - 2) \cdot \dots \cdot (m_{N-1} - N)$$

(kde  $N$  je počet gardistů) možností, jak si gardisté mohou rozdělit kopř.

máme  $F[i]$ , pak výpočet  $F[i + 1]$  odpovřřá spuřření automatu na slovo délky  $i$  a při tom bndeme zpřetnou funkci pořbovat jen pro stavy délky  $i$  nebo menšř, pro které ji již máme hotovo.

Navřc nemůžeme pro jednotlivé přřky spouřřet výpočet vždy znovu od začátku – ( $i + 1$ )-ř přřfx je přeci prodloužením  $i$ -řtého přřfxu o jeden znak. Stačř tedy spustit algoritmus na jehlu bez prvního znaku a sledovat, jakými stavy bnde procházet – to budou přesně hodnoty zpřetné funkce.

Vyrovnání zpřetné funkce se nám tak nakonec zrethkorlo na jedné vyhledávání v textu o délce  $J - 1$ , a proto poběžř v čase  $O(J)$ . Časová složitost celého algoritmu tedy bnde  $O(S + J)$ . Dodáme už jen, že tento algoritmus poprvé popsali panové Knuth, Morris a Pratt a na jejíř počest se mu řřka KMP. Naprogramovaný bude vypadat následovně (čtenř vstupuj jsme si odpušřili):

```
jehla = "INSTINKT"
semo = "INSTINKTINSTINKT"
J = len(jehla)
S = len(semo)
F = [None] * J # Zpřetná funkce
def krok(i, znak):
    if 1 < i and jehla[i] == znak:
        return(i + 1)
    elif i > 0:
        return krok(F[i - 1], znak)
    else:
        return 0
```

```
# Konstrukce zpřetné funkce
F[0] = 0
for i in range(1, J):
    F[i] = krok(F[i - 1], jehla[i])
# Procházení textu
stav = 0
for i in range(S):
    stav = krok(stav, semo[i])
    if stav == J:
        print(i - J + 1, "az", i)
```

```
for i in range(S):
    stav = krok(stav, semo[i])
    if stav == J:
        print(i - J + 1, "az", i)
```

### Poznámky

- Pro anglický nebo český text je použitř takto softsitkovaná něho algoritmu skoro škoda, protože v obou jazycích se sířva jen málokdy, že bychom měli několik slov spojenců dohranady. Praktický bude stačit i na začátku zmíněný natvř algoritmus. Na soutěžích a olympiádách ale přřše raději algoritmus KMP.

- Heřování lze použřt i na vyhledávání řetězce v textu. Obzvlášř vhodné jsou na to *rolling hash functions* (neboli „oběknové heřovací funkce“), které umř v konstantním čase přřpocítat heř, ubereme-li nějaký znak na začátku a přřdáme-li jny na konci – jako kdybychom se dřvali na text skřz posouvájř se oknětko.

### Čvičení

- Rozmyslete si, že když vyhledáváme více slov, ne jen jedno, a algoritmus musř vypsat všechny výskřty na výstup, mžžeme se dobrat výšř než lineární složitosti v závislosti na vstupu. Na čten potom taková časová složitost také záleží?
- Vymyslete nějakou vhodnou okněnkovou heřovací funkci pro vyhledávání jedné jehly.

### Vyhledávání jehluřku

Co kdybychom neměli jen jednu jehlu/hledané slovo, ale celř jehluřek, čili seznam hledanců slov? I to lze řešř podobnou metodou, jakou jsme hledali jedno slovo. Tento algoritmus se nazývá po tvůrciř *algoritmus Aho-Corasickovř* a spočřva v tom, že jednoduchř spojový seznam nahradíme třři a to tre opět přřdáme zpřetné hrany.

Bndeme postupovat podobně jako u KMP. Nejprve naskládáme jehluřek do třře. Pro příklady v této kuchařce použijeme jehluřek **ARAB, ARARA, ARARAT, BAR, BARA, BARABA, BA a RAB**.

Dašřm krokem v KMP bylo sestrojení zpřetných hran. Nejprve jsme sestrořili zpřetnou hranu pro první znak slova, pak pro druhý atd. V třři to bude o něco složitějšř.

Na první pohled se mžže zdát, že bychom mohli automat sestrořit tak, že bychom vyrobili KMP pro první slovo, pak KMP pro druhé slovo s využitřm struktury prvního atd., ale to má háček.

Zpřetné hrany totiž nemusí vést do přředa. Napřřklad pro slovo **BARAB** povede zpřetná hrana do slova **ARAB**, z toho do slova **RAB** a z toho do **B**. Kdybychom ale zkonstruovali automat výše popsaným způsobem (a začali slovem **BARAB**), nebudě existovat v třři ani **ARAB**, ani **RAB**, takže bychom vedli zpřetnou hranu chyběně do **B**.

Mžžeme se ale opřřt o stejný trik, jako při konstrukci KMP. Bndeme opět vyhledávat nejdlšř vlast-ní suřfx. Kani dojde výpočet po jeho vyhledání, tam povede zpřetná hrana.

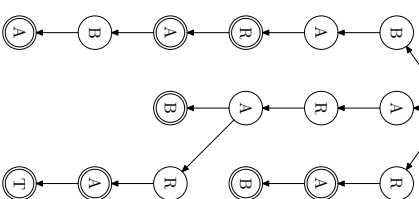
Zkusíme tedy nejprve sestrořit celou třři a pak postupně vyhledat nejdlšř vlastní suřfx pro každé ze slov. Ouhá, to ale také nefunguje. Když zkonstrueme slovem **BARABA**, a bndeme tedy vyhledávat **ARABA**, nalezneme v třři úspěšně přřfx **ARAB**, ale **ARABA** již v třři není. Měli bychom přejřt ze slova **ARAB** po zpřetné hraně, ale tu ještě nemáme zkonstruovano.

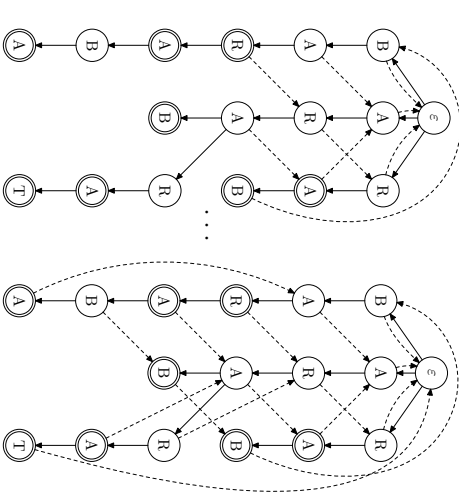
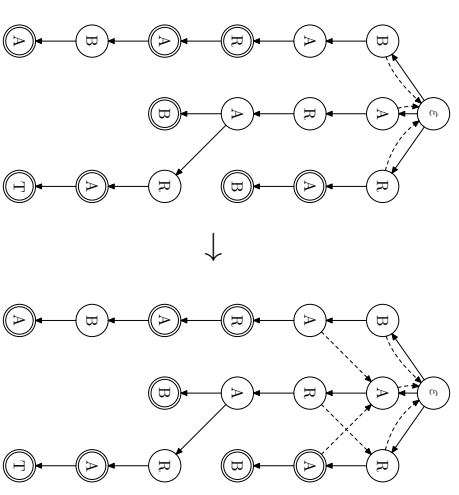
Rozdělme si třři na *vrstvy* – první znaky slov budou první vrstva, druhé znaky budou tvřit druhou vrstvu atd., až  $i$ -té znaky slov budou tvřit  $i$ -tou vrstvu.

Zpřetná hrana jistě povede do kratšřho slova. Z  $i$ -té vrstvy tak povede do vrstvy s nižšřm pořadovým číslem. Pokud tedy bndeme zpřetné hrany konstruovat po vrstvách, dojdeme křžněžno výsledku.

Jestř zbyřv otázka, jak konstruovat zpřetné hrany elektrně, když je musíme vyrábět po vrstvách. Mohli bychom prostě vzřt slovo, pro které hledáme zpřetnou hranu, utřhnout mu první znak a vyhledat. Jenže to budeme dělat spoustu práce zbytečně.

Napřřklad pro slovo **BARABA** bychom mohli vyhledávat **ARABA** v jřž zkonstruované části automatu, ale proč to dělat celé, když jsme při konstrukci předložili vrstvy vyhledávání **ARAB** při konstrukci zpřetné hrany pro **BARAB**?



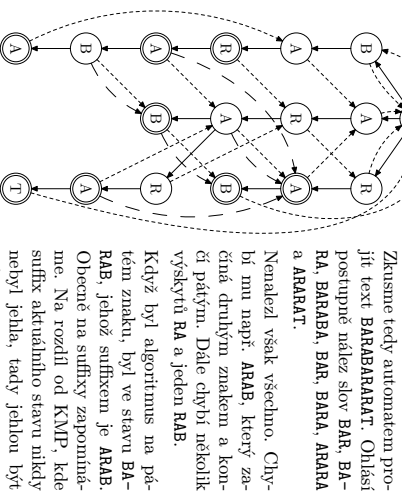


Při konstrukci další zpětné hrany tedy najdeme akordiát, kde jsme minule skončili, a odkamtam vrátnou tam přece vedle zpětná hrana. Takže nutíme postup shrnout do bodů:

1.  $c =$  poslední znak slova (znak stavu  $P$ , pro který hledáme zpětnou hranu);
2. přejdeme se do otce;
3. přičteme se po zpětné hraně;
4. dokud neexistuje syn se znakem  $c$ , nahájeme do kořene, přisouvráme se po zpětných hranách;
5. pokud existuje syn se znakem  $c$ , nahájeme do něj zpětnou hranu z  $P$ , jinak ji nahájeme do kořene.

Automat je zkonstruován. Časová složitost konstrukce sestává z konstrukce trie v  $O(J \cdot |\Sigma|)$ , resp.  $O(J \cdot \log |\Sigma|)$  (pokud použijeme binární strom ve vrcholcích) a z předpočítání zpětných hran. Při předpocítávání uděláme nějaký konstantní počet operací pro každý vrchol (celkem tedy  $O(J)$ ) a také paralelně vyhledáváme všechny jehly z jediného konstantního vyhledání nás slojí  $O(J)$ , resp.  $O(J \cdot \log |\Sigma|)$ .

Tedy konstrukce trvá celkem  $O(J \cdot |\Sigma|)$ , resp.  $O(J \cdot \log |\Sigma|)$ , paměťová náročnost je stejná jako u trie  $- O(J \cdot |\Sigma|)$ , resp.  $O(J)$ , přidali jsme jen  $O(J)$  zpětných hran.



Zkusme tedy automatem projit text **BARABARARAT**. Ohlášejí postupně názvy slov **BAR**, **BARA**, **BARABA**, **BAR**, **BARA**, **BARA**, **BARA** a **ARARAT**.

Nemalozá však všechno. Chybí mu např. **ARAB**, který začíná pátým. Dále chybí několik výskytů **RA** a jeden **BA**.

Když byl algoritmus na pátejn znaku, byl ve stavu **BARAB**, jehož suffixem je **ARAB**. Obecně na suffixy zapomínáme. Na rozdíl od KMP, kde sufixe aktuálního stavu nikdy nebyl jehla, tedy jehlou být může.

V každém stavu bychom tedy měli projít všechny suffixy a zkontrolovat je, jestli náhodou nejsou jehlami. Jak najdeme všechny suffixy? Projdeme postupně po zpětných hranách až do kořene. Má to jen jeden problém – je to pomalé.

Představme si například slovo obsahující  $A$  a  $AAAA \dots A$  (délky  $J - 1$ ). Budeme-li jim vyhledávat v textu  $AAAA \dots A$  délky  $S > J$ , projdeme prakticky pro každý znak až  $J - 1$  zpětných hran, čímž složitost naroste až na nepoužitelných  $O(S \cdot J)$ .

Všimneme si však, že většinou zpětných hran jsme prošli úplně zbytečně. Předpocítáme si tedy *zkratky* – z vrcholů vede zkratka do nejdelšího jeho suffixu, který je jehlou. Na obrázku jsou vyznačeny dlouze čárkovanými šipkami.

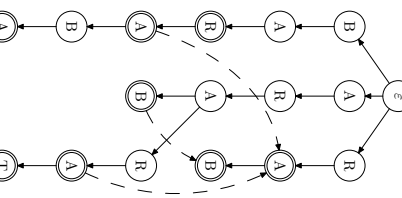
Předpocítání zpětných hran časovou složitost konstrukce automatu jistě nezhorsí, neboť vyžaduje v nejlhorším případě projít všechny zpětné hrany ještě jednou.

Potřebujeme-li ohlášejí všechny výskyt slovy včetně pozice, kde se nacházejí, jsme hotovi. Výsledná časová složitost prohledávání bude  $O(S + O)$ , resp.  $O(S \cdot \log |\Sigma| + O)$ , kde  $O$  je velikost výstupu – počet výskytů všech slov.

Celková časová složitost prohledávání včetně stavů automatu tedy bude  $O(O + S + J \cdot |\Sigma|)$ , resp.  $O(O + (S + J) \cdot \log |\Sigma|)$ .

Jak velký může být výstup? Obecně až  $S^2$ . Extrémně velký výstup je možné vygenerovat třeba slovním obsahujícím všechny prefixy slova  $AAAA \dots A$  délky  $S$  a sečením takřez  $AAAA \dots A$  délky  $S$ . Automat pak hlásí výskyt pro každé podřetězení, kterých je řádově  $S^2$ .

Pokud nám stačí u každého slova jen počet výskytů, nemusíme zohlednit – závislost na počtu výskytů umíme odstranit. Použijeme trik – na každé pozici započítáme pouze nejdelší jehlu, která tam končí (u každé jehly si budeme udržovat čísla). Nebudeme tedy v každém kroku poskakovat po zkratkách až do aleluja, ale

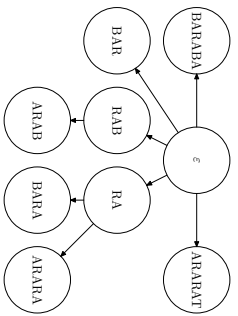


maximálně jednou. Délky tomu nám z časové složitosti zmizí velikost výstupu.

V našem příkladu se seštem **BARABARARAT** tedy na konci budeme mít uloženo, že **ARAB** se vyskytl  $1 \times$ , **ARARA**  $1 \times$ , **ARARAT**  $1 \times$ , **BAR**  $2 \times$ , **BARA**  $2 \times$  a **BARABA**  $1 \times$ . **RA** a **RAB** nemají hlášený žádný výskyt.

Nyní si zkonstruujeme strom jenom ze zkratek a pro každý vrchol spočítáme součet celého jeho podstromu.

Tedy po přepočítání bude mít **RA** tři výskytů a **RAB** jeden výskyt; celkový počet výskytů pak bude 12.



## Vzorová řešení první série dvacátého devátého ročníku KSP

### 29-1-1 Připletené placký

Než se pustíme do samotného řešení, tak si ulehčíme trošku úpravou. Namísto hromady plackiček otočovaných tím či oním způsobem budeme otáčet pole jedniček a nul.

Jednička bude reprezentovat plackičku otočenou správně (tedy spálením dolů), nula otočenou špatně. Jednou dovolíme operaci je přelápnout jedničky na nulu a nuly na jedničku všech hodnot od první až do nějaké  $i$ -té (včetně). Tedy vlastně negace prefixu (souvisle podčásteči začínající na levém okraji) pole.

Hned na začátku si je třeba uvědomit jednu velmi důležitou věc ze zadání. Nediceme o žádnou hodnotu, tedy překlápění plackiček, opravdu provést. Stačí nám říct, kde všude by to bylo třeba.

Toto pozorování způsobí, že naše řešení bude nakonec téměř a ne kvadratické, jak by se na první pohled mohlo zdát.

Pro začátek se pokusíme na pole hodnot podívat jako na soustavu stejnorodých úseček jedniček a nul. Snadno si za chvíli dokážeme, že předklápnutí hodnotu má smysl jen u prefixů, které končí právě v přechodech mezi sousedními úsečkami jedniček a nul (či obráceně).

Určitě platí, že v konkrétním stavu, tedy když jsou všechny hodnoty jedničky, nejsou v poli žádné přechody. Člá věž je totiž jeden stejnorodý úsek.

Otočením hodnot tiskem končíme na rozhraní určité zrušit je jeden přechod. Změníme totiž hodnotu těsně před rozhraním a beze změny ponecháme tu těsně za ním. Což, pokud máme jen dvě možné hodnoty, které se před tím lišily (a vytvářely tak přechod), znamená, že nyní musí být stejné.

Operaci také můžeme neovlivníme jiné úsečky před nebo za původním rozhraním. Teď za rozhraním se totiž ani nedotkneme a těm před pouze otočíme hodnoty.

### Poznámky

- Dalšími kroky po KMP a Aho-Corasickové jsou konkrétní automaty a regulární výrazy, o kterých jsme měli seriál ve 23. ročníku.
- Není moc rozumné snažit se implementovat Aho-Corasickovou v rozzápné době například při soutěži, pokud tento algoritmus nemáte opravdu pod kůží. Raději zkuste použít řešení, pokud budete něco takového potřebovat.

### Čvičení

- Redukci o velikost výstupu můžeme provést i pro případ, kdy výstup nebudeme vypisovat, ale stačí nám mít jej uloženy v paměti. Vynyslete vhodnou úpravu tržku s čískem.
- Zkusíte si implementovat Aho-Corasickovou vlastnoručně ve svém oblíbeném jazyce, alyste si být jisti, že doopravdy chápete všechny záležitosti tohoto algoritmu.

*Martin Bohm, Jan Matějka, Martin Mareš a Petr Skoda*