

Milí řešitelé, řešitelky a řešitelčátá!

Je to tady! Poslední série hlavní kategorie KSP byla doručena expresní hroší poštou až k vám. Pojdte zapřemýšlet třeba nad tím, jak efektivně šifit populašné zprávy nebo jak se vysílají konkrétné signály na UNIXu. A abyste se v textu neztratili, máme pro vás klučáček o vyhledávání v textu. To celé doprovozené příhodami například o roku 2000 či o zapomeném dluhu.

Připomínáme, že každému řešiteli, který v tomto ročníku z každé série dostane alespoň 5 bodů, darujeme propisku, blok a tužku. To vše s logem KSP.

Za úspěšné řešení KSP je možno být přijat na MFF UK bez přijímacích zkoušek. Úspěšným řešitelem se stává ten, kdo získá za celý ročník alespoň 50 % bodů (tedy alespoň 150 z maxima 300 bodů).

Termín série: Pondělí 1. června 2015 v 8:00 SELČ

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

Informace: Další podrobnosti o fungování KSP naleznete na <http://ksp.mff.cuni.cz/>. Pokud budete mít jakkoliv otázku, neváhejte se nás zeptat na ksp@mff.cuni.cz nebo na našem fóru. Přijeme hodně šéští! :-)

Odměna série: Čokoládu pošleme každému, kdo z úloh této série získá alespoň 42 bodů.

Pátá série dvacátého sedmého ročníku KSP

O čem bude průběh této série? Opět o nějaké velké programátorské chybě, která stála desítky miliard dolarů, nebo alespoň desítky biliónů? Ne, tentokrát ušetříme. Dočtete se kromě jiného o chybě nečísle, problému naprobém, jaký byl problém s rokem 2000 nebo! Y2K. Přesunme se nyní do Prvny krůtice před začátek druhého mléčání.

Vážení pane,
dovolujieme si Vás znovu upozornit, že dlužíte za pojištění 0,00 Kč. Tuto částku jste měl uhradit již před třemi měsíci a nečinil jste tak ani po třetí upomínce.

Zadáme Vás, abyste dlužnou částku zaplatil do 7 dnů. Pokud tak nečiníte, bude na Vás vymáhána soudní cestou.

Tento dopis dostal Karel od pojišťovny, když se jednoho dne vrátil z práce.

Co s ním měl dělat? Zahodit, jako ty tři upomínky předtím? V té pojišťovně asi nevědí, že zaplatit nula korun vyjde nastěpno jako nic, nezaplátit, což udělal. Vyhružka soudem přece není jen tak. Mohl by tam zajít osobně, ale je to docela daleko. Nebo to zkusit zaplatit...

A tak se stalo. Karel se rozhodl, že půjde na poštu a zaplatit to. Ještě si ale prohlédne ostatní dopisy. Hurá, má si uzavřednout baňk od manžičky z Anglie!

27-5-1 Šifení populašné zprávy 10 bodů

Firma, ve které pracuje Karlova matka v Anglii, má spoustu kančelářů a místnosti v jedné velké budově.

Jednou za čas tam testují populašný systém. Z centrály mohou zavolat telefonem naráž do dvou kančelářů poplašenou zprávou „Hoří, všichni rychle opusťte budovu!“ Naro z dalších kančelářů vyběhnou zaměstnanci do všech sousedních kančelářů a tam jim předají tuto zprávu. Všichni doběhnou ve stejném okamžiku, tomi celými můžeme říkat třeba jeden takt. V dalším taktu i z těchto kančelářů vyběhnou zaměstnanci do všech sousedních kančelářů, ve kterých se ještě o zprávě nedozvěděli, a tak dál.



Při poplačdu se ale smějí používat pouze chodly a schodiště označené symbolem zeleného uřičá, jeho panáčka, to jsou ty, které vedou k únikovému východu. Tyto chodly tvoří neorientovaný strom.

Nás by zajímalo, do kterých dvou kančelářů máme zavolat, aby se zpráva rozšířila co nejrychleji po celé budově.

V Anglii je mnoho velkých firem, které používají svůj starý informáční systém z 50. let napsaný v COBOLu. V té době byla paměť počítačů drahá, a tak se jí šetřilo, co to šlo. Proto například místo letopočtu ukladali pouze jeho posledních dvě číslice, protože ta 19 na začátku je přece všude stejná, to si každý domyslí.

To ale nepočítali s tím, že ty programy budou chtít používat i po roce 2000, a teď se bojí, že jim to přestane fungovat. Ale místo toho, aby si nechali naprogramovat svůj software znovu úplně od začátku, což stojí spoustu peněz a času, si našli najnovou programátora, který umí COBOL, a ty části, kde se používá datum, jim opraví. Za to mu dají spoustu peněz, hlavně aby měl jistotu, že to bude fungovat i dál.

Popytávka po takovýchto lidech je ale nejednou větší než nabídka. A tak jednou jeden pán z Anglie zanoval i Karlovi mamince. Ta umí programovat v COBOLu, protože byla sekretářka a za jejího mládí se COBOL i v Československu používal pro hromadné zpracování dat, tedy téměř pro všechny věci, jako dnes Excel. Zeptal se jí, jestli by pro ně pár měsíců nechtěla pracovat a přitom si pěkně vydělá. A tak je teď v Anglii.

Cestou na poštu šel Karel kolem autobusové zastávky u jednoho velkého supermarketu. Byla tam spousta lidí, kteří se připravovali na to, že 1. ledna 2000 přestanou fungovat všechny počítače na světě a začalá se naprosit chaos, nбудou lélat letadla, obchody budou zavřené nebo upravené, nastane třetí světová válka kvůli (pauočné) omylem odpáleným rakétám a tak dále. Říká se, že to všechno nastane kvůli tomu, že všechny počítačové programy na celém světě přestanou fungovat, protože nebudou umět zacházet s leto-počtem 2000.

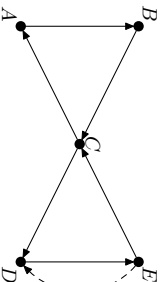
27-5-2 Survivalisté 12 bodů

Tito lidé (říká se jim survivalisté) si v supermarketu nakoupili velké zásoby trvanlivých potravin, aby v následující krizi přežili co nejlépe. Teď si krátí čas při čekání na autobus tím, že si navzájem ukazují, co si kdo koupil. Postávají v hloučcích, podávají si konzervy, balíčky nebo lahve, a ochutnávají.

Žádný z nich ale nechce podat svou věc úplně každému. Například jsou mezi nimi lidé, kteří se neznají, kteří odmítli ochutnat nabízenou věc nebo patří k jiné skupině, než k té, se kterou budou trávit konec světa v bunkru a následně obnovovat lidskou populaci. Takoví lidé se spolu nebaví a věc si nepřijímají. Survivalisty tedy můžeme popsat orientovaným grafem, z každého survivalisty vede hrana do těch, kterým je ochoten věc podat.

Teď každý člověk vytáhne z tašky jednu věc a předá ji někomu jinému, aby ji ochutnal. Zářovet ale chce, aby někdo jiný podal nějakou věc jemu. Na vás je, abyste rozhodli, zda je toto možné.

Například survivalisté na následujícím obrázku si takto věci předávali nemohou. Pokud bychom však přidali hrana $E \rightarrow D$, už by předávání mohlo proběhnout.



Na posmě byla hodně dlouhá fronta. Karel vyphál složenku na 0 korun a ztratil se.

27-5-3 Čekání na poštu 9 bodů

Když je na posmě tolik lidí, že není snadné poznat, kde fronta začíná a kde končí nebo kdo je za kým na řadě, rozmístí na podlahu sloupky a natáhnou mezi nimi páska. Na začátku a na konci fronty je samozřejmě mezera, kterou se prochází, ale pro jednoduchost si představa, že i tam je páska, i když poněkud. Sloupky a páska tak tvoří mnohoúhelník, který neprotíná sám sebe. Lidé čekají uvnitř tohoto mnohoúhelníka.

Váš program dostane na vstupní rozmístění sloupků, tedy N bodů v rovině (nebudou tvořit přímlu). Na vás je, abyste z nich vytvořili mnohoúhelník, který neprotíná sám sebe.

Formát vstupu: Na první řádku dostanete číslo N , tedy počet bodů. Poté následuje N řádků, kde na každém řádku mezerou oddělená x -ová a y -ová souřadnice jednoho bodu. Pro všechny vstupy platí $3 \leq N \leq 250\,000$, $0 \leq x, y \leq 1\,000\,000$.

Formát výstupu: Výstupem programu bude jediný řádek obsahující čísla bodů oddělená mezerami setřazená podle výskytu na obvodu mnohoúhelníka. Body jsou číslovány od nuly podle pořadí na vstupu.

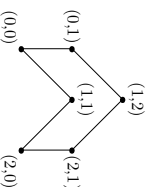
Pozor, výstup nemusí být jednoznačně určený (vypíše libovolný konekt).

¹ <http://ksp.mff.cuni.cz/viz/codex>

Ukázkový vstup:

Ukázkový výstup:

```
6
0 0
2 0
1 2
0 1
2 1
```



Toto je praktická open-data úloha. V odevzdávacím systému si nechtáte vygenerovat vstupy a odevzdaté příslušné vstupy. Záleží jen na vás, jak vstupy vyrobíte.

Konečně se Karel dostal na řadu.

„Dobrý den, já bych chtěl poslat tuto složenku.“ Podal ji pošťáček, ta se na ni podívala, vyfukla ze stojanky průzkou a podala mu ji se slovy: „Tady máte novou složenku, vyphíte to prosím znovu, napsal jste tam nula korun.“

„Paní, to není vyphné špatně, já opravdu chci poslat nula korun.“

„Ale to nechte. Nula korun se nedá poslat. Nejméně, co můžete poslat složenku, je 10 haléřů. Ale k tomu zaplatíte ještě 10 korun poštovně.“

„Ale to musí být, vždyť mně přišla upomínka na nula korun od pojišťovny, podívejte se. Když ty peníze nepošlu, dáji mně k soudu.“

„Bohužel, nejde to. Chcete poslat 10 haléřů? Nebo nic?“

„Vy mi tu složenku nepošlete? Tak mi prosím zavolejte ředitelce této pobočky.“

A tak dále. Karel se rozhodně nenechal odbyt. Ředitelce sice nezavolał, protože již nebyl v práci, místo toho ale slábněl nějakého zaměstnance, který prošel speciálním školením o poštovním řádku.

27-5-4 Školení zaměstnanců 10 bodů

Pošta, jako každá firma, má hierarchickou strukturu zaměstnanců, a tato struktura tvoří strom.

Manažeři ze školního oddělení čas od času vysílají některé zaměstnance na speciální školení, kde se naučí, jak se vypořádat s určitými situacemi, které někdy mohou nastat, ale nejsou tak obvyklé, aby školení využil každý zaměstnanec. Navíc zjistili, že pro zaměstnance je jednodušší se nové věci dozvědět od svých kolegů než se ptát úplně cizího odborníka. A počítají i s tím, že vyskolení se budou chlubit novými znalostmi kolegům, a tak se informace snadno rozšíří po celé firmě.

Konkrétně si řekli, že by bylo dobré, aby se každý zaměstnanec mohl obrátit k vyskolenému zaměstnanci, který prošel například týdenním kurzem razítkování dopisů, přes méně než K spolupracovníků. To znamená, že každý vrchol stromu musí ležet do vzdálenosti K od nějakého zaměstnance, který prošel školením.

Na vás je, abyste pro zadaný strom a číslo K vybrali co nejméně vrcholů tak, aby každý vrchol stromu ležel do vzdálenosti K od jednoho z vybraných vrcholů.

Tato úloha je praktická a řeší se ve vyhodnocovacím systému CodeX.¹ Přesný formát vstupu a výstupu, povolené jazyky a další technické informace jsou uvedeny v CodeXu přímo u úlohy.

Po důkladném prostudování poštovních předpisů a po čtyřech minutách převádění, mezitím, co se lidé z fronty

S programátorskými chybami ale bohužel konec jenom tak nebude. Chyby dělájí lidé, na zlomyslné počítače, a hláskou chyba nemůže uplně vyloučit, dokud z programování nevyhodíte hláskový faktor.

Dalšížte je se z chyb poučit. K podobnému problému by mohlo v budoucnu dojít ještě jednou. Mnohé programy používají znakový 32-bitový typ, který reprezentuje čas jako počet vteřin od 1. 1. 1970, a ten přeteče 19. ledna 2038. Ale pokud do té doby začnou používat 64-bitový typ, přesovová se problém do roku 292277 026 596.

To ale neznamená, že nemáte řešit KSP. S chutí do toho!

Přibít vyprávěl

Dominik Macháček

27-5-7 Shellová automatizace 15 bodů

Poslední díl letošního seriálu o UNIXu a jeho příkazovém řádku se posere v duchu automatizace tímto a lepšího provázání našich skriptů se systémem. Ukážeme si třeba, jak lze spustit stejný příkaz na všech souborech nějakého typu, jak v nějakém složitějším procesu zpracování dat (nebo třeba kompilace programu) zpracovávat jen to, co ovlivní změněné soubory, a také jak například zajistit, aby děle běžící skript po sobě uklidil, pokud se ho rozhodneme ukončit v průběhu práce.

Všechno to jsou drobné pomůcky, které nám krásně zapadnou do všeho ostatního, co jsme se přes rok naučili, a pomohou vám ještě lépe využívat sílu shellu. Pokud se tedy nebojíte, raději vstoupíte.



Hledání souborů

V minulých dílech jsme vám ukázali, jak třeba pomocí wildcardů vybrat všechny soubory v aktuální složce s příponou .pdf. To pomoci můžeme jednoduchše, horší to ale začít být ve divilí, kdy chceme prohledat rekurzivně třeba i všechny podsložky včetně jejich podsložek a tak dále.

Asi by se dal napsat nějaký skript, který by si nedal vypsat příkazem ls všechny složky a na nich by se zavola znovra, ale existuje mnohem snadnější řešení. Zkusíte si třeba ve svém domovském adresáři spustit příkaz find

```
./poznanky.txt
./obrazky
./obrazky/Kevin.jpg
./obrazky/Sara.jpg
./obrazky/Peir.jpg
./Reseni
./Reseni/KSP
./Reseni/KSP/27-5-7.pdf
...
```

Jak vidíte, find vám vypsal uplně všechny složky a soubory ležící ve stromě souborů pod aktuálním umístěním. Pokud mu totiž nezadáme, jakou složku má prohledávat, tak použije aktuální adresář. (a také ho vypíše jako první

² Dokonce může být vnější filtrování pomocí grepu výrazně pomalejší, protože se mezi těmito dvěma příkazy musí přemést skrz rovnou velké množství věcí.

prohledat a připojí ve výpisu před všechny nalezené složky a soubory). Když bychom ve stejném umístění spustili třeba příkaz find Reseni, výpis by pak vypadal takto:

```
Reseni
Reseni/KSP
Reseni/KSP/27-5-7.pdf
```

To je pěkné, na takovýto výstup bychom už mohli použít třeba příkaz grep a vyfiltrovat z něj s trochu práce třeba jen PDF soubory. Ale find to umí sám a ještě spoustu věcí navíc.²

Než se pustíme do dalšího experimentování, tak se na příkaz find podíváme i se všemi jeho skupinami parametrů:

```
find <kde hledat> <kritéria> <prováděný příkaz>
```

• První skupina je asi jasná, určuje místo, kde má find začít své prohledávání. Je možné předat i více umístění, find je prohledá všechna. Pokud není žádné umístění zadáno, tak se použije aktuální adresář. ., jak jsme ukázali výše.

• Druhá skupina parametrů nastavuje různá kritéria omezení výběr souborů a složek. Pokud není nic nastaveno, nefiltruje se nic a vypisují se všechny nalezené složky a soubory. Tím se budeme zabývat zápatí.

• Třetí skupina parametrů určuje, co se má pak s nalezenými názvy souborů a složek dít. Pokud nevolíme žádnou akci samu, tak find použije akci -print a vše jen vypíše na výstup. Můžete si zkusit ho s touto akcí na konci jeho seznamu parametrů spustit.

Mezi další jako akce patří pak třeba formátovaný výpis nebo spouštění nějakého příkazu. Tomu se budeme věnovat po kritériu výběru.

Mocnější find

Co kdyby nás zajímaly všechny README soubory třeba ve složce /etc? V tu chvíli nám stačí použít kritérium -name a spustit příkaz:

```
find /etc -name README 2>/dev/null
```

Protože složka /etc obsahuje pravděpodobně několik podsložek, na jejichž čtení nebudeme mít práva, může nám find vynadat několik chybových hlásekami (jednou za každou nepřístupnou složku). Aby se nám výstup mezi těmito hláskami neztratil, je dobré vzpomnout si na minulé díly seriálu a chybový výstup „odfiltrovat“ jeho přesměrováním do /dev/null, jak jsme v příkladu výše rovnou udělali.

Můžete dokonce použít i shellové wildcardy ke specifikování názvu. Jen pozor na to, že wildcardy se musí dostát až k findu, nesmí je tedy zpracovat už samotný shell a tedy je nutné je buď escapovat, nebo zabalit do úvozovek:

```
find -name "*.pdf"
find -regex ".*[/\[*\]\.pdf"
```

Pokud by vám nestačily shellové wildcardy, je možné podobným způsobem použít i regulární výrazy, ale už s jiným přepínačem. Následující příkaz najde všechny PDF soubory zakínající od písmene a:

```
find -type f -regex ".*"
```

Mezi další zajímavá kritéria patří například specifikace typu souboru pomocí přepínače -type (-type d je složka, -type f běžný soubor, více v manuálové stránce). Věhni

Výsledková listina čtvrté série dvacátého sedmého ročníku KSP

řezitel	škola	ročník	série	4-1	4-2	4-3	4-4	4-5	4-6	4-7	série	celkem	
1.	Jan Špaček	G Wicht	4	9	10	11	11	12	11	14	59,0	232,9	
2.	Richard Hlařík	GOAMaLaz	2	14	10	11	11	12	11	14	59,0	232,9	
3.	Stanislav Lukáš	GPJančekPH	2	5	10	7	11	12	11	9	52,3	212,6	
4.	Marek Černý	G Chrdim	4	9	10	7	11	12	11	2	44,9	206,6	
5.	Marlín Schenbrem	G Mňahn Třb	3	4	4	7	11	12	11	8,5	36,4	183,7	
6.	Štěpán Hudeček	G Lioveš	3	4	4	10	11	11	12	11	50,8	176,2	
7.	Václav Volhejn	GKexpleraPH	2	14	10	11	10	11	12	11	53,6	162,8	
8.	Jakub Tětek	CirkG Pizeň	1	4	2	2	4	11	12	11	23,3	156,9	
9.	Michal Števráhl	GKlatovy	2	2	4	8	10	7	12	11	51,0	153,0	
10.	Přemysl Štrastný	GZanbek	2	7	2	2	2	4	12	12	21,7	143,0	
11.	Václav Štraier	GČeskolipPH	2	4	4	9	10	11	12	6	15,6	142,6	
12.	Jan Tománek	GPelhmov	4	4	4	4	6	7	8	0	8	34,4	130,7
13.	Michal Töpfer	G D-JPaMB	2	4	4	6	7	8	10	0	8	34,4	130,7
14.	Lukáš Uhlíř	SSSVTPraha	4	3	3	3	3	3	3	3	0,0	114,5	
15.	Jan Kočur	G Wicht	4	3	3	3	3	3	3	3	0,0	113,4	
16.	Adrián Goga	SPSvitra	4	3	3	3	3	3	3	3	0,0	110,8	
17.	Jan Knížek	G Strakon	4	4	16	10	11	12	11	8	10,6	68,1	
18.	Jakub Zárubnický	GTomkovaOL	4	4	8	10	11	12	11	11	11,0	33,0	
19.	Anna Gařdlová	GFPeValMaz	4	4	4	4	4	4	4	2	0,0	46,2	
20.	Róbert Selvek	G KošiceS	3	3	3	3	3	3	3	3	0,0	45,2	
21.	Pavel Turinský	G Brandýs	2	2	2	2	2	2	2	2	0,0	45,2	
22.	Jiří Sejkora	GVošarPH	3	3	3	3	3	3	3	3	0,0	45,2	
23.	Jiří Gocník	GŠkodyPR	3	3	3	3	3	3	3	3	0,0	45,2	
24.	Václav Rozhoň	GJihsláčCB	4	4	9	9	9	9	9	9	55,0	68,6	
25.	Jiří Vozár	G UherBrod	3	4	4	4	4	4	4	4	10,6	68,1	
26.	David Cholewa	GMatOS	4	2	2	2	2	2	2	2	0,0	46,2	
27.	Jan Bourček	GKexpleraPH	2	3	3	3	3	3	3	3	12,0	45,2	
28.	Václav Komrtdký	GSOS Fmls	4	4	4	4	4	4	4	4	10,2	43,1	
29.	Jan Pokorný	G BrnoVice	3	6	9	9	9	9	9	9	9,5	41,0	
30.	Martin Zoula	GNačkavaPH	3	3	3	3	3	3	3	3	0,0	38,6	
31.	Barbora Sedláková	GKonštanPv	4	3	3	3	3	3	3	3	8,5	35,9	
32.	Jan Soukup	GKlatovy	4	3	3	3	3	3	3	3	11,0	33,0	
33.	Eva Matoušková	G_Skolov	4	1	8	8	8	8	8	11	23,4	23,4	
34.	Filip Bičals	GOpatovPHA	2	5	5	5	5	5	5	5	23,4	23,4	
35.	Vít Macura	GOAMaLaz	2	2	2	2	2	2	2	2	0,0	20,0	
36.	Jakub Lukáš	GNAlejPH	4	1	1	1	1	1	1	1	0,0	18,9	
37.	Dalimil Hájek	GKepleraPH	2	1	1	1	1	1	1	1	0,0	14,0	
38.	Marlín Kuběša	GŠkodyPR	3	1	1	1	1	1	1	1	0,0	13,4	
39.	Jakub Matěna	GČeskolipPH	3	2	2	2	2	2	2	2	0,0	13,4	
40.	David Jurtča	GNadStbOPH	2	1	1	1	1	1	1	1	1,0	11,0	
41.	Jan Kaifer	GČesBrod	-1	2	2	2	2	2	2	2	0,0	10,9	
42.	Matěj Konečný	GJhrovoCB	4	1	2	2	2	2	2	2	0,0	10,6	
43.	Jan Barča	G_Holice	1	1	1	1	1	1	1	1	0,0	10,0	
44.	Václav Steinhäuser	GDačice	1	1	1	1	1	1	1	1	0,0	9,0	
45.	Josef Vávra	Slec	4	1	1	1	1	1	1	1	0,0	7,9	
46.	Jan Mláz	G_Holice	1	1	1	1	1	1	1	1	0,0	5,7	
47.	František Dostál	VSPSeOC	4	1	1	0	0,5	1	0	1	4,4	4,4	
48.	Michal Novák	SSSVTPraha	4	1	1	1	1	1	1	1	0,0	4,0	
											0,0	2,0	

Chcete-li s námi komunikovat, sdílejte své zkušenosti, můžete si otevřít náš HTTPS certifikát – jeho SHA1 fingerprint je: OE:D9:B6:EE:6F:B0:51:D9:66:EB:99:99:5F:5F:99:D6:FD:A3.

číst ze vstupu, jejím spuštěním na pozadí ji okamžitě pozastavíme znovu. Příkaz `bg` ale má stále své využití. I úlohy běžící v popředí totiž můžeme z terminálu snadno pozastavit. Věšinou stačí značkou `Ctrl+Z`. Hned si vysvětlíme, jakým způsobem toto pozastavování funguje.

Připomeňme si ještě, co bylo o UNIXových procesech vime. Letmo jsme se s nimi seznámili ve druhém díle. Po vědětli jsme si, že systém si pro každý proces pamatuje jeho identifikční číslo PID, stav, seznam otevřených souborů, uživatele, s jehož právy běží, a mnoho dalších informací. Seznam všech běžících procesů dostaneme příkazem `ps ax`. Tyto znalosti využijeme v další části.

Pokud bychom chtěli počkat na dokončení některého procesu běžícího na pozadí, poslouží nám k tomu interní příkaz `wait` PID. Pokud žádný parametr nedostane, počká jednoduše na všechny podprocesy.

Signály a meziprocesová komunikace

Možná se ptáte, jakými prostředky spolu výhbec mohou procesy komunikovat. Používali jsme již routy (a to jak nepojmenovanou, tak pojmenovanou vzhledem příkazem `mft1-f0`). Z shellu si ale můžeme snadno vytvořit ještě jeden komunikační kanál, jsou jim signály.

Signál si můžeme představit jako takové malé štronchnutí. Procesy si je mohou mezi sebou navzájem posílat a tím si vlastně předávat informace. Rozhodně to není způsob, kterým byste měli přenášet kilobyty dat (načto více). Signály slouží spíš k upozornění na asynchronní události. Na běžném dřemím Linuxu jich najdeme 64, na OpenBSD jen 32. Každý signál má své jméno a číslo. Pozor na to, že různé operační systémy mohou mít přiřazené čísel signálů různé.

A k čemu se signály hodí? Například při vypnutí počítače by bylo dobré dávat všem programům vědět, že se mají vypnout a případně uložit rozdělanou práci. K tomu se používá signál `SIGTERM`. Pokud by na to program nereagoval a nechtěl se vypnout, můžeme jeho činnost ukončit natvrdo signálem `SIGKILL`.

Podobně existují také signály `SIGTSTP` a `SIGSTOP`, které způsobí zastavení běžícího procesu a naopak `SIGCONT`, pro opětovné spuštění. `SIGTSTP` posleme procesu právě pomocí stisků `Ctrl+Z`. Naopak `Ctrl+C` posílá `SIGINT`.

Za zmiňku stojí ještě `SIGHUP` a `SIGCHLD`. Kdykoli nějaký proces skončí, je na to upozorněn jeho rodič signálem `SIGCHLD`. Stejně upozornění přijde i v případě, že je synovský proces pozastaven, případně znovu spuštěn. `SIGHUP` má hned několik významů. Tento signál obdrží programy, pokud jim zavřeme terminál, ve kterém běží. Věšina aplikací se proto ukončí. U *daemonů*, programů určených k tomu, aby běžely na pozadí a s člověkem nekomunikovaly pomocí terminálu, `SIGHUP` obvykle způsobí znovuvytváření konfigurace z konfiguračních souborů.

Chceme-li proces poslat signál, stačí z shellu zavolat `kill [-signál] PID`. Pro oznáčení signálů můžeme použít jak číslo, tak název. Pokud signál nespecifikujeme, posílá se `SIGTERM`. Jako identifikátor procesu můžeme zvlášť 1-4. Potom je daný signál poslán úplně všem procesům, kterým můžeme nějaký signál poslat. Signály totiž můžeme posílat pouze svým vlastním procesům. (Jenom root má výjimku a umí signalizovat všem.)

Signály mají za sebou bohatou historii a jejich význam se průběžně trochu měnil. Například někdy naznačme na

to, že signály nezakládají na SIG. Máme pak `INT`, `TERM`, `TSTP`, ... Pokud by vás zajímaly detaily, podívejte se do manuálových stránek `man 7 signal`. V zásadě můžeme signály dělit podle toho, jestli proces ukončí, ukončí a vytvoří obraz jeho paměti (tzv. core dump), pozastaví znovu spuštění, případně jestli jsou ignorovány a nedělají nic.

U větších signálů si proces může výchozí chování přenastavit. Jednou výjimkou jsou `SIGKILL` a `SIGSTOP` – ty vždy znamenají to samé. Diky tomu jde každý proces ukončit, případně zastavit. Typicky chceme některé signály ignorovat, nebo odplyt vlastní funkcí. Pokud pak náš proces obdrží signál, operační systém přemění aktrnálně vykonávanou práci a spustí naši funkci. Signály můžeme oddělyávat i v shellu.

Čináme na signál

Abychom mohli signály zachytit, potřebujeme na něj nejprve nalít past. Jednoduše pomocí `trap` řekneme, co se má provést v případě, že daný signál přijde. Například po zavolání `trap "echo bar" SIGUSR1` shell vypíše na svou výstupní napís „bar“ vždy, když obdrží signál `SIGUSR1` nebo `SIGINT`. Zavolání `trap bez parametru` vypíše, jaké příkazy bude shell při kterém signálu provádět.

Úkol 4 [3b]: Napište skript, který vám pomůže se snyšaplým využitím desítek víceřádkových počítáčů naplo. Váš skript dostane jediný parameter `N`. Na standardním vstupu pak bude číst podobně jako shell příkazy a bude je spouštět. Řádky by však měl provádět paralelně. Vždy se smí provádět nejvýše `N` řádek současně.



Proč / proč?

Aby nebylo nutné vytvářet pro zjištění všech možných informací specializovaná systémová volání, vznikl v Linuxu virtuální filesystém `/proc`. V `proc` najdeme pro každý spuštěný proces adresář s celou řadou zajímavých souborů. Například v `/proc/PID/fd` najdeme jako symboliku seznam všech souborů otevřených daným programem. Podrobnosti najdete v manuálu `man 5 proc`.

Úkol 5 [2b]: Napište malou náhradu programu `ps ax`. Na výstupu vašeho skriptu by se měl objevit o každém procesu jeho PID, nějaký identifikátor uživatele, příkaz i se všemi jeho parametry a alespoň jeden další údaj podle vašeho výběru. Můžete si vybrat cokoliv, co se vám bude zdát aspoň trochu zajímavé či užitečné. Všechna potřebná data čtete přímo z `/proc`.

Make – základy

Poslední velký pomocník, kterého si letos ukážeme, je příkaz `make`. Ten se stará o automatizaci procesu nějaké kompilace, překladu nebo třeba jen zpracování dat. Věšinou je třez souborem s názvem `Makefile` (všimněte si velkého prvního písmena) v adresáři, kde `make` zavoláme.

A co že přesně umí? Základem celého procesu je, že se `make` pokouší spustit nějaké `Makefile` definované cíle, což většinou znamená vytvořit nějaké soubory. Pokud `Makefile` zjistí, že požadovaný cíl ještě neexistuje, zkuste ho podle pravidel uvedených v `Makefile` vytvořit. Pravidla pro výrobu cíle

Úkol 1 – Náhodné dvojice

Věšina řešitelů se shodla na základním postupu: vstoupí soubor zamícháme (pomocí `shuf`) a poté spárujeme sousední dvojice řádků. Ukážeme si hned několik způsobů, které se mezi řešeními objevily.

Asi nejjednodušší na vynyšlení je použití bashového cyklu a `read`:

```
shuf | while read a; do
    read b
    echo "$a:$b"
done
```

O něco elegantnější a stále jednodušší se dá využít `sed`:

```
shuf | sed -re 'N; s/\n/;/'
```

Jak již víme, příkaz `N` načte další řádek a přilepí ho do paternu space odděluje znakem `\n`. Ten stačí nahradit dvojitě-konou a jsme hotovi. Při příští iteraci se pokračuje nejbližším ještě nepracovaným řádkem, tedy třetím, pak pátým, atd. Velmi podobná věc se dá udělat i pomocí `awk` (inspirováno řešením Štěpána Hudečka):

```
awk '{ print "%s: ", $0, getline; print; }'
A na závěr jedno mile bláznivé řešení podle Jakuba Těhla:
shuf | awk '{ORS = NR%2 ? "\n" : "\n\n" } 1'
```

Zkusíte schválně vynyšlet, co dělá ta jednička na konci :-).

Tato úloha byla původně zamýšlena jako cvičení na `paste`, leč ukázalo se, že naše řešení je výrazně složitější než ta ukázaná výše. Základem myšlenka je zamíchat `vstup`, rozdělit do dvou souborů jeho první a druhou polovinu, a ty poté spojit do dvojice právě pomocí `paste`:

```
shuf > jmena.rand
half=$(( $(wc -l jmena.rand) / 2 ))
head -n $half >prni.tmp
tail -n $half >druzi.tmp
paste -d: prni.tmp druzi.tmp
```

Úkol 2 – Převracená slova

Nejprve vyřešíme nalezení vyhovujících slov. Výběr nejdlešho doplníme na konci. Dobrym začátkem určitě bude vytvořit si soubor s převracenými verzemi všech slov ze slovníku, což nám zaručí onen podvrtý příkaz `rev`. Snadno si rozmyslíte, že hledaným řešením jsou právě slova, která se vyskytují jak v původním slovníku, tak v tomto pomocném souboru. K hledání přímků dvou souborů přimocíte posloupí příkaz `comm` – jen si nejdříve oba musíme seřadit. Řešení by mohlo vypadat takto:

```
sort -n slovník >slovník.sort
rev slovník | sort -n >slovník.rev
comm -12 slovník.sort slovník.rev
```

Některé řešitelé přišli s alternativním postupem, který namísto `comm` používá příkaz `uniq`. Hledání přímků můžeme provést taky tak, že oba soubory (původní slovník a jeho reverz) slepíme do jednoho a z něj vybereme všechny duplicitní řádky. K tomu lze (po seřídění) použít příkaz `uniq -d`, který příkazuje vypisovat pouze řádky s více než jedním opakováním.

Alternativně lze použít `uniq -c -a` z výsledku odgropovat všechny řádky s technickou na místě počtu opakování. Za předpokladu, že původní slovník neobsahoval duplicity, jde o výskyt duplicitního řádku musel pocházet z jednoho souboru a druhý z druhého, tedy patří do přímků. A při-

padných duplicit se snadno na začátku zbytečně příkazem `sort -u`.

```
{ sort -u slovník; rev slovník | sort -u; } \
| sort | uniq -d
```

Nyní k výběru nejdlešho. To snadno vyřešíme jednodušším skriptem pro `awk`, který bude v nějakých proměnných průběžně udržovat dosud nejdleší slovo a jeho délku a na konci ho jen vypíše:

```
awk '{
    if (length > maxlen) {
        maxlen = length;
        word = $0;
    }
}
```

```
END { print word; }'
```

Pokud vám přijde toto řešení příliš „céčkové“, dá se postporovat i jinak. Nedaíme `awk` jen připsat ke každému řádku jeho délku a pak využijeme příkaz `sort`. Takové řešení má síce složitost $O(N \log N)$ namísto lineární, ale to nás v praxi příliš netrápí.

```
awk '{print length,$0}' | sort -nr | head -n 1 \
| cut -d' ' -f2-
```

Úkol 3 – Seriaty

Tato úloha byla spíš technickým cvičením a asi se nedá vynyšlet úplně hezky a elegantně řešen. K úloze lze přistoupit ze dvou stran:

a) Projít soubory v číselném adresáři z každého z nich zkusit vytáhnout správná čísla, najít informace v seznamu epizod a přejmenovat ho.

b) Projít stažené seznam epizod, pro každou se pokusit najít odpovídající soubor a přejmenovat ho.

Věšina řešitelů se přiklonila k první variantě. Ta se ale ukázala poměrně nebezpečnou, neb téměř nikdo nečísil, co se stane se soubory, je-li-liž název není ve správném formátu, případně pokud příštná epizoda není nalezena v seznamu. Ve většině případů to skončilo tím, že se čísla seřít a epizody napsávaly jako prázdne řetězce, název se v seznamu nemasá a z toho vznikl název typu `SO0E00_ -_av1`. Pokud by bylo v adresáři nerozpoznávaných souborů více, všechny byly přejmenovány na tento stejný název, čímž se navzájem přepsaly. Například věšina skriptů byla obohacena takto přeměnovat i sebe sama, pokud si je člověk uloží do stejného adresáře, připadne i své pomocné soubory.

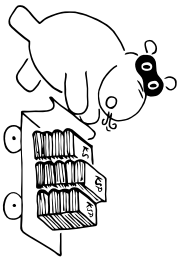
Tento úkol měl přesvědčit, že opravování KSP je riziková činnost. Jeden z doských skriptů totiž začal takovýmto způsobem přeposlat soubory v mém domovském adresáři. Za to mohl i společnost s problémem popsaným výše) neměnně vypadáající příkaz na začátku skriptu:

```
cd $1
```

Pokud takový skript omylen spuštěte bez parametru, dle pravidel bashové expanze se `$1` zcela zahoří a spustí se příkaz `cd` bez parametru, který skóčí do domovského adresáře uživatele.

Samozřejmě všechny tyto problémy se dají osčřit (testováním, zda soubory mají správný formát názvu, přidáním uzovozů na vhodné místa, použitím `mv -i`, atd.), ale je to docela otrava a není jednoduché na mě nezapomenout. Proto si raději ukážeme druhý způsob, který těmto neduhům netří.

Jak jste možná poznali, přесonování kotoučů ze zadního odpovídá slavné úloze o Hanojských věžích, kterou poprvé popsal francouzský matematik Édouard Lucas, ⁶ Pro úplnost dodáme, že v původní podobě úlohy se mistra A, B a C nezývají výčemi, kotoučů je celkem čtyřnáctdesát, všech my kotouče jsou z ryžlo zláta a na začátku jsou umístěny na jedné tyči. Michalové každých den jeden kotouč přесonou z tyče na tyč (skladě plati, že nesmí být větší kotouč položen na menší). V okamžiku, kdy se jim podaří přесonit všechny kotouče na třetí tyč, nastane podle legendy konec světa.



Než se pustíme do samotného uspořádávání fotografií pořízených při střehování dat, potřebujeme porozumět tomu, jak se kotouče spravně přесonují. Konkrétně si ukažeme, jak lze přесonit všechny kotouče na co nejmenší počet tařhů, a společně s tím dokážeme, že se jedná o jediný možný postup s nejmenším počtem tařhů. To nám dá jistotu, že tenatý postup použijí i pracovníci vystupující v zadání úlohy. Celý postup popíšeme rekurzivně (pokud si s rekurzí příliš nerozumíte, můžete v případě nesmáři nahlídnout do naší kuchárky o základních algoritmech). ⁷ Celkové máme N kotoučů, očísujeme si je od nejmenšího po největší čísly 1 až N.

Funkce *presun(k, vychozí, cilová, pomocná)*:

1. Pokud $k = 0 \rightarrow$ skonči
2. Zavolej *presun(k-1, vychozí, pomocná, cilová)*
3. Přесun kotouč k z tyče *vychozí* na tyč *cilová*
4. Zavolej *presun(k-1, pomocná, cilová, vychozí)*

Funkce *presun(k, vychozí, cilová, pomocná)* předpokládá, že na tyči *vychozí* se nachází kotouče 1, 2, ..., k , a začítí přесun všech těchto kotoučů na tyč *cilová*. Nyní pomocí matematické indukce dokážeme, že nejmenší možný počet tařhů pro přесunuti všech N kotoučů je $2^N - 1$ a lze jej dosáhnout pouze pomocí nashého postupu.

Když nemáme žádné kotouče, přесuneme je pomocí nula tařhů a je to jediný možný způsob přесunu (nebo spíše ne přесunu). Platí rovnost $2^0 - 1 = 0$ a základí indukce je ověřen.

Předpokládejme nyní, že jsme již toto tvrzení dokázali pro $N = k - 1$. Dokážeme jej pro $N = k$. Klíčové je následující pozorování: v okamžiku, kdy je přесonová nejvíce kotouč mezi dvěma výčemi, musí být všech $k - 1$ ostatních kotoučů umístěno na třetí tyči. Z předpokladu víme, že jediný způsob, jak přесunout $k - 1$ kotoučů z výchozí tyče na pomocnou tyč na co nejmenší počet tařhů, je použít náš postup, který potřebuje $2^{k-1} - 1$ tařhů. V dalším kroku je přесun největší kotouč na cílovou tyč a zbývá na tuto tyč přесunout nět zbývajících $k - 1$ kotoučů. Opět díky předpokladu víme,

že na nejmenší počet tařhů toho docílíme jedine s použitím nashého postupu. Ukázali jsme tak, že v případě $N = k$ náš postup potřebuje $2(2^{k-1} - 1) + 1 = 2^k - 1$ tařhů a každý jiný postup by tařhů potřeboval víc.

Když už rozumíme tomu, jak kotouče přесonovat, můžeme se pustit do samotného řazení fotografií. Postupně budeme rekonstruovat pro zadané fotografie, v jaké řázi přесunmu se nacházíme. Na začátku si fotografie můžeme rozdělit do dvou přhřádek podle toho, kde se nachází největší kotouč. Pokud se nachází na výchozí tyči, znamaná to, že teprve přесonovává menší kotouče z výchozí na pomocnou tyč; tařkové fotografie určí předkládá všechny fotografie, u kterých se největší kotouč nachází na tyči cílové, tam už jsme ve řázi přесunu menších kotoučů z pomocné na cílovou tyč. Podobný postup můžeme zopakovat pro uspořádání každé z hřmádek, když se podíváme na polohu druhého největšího kotouče. Tento postup lze opakovat, dokud neseřadíme všechny fotografie. Celou proceduru shrnuje pseudokód.

Funkce *usporádej(k, fotografie, v, c, p)*:

1. Pokud $k = 0 \rightarrow$ skonči
2. Pro každou fotografii z pole *fotografie*:
3. Pokud se k -ty kotouč nachází na tyči v , umísti fotografii do pole h_1
4. Jinak umísti fotografii do pole h_2
5. *serazení* = *usporádej(k-1, h_1, v, p, c)*
6. *serazení* = *usporádej(k-1, h_2, p, c, v)*
7. *Výstup*: Zřetězení *serazení* a *serazení*?

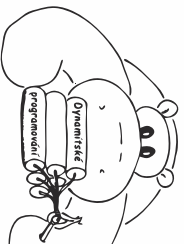
Uspořádání všech fotografií docílíme následujícím zavoláním: *usporádej(N, fotografie, A, B, C)*, kde *fotografie* označuje pole všech fotografií.

Celkový počet fotografií označme jako F . Každá fotografie je v průběhu řazení zpracována celkem N -krát. Pokud bychom proceduru implementovali doslova podle pseudokódu, dosáhli bychom složitosti $O(FN^2)$, protože bychom potřebovali čas $O(N)$ na zpracování jedné fotografie, například při umístování do přhřádek. Nám ovšem stačí pracovat pouze s odkazy na fotografie, což nám dá výslednou časovou složitost $O(FN)$. Paměťová složitost je přirozeně taktéž $O(FN)$.

Program (C):

<http://ksp.mff.cuni.cz/viz/27-4-6.c>

Lukáš Poluvarský



27-4-7 Nástroj pro zpracování textu

Seriál opět patří mezi vaše oblíbené úlohy, patří mezi tři úlohy s nejvíce odezvami. Proto doufáme, že se řešení naučíte třeba nové triky a že s námi zůstanete i v poslední seri.

se v něm zapisují odkazová tahnátorum pod názvem cíle. Takový jednoduchý Makefile tedy může vypadat takto:

```
datum.txt:
    echo Datum: > datum.txt
date >> datum.txt
```

Pokud tedy ve složce s tímto Makefile zavoláme příkaz `make datum.txt`, tak se provedou příkazy definované výše a vznikne nám tento soubor. Pokud ale zkusíme teď stejný příkaz zavolat znovu, tak se už nic neprovede – `make` už totiž ví, že je tento cíl splněn (soubor existuje) a že tedy není potřeba nic vytvářet.

Make a závislosti a virtuální cíle

Nejlepšími zbraní, kterou ale `make` disponuje, jsou závislosti. Dá se prohlásit, že cíl závisí na určitých souborech, například že vytvářený soubor se statistkou závisí na zdrojových datech měření. Když je pak `make` požádán o splnění nějakého cíle, tak se nejdříve pokusí splnit všechny závislosti. Pro každou závislost se podívá, jestli neexistuje stejně pojmenovaný cíl, a zkusí ho také splnit. Taktó může rekurzivně postupovat prakticky do neomezené hloubky.

Poré, co má nějaký cíl splněné všechny své závislosti, tak se `make` ještě rozhodne, jestli je nutné provést i tělo tohoto cíle. Pokud soubor stejného jména, jako je jméno cíle, ještě neexistuje, není potřeba váhat a tělo se provede. Pokud ale takový soubor už existuje, je to zajímavější. Pak se `make` podívá na čas modifikace všech souborů, na kterých altnahá cíl závisí, a porovná je s časem modifikace již existujícího souboru.

Když zjistí, že již existující soubor je novější než všechny jeho závislosti, tak není potřeba dělat nic. Pokud se ale alespoň jedna závislost s novějším časem modifikace, než má již existující soubor, tak se tělo cíle provede.

Abychom si to ukázali na příkladě, tak uvažujme následující Makefile používaný ke generování seznamu kapitol a nějakých statistik ze sepsanové knihy:

```
all: kapitoly.txt statistika.txt
kapitoly: knizka.txt
    grep "kapitola:" <knizka.txt >kapitoly
statistika.txt: knizka.txt kapitoly
    echo Počet řádků > statistika.txt
wc -l knizka.txt >> statistika.txt
echo Počet kapitol >> statistika.txt
wc -l kapitoly >> statistika.txt

clean:
    rm -f statistika.txt
    rm -f kapitoly
```

Jako první jste si asi všimni podivných cílů `all` a `clean`. Oba jsou to takzvané virtuální cíle, tedy jim neodpovídají žádné skutečné vytvářené soubory, ale slouží pro speciální účely. `all` `make` nebyl zamaten, kdyby se přelepe jen objevily soubory tohoto jména, tak mu to raději sdělíme mině magickou formou: `PHONY: na konci ukázký - ta zajistí to, že make bude stále se jmenující soubory ignorovat a vyjmenovaný cíle brát vždy jako virtuální a vždy se provedou její těla, jako by soubory neexistovaly.`

Nyní můžeme spustit příkazy `make all` pro výrobu všeho, na čem cíl `all` závisí, nebo `make clean` pro odstranění pra-

covních souborů. Cíl `all` je navíc uveden jako první proto, že když zavoláme jen `make bez cíle`, provede se první cíl.

Když jsme teď ukázali pár nových triků, pojďme se vrátit k závislostem. Rekněme, že jsme už provedli `make all` a máme tedy v adresáři všechny tři soubory. Když teď zavoláme `make kapitoly` nebo `make statistika.txt`, tak se nic nestane. Při volání `make statistika.txt` se `make` rekurzivně zavolá na splnění cíle `kapitoly`, ale protože ten nevygeneruje žádné nové závislosti, tak se neprovede ani tělo cíle `statistika.txt`.

Pokud ale nejdříve změníme obsah souboru `knizka.txt`, začne to být mnohem zajímavější. Po opětovném zavolání `make statistika.txt` se `make` znovu pokusí obstarat všechny jeho závislosti. Pro soubor `knizka.txt` žádný cíl nemá a tedy ho bere tak, jak je, ale pro cíl `kapitoly` cíl existuje a tak se ho pokusí rekurzivně splnit.

Při plnění cíle `kapitoly` zjistí, že jsou závislosti novější, než je vygenerovaný soubor, a tedy spustí tělo příkazů `kapitoly` a vygeneruje tak nový soubor. Cíl `statistika.txt` pak zjistí, že dokonce obě jeho závislosti (soubory `knizka.txt` i `kapitoly`) jsou novější, než vygenerovaný soubor, a proto taky spustí své tělo a vygeneruje nový obsah souboru `statistika.txt`.

Zde je vidět, že `make` vždy dělá jen tu nejmenší nutnou práci, spouští jen těla těch cílů, jež už ztratilo, na kterých závisí, se změnilo. U malého projektu je nám to asi jedno, ale když bychom třeba `make` používali k překladařům Unixového jádra a provedli jsme změnu jen v jednom zdrojovém souboru, budeme rádi, že `make` během několika málo sekund přegeneruje jen to, co potřebuje, a nemusíme tak čekat dlouhé minuty, než by se provedlo přegenerování úplně všeho, i když to nebylo potřeba.

Make a speciální proměnné

Mozná vám přijde, že třeba přejmenovat soubor se statistikami v příkladu výše by bylo docela pracné a mále pravdě. Znamenal by to na spoustě míst přepsat jeho název na něco jiného, hlavně v těle cíle, jež ho vytváří. Na to má ale `make` docela pěknou ležbu v podobě speciálních proměnných.

V tělech cílů se dají používat třeba tyto tři základní:

- `$(0)` zástupuje název cíle (jméno vytvářeného souboru)
- `$(*)` zástupuje název první závislosti
- `$(*)` zástupuje názvy všech závislosti

Klíčové cíle z dřívější ukázký by se tak daly přepsat třeba takto:

```
kapitoly: knizka.txt
    grep "kapitola:" < $(*) > $(0)
statistika.txt: knizka.txt kapitoly
    echo Počet řádků > $(0)
wc -l knizka.txt >> $(0)
echo Počet kapitol >> $(0)
wc -l kapitoly >> $(0)
```

Velmi mocným nástrojem je také konstruování obecných cílů. Pokud bychom třeba chtěli mít obecná pravidla, která nám pro každý textový soubor (třeba pro každou knihu, kterou jako úspěšný spisovatele sepsjeme) vygeneruje statistiku jako výše, můžeme k tomu právě tyto obecné cíle využít.

Obecný cíl vytvoříme tak, že v jeho názvu použijeme zástupnou část reprezentovanou znakem `%`. Za tu pak `make` při

⁶ Zadání úlohy z pera autora si můžete přečíst i nyní. Je sepsáno ve třetím svazku jeho díla *Recreations mathematiques* na straně 55, dostupno on-line: <https://archive.org/details/recreationsmathe005lucarc/h> <http://ksp.mff.cuni.cz/viz/kuchariky/zakladni-algoritmy>

hledání cili může doplnit cokoliv chce, ale zástupný znak můžeme použít jen na jednom místě v názvu cíle (nelze například vytvořit cíl `KSP-%-5-%.pdf`). Můžeme jej pak ale libovolně používat v závíslostech (lze tak třeba vyjádřit, že libovolný předložený program zavíší na svém zdrojovém jazyce C a podobně).

Když zástupný znak použijeme, tak se nám pak pro použití v těle cíle přidává ještě jedna speciální proměnná:

- `$(*)` obsahuje to, co se dosadilo za zástupný znak `%`

Hlavní část tohoto obecného Makefile by tedy mohla vypadat třeba takto:

```
%-kapitoly: %.txt
grep "%kapitola:" < $* > $@

%-stat.txt: %.txt %-kapitoly
echo Počet řádků > $@
wc -l $*.txt >> $@
echo Počet kapitol >> $@
wc -l $*-kapitoly >> $@
```

Ted můžeme volat třeba make `detektivka-stat.txt` nebo třeba make `roman-stat.txt` a stačí nám jen, aby `detektivka.txt` a `roman.txt` existovaly.

Někdy se nám dokonce může hodit mít nějaké obecné pravidlo a pak pro jednotlivé soubory přidávat zavislosti navíc. Pokud uvedeme jen hlavní cíle bez těla obsahujícího příkazy, tak se pro tento konkrétní cíl jen přidají zavislosti. Třeba příklad níže má pro soubory `prvni.txt` i `druby.txt` definovaný stejný příkaz, ale díky speciálním proměnným se pro každý zpracuje jiný zdrojový soubor:

```
%.txt:
grep "body" <$* >$@

prvni.txt: KSP_serial.txt
druby.txt: KSP_serial2.txt
```

Poslední věcí, kterou v spojení s příkazem make zmíníme, je definování a použití vlastních konstant. V podstatě to jsou proměnné, ale doporučujeme vám používat je skutečně jako konstanty (tedy do každé příhrádky pouze jednou), jinak se to celé zamotává. Typické použití je třeba v Makefile pro překlad zdrojů jazyka C, kde se jednou globálně nastaví všechny přepínače kompilátorem, a pak se používají. Hodnotu v proměnné použijete zapsáním `$(nazev-promenne)`.

Ukázkový Makefile využívající definované proměnné, speciální proměnné i obecné cíle vidíte níže. Používá se jak verze `$(promenna)`, tak i verze `$(promenna)`.

```
PROG=mujprogram
OBJ=$(OBJ)
CFLAGS=-Wall -c
LDFLAGS=-lm -lpthread
CC=gcc

%.o: %.c
$(CC) $(CFLAGS) -o $@ $*
$(PROG): $(OBJ)
$(CC) $(LDFLAGS) -o $@ $*
```

K tomu ještě dodáme to, že make zpřístupňuje i proměnné prostředí (tedy proměnné `ze shellu`). Tak si třeba v shellu

můžeme nastavit proměnnou, která nám ovlivní provádění příkazy. Když make žádnou definovanou proměnnou daného jména nemá (ani interní, ani v prostředí), dosadí místo ní prázdňý řetězec.

Úkol 6 [20]: Představte si, že máte tři zdrojové soubory `data: A.data, B.data` a `C.data`. Dále máte program `generuj <vstup>`, kterému jako parametry můžete předložit libovolný počet vstupních souborů a on na svůj standardní výstup vypíše nějaký vygenerovaný obsah. Ten by se měl ukládat do souborů, které budeme označovat jako `A.AB` a podobně (podle jejich vstupu).

Pak máte tři finální soubory, které se vytvářejí obdobným příkazem `finalizuj <vstup>` (opět vezme libovolně mnoho vstupních souborů a vygenerovaný obsah vypíše na svůj standardní výstup). Soubor `FINAL` se vytváří ze souborů `A.AB` a `B.BC` a `C.A`.

Septe Makefile, který bude odpovídat těmto zavislostem (zkuste co nejvíce pravidel néžak šikovně seskupit), a pak se zamyslete (a odpovězte), co vše se přegeneruje, když postupně provedeme následující příkazy:

```
make FINAL
touch A.data
touch B.data
touch C.data
make FINAL
```

Úkol 7 [30]: V KSPřku používáme make i na překlad zdrojů v TeXu. Ale když chceme někde použít automaticky generované obsahy, nastává problém – TeX totiž generuje soubor s obsahem během generování svého výstupu, a je tak potřeba často spustit první překlad TeXu „naprázdno“ a teprve při druhém překladu použít již vygenerovaný soubor s obsahem, který se vloží na správné místo zdrojů. Protože se soubor s obsahem (označme ho `toc`) vkládá do zdrojů (`tex`) a z něj se pak generuje pdf soubor, má by zavislost být `toc ← tex ← pdf`. Jenže čas změny `toc` souboru se znění vždy s překladem nového pdf a tedy vzniká vlastně cyklus.

Zkuste vymyslet řešení a sepsat základní Makefile, který bude toto řešení ilustrovat. Může se vám hodit třeba příkaz `diff`.

Závěr

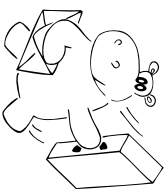
Dnes jste se dozvěděli další konkrétní velké skladačky o UNIXu a jeho příkazovém řádku. Pověděli jsme si o příkazech `find` a `xargs`, podívali se na zpracování signálů a obsah `/proc` a nakonec jsme si ukázali mocného pomocníka v podobě příkazu `make`.

Existuje ještě spousta konusů této skladačky, o kterých jsme neměli čas se zmínit, ale doufáme, že jsme vám poskytl aspoň to, co jsme považovali za základ, a pomohl vám třeba k tomu, abyste se o UNIX začali zabývat sami. Děkujeme, že jste s námi a našim seriálem přes celý rok vydrželi. :-)

*Poslední díl seriálu pro vás připravili
Jirka Scheiřka & Jenda Hadavna*

27-4-4 NP-úplný hlavolam

Jak dobře víte z kuchařky, důkaz NP-úplnosti obvykle sestává ze dvou kroků – důkazu toho, že problém leží ve třídě NP, a převodu některého problému, o kterém již víme, že je NP-úplný, na tento náš problém.



První krok je v našem případě velice snadný. Jako certifikát použijeme seznam sloupců, pod které jsou zasunuty barevné proužky. Ověření certifikátu učíte v polygomiálním čase zvládneme, stačí totiž pro každý řádek přimocovat spočítat, kolik barevných polí je vidět.

Druhý krok lze provést více způsoby. Pokud jste pečlivě prohlédli informatickou literaturu (či Wikipedii), mohli jste zjistit, že náš problém (přirozeně trochu jinak formulovaný) lze najít už ve slavné knize *Computers and Intractability: A Guide to the Theory of NP-Completeness* panů Gareyho a Johnsona z roku 1979 pod kódem LO4.

Zde si přečteme převod z problému Trojbarvnosti grahu. Jeho znění pro jistotu připomínáme.

Název problému: Trojbarvnost grahu
Vstup: Neorientovaný graf

Problém: Lze vrcholy tohoto grahu obarvit třemi barvami tak, že každá hrana sousedí s vrcholy dvou různých barev? (V obarvení musí mít každé dva sousední vrcholy různou barvu.)

Necht $G = (V, E)$ je neorientovaný graf, jehož trojbarvnost máme rozhodnout. Značíme $n = |V|$ a $m = |E|$. Naším cílem je konstruovat hlavolam, který má řešení právě tehdy, když je graf G trojbarvený. Tento hlavolam bude mít celkem $3n + 5m$ sloupců. Pro každý vrchol $v \in V$ mějme tři sloupce \tilde{c}_v, m_v, a_v . Pro každou hranu $e \in E$ a každou barvu $b \in \{\tilde{c}_v, m_v, z\}$ sloupec $s_{e,b}$.

Nejprve zajistíme, že každé řešení hlavolamu vůbec reprezentuje nějaké obarvení grahu. Pro každý řádek přídáme řádek, který vymocuje, že právě jeden ze sloupců \tilde{c}_v, m_v a a_v je barevný. Jak poznamenejme, jiskou barvu má vrchol v ? Podtrváme se, který ze sloupců \tilde{c}_v, m_v a a_v je barevný.

Zbývá zajistit, aby každé řešení hlavolamu reprezentovalo obarvení, ve kterém mají sousední vrcholy různé barvy. Pro každou barvu $b \in \{\tilde{c}_v, m_v, z\}$ a každou hranu $e = \{u, v\}$ přidáme řádek, který říká, že právě jeden ze sloupců b_u, b_v a $s_{e,b}$ je barevný. Podotkneme, že sloupec $s_{e,b}$ má úlohu ryze pomocnou – umožní nám říci „nejvyšší jeden ze sloupců b_u a b_v je barevný“.

Každé řešení hlavolamu tudíž odpovídá správnému trojbarvení. Naopak i pro každé správné trojbarvení, jak si snadno rozmyslíme, lze nalézt odpovídající řešení zadaného hlavolamu.

Náš převod je tedy korektní a zřejmě je možné jej realizovat v polygomiálním čase. Důkaz je hotov.

27-4-5 Večere pro opraváře

V úloze procházíme čtvercovou síť a cílem je najít nejkratší cestu, která vede přes všechny vyznačené hospody a restaurace, tedy je maximálně $n \leq 20$. Na zadání se

podíváme jako na úplný obdohocovaný graf s n vrcholy reprezentujícími hospody. Obdohocení získáme jako zvládnutost mezi hospodami ve čtvercové síti. Když můžeme získat spustěním přechodu do síťky z každé hospody zvláště. Určité se nám nevyplácí používat jiné než tyto nejkratší cesty, protože mezi dvěma hospodami se vždy vyplatí jít přímo bez jakéhokoli obdohocování. Toto převedení zvládneme v čase $O(n \cdot WH)$, kde W a H značí šířku a výšku mapy.

Ted stačí najít nejkratší cestu v úplném obdohocovaném grafu, která každý vrchol navštíví právě jednou. Tento problém je známý pod jménem „Problém obdohodního cestujícího“ (TSP) a je NP-úplný. To znamená, že není známý žádný polygomiální algoritmus, který by jej řešil. My si ukážeme algoritmus, který jej řeší v čase $O(n^2 2^n)$ a v prostoru $O(n 2^n)$, což je pro $n \leq 20$ dostatečný.

Úloha by se dala řešit obyčejným backtracem. Víme, v jakém vrcholu aktuálně stojíme (pro počáteční vrchol zkusíme všechny možnosti) a které vrcholy jsme již navštívili. Pro další v pořadí postupně zkusíme všechny možnosti navštívených vrcholů, přisumujeme se tam a pokračujeme v backtracem z nich. Takové řešení by ale mělo časovou složitost až $O(n!)$, protože vlastně postupně zkoušíme všechny permutace vrcholů.

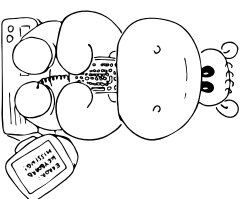
My ale toto řešení dokážeme vylepšit. Stačí nám do backtracem přidat ještě myšlenku dynamického programování. Při každém volání backtracové funkce se stejnými parametry (aktuální vrchol, množina navštívených vrcholů), musíme nutně dostat i stejný výsledek. Tedy jakmile jednou výsledek pro konkrétní parametry spočítáme, tak si jej uložíme do paměti a při dalším volání jej v konstantním čase vrátíme. Počet možností pro aktuální vrchol je 2^n (každý vrchol v ni buď je, anebo není). Celkově tedy backtracová funkce proběhne maximálně $O(n 2^n)$ krát.

Zbývá určit čas, který nám zabere jedno volání backtracové funkce. Ta jen pro každý z n vrcholů zkontroluje, zda je v množině již navštívených vrcholů a pokud ne, tak v něm rekurzivně zavolá backtracem. Pokud testování, zda vrchol je součástí množiny, zvládneme v konstantním čase, dostáváme se na časovou složitost $O(n^2 2^n)$ a paměťovou $O(n 2^n)$.

Jelikož vrcholů je maximálně 20, můžeme si jako konvolci pod-množin vrcholů reprezentovat jako n -bitové číslo. Pak, pokud i -ý bit má hodnotu 1, vrchol je v podmnožině a pokud 0, vrchol v ni není. Přidávání do množiny a zjišťování, jestli v ni vrchol je, tedy zvládneme v konstantním čase pomocí bitových operací, vizte vzorový kód.

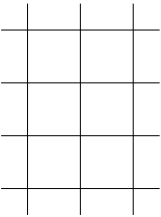
Program (C++):
<http://ksp.mff.cuni.cz/viz/27-4-5-cpp>

Karel Tesar



stále dokola. Takto provádíme všechny skupiny na sebe kolmých rovnooběžek, a to v lineárním čase. Celkové nám to zabere čas $O(N \log N)$ (kvůli řídění).

Pokud by vždy byly tři rovnooběžky, bylo by otázkou, jestli tvoří čtverec, triviální (jen bychom zkontrolovali jejich vzdálenosti, jestli jsou stejné). Za takové řešení jste mohli získat větší bodů, ale ne všechny. Pro plný počet bodů bylo potřeba zanalyzovat se i nad situací, kdy se vyskytne velké množství na sebe kolmých rovnooběžek – třeba v případě přímké uspořádaných v pravouhlé síti (jako na obrázku ze zadání níže).



Co dělat v takové situaci? Bylo by možné vyzkoušet každou dvojici rovnooběžek z první skupiny s každou dvojicí rovnooběžek z kolmé skupiny. Ale to by nám vlastně zdegenerovalo až na naše původní triviální řešení v čase $O(N^4)$. Nás však nezajímají jednotlivé příklady, ale jen vzdálenosti mezi nimi. V každé skupině si tedy vezmeme všechny možné vzdálenosti mezi přímkami (těch je K^2 pro K přímek, tedy $O(N^2)$ pro nejhorší případ), ty si v čase $O(N^2 \log N^2) = O(N^2 \log N)$ utřídíme a pak tyto dvě seřazené posloupnosti porovnáme.

Stačí nám pro každou vzdálenost, která se vyskytne v jedné z posloupností, spočítat součin počtu výskytů této vzdálenosti v první posloupnosti a počtu výskytů této vzdálenosti ve druhé posloupnosti (součin proto, že čtverec tvoří každé dvě dvojice). Všechny tyto součiny sečteme a dostaneme tak počet vytvořených čtverců.

Pro příklady, kdy se na vstupu vyskytne mnoho rovnooběžných přímk, umíme dosáhnout času $O(N^2 \log N)$ s paměťovou složitostí $O(N)$. Pro příklady, kdy bude rovnooběžných přímk málo, se bude čas běhu našeho postupu blížit spíše $O(N \log N)$.

Jirka Šedivka

27-4-3 Vysoké napětí

Letci varianta

Věšina z vás si všimla, že máme-li nějaké korektní řešení, můžeme vytvořit další správné řešení tím, že prohmotíme všechna 100kV napětí v uzlech za 0kV a 100kV za 0kV. Z toho speciálně dostáváme, že si můžeme vybrat libovolný uzel a přičíst mu hodnotu 0kV. Existuje-li totiž nějaká korektní řešení, kde tento uzel má napětíovou hladinu 100 kV, existuje korektní řešení, kde má hladinu 0kV.

Pak si stačí všimnout, že napětíová hladina v prvním uzlu jednoznačně určuje napětíové hladiny v sousedních uzlech, ty zase ve svých sousedech a tak dále, až se jednoznačně určí celý graf. Toto řešení můžeme tedy nazvat jednoduchým příkladem například do hloubky.

Vždy když zkusíme nějaký uzel, který ještě nemá určenou napětíovou hladinu, určíme ji podle napětíové hladiny předchozího uzlu a rozdíl napětí na vodiči, jež tyto uzly spojuje – v případě, že rozdíl napětí byl 0kV, bundou napětí v uzlech stejna, pokud byl rozdíl 100 kV, bundou napětí opakata. Pak začneme prohledávat všechny jeho sousedy.

Pokud zkusíme uzel již má určenou napětíovou hladinu, jen zkontrolováme, jestli tato hladina souhlasí s hladinou, kterou bychom ji jinak přiřkli. Jestliže nesouhlasí, dostáváme spor a graf nelze obhospodit. Pokud souhlasí, tak je vše v pořádku (jeho sousedy již prohledávat nemusíme, neboť jsme je prohledali při první návštěvě tohoto uzlu).

Musíme si pamatovat celý graf a u každého uzlu a hrany konstantní množství informací, paměťová složitost bude tedy $O(n + m)$, kde n je počet uzlů a m počet hran. Celý graf musíme načíst a pak pro každou hranu provedeme konstantní množství kroků (zpracování jednoho uzlu), dostáváme tedy opět $O(n + m)$.

```
Program (C++):
http://ksp.mff.cuni.cz/viz/27-4-3a.cpp
```

Tři napětíové hladiny

Dobry nápad je využít řešení jednoduchší varianty. Vezmeme náš graf a odstraníme z něj všechny vodiče s rozdílem napětí 100kV. Tímto se nám graf rozpadne na několik komponent. Smažeme ty, které neobsahují žádný vodič s rozdílem 200 kV. Zbývající komponenty vyřešíme podobně jako jednodušší variantu. Jen musíme místo 100kV pracovat s 200 kV a také nám tu nastává ten problém, že máme více komponent. Musíme tedy prohledávat z jednodušší varianty postupně spouštět ze všech uzlů. Rozmyslete si, že to nám nijak nezhorší asymptotickou časovou složitost (většinou se spustíme na uzel s již určenou hladinou, a tedy hned skončíme).

Máme tedy určené napětíové hladiny všech uzlů sousedících s vodičem s rozdílem 200kV a uzlů, které jsou z těchto uzlů jednoznačně určité. Zbývá tedy určit napětíové hladiny uzlů, které s žádným 200kV vodičem nenesou, a zkontrolovat, jestli hladiny, které jsme určili, souhlasí s uzly k nim připojenými 100kV vodičem (0kV vodiče řeší neustále, neboť ty jsme řešili již v prvním kroku).

Rozmysleme si, že obhodnocené komponenty jsou navzájem spojeny cestičkami z vodičů s rozdílem 100 nebo 0kV, kde navíc první a poslední vodič v každé cestičce je 100kV. Uzel na druhém konci tohoto vodiče musí mít hladinu 100kV, protože vyčtoží uzel má hladinu při 200kV, nebo 0kV. Kromě toho musíme ještě zkontrolovat, že toto přiřazení nenarušilo dosavadní přiřazení hladin (to by nastalo v případě, že by byly dvě komponenty spojeny právě jedním vodičem), pak by řešení neexistovalo.

Nyní si stačí doceně odmyslet již obhodnocené komponenty (až na ty poslední uzly s hladinou 100 kV). Zbudou nám tedy pouze uzly s hladinou 100 kV a vodiče s rozdíly napětí 0kV a 100kV. Nabízí se tedy opět řešení jednoduchší varianty. Spor s dosud obhodnocenými komponentami nám určitě nenarouká, neboť s nimi jsme spojeni pouze prostřednictvím 100kV uzlů. Takto tedy dojdeme ke sporu a zjistíme, že řešení neexistuje, nebo nalezneme korektní řešení.

Opět si stačí pamatovat graf, takže paměťová složitost bude $O(n + m)$. Co se týká časové, tak nejprve provedeme jednou jednodušší variantu na obhodnocené části grafu, poté v lineárním čase označíme nějaké 100kV uzly a poté znovu provedeme variantu jednoduššího algoritmu. Opět tedy dostaneme časovou složitost $O(n + m)$.

```
Program (C++):
http://ksp.mff.cuni.cz/viz/27-4-3b.cpp
```

Dominik Smrz

Recepty z programátorské kuchyně: Hledání v textu

Řetězec je v podstatě jakákoliv posloupnost symbolů zapisaná za sebou a s nimi budeme v této kapitole pracovat. Každěmu napadne „vyhledávání v textu“ nebo „hledání jmen v telefonním seznamu“ ale řetězec najdeme i na míšičkách tvůrčích informatických. Například celé číslo zakódované v binární soustavě, které dostaneme na vstupu programu, je také jen řetězec nult a jedniček.

Jiný příklad použijí řetězci (a jejich algoritmy) najdeme v biologii. DNA není o mnoho více, než čtyřtří uložení posloupnosti čtyř znaků/nukleových bází – a chceme-li hledat vzory anebo konkrétní podpodobnosti, bude se nám hodit znalost základních algoritmů pro práci s řetězci.

Nemáme hohlzákl 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 kuchyňce, v této se budeme soustředit na algoritmy, které se objevují spíše v práci s textem.

Kromě úvodu popíšeme dva *stavební kameny* textových algoritmů, což bude jedna datová struktura pro slovníky (tree) a jedno vyhledání v textu s předzpracováním hledaného slova (a jeho rozšíření pro více slov). S jejich znalostí se pak mnohem snáze vyřídí řešení složitějších, realističtějších problémů.

Jak řetězec chápat

Když programátor dělá první křížky, často moc netuší, co s těmi řetězci vlastně má dělat a nemít dělat. V programování jazyce to je jasné – něco mu jazyk doví a na něco nejsem prostředky. Ale jak to je na úrovni více teoretické?

Jak jsme si řekli na začátku, řetězec bude posloupnost nějakých symbolů, kterým říkáme *znaky*. Tyto znaky jsou z nějaké množiny, která říkáme *abeceda*. Abeceda může být jen {0, 1} pro čísla v binárním zápisu, klasické {A-Z, a-z} pro anglickou abecedu anebo plný rozsah univerzální znakové sady Unicode, která má až 2^{21} znaků. Nezapomínáme, že nejmenším 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 $|Z|$. Abeceda sama se v textech o řetězích často značí řeckým Σ .

O znacích samotných 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 chápat řetězec jako pole znaků, nebo jako spojový seznam? Šalounůvská odpověď: můžeme s ním pracovat tak i tak. Když budeme potřebovat přehled řetězce na spojový seznam (protože se nám hodí rychlé přepojování řetězci), tak si jej převedeme. 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 uložen v poli někde v paměti (jíz od začátku algoritmu), takže ke každému znaku můžeme přistupovat v konstantním čase.

Jelikož jsme řetězec dehmovali jako posloupnost, nesmíme zapomenout ani na *prázdný řetězec* ϵ . A když už máme řetězec, určitě máme i *podřetězec* – souvislou podposloupnost

znaků jiného řetězce. Například BAR, BET, 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*, kterému říkáme *prefix* (česky předpona), a pokud odstraníme nějaký souvislý úsek ze začátku, dostaneme *suffix* neboli příponu. BET je suffix slova KABARET, KABA je zase jeho prefixem.

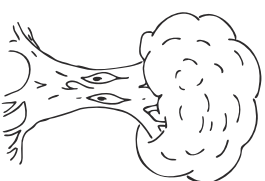
Terminologie dovoluje zepřehdí i zezadu odstranit prázdný řetězec – to znamená, že slovo je samo sobe prefixem i suffixem. Pokud chceme mluvit o prefixech, suffixech nebo oboecé podřetězci, kde jsme mysleli alespoň jeden znak odtrhnout, označíme takové podřetězce jako *vládní*.

Pro některé použijí řetězci je důležitě, abychom je mohli porovnávat – když máme řetězec R a S , chceme rozhodnout, který je menší, a který je větší. Jaké přesné toto uspořádání bude, závisí na naší aplikaci, ale mnohdy se používá tzv. *lexikografické uspořádání*.

Pro lexikografické uspořádání potřebujeme nejprve zadane lineární uspořádání na znacích. Tím se myslí takové, kde jsou všechny prvky navzájem porovnatelné a v podstatě to znamená, že jsou uspořádány 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 nějaký uzel jednoho z řetězci dojdou dříve, prohlásíme tento řetězec za menší.

Přati tedy třeba $\epsilon < A < AUTO$ < AUTOBIS < AUTOGRAF < AUTOR < BAWELTKA < BARNABAS < Z.



Adresář pomocí tree

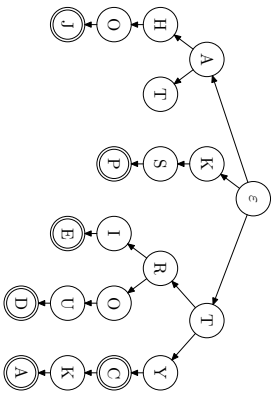
Typický „textový“ problém je uložování množiny řetězci – můžeme si představit třeba slovník. Slova ve slovníku si chceme šikovně předzpracovat, abychom pak mohli efektivně odpovídat na otázky typu: „Je slovo S obsaženo ve slovníku?“. Můžeme také po předzpracování chtít přidávat nové položky, nebo i odebrat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

Pokud bychom nemohli odebrat slova, můžeme použít hesování, které je rychle a účinné. Více o něm najdete v hesovací kuchyňce.³ Má však tu nevýhodu, že při velkém zaplnění se začne chovat pomaleji a mírně nepřehledně.

Ukážeme si jiné řešení, které je také asymptoticky rychle a není ani příliš náročné na paměť. Využívá stromové struktury a říká se mu *trie* (vyslovujeme česky „trýý“ a anglicky jako část slova „retrieve“, z něhož slovo trie vzniklo). V české tině se občas používá také označení „písmennkový strom“.

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

Obrázek vydá za tisíc definic, pojďme se podívat, jak vypadá trie pro slova AHQJ, AT, KSP, TRIE, TROUD, TYG, TYCKA. 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 zadanych slov. Například prefixy délky 2 jsou AH, AT, KS, TR a TY. Hrana mezi prefixy vede právě tehdy, lze-li jeden z druhého získat připsáním písmene na konec.

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

Z takto popsané trie bychom nepoznali, 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 ne (jak je to naznačeno dvojitými kruhy v obrázku), anebo si rozšíříme abecedu o speciální znak, který se v ní předtím nevykřetoval – třeba \$ – a pak všem slovům přilepíme tento \$ na konec.

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

Jestli jsme si nerozumysleli, 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 vyšší 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 oždet konstantní rychlost dotazu

⁴ <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}. Tehdy nahradíme každé slovo dvojkové soustavě. Tím se časová složitost konstrukce zlepši na $O(|D| \cdot \log |Σ|)$ a časová složitost dotazu na slovo délky L zhoší na $O(L \cdot \log |Σ|)$.

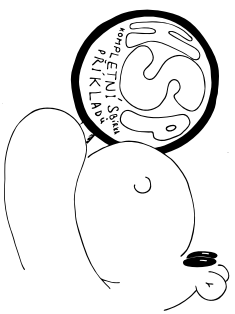
A jsme hotoví! S trii můžeme v lineárním čase odpovídat na dotazy „Vyskytuje se dané slovo ve slovníku?“, přidávat a odebrat další položky za běhu a nejen to – víc o tom ve cvičeních.

Poznámky

- Chcete-li algoritmus konstrukce trie vidět napsaný v Pascalu, podívejte se do knihy *Algorithm a programovací techniky*.
- 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 před postavit slovník ze všech slov v daném textu, a pak procházet příslušnou trii. Má to ale pár háček: jednak je často hledány řetězce krátké, ale text se nevejde do paměti. Druhák, pokud bychom použili jako oddělovací mezery, mohli bychom hledat jen jednotlivá slova, a nikoli jejich konce nebo delší kusy věty.
- Asi se po posledním poznámce ptáte – existuje nějaká modifikace trie, která umí hledat libovolnou část textu? Ano, jmenuje se *suffixový strom* a jsou s ní další spousty krásných kousků. Říká se, že každou řetězovou úlohu lze řešit v lineárním čase pomocí suffixových stromů. Víc se o nich dočítete třeba v knize *Krajinnou grafových algoritmů*.⁴

Cvičení

- Řekneme, že chceme slovník na vstupu setřídit v lexikografickém pořadí (definovaném v sekci „Jak řetězce dle pat“). Problémem pro klasické třídící algoritmy je to, že porovnání dvou řetězců není bohužel konstantně rychlé. Vymyslete způsob, jak setřídit takový slovník rychle pomocí trie.
- *Kompresce trie*. Co kdybychom chtěli odstranit přebytečné nevrcholy trie, tedy ty, v nichž se slova nevětví? Rozmyslete si, jestli by něčím vadilo místo takového cest mít jen jednotlivé hrany. Zesložte se konstrukce nebo vyhledejte řešení? Mimoходом, je celkem jasné, že takováto *kompresovaná trie* přinese jen konstantní zrychlení dotazů i prostoru, a tak na soutěžích apod. stačí použít základní variantu.



Vzorová řešení čtvrté série dvacátého sedmého ročníku KSP

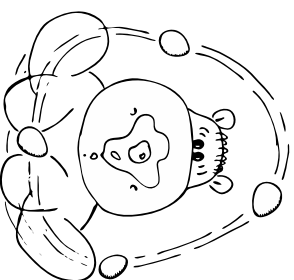
27-4-1 Zadávání úkolů

Například uděláme jedno jednoduché pozorování. Kdybychom věděli, kolik úkolů přes konkrétní zaměstanec projde, je snadné určit, komu bude posílat další. Pokud prošel sudý počet úkolů, bude to stejný podřízený jako na začátku, v jiném případě to bude ten druhý.

Stejně tak lze snadno určit, kolik úkolů takový vedoucí předá kterému podřízenému. Pokud je počet sudý, oba podřízené dostanou polovinu. V jiném případě dostane ten zatímající o jeden úkol více.

Jak tedy určit, kolik úkolů skrz kterého vedoucího projde? Zapišme si řetěze. Ten svůj počet úkolů v_i máme ho tedy vyřešeno. Jeho příjím podřízeným přidáme úkolý. Poté si vybereme některého zaměstance, který už všechny úkolý od svých nadřízených dostal: už tedy také v_i , co ho čeká a nemáme. Opět rozdělíme jeho úkolý a opět si vybereme někoho, kdo má všechny své nadřízené vyřešené. To děláme tak dlouho, dokud ještě někdo zbyvá. Zjevněšně, vedoucí začne rozdělávat, až když dostane všechny své úkolý.

A proč se nám nemůže stát, že siče máme ještě nějaké nevyřešené zaměstance, ale všichni mají nějakého svého nadřízeného ještě nevyřešeného? Pro spor si představme, že se nám přesně taková nepřijemná věc stala. Vezměme tedy libovolného nevyřešeného zaměstanca. A z něj postupně do některého jeho nevyřešeného nadřízeného – takový musí z definice této nepřijemné situace existovat. A u toho nadřízeného si zase vybereme některého jeho nevyřešeného nadřízeného. Takto budeme pokračovat „nahoru“ v hierarchii, ale nikdy neskončíme. Zaměstanec „nahoru“ v hierarchii počet (zřejmě v dané společnosti mohou pracovat maximálně všichni lidé na Zemi), musíme tedy jednou navštívit některého vícekrát. To ale znamená, že je v grafu cyklus, a ten máme zadáním zakázáný.



Nyní nám tedy zbyvá rozmyslet, jak to napsat s co nejlepší složitostí. Napřed si spočítáme, kolik má který zaměstanec nevyřešených nadřízených. Všichni začnou na nule, projdeme všechny vedoucí a každému podřízenému vždy přičteme jednotku. To zvládneme v konstantním čase na každého vedoucího. Založíme si skladiště na zaměstancace bez nevyřešeného nadřízeného (třeba zásobník, ten je příjemně jednoduchý) a při dalším příchodu přes zaměstancace do něj nasbíráme všechny, kteří mají nulu (v době

fungující společnosti by to měl být jen řetězec). To opět zvládneme v celkové lineárním čase. Nyní opakovaně vyndáme zaměstancace ze skladiště, vyřešíme ho a oběma jeho podřízeným odečteme jedničku. Pokud číslo u některého z nich (nebo obou) klesne na nulu, přidáme ho do skladiště také. Vyřešení jednoho nám opět bude trvat konstantní čas.

Celkové řešení zvládneme celý výpočet v lineárním čase. Lěpe to neuplyne, protože jen nastavení výsledků každého zaměstancace bude trvat takovou dobu.

Co se týče paměťové složitosti, potřebujeme si ke každému zaměstancaci zapamatovat konstantně mnoho informací a potřebujeme skladiště, do kterého uložíme každého zaměstancace maximálně jednou. Tedy si vystačíme s lineární paměťovou složitostí.

A pro znalce: ano, je to topologické třídění.

Program (C):

<http://ksp.mff.cuni.cz/viz/27-4-1.c>

Michal „norma“ Vancr

27-4-2 Čtverce v síti

Nějaké se zamysleme nad triviálním řešením. Můžeme zkoušet vřít každou čtvercív přímk a podívat se, jestli náhodou netvoří čtverec. To pro každou čtvercív zkontrolujeme snadno – ověříme, že dvě a dvě z nich jsou rovnoběžné, tyto dvojice rovnoběžek jsou na sobe kolmé, a navíc jsou od sebe stejně daleko. Pokud si potřebujete trochu připomenout základy analytické geometrie, nahleďte do naší kuchařky o geometrii.⁵

Takový přístup nám zabere celkové čas $O(N^4)$. Neumíme to ale lépe. Již jsme si všimni toho, že čtverce tvoří vždy dvě dvojice rovnoběžek, toho jistě můžeme nějak využít. Tím nám ale vzniká nová otázka, a to jak rychle najít rovnoběžky mezi N přímkami v rovině.

Pokud bychom místo přímk měli třeba reálná čísla, stačilo by nám je v čase $O(N \log N)$ setřídit, čímž by se stájně hodnoty dostaly k sobě, a pak bychom jedním lineárním příchodem našli všechny duplicity. Přímký siče nejsou reálná čísla, ale můžeme si je jimi popsat.

V této fázi nás zajímá jen směr přímk, ne jejich poloha, takže nám stačí namísto popisu celé přímký vřít její *směrový vektor* (tedy dvojici čísel (x, y) vyjadřujících směr přímký vůči souřadným osám). Směrové vektory bychom museli ještě „znoromalizovat“, tedy upravit je všechny tak, aby třeba $x = 1$, čímž z nich vlastně získáme jedno reálné číslo, a to jejich *směrnici*.

Nyní si tedy můžeme všechny přímký v čase $O(N \log N)$ setřídit, dostat tak rovnoběžky k sobě a pak takto setřídné přímký lineárně projít. Pro každou nalezenou skupinu rovnoběžek budeme chtít nalézt rovnoběžky na ně kolmé. To můžeme pokaždé dělat binárním vyhledáváním v čase $O(\log N)$ na dotaz, nebo na to můžeme jít čtyřtříjí.

Stačí nám držet si v setříděném seznamu přímk dva ukazatele posunuté od sebe o 90° a posouvat je oba najednou. Pokud nám druhý ukazatel skočí až na skupinu rovnoběžek, které ve směru otáčení svraží s rovnoběžkami na první pozici úhel větší než 90°, posuneme zase první ukazatel a tak

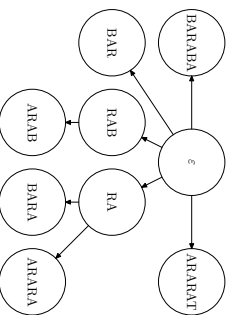
⁵ <http://ksp.mff.cuni.cz/viz/kuchařky/geometrie>

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

V našem příkladu se seamen **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 hlásky žádné vyskytly.

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

Tedy po přepočítání bude mít **RA** tři vyskytly a **RAB** jeden vyskyt, celkový počet vyskytů pak bude 12.



Poznámky

- Dalšími kroky po KMP a Aho-Corasickové jsou konečné automaty a regulární výrazy, o kterých jsme měli seriál ve 23. ročníku.

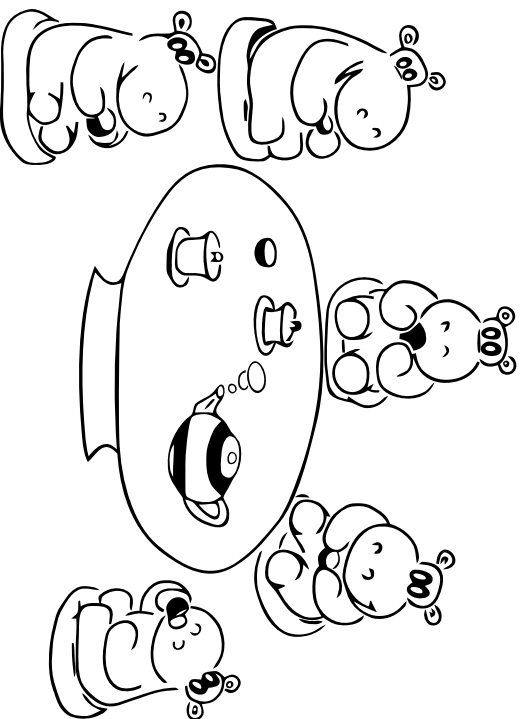
- Není moc rozumné době například při součtu, pokud tento algoritmus nemáte opravdu pod kůží. Radši zkuste použít řešení, pokud budete něco takového potřebovat.

Výčtení

- Redukci o velikost výstupu můžeme provést i pro případ, kdy výstup nebudeme vypisovat, ale stačí nám mít jej uloženy v paměti. Vymyslete vhodnou úpravu tržku s čítkem.

- Zkuste si implementovat Aho-Corasickovou vlastnoručně ve svém oblíbeném jazyce, abyste si byli jisti, že doopravdy chápete všechny záležitosti tohoto algoritmu.

Martin Bohm, Jan Matějka, Martin Mareš a Petr Škoda



Vyhledávání v textu

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 nahlásieme jeden nebo všechny výskytly slova. Zajímají nás při tom i výskytly, které se nazývají překryvaty: v textu **MAAANA** se slovo **MAA** vyskytuje dvakrát. Často se hovoří o „hledání jehly v kusep sena“, pročez se textu přezdvává *sema* a hledanému slovu *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 kdybydrom v textu narazili na slovo **INSTINKTIT**?

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

Sice je předlehozí popis skutečně v nejlhorším případě složitý $O(S \cdot J)$, avšak stačí malá úprava a složitost přijde na lineární $O(S + J)$. Ve skutečnosti algoritmus nepomalovalo vracení se – za špatnou složitost mohl fakt, že jsme se vraceli *příliš zpátky*.

Třeba v našem příkladu s textem **INSTINKT** se nemusíme vracet ve spojovém seznamu na začátek, jakmile nachíme **INSTINS**. Mohli jsme se vrátit jen na druhý znak, tedy do prvního **I**, a pak kontrolovat, jaký znak pokrácuje 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 skočíme, 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 pole, 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 máte rádi grafovou terminologii, můžete se na náš spojový seznam dívat jako na graf a hovořit o *zpětných hranách*.

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

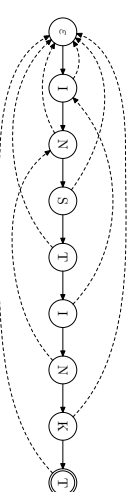
Abychom náš požadavek upřesnili, zamysleme se nad zpětným políčkem pro jiné slovo. Co kdyby jehlou bylo slovo **ABABABC** a my určovali zpětné políčko pro **ABABAB**? Kdybychom ukázali na první písmenko **B**, nebylo by to správné,

protože pak bychom pro text **ABABABABC** nezahášíli výskyt jehly, což je jasná chyba. Musíme se vrátit už na **ABAB**!

Zajímá nás tedy ne lhovolyý sufix, který je stejný jako začátek, ale nejdlejší takový konec/suffix. A ještě navíc ne jen ten nejdleší, ale nejdleší „ netriviální “ – slovo **INSTIN** je samo sobě prefixem a suffixem, ale zpětná funkce pro **N** by se neměla cyklit, měla by vesti zpátky.

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

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

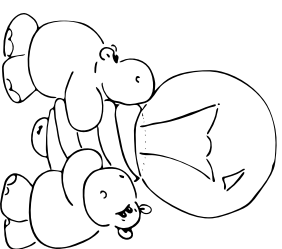


Nyní vystavíme dvě otázky: Jakou to má celé časovou složitost? A jak spočítat zpětnou funkci?

Poporne se nejdříve s tou první. Pro každý znak vstupního textu můžeme nastat dva případy: Bud znak rozšiřuje aktuální prefix, nebo umisíme 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 řešení 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 přečetli znaků textu. Celkem je tedy počet kroků automatu lineární v délce textu, $O(S)$.

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



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

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

máme $F[i]$, pak výpočet $F[i+1]$ odpovídá spuštění automatu na slovo délky i a při tom budeme zpětnou funkci porovnávat jen pro stavy délky i nebo menší, pro které již máme hodnotu.

Navíc nemůžeme pro jednotlivé prefixy spouštět výpočet vždy znovu od začátku $-(i+1)$ -ní prefix je přesí předložkami i -tého prefixu o jeden znak. Stačí tedy spustit algoritmus na jehlu bez prvního znaku a sledovat, jakými stavy bude procházet – to budou přesné hodnoty zpětné funkce.

Vyrovnání zpětné funkce se nám tak nakonec znehodnotilo na jedné vyhledávání v textu o délce $J-1$, a proto poběží v čase $O(J)$. Časová složitost celého algoritmu tedy bude $O(S+J)$. Dodáme už jen, že tento algoritmus poprvé popsal panové Knuth, Morris a Pratt a na jejich počest se mu říká KMP. Naprogramovaný bude vypadat následovně (tém vstupu jsme si odpustili):

```

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

```

```

# Poznámky
• Pro anglický nebo český text je použit takto softlkovaneho algoritmu skoro škoda, protože v obou jazycích se stává jen málokdy, že bychom měli několik slov spojených dohromady. Prakticky bude stačit i na začátku zmíněný nativní algoritmus. Na soutěžích a olympiádách ale přiče raději algoritmus KMP.
• Hesování lze použít i na vyhledávání řetězce v textu. Ohzvláště vhodné jsou na to rolling hash functions (neboli „oběhové hesovací funkce“), které umí v konstantním čase připočítat hes, ubereme-li nějaký znak na začátku a přidáme-li jiný na konci – jako kdybychom se držali na text šláz posouvající se okénko.

```

- Rozmyslete si, že když vyhledáváme více slov, ne jen jedno, algoritmus musí vypsat všechny výskytů na výstup, mžžeme se dobřat vyšší než lineární složitosti v závislosti na vstupu. Na čten potom taková časová složitost také záleží?
- Vynyslete nějakou vhodnou okénkovou hesovací funkci pro vyhledávání jedné jehly.

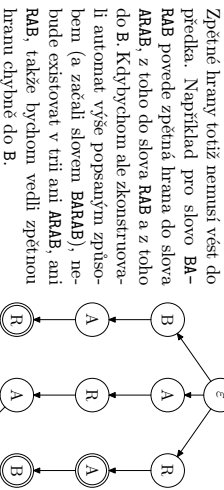
Vyhledávání jehleňku

Co kdybychom neměli jen jednu jehlu/hledané slovo, ale celý jehleňček, čili seznam hledaných slov? I to lze řešit podobnou metodou, jakou jsme hledali jedno slovo. Tento algoritmus se nazývá po tvůrcích *algoritmus Aho-Corasickovi* a spočívá v tom, že jednoduše spojový seznam nahradíme tři a do tří opět přidáme zpětné hrany.

Budeme postupovat podobně jako u KMP. Nejprve naskládáme jehleňček do trie. Pro příklady v této kuchařce použijeme jehleňček **ARAB, ARABA, ARARAT, BAR, BARA, BARABA, RA a BAB**.

Dalším krokem v KMP bylo sestrojení zpětných hran. Nejprve jsme sestrojili zpětnou hranu pro první znak slova, pak pro druhý atd. Ve trii to bude o něco složitější.

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



Zpětné hrany totiž nemusí vést do předka. Například pro slovo **BARAB** povede zpětná hrana do slova **ARAB**, z toho do slova **BAB** a z toho do **B**. Kdybychom ale zkonstruovali automat výše popsaným způsobem (a začali slovem **BARAB**), nebudou existovat v trii ani **ARAB**, ani **BAB**, takže bychom vedli zpětnou hranu chybě do **B**.

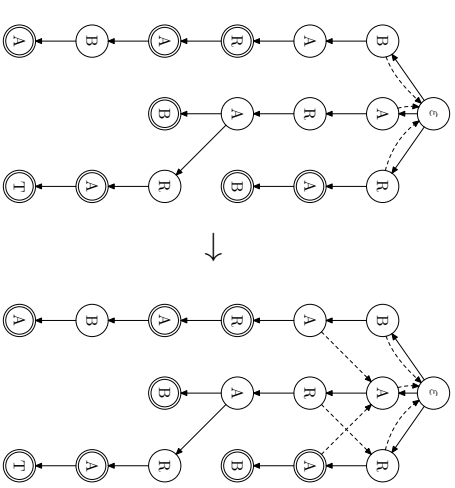
Mžžeme se ale opřít o stejný trik, jako při konstrukci KMP. Budeme opět vyhledávat nejdelší vlastní sufix. Kam dojde výpočet po jeho vyhledání, tam povede zpětná hrana chybě do **B**. Zkusíme tedy nejprve sestrojiti celou trii a pak postupně vyhledat nejdelší vlastní sufix pro každé ze slov. Ouhá, to ale také nefunguje. Když zkontrolujeme slovem **BARABA** a budeme tedy vyhledávat **ARABA**, nalezneme v trii úspěšně prefix **ARAB**, ale **ARABA** již v trii není. Měli bychom přejít ze slova **ARAB** po zpětné hraně, ale tu ještě nemáme zkonstruovanou.

Rozděláme si trii na *vrstvy* – první znaky slov budou první vrstva, druhé znaky budou tvořit druhou vrstvu atd., až *i-té* znaky slov budou tvořit *i-tou* vrstvu.

Zpětná hrana jisté povede do kratšího slova. *Z i-té* vrstvy tedy povede do vrstvy s nižším pořadovým číslem. Pokud tedy budeme zpětné hrany konstruovat po vrstvách, dojde-me kýžnému výsledku.

Jestě zbývá otázka, jak konstruovat zpětné hrany efektivně, když je musíme vyrábět po vrstvách. Měli bychom prostě vzít slovo, pro které hledáme zpětnou hranu, utrhnout mu první znak a vyhledat. Jenže to budeme dělat spoustu práce zbytečně.

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



Při konstrukci další zpětné hrany tedy najdeme akorát, kde jsme minule skončili, a odkamtud pokračujeme dál. Jak to najdeme? *Z otce* našeho vrcholu tam přece vede zpětná hrana. Takže můžeme postup shrnout do bodů:

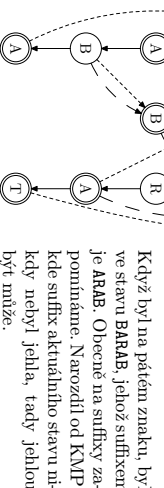
1. c = poslední znak slova (znak stavu P , pro který hledáme zpětnou hranu);
2. přisuneme se do otce;
3. přisuneme se po zpětné hraně;
4. dokud neexistuje syn se znakem c nebo nejíme v kořeni, přisunováme se po zpětných hranách;
5. pokud existuje syn se znakem c , natáhneme do něj zpětnou hranu z P , jinak ji natáhneme do kořene.

Automat je zkonstruován. Časová složitost konstrukce sestává z konstrukce trie v $O(J \cdot |Σ|)$, resp. $O(J \cdot \log |Σ|)$ (pokud použijeme binární strom ve vrcholech) a z předpočítání zpětných hran. Při předpočítání nděláme nějaký konstantní počet operací pro každý vrchol (celkem tedy $O(J)$) a také paralelně vyhledáváme všechny jehly z jehleňku, jejíž délka vyhledání nás stojí $O(J)$, resp. $O(J \cdot \log |Σ|)$.

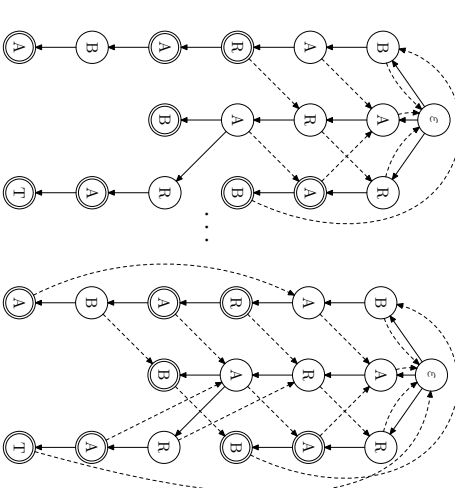
Tedy konstrukce trvá celkem $O(J \cdot |Σ|)$, resp. $O(J \cdot \log |Σ|)$, paměťová náročnost je stejná jako u trie – $O(J \cdot |Σ|)$, resp. $O(J)$, přidali jsme jen $O(J)$ zpětných hran.

Zkusíme tedy automatem projit text **BARBARARAT**. Ohláší postupně názvy slov **BAR, BARA, BARABA, BAR, BARA, ARARA a ARARAT**.

Nenalezl však všechno. Chty-bí mu např. **ARAB**, který začíná druhým znakem a končí pátým. Dale clybhi několik výskytů **RA** a jeden **BAB**.



Když byl na pátém znaku, byl je stavu **BARAB**, jehož suffixem je **ARAB**. Obecně na suffixy zapomínáme. Narozdil od KMP, kde suffix 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 slovník obsahující **A a AAAA...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. Předpočítáme si tedy *zkratky* – z vrcholu vede zkratka do nejdelšího jeho suffixu, který je jehlou. Na obrázku jsou vyznačeny dlouze čárkovanými šipkami.

Předpočítání zpětných hran časovou složitost konstrukce automatu jisté nezhorsí, neboť vyžaduje v nejlouším případě projít všechny zpětné hrany ještě jednou.

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

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

Jak velký může být výstup? Obecně až S^2 . Extrémně velký výstup je možné vygenerovat třeba slovníkem obsahujícím všechny prefixy slova **AAAA...A** délky S a se-nem takéž **AAAA...A** délky S . An-tomat pak hlásí výskyt pro každé podслово, kterých je řádově S^2 .

Pokud nám stačí v každé slova jen počet výskytů, nemůžeme zohat – závislost na počtu výskytů umíme odstranit.

Použijeme trik – na každé pozici započítáme pouze nejdelší jehlu, která tam končí (i u každé jehly si budeme udržovat čítač). Nebudeme tedy v každém kroku poskakovat po zkratkách až do alehňja, ale

