

na jedno další místo C – všude je však museli skládat na sebe za dodržení uvedených pravidel.

Vysvětlivatel získal fotografie vzniklé při tomto stěhování a chce se ujistit, že během něj nedošlo k manipulaci s daty. Na vstupu dostane každou fotografii popsanou jako řetězec znaků A, B a C, kde A-ty znak určuje polohu A-tého nejmenšího disku na fotce (protože jsou na každém místě disky seřazeny od největšího po nejmenší, je tím jejich poloha určena jednoznačně). Na obrázku vidíte rozmístění kotoučů odpovídající řetězci BACB a také rozmístění odpovídající řetězci AOCB, které vzniklo z prvního rozmístění přesunutím nejmenšího kotouče na místo A.



Víme, že pracovníci přemístili všechny kotouče z místa A na místo B zpětně, při kterém byl počet přesunů kotoučů mezi místy nejmenší možný. Seřadte zadané fotografie podle času jejich pořízení.

| | |
|------------------------|-------------------------|
| <i>Ukazkový vstup:</i> | <i>Ukazkový výstup:</i> |
| CCA | AAA |
| CCB | BCA |
| BCA | CCA |
| AAA | CCB |

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.

Podotčení, že počítací byl napaden úřem, se ze záznamů nepobírá. *Hardwarové schéma bylo téměř vyloučeno, a proto nešlo nic jiného, než přikročit ke zkoumání zdrojového kódu. Tím programátorní se několik týdnů pokoušeli simulovat různé situace a testovat program na neplněných nastupech. A nakonec ji našli – tenkrát nevítanou chybou projevil se jen při kombinaci různých podmínek, které se právě toho nečkáštrádo odpoledne spánky.*

Výstrážný systém pro tazání údajů používal datovou strukturu, do níž mohlo simultánně zapisovat více vláken, tedy procesů, které spolu sdílejí data. V takovém případě se používají různé mechanismy zápisů, aby se taková struktura nerozbitla. Přestože systém toto obvykle zadržel dobře, při splnění podmínek – jednou z nich bylo i velké množství přidávacích výstrah – mohla vláknem zapisovat do stejného místa ve strukturně současně. Takový chybový stav se obecně nazývá race condition, a aby ve vícevláknové aplikaci nemohlo nastat, musí být každá věnována velká pozornost. Pochopení dat vedlo k chybám při přidávání nových zprávek, procesy se chodily nekorektně, zachytily se a přibíhaly server následně spadl. Záložnímu serveru však nemětená data předkl, k chybovým situacím docházelo nadále a i ten nakonec přestal fungovat.

Ve FirstEnergy bohužel nefungovalo nic, co by dispečerů na vypnutí narovněl účas upozornilo. Kvůli tomu je ani nemogoulo sledovat sř pomocí jiných prostředků a na to, co způsobili, začali přicházet už po samotném kolapsu.

Člověk údajost přečpřivatel

Jakub Maroušek

¹ <http://ksp.mff.cuni.cz/viz/12/codex>

² Těba programek jménem conry. Ten sice počásá už umi zobrazovat sám, ale i tak je zajímavé cvičení ho to naučit po svém. http://www.planetarium.cz/meteo/PL_meteo.htm

27-4-7 Nástroj pro zpracování textu 14 bodů

Již v prvním dílu jsme si řekli, že Unix byl zpočátku z politických důvodů oficiálně vyvíjen jako „nástroj pro zpracování textu“. Je pomalu nčasé představit si jeho schopnosti v této oblasti. Ale nejprve motiváční příklad.

Příklad: Počasí

Svého času bývalo populární nechat si zobrazovat na ploše různé aktuální informace: například o počasí. Samozřejmě existovala spousta programů (zvláště na Windows), která dělala právě tohle: zobrazila na ploše informace o počasí. To je sice užitečné, ale jen omezené. Co když si vedle toho chcete zobrazit aktuální zprávy, čas do odjezdu nejbližšího autobusu, aktuálně přehrávanou písničku, čas výchozu a západu slunce, kurz dané koruny nebo cokoliv jinho? A co když tyhle všechny věci chcete naozac zobrazit někde jinde, řekneme třeba na nástěnných hodinách?

Učtitě chápete, že na každou z takových věcí vám samozřejmě program nikdo nenapíše. Unixový svět se k tomu postavil trochu jinak. V něm najdeme programy, které umí zobrazit na ploše cokoliv,² a jak asi tušíte, tím cokoliv je v tomto případě výstup nějakého jiného programu či skriptu. To nám nabízí poměrně bohaté možnosti, neb na rozdíl od programů kreslících na plochu, je snadné vytvořit šablony skriptu, který něco vypíše.

Zkusme si právě to: napsat skript, který na svůj výstup vypíše aktuální teplotu v nějakém městě. Pokud přidáte v Praze, dobrým zdrojem informací o počasí je meteostránice Planetařia a Strumovec.³ Relevantní kousek HTML kódu této stránky vypadá následovně:

```
<TD VALIGN="top" WIDTH="200">
<FONT [ ... ]>aktuální teplota vZadach</FONT>
</TD>
<TD VALIGN="top" WIDTH="150">
<FONT [ ... ]>-2.3°C</FONT>
</TD>
```

Pokud jste nikdy o HTML neslyšeli, zkonkultujte např. Wikipedia. Pro účely seriálu bude stačit vědět, že každá webová stránka je ve skutečnosti textový soubor, který popisuje, co se má užívateli zobrazit, a právě tento textový soubor vám curl změníme níže vypíše.

Podrženy je dáaj o teplotě, který bychom rádi získali. [...] jsou vynechané nezájímavé části. Jen vás upozorníme na obitř se zpracováním speciálních znaků (zvláštěých například na kódování, jako jsou diakritická písmenka, nebo třeba znak stupňů). Až s nimi budeme pracovat, dále v různých příkazech, budou typicky nahrazeny otázkami. Potřebovali bychom stránku stáhnout, vybrat z ní správný řádek a z něj oddělit číselnou hodnotu teploty. A nepřítř překvapivě na každou z činností poslouží jiný nástroj.

Všechny ty internety...

Již známe spoustu zajímavých unixových nástrojů pro zpracování informací (a na konci tohoto dílu budeme znát ještě více), ale zatím jediné, s čím umí pracovat, jsou soubory na našem disku. Když bychom jim tak dokázali „podstrčit“ data získaná z Internetu, otevře se nám nepřeberné množství nových možností, stále s tím stejnými nástroji.

27-3-7 UNIXové děja vu

S třetím dílem seriálu o UNIXu jste se poprali statečně. Nejvíce zážitelné bylo paradoxně v nejméně hodnotném prvním úkolu, jehož vzorové řešení si necháme až na konec, abyse se nechtli a nechteli hned na začátku. Trochu jste také zápsali se skriptovacími úkoly, kde se projevily jednak nepostarek zkušenosti s elektřinými kombinovanými šiflovacími utilit, jednak znalost získané jinde. Ty některým z vás dovolily napsat neporablní kód, který by v jiném shellu než Bash neobžel. Používající shell částěji, zkuste si s jeho pomocí šetřit práci, hraje si s nimi. Stručně, elegantní a efektivní vyjadřování vám v něm pak přijde snáz a nebudete už to tolik sváček k používání cizíchových konstrukcí v shellu.

Úkol 2 – Hardlinky na adresáře

Pohlédneme do man 1s nebo do předchozho dílu seriálu můžeme zjistit, že 1s -1d adresar vypíše podrobnosti o daném adresáři, nikoliv o souborech v něm umístěných, jako by to udělal bez přepínače -d. Počet linků je první číslo vypavo od práy. Experimentální pozorování na všech slusných systémech prozradí, že toto číslo je vždycky aspoň 2. Porovnáváním adresářů s nízkým a vysokým počtem linků přijde na to, že se o jedničku zvešší za každý podadresář. Po troše přemýšlení a přečetání výkladu v seriálu by mělo být jasné, že jeden link vede z nadřazeného adresáře, druhý z adresáře samotného (jmenovitě se . – tečka) a zbytek jsou dvě tečky (..) z podadresáře.

Mezi slusné systémy tu nemůžeme počítat Cygwin, který o každém adresáři tvrdí, že má jediný link. Tečka a dvě tečky implementuje úplně speciálně, a to hlavně kvůli podpoře windowsových souborových systémů. Při hlubším pohledu do POSIXu jsem s podivem zjistil, že norma takové chování dovoluje – konkrétně v definici prážného adresáře.

Tenhle úkol se ukázal být celkem jednoduchým na slusných i neslusných systémech.

Úkol 3 – Mazání obsahu a linku

Příklad cat </dev/null >soubor zkrátí délku souboru na null. Pokud před jeho zavoláním soubor smazáme příkazem rm (té operaci se říká také unlink – maže se jason jeden link), pítvodi data se nemění a vytvoří se nový prážný soubor. Důležitkem je, že pítvodi data jsou pořád přístupná přes zbývající linky, pokud nějaké existují. Pokud se soubor předeem nemalinkuje, přepíše se pítvodi obsah a změna se projeví při přístupu přes libovolný link.

Na tomto místě bychom rádi připomenli, že inođe je datová struktura, která si pamatuje metadáta souboru a odkazy na jeho datové bloky. Link tkezuje na inođe. Je to záznam v adresáři, který má jmeno (jmeno souboru) a číslo inođe souboru. Data adresáře se vlastně skládají jen z takových záznamů. V técho pojmech jste měli v řešení dost znaků.

Úkol 4 – Nucené privátání home

Při klasickém umístění domovských adresářů, tedy v /home/kazrvae12, je úloha představenými prostředky netešitelná. Vlastník souboru totiž vždycky může měnit jeho práva, takže bychom museli buď měnit práva adresáře home a odfitnout tak přístup do vlastních domovských adresářů ostřinými vlastnosti, nebo používat nějaký speciální mechanismus. Norma POSIX ale žádá nepokouřuje, jen na souborovém systému ext2 (a vyšších) jste objevili příkaz chattr, který si dovolim ignorovat jako řetězec neporablní.

Nemusíme se ale tak moc omezovat, v /etc/passwd je možné cestu k domovskému adresáři nastavit. To je dobře vědět, abyse se nikdy na klasické umístění domovského adresáře nespofetali a vždycky ve skriptech používali expanzi vlnky („broch...“) nebo obsah proměnné prostředí HOME (\$HOME expanduje na cestu k domovskému adresáři aktuálního uživatele).

Potřebný nápad je použítí nadřazeného adresáře pro kontrolu přístupu. Root si může přiřadit adresář nadřazený domovskému adresáři brocha, sebrat práva stětu (others), jako skriptového Vlastnika nastavit skriptu obsahující jen brocha a této skriptu dát právo k prohlédávání (x). Nikdo jiný než root a broch tak nemá možnost se k brochovu domovskému adresáři dostat, ať se snaží sebrat.

Budeme tedy potřebovat skriptum, ve které je broch sám. Často taková skriptina již existuje a jmenuje se broch. Pokud nekustuje, snadno ji vytvořime:

```
groupadd broch
usermod -ag broch broch
```

Ted už jen vytvořim „reprehlitný“ adresář a přesunout do něj brochův stávající domovský adresář:

```
usermod -md /home/broch_priv broch
chown root:broch /home/broch_priv
chmod 750 /home/broch_priv
```

Úkol 5 – Preprocessing

Od pohledu chceme použít jeden cyklus přes řádky a jeden cyklus přes tokeny na řádku. Když na řádku najdeme token COMMENT, zahodime všechno až do konce řádku – tomu odpovídá příkaz break ve vnitřním cyklu. Token BYE odpovídá vnitřnímu return – funke v ten okamžik má skončit a víc se vstupu nečít.

Budeme číst zdroják, takže se nám bude hodit prepnačr-příkazem read. Jinak by zpětná lomítka dělala neplehnu. Při volání funkce se hodí možnost mít lokální proměnné. Těba když měnime proměnnou IFS, není pěkné ji nechat změněnou po návratu z funkce, jiné kusy kódu se pak mohou chovat neočekávaně. POSIX v tomto ohledu nabízí dvě řešení, obě dost španná.

- Můžeme uložit starou hodnotu IFS do jiné proměnné, změnit IFS a nakonec vrátit IFS původní hodnotu. Na to se ale snadno zapomná, funkce bývá možné opustit více způsoby a po návratu z funkce zůstane nastavená zbývečná proměnná.
- Také můžeme celé volání funkce prováčet v subshell – stačí tělo funkce obalit do kulatých závorek a místo return volat rovnou exit. Nevýhody jsou nasmadě: pouští se další procesy, nejde snadno nastavovat melokání proměnné nebo ukončit shell.

Bash nabízí třetí řešení, které jiné shelly nemají. Tím je vestavěný příkaz local, o kterém se více dozvíte v help local. Používá se podobně jako export a příklad použití najdete ve vzorovém řešení.

Už ve druhém dílu seriálu jste potkali konstrukce for, case a podmínky. Podmínky jste se učili i pomocí ||. Přidějme k tomu obsah třetího dílu a řešení je na světě.

Program (shell):

```
http://ksp.mff.cuni.cz/viz/27-3-7-preprocess.sh
```

Zvolil jsem nejjednodušší realizaci smyčky přes tokeny. Nedal jsem shell pomocí IFS rozdělň řádku vstupu podle mezer a přes výslechnou množinu slov jsem psal for-cyklns.

Štaří tedy z počátečního stavu prodáváček do šifry. Tím máme zaručeno, že projdeme všechny dosažitelné stavy a navíc je problémy v tom pořadí, v jakém násilnějí při průlomu sou- dy. Můžeme tedy přiblížit ve všech stavech počítat největší získatelné množství energie.

Při přechodu z (a, b) do $(a+1, b)$ odebereme energii rovnu vzdálenosti mezi misy a a $a+1$. Pokud však přicházíme přes prošlou oblast, tedy například z (a, b) do $(b-1, a)$, musíme odečíst i součet vzdálenosti mezi všemi již prošlými misy. Nakonec resimne zapomenout přičíst energii získanu na nové pozici.

Čelý přechod je nakonec velmi jednoduchý, jenom se nesmíme ztratit v indexování a přičítání ± 1 . Na detaily se podívejte do zdrojového kódu. Casová složitost celého algoritmu je $O(N^2)$, protože každý z kvadraticky mnoha stavů znovu musíme nejdříve jednout. Díky tomu, že nepotřebujeme optimální cestu soudu zkonstruovat, vystačíme si s lineární pamětí.

Program (C):
<http://ksp.mff.cuni.cz/viz/27-3-4.c>

Jenda Hadwara

27-3-5 Komprese obrázu

Protože na vstupu dostaneme vždy obrázek o rozměrech $2k \times 2k$, v každé úrovni rekurzivního komprimování pracujeme vždy se čtvercem. Pokud má na aktuální úrovni čtverce rozměry 1×1 , nic již nekomprimujeme a jen vrátíme jeho barvu.

Co přesně máme v každém kroku kvadrantistické komprese dělat? Potřebujeme vybrat dvě čtverci uvažovaného čtverce a každou „odbyť“ jednou barvou, respektive jednou pro- blášt za celobloku a druhou za celoblokem. Řekněme, že cenou výsledné komprese je počet pixelů, které smysly být přepravky.

To nás může svádět k řešení, ve kterém pro každou čtverci spočítáme počet čtverců a jejich pixelů. Následně „nejčert- nější“ z nich obarvíme černě, nejbližší bíle a zbylé čtverci pak zpracováváme rekurzivně, přičemž rekurzivní funkce bude vracet cenu zakomprimování dané části. Toto řešení má dva problémy – jednak není úplně jasné, jak vybrat nejbližší a nejbližší čtverci (může jich být více se stejným počtem pixelů/čtverců pixelů), a navíc nemusí vést k optimálním řešením.

Mohlo by se totiž stát, že sice v jedné čtverci je hodně čtverců pixelů, více než v ostatních, ale díky rozložení čt- verců pixelů se dobře kvadrantisticky komprimuje. Pak by ji- temo algoritmus obarvil celobloku a namsto toho se kvad- rantisticky snažil komprimovat čtverci, která má pixely rozložené dost nepravdělně.

Zkusme to trochu jinak – na každé úrovni se pro každou čtverci započítáme naši rekurzivní funkce, kolik by stálo její zakomprimování. Budeme si pamatovat i počet čtverců a bi- lých pixelů v každé čtverci. Pokud známe tyto informace, můžeme prostě vyzkoušet všech 12 možností, jak vybrat ce- ločíslo na celobloku čtverci, a seskládat pro každou z mož- ností její cenu. Ze všech dvacícti cen si pak vybereme tu nejmenší.

Kolik na takový výpočet potřebujeme času? Povšimněme si, že rekurzivní funkce se nikdy nevolá znovu pro stejný čtverec, protože uvažované čtverce stejné velikosti se nikdy nepřekrývají. Víme, že na každém čtverci velikosti $2^m \times 2^m$

strávíme $O(4^m)$ času, a takovýcht čtverci uvažujeme na kaž- dé úrovni rekurze $\frac{4^k}{4^m}$. Celkem tedy potřebujeme $O(4^k \cdot K)$ času i paměti.

Toto řešení by šlo ještě zrychlit: určit počet bílých pixelů pro každý čtverec lze v čase $O(1)$, protože stačí sešit- ť výsledky ze čtyř menších čtverci. A pokud si vybereme „od- byté“ čtverci budeme uvažovat trochu lépe, dostaneme řešení, které si vystačí s lineárním časem i pamětí v počtu pixelů.

Program (C++):
<http://ksp.mff.cuni.cz/viz/27-3-5.c++.cpp>

Ondřej Hübsch

27-3-6 Ukládání přepravky

Některí z vás poznali, že jde jen o jiné zadání daleko zna- mější úlohy, totiž barvení intervalových grafů. Tato úloha se dá ve zkratce zadat jako hledání nejmenšího počtu po- sluháren, které potřebujeme pro přednášky zadané svými časy začátku a konce.

Tato úloha se dá řešit hladově, a to ve velkém množství variant, proto byla většina řešení funkcí. Hlavní část řešení je ovšem obhájit, že hladově řeší opravdu jde.

Začneme tím, že si přepravky seřadíme podle vnějšího pří- měru sestupně. Na vnitřním průměru tedy nezáleží, jak později uvidíme. Budeme si udržovat seznam potřebných komuniků, jakési přibížeře řešení. Jedná hodnota, kterou z každého komuniků potřebujeme, je vnitřní průměr nejmen- ší přepravky.

Přepravky budeme postupně probírat a snažit se je vložit do nějakého komuniků. Protože víme, že žádná větší přepravka už nepřijde, stačí se seznam vybrat maximum. Bude se nám tedy hodit reprezentovat komuniků maximum haldu. Pokud je maximum menší než vnější průměr aktuální pře- pravky, pak žádný vhodný komunik neexistuje a proto při- dáme nový. Na konci běhu algoritmu prostě jen spočítáme počet komuniků ve stroju.

Konečnost algoritmu je zřejmá, podíváme se na správnost. Dokážeme ji indukci podle velikosti průběžného řešení. Do- kad máme jen jeden komunik (případně žádný), nemůže existovat lepší řešení.

Mějme tedy k komuniků. Jediné, co může k zvýšit, je pak situace, kdy neexistuje pro nějakou přepravku vhodný ko- munik. V tu chvíli ale existuje alespoň k přepravka (máji vnější průměr alespoň takový jako aktuální přepravka (dlky setřezení), a vnitřní průměr ostře menší (neexistuje vhodný komunik). S aktuální přepravkou je to tedy $k+1$ přepra- vka, ze kterých žádná dvě nejdou vložit do sebe. Optimální řešení tedy nemůže používat menší než $k+1$ komuniků.

Jaké vlastnosti algoritmus má? Už kvůli prvotnímu seřá- zení vídme, že nepoběží rychleji než v $O(N \log N)$. Další kroky algoritmu závistí na počtu komuniků, tedy na řešení. Ten nemůžeme nijak rozumně odhadnout – počet komuniků může být až N . Každý krok tedy zabere až $O(\log N)$ času. Tedy nakonec se opravdú do $O(N \log N)$ vejde. Paměti spotřebujeme $O(N)$, více nepotřebujeme.

Vzorový kód je napsaný v C++ a využívá standardní imple- mentaci haldu.

Program (C++):
<http://ksp.mff.cuni.cz/viz/27-3-6.c++.cpp>

Ondra Hlavatý

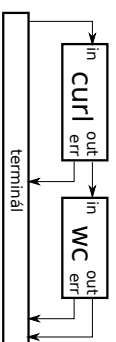
To můžeme zařídit programkem curl, který na spoustě uni- xových systémech najdeme ve výchozí instalaci, případně si jej lze snadno doinstalovat (i v Cygwinu). Jeho použití je přímocare: jako parametr dostane URL webové stránky (či stránku určitého souboru) a na standardní výstup vypíše je- ji obsah (HTML kód). Odhul jej můžeme přeměňovat do souboru či poslat kamkoli rovnou, jak jsme zvyklí.

Zkusme si spustit například
`curl http://ksp.mff.cuni.cz/ | wc -c`

Asi vás překvapí, že kromě očekávaného výstupu se obje- vi několik podivných řádek, které vypadají jako informace o průběhu stahování. Jak je to možné, když je výstup při- kazu curl přeměňován?

Ve skutečnosti má každý proces kromě svého standardní- ho vstupu (též státní) a výstupu (státní) ještě třetí „datový kanál“, takzvaný *standardní chybou výstupu (stderr)*. Ten směřuje na terminál, i když je výstup příkazu přeměňován do souboru či rovnou. Jak již název napovídá, slouží k tomu, že když za běhu příkazu nastane chyba, uvidí ji uživatel na- místo toho, aby se zapsala doprostřed výstupního souboru. Ale nepoužívá se jen pro chyby a varování, nýbrž i pro infor- mace o stavu a průběhu programu. A právě curl na něj vy- píše informace o průběhu stahování, díky čemuž když sta- huje velký soubor příkazem curl `http://adresa.>sou- bor`, vídíte, kolik již je stáženo a kolik času ještě zbývá.

Náš *pipelina* (tímto pojmem se označuje řada příkazů po- spojitelných rovnou) ve skutečnosti z pohledu operačního systému vypadá takto:



Tedy už je jasné, když se ona hlášení na terminálu dostanou. I `stderr` se dá přeměňovat, a to operátorem `2>`. Nejčastěji se to používá k umlčení všech hlášení, například takto: `curl http://...>2>/dev/null | wc -c`. O speciálním souboru `/dev/null` byla řeč v minulém díle. Je třeba dávat pozor, že křetením příkazů přeměňování patří. Pokud byste ho na- psali na konec, křženošlo elektu nedosáhnete, neb přemě- nujete `stderr` procesem `wc`, nikoli `curl`. V případě `curl-2` ale přeměňování používat nemůžete, neboť nabízí přepínač `-s`, který stavové zprávy utiší.

Na závěr dodáme, že ke stejnému účelu jako curl lze použít i příkaz `wget`, který ovšem ve výchozím nastavení nkládá sřazenou stránku do souboru. K vypsobání na `stderr` ho přiměňte parametrem `-0 -`, hlášení o průběhu umlčíte `-q`. Více v manuálové stránce.

Filtrování řádek: grep

Občas by se nám hodilo z textového souboru vybrat řádky splňující nějaké kritérium. Nejjednodušším nástrojem, kte- rý takovou věc dělá, je `grep`. Ten jako parametr přijímá slovo (kus textu), čte postupně řádky ze svého vstupu a na výstup vypisuje jen ty, které zadané slovo obsahují (mysle- mo obshnují jako podřetězec, nemusí být oddělené meze- rami). S parametrem `-v` (in Vert) naopak vypisuje jen řádky *neobshnující* dané slovo. Pokud hledané slovo obshnuje ně- jaké „divné“ znaky, je třeba `grep-n` dát ještě parametr `-F`.

Které přesné znaky jsou divné a proč, si vysvětlíme později, prozatím můžeme za bezpečná považovat písmena a číslice.

S parametrem `-i` ignoruje `grep` při hledání slova velikost písmen (v závislosti na nastavení systému nemusí fungovat pro české znaky).

Pokud dáte `grep-n` více parametrů, další jsou brány jako názvy souborů, ve kterých se má hledat. Tedy `grep slo- vo >soubor1 >soubor2` se slova podobně jako `cat >soubor1 >soubor2 | grep slovo`, s tím rozdílem, že pokud je sou- bor více než jeden, `grep` před každý vyhovující řádek napíše název souboru, ze kterého pochází (to se dá vypnout přepí- nácem `-h`). Například když bych chtěl zjistit, zda jsme v ně- kterém díle seriálu zmínováli proměnnou `$RANDOM`, použiji:

```
$ grep -F '$RANDOM' serial1*
[...].text:3:tex:3:$(($RANDOM % 100))
[...]
```

a vídím, že to bylo ve třetí sérii. Pokud nás zajímají jen názvy souborů, ve kterých se slovo vyskytuje, můžeme po- užít přepínač `-l` (dobře se pamatuje podle 1s, které taky vypisuje názvy souborů), případně `-l` pro operátory opera- či (výpis souborů neobshnujících ani jednou dané slovo). S přepínačem `-r` můžeme `grep-n` předávat za parametry i názvy adresářů: v těch pak prohlédá všechny soubory re- kurzivně. `grep -lr bagr` najde v domovském adresáři všechny soubory obsahující slovo `bagr`. Takového hledání může divití trvat.

Přepínačem `-C K` (Context) řeknete `grep-n`, že má kromě vyhovujících řádek vypsat i K řádek okolo každé z nich. Pokud místo `-C` použijete `-B` (Before) či `-A` (After), bu- de vypsat jen kontext jednostranný (K řádek před, resp. za každým vyhovujícím). Pokud by se kontexty překřýva- ly, jsou šlhy do jednoho bloku, žádné vstupní řádky se na výstupu neobjeví dvakrát. Tedy např. `grep -C 999 : na /etc/passwd` vypíše to samé, co `cat`. Pokud na sebe sou- sední kontexty nenavazují, jsou odděleny řádkem s dvěma pomlčkami pro snazší vizuální orientaci.

Pokud používáte `grep` ručně, mohl by vás zajímat přepínač `--colort`, který výskry hledaného slova barevně zvýrazní. Ve skriptech by se vám naopak mohl hodit přepínač `-q`, kte- rý zprůsobí, že se na výstup nevypíše nic (jako `/dev/null`). K čemu je taková věc dobrá? Zatím jsme vám zatajili, že `grep` vrací nulovou návratovou hodnotu, pokud našel ales- poň jeden vyhovující řádek, jinak nenulovou. Takže můžete použít `if $(grep -q slovo soubor` pro test, zda je v sou- boru obsaženo slovo.

Tedy už máme asi vše potřebné k výběru správného řádku:

```
$ curl -s ... | grep -a -A 3 "teplota vzduchu" \
    | head -n 4 | tail -n 1
<FONT COLOR="#FF000" FACE="Arial">0.47</FONT>
```

Od začátku psaní se trochu otevřelo. :-)

První `head` je potřeba, protože se text na stránce vysky- tuje vícekrát (a kvůli podivné kódováním českým znakům nemůžeme hledat slovo „aktální“). Přepínač `-a` slouží k to- mu, aby `grep` vypsal obryklý výstup, i pokud považuje sou- bor za binární (obshnuje nějaké neobvyklé znaky).

Jestě snad dodáme, že zpětné lomítka na konci prvního řádku znamená, že příkaz pokračuje na řádku následujícím. Do skriptu to můžeme napsat přesně takto, jen je třeba dát si pozor, aby za lomítkem nebyla žádná mezera. Pokud by jste chtěli příkaz spustit v terminálu, je lepší napsat vše na jeden řádek (a lomítka vynechat).

Sekání a slopování řádek: cut, paste a spoli.

Už umíme vybrat ze souboru celé řádky. Občas by se nám mohlo hodit získat i jejich části. K tomu nám poslouží příkaz cut. Ten se nejčastěji používá ve souboru „tabulka.csv“, charakteru (např. /etc/passwd), kde každý řádek je rozdělen na několik „slopoveček“ nějakým oddělovačem (v Unixu typicky dvojtečka či libovolná posloupnost bílých znaků – s touti si ale cut neporadí). Důležitě přepne jsou -d, který nastavuje oddělovač (má být jednoznačný) a -f, jenz říká, které slopovečky chceme na výstup. Jeho parametrem může být jedno číslo, rozsah čísel (3-5), případně seznam čísel a rozsahů oddělovačů čárkami (1,3,5-10). Slopovečky jsou číslovány od jedničky.

Kombinací příkazů cut a grep můžeme nyní třeba zjistit ID uživatele hroch:

```
$ grep hroch /etc/passwd | cut -d: -f3
4242
```

To ale není úplně spolehlivé, například se to rozdíje, pokud se někdo bude jmenovat drubhroch: později si ukážeme lepší způsob. Stejně jako sponsta, jiných příkazů, pokud cut dostane jako parametr jeden nebo více názvů souborů, čte z nich, jinak čte ze standardního vstupu. Zapsisjte vždy na standardní výstup. Slopovečky nejdle prohazovat: -f 3, 1 je to samé, jako 1, 3.

cut má ještě alternativní režim, kdy místo slopovečky vysekává z řádků jednotlivé znaky; např. cut -c 1-3 vybere z každého řádku první tři znaky.

Inzerovní operaci ke cut je paste, který dostane za parametry několik souborů, které považuje za jednotlivé slopovečky. Pak vezme první řádek z každého souboru a všechny spojí zadaným oddělovačem (-d), čímž vznikne první řádek výstupu. A tak dále pro další řádky. S ním bychom mohli, byť trochu neobratně, zarídit například zmiňované prohlazení slopoveček:

```
$ cut -d: -f1 /etc/passwd > jmena
$ cut -d: -f3 /etc/passwd > uid
$ paste -d: uid jmena
0:root
4242:hroch
[...]
```

Formát „slopoveček“ oddělených dvojtečkami je sice příjemný pro strojové zpracování, ale člověku se čte o dost hůře než opravdové slopovečky, kde jsou odpočítávají si hodnoty zarovnané pod sebou. Takovýto oku přívětivý formát lze vytvořit programem column.

Ten pracuje ve dvou režimech. Prvím z nich je tabulkaový (column -t -s oddělovač), ten načte ze vstupu soubor v „oddělovačovém“ formátu a na výstup ho vypíše jako hezkon tabulku:

```
cut -d: -f 1-4 /etc/passwd | column -s: -t
root      x      0      0
hroch     x      4242   4242
[...]
```

Slopovečný režim očekává na vstupu jednořádkový seznam položek, které vypíše na výstup ve víceřádkové sestbě, podobně, jako to dělá ls. Tím se dá šetřit místo na obrazovce, pokud jednotlivé vstupy řádky jsou krátké. Například kdybychom chtěli vypsat seznam všech uživatelských jmen v systému:

```
$ cut -d: -f1 /etc/passwd | column
root      hroch     guest
hroch     hacker    nobody
```

Počet slopoveček je zvolen automaticky, dá se ovlivnit přepínači, stejně jako se dá zarídit, aby se hodnoty vyhlazovaly po řádkách namísto po slopovečkách (použijte jen pokud víte, co děláte, čte se to hrozně).

Nahrávání znaků: tr

V nejjednodušším použití tr nahradí všechny výskytý jednoho znaku (určeného prvním parametrem) na svém vstupu za jiný (druhý parametr) a výsledek vypíše na výstup. Každý z parametrů může být i celým seznamem znaků, zadaným buď vymeňováním těsně za sebou, rozsahem (např. a-z) nebo kombinací obojího. Pak platí, že každý výskyt nějakého znaku z prvního seznamu se nahradí za odpovídající znak z druhého seznamu.

Například tr a-z A-Z převede svůj vstup na velká písmena (funguje jen pro znaky anglické abecedy, ostatní nechá beze změny), tr a-zA-Z A-Za-z prohodí velikost písmen (z velkých udělá malá a z malých velká). Pokud některý ze seznamů je kratší, je doplněn zopakováním posledního znaku. Např. tr aeiouy e nahradí všechny (malé) samohlásky v textu za e.

S přepínačem -d očekává tr jen jeden seznam znaků a všechny znaky na tomto seznamu ze vstupu smaže. Přepínač -c vezme místo prvního seznamu znaků jeho doplněk. Nejstříhší se používá spůl s -d pro smazání všech znaků kromě zvolených nebo s jednoznačným prvním seznamem za tr -c a-zA-Z - nahradí všechny znaky kromě písmen za podzřítku.

Řešení: Počasí

Nyní už máme všechny střípky k vyřešení našeho úvodního příkladu:

```
$ curl -s \
http://www.Planetarium.cz/meteo/PL_meteo.htm \
  | grep -a -A 3 "teplota vzduchu" \
  | head -n 4 | tail -n 1 \
  | cut -d'>' -f2 | cut -d'<' -f1 \
  | tr -cd '0-9.-'
-0.7
```

Takovéto číslo si pak můžeme nejen nechat někde zobrazit, ale také s ním libovolně dále pracovat. Například by nebylo těžké napsat skript, který se před vypnutím počítače podívá na aktuální teplotu, a pokud je méně než 15, zobrazí upozornění „Vezmi si bundu!“

Rest: seq

Při povídání o for-cyklech jsme vám zatajili velmi užitečný příkaz seq. seq A B na svůj výstup vypíše všechna čísla od A po B (obojí včetně, A lze vyměnit, pak se vypisuje od jedničky), každé na samostatný řádek. Spolu s for-em se dá použít pro zopakování nějakých příkazů n -krát.



(Kdybychom se rozhodli slova nevíkládat do tří, ale řádit, zvládneme to také lineárně – díky pevné velikosti abecedy to jde pomocí RadixSortu.)

Za zámku možná stojí, že v řetěze uloze občas realita správně asynchronizace a pošle ji stydět se do kouta. I nepřilíš optimalizovaná řešení se složitostí $O(\log l)$ mlžou pro to znume velké vstupy doobhnout rychleji než řešení lineární. Souvisí to mimo jiné s tím, že tře je budování pomocí ukazatelů, a s efekty keší... ale to už by byl úplně jiný příběh.

Program (C):

```
http://ksp.mff.cuni.cz/viz/27-3-2.c
```

Karolina „Karyyanna“ Burešová

27-3-3 Výběr vysílaců

Ulohou je vlastně obarvit strom barvami 1, 2, 4, 8, ... tak, aby sousední vrcholy měly různé barvy a součet hodnot barev přes všechny vrcholy byl co nejmenší. Připomeneme, že bez požadavku na minimální součet lze každý strom korektně obarvit dvěma barvami. To můžeme udělat například přichodem stromu do hloubky. Ten může vypadat například takto:

1. obarví (v, b):
2. barva[v] = b
3. for (u in soused(v)):
4. if (barva[u]==0): obarví(u, 3 - b)

Negativně obarvime libovolný vrchol barvou 1, pak víme, že jeho sousedé musí mít barvu 2, jejich sousedé zase barvu 1, a tak dále. Tento jednoduchý algoritmus má časovou složitost $O(N)$ a budeme na něm dále stavět.

Tím, že strom umíme obarvit barvami 1 a 2, dostáváme obarvení se součtem maximálně 2N. Tedy hledané obarvení s minimálním součtem určitě nepoužije barvy větší než 2N. To zároveň znamená, že barev použijeme maximálně $\log N + 1$.

Nyní pomalu přejdeme k popisu algoritmu. Strom si zakoreníme v libovolném vrcholu a z něj postupně prohledáváme do hloubky. To vždy nejříve spočítá řešení pro podstrom tvořené syny vrcholu a z těchto řešení složí řešení pro daný vrchol. Tento postup je takovým standardním dynamickým programováním na stromě. A co přesně v podstromech budeme počítat?

Pro každou barvu $c = 2^i$ pro $i = 1, \dots, \log N + 1$ spočítáme nejlepší možné obarvení podstromu, ve kterém je kořen obarven barvou c . Abychom zjistili nejlepší obarvení pro barvu c , tak stačí pro každého syna vybrat nejlepší barvu jinou než c . To pro konkrétního syna vždy bude buď jeho nejlepší barva, anebo jeho druhá nejlepší barva. Třídá nám stačí si pro každý podstrom pamatovat pouze dvě nejlepší možnosti.

Po zavolení výpočtu v kořeni stromu dostaneme výsledek. Během výpočtu pro každý vrchol zkonštruje $\log N + 1$ barev plus ($\log N + 1$)-krát n něj vybráme jednu ze dvou barev pro otce. Tedy časová složitost algoritmu je $O(N \log N)$. My ale algoritmus ještě vylepšíme.

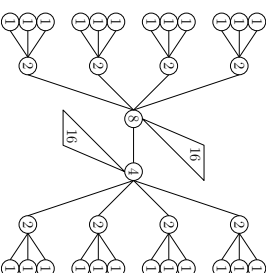
Všimneme si, že pokud se má vyplatit vrcholu přiřadit barvu 2, tak všechny barvy z $2^1, \dots, 2^{i-1}$ musí být zastoupeny v jeho synech. Tedy v vrcholu s k synů nám pro získání dvou nejlepších možností stačí vyzkoušet $k+2$ barev (druhá nejlepší možnost pořadí může mít barvu $k+2$).

Další věc, kterou na předchozím algoritmu můžeme vylepšit, je operativně vybrání z nejlepší a druhé nejlepší barvy synů. Vime, že vždy vybereme tu nejlepší kromě případu, kdy pro kořen zvolíme stejnou barvu, pak vybereme druhou nejlepší. Stačí nám spočítat součet nejlepších možností synů a pro každou barvu $c = 2^i$ pro $i = 1, \dots, k+2$ si předpochat, o kolik se dohoromady liší nejlepší a druhé nejlepší možnosti synů, kteří jako syny nejlepší barvu mají c . To pro vrchol s k synů můžeme jednoduše udělat v čase $O(k)$.

Jelikož každý vrchol je synem právě (maximálně) jednoho jiného vrcholu, dostáváme se na časovou složitost $O(N)$.

Pro zajímavost ještě řekneme, že stejný algoritmus funguje i pro verzi úlohy, kde strom barvime barvami 1, 2, 3, ..., N. Sami si můžete rozmyslet proč.

A kdo je zvědavý, tak strom, pro který jsou potřeba alespoň čtyři barvy, může vypadat například takto:



„Šestnáctkové“ trojhlubky zastupují šestnáct přímých potomků uzlu, kde je každý obarvený jedničkou.

Karel Tesar

27-3-4 Doplnování energie

První, co nás může napadnout, je vyzkoušet všechny možnosti příchodů. Těch je však velmi mnoho, protože nám nic nebrání chodit přes některá místa vícekrát.

Ujasněme si nejprve, jak bude vypadat přetek, který může získat maximální množství energie. Protože nemáme žádné omezení na uvolňovanou energii, vyplácí se nám vždy při prvním navštívení nějakého místa z něj všekrou energeti vyčerpát. Nemá tedy smysl si ji zde šetřit na později. Stejně tak je zbytečné nějaké místo při přeletu vyměchat a nezastavovat se na něm.

Díky tomu můžeme předpokládat, že při každém teprání energie máme pro jítom souvislou oblast, a navíc se nacházíme na jejím kraji. Tím jsme výrazně snížili počet možností, které dceeme vyzkoušet. Nyní máme pouze N levých a N pravých konců, celkem N^2 možných koncových stavů. Stav je tedy dvojice čísel (a, k) , kde a je aktuální pozice a k druhý konec proše oblasti.

Jak ve stavech spočítat optimální množství energie? Pro stav s oběma konci intervalu na startovní pozici M energii určime snadno. Bude rovna právě energii, kterou z této pozice můžeme teprát.

Do ostatních stavů se můžeme dostat nejryše ze dvou předchozích. Konkrétně pro $a > k$ se do stavu (a, k) dostaneme ze stavu $(a-1, k)$ nebo $(k, a-1)$. A pro $a < k$ to jsou stavy $(a+1, k)$ a $(k, a+1)$. Do nich se umíme dostat zase ze dvou předchozích, a tak dál. V každém případě můžeme začít pouze ve stavu (M, M) .

Náš algoritmus počítají s délkami intervalů a plochami čtverců jako s normálními čísly. Tato čísla ale rostou exponenciálně s hloubkou stromu, takže by se mohlo stát, že se nevejdou do běžné celočíselné proměnné. Zadáni o této možnosti nemáme námitky, ale zkusme ji aspoň na chvíli připsát: Dřívě si kloboňky...

První z algoritmů na konstrukci grafu s délkami ani plochami nepočítá, takže ho můžeme nechat.

Druhý algoritmus používá délky při propojování seznamů intervalů. Můžeme si pomoci takto: podíváme se na intervaly x a y na začátku seznamu. Pokud jsou stejného řádu, nic se nemění. Pokud je $\text{B}(\text{UNO}) x$ vyššího řádu, rozdělíme ho na dva intervaly o 1 nižšího řádu a postup opakujeme. Části vzniklé rozdělením zředí informace o vrcholech grafu a o barvě z původního intervalu. Jako cvičení ponecháváme dokázat, že vytvořené nejvyšší lineární množlo nových intervalů (nápodobě představitel si strom popíšijte postupně dělení intervalů).

Hledání komponent souvislosti získáme bez zmatků, ale pak potřebujeme počítat plochy komponent. Nejprve každému čtverci přiřadíme číslo komponenty. Poté projdeme celý strom po hládnatých shora dolů, a kdykoli narazíme na jednohrotový čtverec, přidáme ho k přislíbené komponentě. Pro každou komponentu tudíž vznikne seznam jejích čtverců uspořádaný sestupně podle řádu.

Nyní v tomto pořadí budeme sčítat plochy čtverců. To jsou potenciálně ohrovenská čísla. Budeme je zapisovat ve dvojkové soustavě a od paměti ukládat jako seznam pozic jedniček ve dvojkovém zápisu, uložené od nejvyššího bitu k nejnižšímu. Každý další čtverec přitom přispěje plochou, která je menší nebo rovna zatím nasčítané ploše. Můžeme tedy připsat jedničkový bit na konec čísla, jen se nám mohlo stát, že už jím jeden bit tohoto řádu mohl být, takže dojde k přenosu. Vyrizováním všech přenosů ovšem strávíme lineární čas, protože s každým přenosem klesne celkový počet jedničekové bitů o 1. (To je podobná úvaha jako orelidy s intervaly.)

Plochy máme spočítány, zbývá z nich najít maximum. Kdykoli při tom porovnáváme dvě čísla, procházíme je od nejvyššího řádu a jakmile se přestanou shodovat, porovnáváme u ukončíte. Čas strávený porovnáváním přitom účtujeme jedničkám v menším z obou čísel, které se už žádádného dalšího porovnávání neúčastní. Takto celkem načítáme každé jednotkové čas $O(1)$ a všech jedniček za celou dobu běhu algoritmu vznikne $O(n)$.

Is ohrovenskými čísly jsme tedy dokázali nahřít lineární složitost algoritmu. Kouzlo se podařilo :-)

Martin „Machdick“ Mareš *et Filip Štědroňský*

27-3-2 Návrhy pro komisi

V álouze se pracuje s řetězci, zkrácený řetězel si tedy přete zadání, zamysli se a řekne si „Ha, triel“. To je speciální strom, ve kterém se hranám přiřazují písmenka a cesta od kořene nějak odpovídá vstřpním řetězcům: více o ní píšeme v knihučce o hledání v textu.¹⁵

Jenže jak přesně tři na naši tlouh použít? Můžeme jednoduše počítat, kolik členů komise sděluje určité slovo – stačí si do vrcholů přidat počítadlo. Když se pak vydáme

od kořene dolů, a budeme tyto členy postupně sčítat, zjistíme nás, jestli v daném vrcholu součet dosahuje alespoň hodnoty K . Pokud ano, jakýkoliv návrh začínající slovem odpovídajícím aktuálnímu vrcholu bude schválen.

Kolik takových návrhů existuje? Když aktuální hloubku označíme jako d , dá se její počet vyjádřit jako 26^{d-1} . Za každý ze znaků, který dýchá do přijatelné délky návrhu D , totiž smíme zvolit libovolný znak anglické abecedy.

Také z uzlu, v kterém se nachrálo alespoň K hlasů, nedreem postupovat níž – všechna možná pokračování budou začínat schvalovaným řetězcem, takže jsme je už zahrnuli. Při dosažení hodnoty K se tedy zastavíme, resp. vrátíme o úroveň výš a nebudeme pokračování postupovat dál.

To zní dobře. Ale nebude to opravdu fungovat? Co když by nějaký člen schvaloval jak slovo „psa“, tak slovo „psal“? Ouhai! To bychom ho započítali vícekrát, což nedreeme – a zadání nám rozhodně nesluhuje, že taková situace nenastane.

Připomenutí si naznačeno myšlenku: pokud člen (či komise) schvaluje nějaké slovo, sděví i jakékoliv další ktere tímto slovem začíná. Když bychom tedy věděli, že člen už schvaluje nějaký prefix aktuálního slova, můžeme právě zpracovávané slovo s kldem zahodit.

Tady se nabízejí dva přístupy: Můžeme slova každého člena lexikograficky seřadit a přidávat je do tree už seřazená. Jen si musíme rozmyslet, jak v tri komise poznat, zda jsme do aktuálního vrcholu už přičítli hlas právě zpracovávaného člena.

Alternativu, se kterou pracuje i vzorový program, představuje vřhndování druhé třídy, tentokrát specifické pro člena. Do ní jednoduše naházíme všechna jeho oblíbená slova (zaráčime si ve vrcholu, jestli tam končí slovo). Následně ji projdeme a zřídí, když narazíme na konec slova, odpovídatí slovo přidáme do tree pro komisi a hned se vrátíme.

Pro zpracování slov všech členů tři projdeme tak, jak jsme popsali na začátku. Teď už hlas každého člena započítáme pro libovolný řetězec maximálně jednon, takže program vřda spřávný výsledek. Zbývá se zamyslet nad složitostí.

Co je vlastně vstup? Kromě parametrů C , K a D to jsou zejména všechna oblíbená slova. Označme si počet všech znaků v nich jako ℓ .

Operace s trii, když máme konstantně velkou abecedu, můžeme bez obav prolášt za konstantní. Do tree pro komisi přidáme maximálně ℓ znaků, čímž vytvoříme maximálně ℓ vrcholů. Při závěrečném průchodu navštívíme každý vrchol nanejvyšší jednon, zpracování tree pro celou komisi tedy trvá $O(\ell)$ času a zabírá $O(\ell)$ paměti.

To samé platí pro trie jednotlivých členů. Zdálnivě tedy máme celkovou složitost $O(C \cdot \ell)$. Jenže trie všech členů nemůžeme mít dohromady víc než ℓ znaků, takže zpracování všech členských trii dohromady nemůžeme zabrat víc než $O(\ell)$.

Započítání jsme všechny operace? Skoro. Neměli bychom zapomenout na mocnění, ac ve všech řetězích jsem jeho zohlednění nevyžádovála. Zatímco násobení za konstantní prolášt můžeme, u mocnění by to bylo poněkud odvážné, zvlášť když by teoreticky mohlo nastat $D \gg \ell$. Třochu si ovšem život usnadníme a jen prohlásíme, že mocnění trvá $O(m)$. Pak má celý náš program časovou složitost $O(\ell \cdot m)$, paměťovou $O(\ell)$.

Přenosování řádků: sort, shuffle, tac

Představíme si skupinu programů, která nějakým způsobem mají pořadí řádků na svem vstupu. Nejdřívejším z nich je `sort`, který seřadí vstup čí soubor do zadáných kritérií. Ve výchozím nastavení třídí řádky abecedně vzestupně. Můžeme použít přepínač `-r` pro sestupné třídění, `-f` pro ignorování velikosti písmen, `-n` pro číselné porovnávání (abecedně) ve zatřídění „100“ před „11“) a `-k` pro třídění podle některého sloupceku (ve stejném významu jako u `cut`-u, oddělovací se nastavuje `-t`, výchozím oddělovačem je whitespce, tedy libovolná posloupnost bílých znaků). Například `sort -t: -k3 -n /etc/passwd seřadí` zřaznamy v `/etc/passwd` podle ID uživatele.

S parametrem `-R` `sort` místo třídění vstupní řádky náhodně zamíchá. To samé dělá `shuf`, jen nabízí nějaké parametry navíc, například vybrat ze vstupu náhodně jen K řádek (`-n K`) nebo povolit vybrat jeden řádek vícekrát (`-r`).

Například pokud byste byli učitel a měli v nějakém souboru uloženo hromadu otázek, ze kterých chcete vygenerovat 30 náhodných písemek po 10 otázkách, můžete napsat:

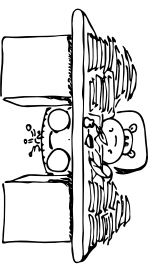
```
for i in $(seq 30); do
  shuf -n 10 otázky.txt >pisemka-$i.txt
done
```

Tac vypíše řádky vstupního souboru v opačném pořadí (od posledního po první) a uvádíme jej zde hlavě, abyse se taky mohli těšit z kreativity, kterou do názvů příkazů rani umňoví programatři vložili :-). Tak trochu dalším příkazem `k` `tac` je `rev`, který pořadí řádků nemění, ale z každých z nich napíše pozpátku. Takže pokud byste chtěli obrátit celý soubor (od posledního znaku po první) `tac | rev` zadáte přesně to.

Třídí soubovy nemusíme jen z vlastního zájmu, ale také proto, že některé příkazy srň vstup seříděný vřadují. Jedním z nich je `uniq`, který ze vstupu vřadí duplicitní řádky, ale pouze, pokud leží všechny duplikáty vedle sebe (což zaráčíte právě seříděním). `uniq` má spousta zajímavých přepínačů, jako např. `-c`, který před každý výstupní řádek napíše, kolikrát se vyskytl na vstupu. Třeba pokud máme dlouhý seznam jmen uživatelů a chcři bychom vědět, která křestní jména se u nich vyskytují nejčastěji, můžeme použít:

```
$ sort jmena.txt | cut -d' ' -f 1 \
| uniq -c | sort -nr | head -n 3
64 Jan
51 Martin
42 Kateřina
```

S parametrem `-d` se naopak vypisují pouze duplicitní řádky (každý jen jednon), `-1` při porovnávání ignoruje velikost písmen a mnoho dalších zajímavých parametrů najdete ve manuálové stránce. Protože `sort | uniq` je tak často kombinace, existuje za ni zkrátka `sort -u`.



⁴ <http://git-scm.org/>

Porovnávání souborů: comm, cmp, diff

Dalším příkazem, kterému se hodí seříděný vstup, je `comm`. Slouží pro porovnání dvou množin reprezentovaných řádky seříděných souborů. Ve výchozím nastavení vypíše výstup do tří sloupečků: v prvním jsou řádky obsažené pouze v prvním souboru, v druhém řádky unikátní pro druhý soubor a ve třetím řádky oběma soubořům společné. To je hezké pro vizuální porovnání, ale ve skriptech víceem zajistit vypsat jen jednoho z těchto sloupečků, ale syntaxe není příliš intuitivní: parametrem `-n` říkáme, že `n`-tý sloupeček nedreeme zobrazit. Takže typické použití je např. `comm -12 soubor1 soubor2` pro zobrazování řádků vyskytujících se v obou soubořech.

Pokud chceme dva soubovy jen porovnat na shodnost (pro to už podstatně nemusi být seříděné), můžeme použít příkaz `cmp`. Ten srovná s nulovou návratovou hodnotou, pokud je obsažen soubor shodný, jinak s nenulovou (dá se tedy použít v rámci příkazu `if`). Při použití ve skriptech doporučujeme přidat parametr `-s`, aby se při neshodě nevypisovalo informativní hlášení.

Dalším nástrojem pro porovnávání souborů je `diff`, který umí zobrazit „v čem“ se dva soubovy liší. Nejčastěji se používá pro porovnání dvou verzí téhož souboru, například když přemýšlíme, proč stará verze našeho programu fungovala a nová už ne. Doporučujeme používat přepínač `-u` (zapne trochu snyšší nejlepší výstupní formát). `S -N -r` lze porovnávat rekurzivně celé adresáře. Zkusme si s ním porovnat.

Výstup programu `diff` se obvykle říká buď „diff“ (mezi ně jakož verzi a nějakou jinou), nebo „patch“. Druhé označení je odvozeno od příkazu `patch`, který s `diff`-em úzce souvisí. Vstupem příkazu `patch` je stará verze souboru a `diff` mezi starou a novou verzí. Výstupem pak zrekonstruovaná nová verze. Například po spuštění

```
diff -u A B >AB.diff
patch -o C A AB.diff
```

bude mít soubor `C` stejný obsah jako `B`. To se dá použítvat k snadnému šíření úprav: pokud třeba udláče malou změnu ve velkém programu a chcete se o ni s někým podělit, nemusíte mu posílat celý kód, stačí jen `patch`.

Jesté větší výhodou je to, že `patch` lze obvykle aplikovat, i když se mezitím původní soubor (`A` v našem příkladu) změnil. Tohle umožňuje několika lidem nezávisle na sobě změnit nějaký soubor a poté všechny změny automaticky sloučit. Stačí, když každý vytvoří patch mezi společenou původní verzí a svou upravenou verzí a nakonec se všechny tyto patche na soubor postupně aplikují. Problém nastane pouze v případě, že se dva lidé pokusí změnit stejnou část souboru: v takovém případě dojde k takzvanému *konfliktu*, který je třeba vyřešit ručně. Ale to se stává překvapivě málo často.

Na tomto principu je založen vývoj mnoha open source projektů. Přispěvatelé si lokálně program mění a testují, a hotové změny posílají autorům ve formě patchů. Ale i vytváření a aplikování patchů a udržování historie vývoje je spousta mrační práce, protože se tyto procesy snaží automatizovat takzvané verzovací systémy, jako např. `git`.⁴

¹⁵ <http://ksp.mff.cuni.cz/viz/kucharky/hledani-v-textu>

Příklad: Spam

Některí jste si již možná všimli, že v KSPřku většinu hromadných mailů posíláme každému z vlastním oslovením, včetně správně vyškolených tvarů slov ale pohlaví adresáta. Jak to děláme? Obvykle k mailu vytvoříme šablonu, prostý textový soubor, který může vypadat třeba takto:

```
Mil [y|a] $osloveni,
```

```
byl[a] jsi vybrán[a] jako $co na soustředění...
```

a pak seznam lidí, kterým se má poslat:

```
ke@example.net:Květoslava Čeřka:M:náhradník
poc@example.net:Pokusná Osoba:F:tčesťstřík
hroch@example.net:Hrochou:M:maskot
```

A zbytek zařídí jednohubý shellový skript. Pojdme si ho zkusit napsat. Šablona by nám vymyslel, jak vytvořit z seznamy mailů pro jednoho konkrétního adresáta, hromadné zpracování už pak snadno zjistíme nějakým while reed. Potřebovali bychom umět v textu nahradit všechny řetězce v nějakém tvaru (např. *Intečco Intečco*) za jiné řetězce, a ještě k tomu v náhradě použít nějaké části řetězce původního.

Regexy

Regex alias regulární výraz je řada písmenek a speciálních znaků, která právě dokáže popsat „řetězce nějakého tvaru“, jako v příkladech výše. O libovolném řetězci lze rozhodnout, jestli danému výrazu *vyhovuje* (na správný tvar), nebo nikoli. Regex tak vlastně popisuje nějakou (potenciálně nekonečnou) množinu řetězců. S něčím podobným jsme se již setkali: byly to wildcardy. Například nahrazování řetězců z příkladu výše bychom se mohli pokusit popsat wildcardm `\\[*|*|*]`. Ale možnosti wildcardů jsou poměrně omezené; s regexy se mají dělat mnohem zajímavější věci.

Níže najdete seznam konstrukcí používaných v regexech. Literál je libovolná posloupnost znaků, které nemají speciální význam (nejsou použity v prvním sloupceku tabulky).

Kulaté závorky lze vynechat tam, kde by se v nich nacházel jediný znak nebo jiný nedělitelný element (např. množina), v případě alternativy, když by měly být kolem celého regexu.

Par příklady:

- `[a-zA-Z]` `[a-zA-Z0-9]` * vyhovuje platný identifikátor, jak je definován většinou programovacích jazyků – tedy neprázdná posloupnost písmen, čísle a podtržítka *nezačínající* čísle.
- `[01]+` `[0-9]` `[12|0-3]` `[0-5]` `[0-9]` vyhovuje čas zapsaný v 24-hodinovém formátu (např. 0:42).

| Syntaxe | Význam | Příklad | Vyhovují | Nevyhovují |
|-----------------------------|-------------------------------------------------------------------------------|----------------------------------------|---------------------------------------------------------------------|---------------------------------------------------------|
| <i>literál</i> | Řetězec speciální s literálem. Escapuje speciální znak (udělá z něj literál). | <code>bagr</code> | <code>'bagr'</code> | <code>'bagrovat'</code> , <code>'1apata'</code> |
| <i>speciální-znak</i> | Libovolný znak. | <code>*</code> | <code>'a'</code> , <code>'%'</code> | <code>'abc'</code> , <code>'\'</code> |
| <i>fmnožina</i> | Libovolný znak patřící do množiny (jako v sblatlových wildcardch) | <code>[a-z-]</code> | <code>'q'</code> , <code>'-'</code> | <code>'A'</code> , <code>'-'</code> , <code>'aa'</code> |
| <i>(regex regex2 ...)</i> | Řetězec vyhovující alespoň jedné z možností. | <code>bro(ch s ik)</code> | <code>'hroch'</code> , <code>'hrosik'</code> , <code>'hrosi'</code> | |
| <i>(regex)*</i> | Nula nebo více opakování. | <code>*</code> | <code>'abc 123'</code> | <code>'baba'</code> |
| <i>(regex)+</i> | Jedno nebo více opakování. | <code>[ab]+</code> | <code>'ab'</code> , <code>'ababab'</code> | |
| <i>(regex)?</i> | Nepovinný prvek (0-1 opakování). | <code>[a-9]?</code> <code>[0-9]</code> | <code>'5'</code> , <code>'42'</code> | <code>'01'</code> , <code>'333'</code> |
| <i>(regex){m}</i> | <i>m</i> opakování. | <code>[A-Z]{2}</code> | <code>'AA'</code> , <code>'ZQ'</code> | <code>'X'</code> , <code>'ERF'</code> |
| <i>(regex){m,n}</i> | <i>m</i> až <i>n</i> opakování (obaji včetně). | <code>[0-9]{1,3}</code> | <code>'4'</code> , <code>'007'</code> | <code>'1337'</code> , <code>'xyz'</code> |
| <code>\$</code> | Začátek řádku. | | | |
| <code>^</code> | Konec řádku. | | | |

Podobnější úvod do regexů najdete v seriálu 23. ročníku⁵ či v deskách internetových tutoriálů.

Neplette si regexy s wildcardy! Sice řeší podobný problém (popis množiny řetězců), ale jinak spolu nemají nic společného. Například pozor na to, že * v regexech je kvantifikátor, který značí opakování toho, co stojí před ním, samostatně stojící * nemá smysl. Obdobou wildcardové hvězdičky je regex `*`. Lísti se také použitím wildcardů interpretuje shell a dají se použít pouze pro hledání názvů souborů, regexy se obvykle předávají nějakým pomocným utilitám a používají se pro hledání kusů textu.

Použití regexu: hledání a nahrazování

Kde lze regexy v Unixu použít? Například v nám době známeň přikazu `grep`. Pokud nám místo přepínače `-F` (Fixed, hledá pevný řetězec) dáme `-E`, hledá řetězce vyhovující nějakému regexu. E znamená Extended a zaplní tabulcovou rozšířenou syntaxi řetězů (tu jsme si vysvětlili v předchozí kapitole). Existuje ještě „zvládnutí“ syntaxe (starší), kterou se používá, pokud `grep`-n neudáte žádný přepínač. Ta je ale matoucí a nekonzistentní, tak ji raději nepoužívejte. Za `grep -E` existuje zkratka `egrep`.

Nezapomente, že `grep` hledá řádky *obsahující* řetězec vyhovující regexu, ne řádky celé *vyhovující* regexu. Například řádek `abcd` nevyhovuje regexu `[a-z]`, ale obsahuje a, které mu vyhovuje, a tedy `echo abcd | grep -E '[a-z]'` joi výše. Pokud bychom chtěli hledat řádky celé vyhovující regexu, stačí použít `rgrep$`.

Slibili jsme spolehlivější zjištění ID uživatele hroch, zde je: `grep -E 'hroch:' /etc/passwd | cut -d: -f3`

Tady hledáme slovo `hroch` pouze na začátku řádku a následované dvojtečkou, neznáme nás tedy uživatele `druny-hroch` či `hrochodlak`, ani někdo, kdo si nastavil jako shell `/bin/hrochsh`.

`grep` má ještě jeden přepínač, který je užitečný až s regexy: `-o`. Ten zajistí, že se nevypisují celé vyhovující řádky, nýbrž jen nalezené výskyt regexu. Pokud je na jednom řádku více výskytů, každý se vyíše na samostatný výstupní řádek. Například takto bychom mohli najít seznam všech identifikátorů (názvů funkcí, proměnných apod.) použitých v nějakém programu:

```
grep -Eio '[a-zA-Z][a-z0-9_]*' program.c | sort -n
```

Pro plnou funkčnost bychom ještě museli odfiltrovat komentáře, stringové literály a klíčová slova. K tomu by nám mohli pomoci nástroj, který si představíme za chvíli.

Zbývá si rozmyslet, jak $S_p(v)$ určit. Pro jednohubost ura-
zíme $s = P$, ostatní směry vyřesíme obdobně. Označme si p roditel vrcholu v . Pokud čtverec odpovídá v leží v levé polovině čtverce p , stačí vzít správného sousedce v . Např. je-li v levým dolním synem p , $S_p(v)$ je pravý dolní syn p . V opačném případě se podíváme na vrcholu $u := S_p(p)$ (p je ve vyšší polovině, tedy jeho S už máme spočítané). Pokud u je list, pak $S_p(v) = u$. Pokud není list, vime, že leží na stejné úrovni jako p (kdyby ležel výš, nebyl by *něžhlubším* sousedem p , neb některý z jeho synů by také sousedil s p a byl hlouběji). V takovém případě $S_p(v)$ musí být jeden ze synů u , ten, který správně přiléhá k v . Například je-li v prvním horním synem p , pak $S_p(v)$ je levý horní syn u . Zvlášť musíme osvětlit případ, kdy v je v daném směru na kraji obrázku.

Určení $S_p(v)$ nás stojí konstantní čas, tedy celý příchod stromu zvládneme v lineárním čase a v jeho průběhu sestojíme graf sousednosti jednohubých čtverců. A už vime, že sestrojili kvadrantového stromu, rozklad grafu sousednosti na komponenty a výběr největší zvládneme rovněž v lineárním čase. Tedy i celý algoritmus si vystačí s $O(n)$ času i paměti, kde n je délka vstupního kódu.

Program (Python):

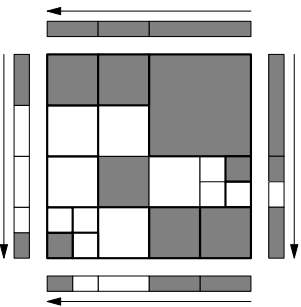
```
http://ksp.mff.cuni.cz/viz/27-3-1-bfs.py
```

Řešení příchodem do hloubky

Předchozí řešení vytvářelo sousednosti po hladinách shora dolů. Ukážeme jiný způsob, který postupuje naopak zdola nahoru a dosahuje stejné časové složitosti.

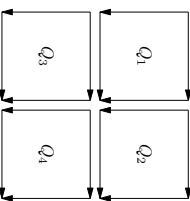
Strom hledáme procházet do hloubky a přiběžné hudovat graf sousednosti. Přesnější řečeno, kdyžkoliv se při procházení stromu hledáme vracet z nějakého čtverce, budou už zaznamenaný všechny sousednosti mezi jeho potomky. Sousednosti vedoucí přes hranici aktuálního čtverce doplníme později.

Proto se nám bude hodit předávat do vyšších pater stromu informace o tom, jaké čeré a bílé podčtverce leží na hranici aktuálního čtverce. To uložíme do čtyř seznamů. *Levý* seznam bude popisovat podčtverce přiléhající k levé hraně aktuálního čtverce, uspořádané shora dolů. Z každého podčtverce nás zajímá jen to, jak se dotýká hranice, což je nějaký *interval* g -ových souradnic. Za každý interval přidáme do seznamu dvojici (v, r) , kde v je příslušný vrcholu grafu sousednosti a r řád podčtverce. Podobně vytvoříme *pravý* seznam (též uspořádaný shora dolů), *horní* a *dolní* (oba zleva doprava).



Vracíme-li se z jednohubového čtverce (listu stromu), ne-
musíme vytvářet žádné sousednosti. Všechny čtyři seznamy
budou obsahovat odkazy na tento čtverec.

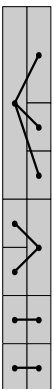
Zajímavější věci se dějí, pokud se vracíme z vícehubové
ho čtverce (vnitřního vrcholu stromu). Někdy Q_1 až Q_4
jsou kvadranty aktuálního čtverce. Z rekurze už známe sou-
sednosti vnitř kvadrantů a také podčtverce na hranicích
kvadrantů. Nyní potřebujeme rozpadnout sousednosti vedou-
cí přes hranice kvadrantů a sestrojiti hranici celého čtverce.



Nejprve vyřesíme hranici: levý seznam celého čtverce zis-
káme spojením levého seznamu kvadrantu Q_1 s levým se-
znamem kvadrantu Q_3 . Podobně získáme ostatní seznamy
slepením seznamů z vnějších hranic kvadrantů.

Nyní sousednosti: urážijme čtverce sousedící přes spoje-
nou hranici kvadrantů Q_1 a Q_2 (ostatní hranice zpracujeme
obdobně), tedy intervaly z pravého seznamu Q_1 a levého
seznamu Q_2 . Kdykoliv se interval z Q_1 překrývá s interva-
lem z Q_2 , chceme propojit hranou příslušné vrcholy v grafu
sousednosti.

Propojování probíhá takto: podíváme se na počáteční in-
terval každého seznamu – tento intervalům říkejme třeba
 x a y . Pokud jsou stejného řádu, vytvoříme hranu. Pokud
se řády liší (bez újmy na obecnosti x má vyšší řád), vezme-
me x a z protilehlého seznamu hledáme odebrat intervaly,
dokud se jejich délky neakšití na délku x . (Všimněte si,
že k tomu musíme dojít přesně na hranici intervalu, neboť
všechny intervaly vznikly postupným plněním délky hrany
celého obrázku.) Kdykoliv odebráme nějaký interval, vy-
tvóříme hranu. Poté pokračujeme zbytkem obou seznamů,
než se oba vyprázdní. Dopadne to například takto:



(Pokud vám na obrázku chybí hrana např. mezi třetím a
čtvrtým intervalem v prvním řádku, uvědomte si, že odpo-
ovídá sousednosti v koloně směru, takže vznikne při propo-
jování jiných seznamů než těchto dvou.)

Uvnime tedy zpracovat každ listy stromu, tak jejich vnitřní
vrcholy. Nakonec se probíhávání vrátí z kořene stromu a
graf sousednosti je hotov.

Zbývá určit časovou složitost. V každém vrcholu stromu
strávíme konstantní čas spojováním seznamů – to je cel-
kem $O(n)$ za celý strom – a nějaký další čas zpracováním
vnitřních hranic. Vnitřní hranice přitom může být poměrně
komplikovaná, ale stačí si uvědomit, že kdykoliv saháme
na nějaký interval, vypadne tento interval ze svého seznamu
a už se nikdy do žádného seznamu nevrátí.

Čas strávěný zpracováním všech hranic dohromady je tedy
shora omezen počtem všech intervalů, které vložíme do se-
znamů. Vkládáme ale pouze v listech stromu: 4 intervaly za
každý list. Tak vznikne celkem $O(n)$ intervalů, takže jejich
odebráním strávíme čas $O(n)$.

Program (C):

```
http://ksp.mff.cuni.cz/viz/27-3-1-dfs.c
```

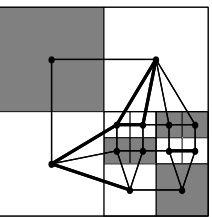
⁵ <http://ksp.mff.cuni.cz/viz/23-1-7>

Vzorová řešení třetí série dvacátého sedmého ročníku KSP

27-3-1 Plocha k přístání

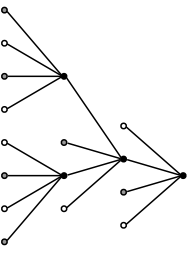
Většina z vás asi ví, jak se největší bílá oblast hledá, pokud jde o obývaný břítanový obrázek. Na obrázek se lze dívat jako na graf, kde pixely představují vrcholy a každá dvojice sousedních bodů je spojena hranou (tedy graf má tvar mřížky). V takovémto grafu bychom chtěli najít „komponenty bílé souvislosti“ – tedy maximální souvislé části grafu, které jsou celé bílé. K tomu můžeme použít klasický algoritmus na hledání komponent souvislosti prohlédáváním do šířky či hloubky (pokud jej neznáte, nahlédněte do naší grafové kuchařky),¹⁴ jen s tím rozdílem, že při prohlédávání zcela ignorujeme černé vrcholy – vzhled je neanavstevujeme. Snadno si rozmyslíte, že takto získáme očekávaný výsledek. Průběžné počítání velikost nalezených komponent a nakonec jen vybereme maximum. Na tento postup jde polhlízet také jako na opakované použití klasického algoritmu *lood fill*.

Zkusme se tím inspirovat. I kvadrantový obrázek má něč jako „pixely“, tedy elementární jednobarevné plochy – jsou to všechny nepodrozdělené čtverce. Jen je každý jinak velký a může mít víc než čtyři sousedy. Navíc není vůbec jasné, jak tyto sousedy najít. Ale pokud by nám někdo dal graf, jehož vrcholy jsou elementární čtverce obdohodocené svou plochou a hrany představují jejich sousednosti, dokázali bychom už jednoduše úlohu vyřešit.



Zbytek řešení bude pojednávat o tom, jak si takový graf pořádit, a to hned dvěma různými způsoby. Na úvod osmšen pověrně pár věd oberna řešením společně. Řešení více-méně každé úlohy nad kvadrantovým kódem začíná tím, že z kódu vytvoříme takzvaný *kvadrantový strom*. Kořen tohoto stromu reprezentuje celý obrázek. Pokud je jednobarevný, pamatuje si svou barvu a nemá žádné potomky; v opačném případě má právě čtyři syny představující kvadrantové stromy pro jednotlivé čtřtiny.

Strom pro kvadrantový kód $O((1010)1(0101)0)10$ (odpovídá obrázku výše) vypadá takto:



Každý vrchol kvadrantového stromu představuje čtverec $2^r \times 2^r$, kde r je takzvaný *řád* vrcholu (čtverce). Kořen má řád h (výška stromu), jeho synové $h-1$, atd.

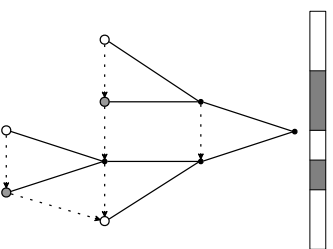
Strom vytvoříme přímočarou rekurzivní funkcí. Pokud se vyhneme zbytečnému kopírování řešícího (například si kód uložíme v globální proměnné a rekurzivním voláním bude předávat jen index znaků, od kterého začít), zvládneme to v lineárním čase. Podrobněji ve zdrojáku.

Řešení přidáním do šířky

Graf souvislosti vytvoříme tak, že pro každý jednobarevný čtverec (tomu odpovídá list kvadrantového stromu) najdeme seznam jednobarevných čtverců s ním v obrázku sousedících. To je ale těžké, neb jeden čtverec může mít sousedů hodně, a to i v doceła vzdálených částech stromu. Nám ovšem stačí, když každou souvislost „objeví“ jen jeden z zúčastněných čtverců, např. ten menší. Pak nám stačí pro každý čtverec hledat jen sousedy stejného nebo vyššího řádu (příslišně listy jsou v stromu stejně vysoko nebo výš). A takový je v každém směru nejvýše jeden. Tím získáme všechny hrany grafu souvislosti orientované směrem od menšího čtverce k většímu, opatřeny směrem doplníme jejich otočením, což můžeme dělat i průběžně. Z toho je také hned vidět, že hran je lineárně mnoho.

Samotné hledání souvislosti provedeme průchodem stromu po pátrech (neboli do šířky od kořene). Když navštívíme nějaký vrchol v , určíme jeho nejbližšího (obrázkového) souseda na každé straně, který ale není hlouběji než on sám. Toho si označíme $S_s(v)$ (kde $s \in \{L, P, H, D\}$ je příslišná strana). Všimneme si, že pokud je v tím hlubším z nějaké dvojice sousedních listů, je $S_s(v)$ právě druhým vrcholem této dvojice, a tedy jsme objevili jednu hranu grafu souvislosti (a po průchodu všech vrcholů objevíme všechny). V opačném případě je $S_s(v)$ nějaký vnější vrchol stromu, který se ale bude dále při hledání hodit.

To by si určitě zasloužilo obrázek. Alysyste se neztrátili v obrovském množství čar, ukážeme si to na příkladu jednorozměrného „kvadrantového“ obrázku. To je dlouhá „mudle“ rozměrně $2^k \times 1$, která je buď celá bílá, celá černá, nebo rozdělena na dvě poloviny rekurzivně splňující stejnou definici. Kvadrantový strom takového obrázku je binární a má tu výhodu, že směr vlevo a vpravo ve stromu odpovídá stejnému směru v obrázku.



Tečkované šipky značí ukazatele S_P (vedou od v k $S_P(v)$).

Zauvysleme se ještě nad jednou věcí. Pokud je v programu například identifikátor `bagr`, nachází se vnitř něj i další plně identifikatory, jako např. `ag`, `grep` – o takovém případě udělá to, co tenkrát vždy chová: z každé množiny překryvných výskytů vypíše pouze ten nejdelší a nejdelší (což bude přesně odpovídat identifikátorům v programu doopravdy použitým).

Při předávání regexů jako parametru příkazů je třeba dát si velký pozor na to, že mnoho regexových speciálních znaků má zvláštní význam i pro shell, a tedy pokud chceme, aby se např. ke `grep`-u dostaly nezmenšené, musíme je před shellem zescapovat. Nejjednodušším řešením je psát všechny regexy do apostrofu, kde se nic escapovat nemusí. Pokud přece jen z nějakého důvodu escapovat budeme, je třeba si uvědomit, že máme co do činení se dvěma úrovnemi escapování. Například pokud napíšeme `grep -*` soubor, nejprve dostane řešec do rukou shell, který ví, že `\` ve skriptování znamená `\` atp., všechny escapy odstraní a `grep`-u předá jako první parametr řešec `*.*` `grep` zase ví, že `\` znamená literál `*`, tedy tento příkaz tedy vybere ze souboru řádky začínající hvězdičkou.

Když už umíme výskyt regexů v textu hledat, hodilo by se také ušetřit je nahrazovat něčím jiným. Řešení si představíme zatím jen jako zaklamadlo `sed -re 's/regex/nahrada/g'`. Nebojte, za chvíli si ho vysvětlíme.

Jak jsme slibovali na začátku, v textu nahradit se lze odkazovat na části původního textu: přesejdi na obsah libovolné kulaté závorky. Obalí závorkou můžeme cokoliv, aniž by se tím změnil význam regexu. Počítají se všechny závorky, včetně těch vnutřních např. kvůli ohraničení skupiny alternativ. Obsah závorky do nahradě vložíme speciální sekvenčí *Věšlá*, kde číslo je pořadové číslo závorky (přesnější řečeno, páry závorek se číslují od jedničky v pořadí jejich otevřících závorek). To mimo jiné znamená, že i v nahradě musíme escapovat zpětná lomítka, pokud je tam chceme vložit doslovně.

Například pokud chceme v textu nahradit slovo `bagr` ve všech tvarech za slovo `kombaajn`, můžeme použít příkaz `sed -re 's/bagr(1|l|e|j|l|t|l|e|c|)/kombaajn1/g'`

Tehle trik samozřejmě funguje pouze, pokud mají slova stejny (pod)vzor a žádné nepravdivosti.

Řešení: Span

Nyní už máme téměř vše, co potřebujeme k řešení spanovacího příkladu. Můžeme si jej zkusit načrtnout. Předpokládáme, že v shelových proměnných `$osloveni`, `$pohlavi` a `$co` máme příslišně údaje k vyplnění.

```
 tvar= $\$(echo $pohlavi | tr MF 23)$ 
```

```
<seblona sed -re "s/\\$osloveni/$osloveni/g" \
| sed -re "s/\\$co/$co/g" \
| sed -re 's/[[[:*]]?([*])\?/\$tvar/g
```

Jak to funguje? První dva `sed-y` jsou přímočaré nahrazení jednoho řešce za jiný, jen pozor na escapování dolárů (od shellu i regexu). Poslední regex hledá řešce tvaru `[[[:*]]?([*])`, kde každé „`[:*]`“ načte do jedné závorky. V proměnné `$tvar` pak máme číslo závorky, kterou chceme vybrat. Parametr `sed-u` bude po expandování proměnných vypadat např. jako:

```
s/[[[:*]]\?([*])\?/\3/g
```

První část (`[:*]`) je nepovinná, pokud mužský tvar není uveden, předpokládá se `žádny`.

Toho řešení skoro funguje, ale ještě ne úplně. Zapomněli jsme totiž upozornit na jednu věc: regexy jsou *žvané*. To znamená, že při hledání vždy vyberou nejdelší kus textu, který jim vyhovuje. Tedy například pro text `MIll[ý]á dčera[hl]lce[el] nebndou nalezeny očekávané dva výskyt-y, nýbrž jeden velký, kde v první závorce skončí „[ý]á dčera[hl]lce“ a v druhé „lce“`. Snadno si rozmyslíte, že to vyhovuje našemu regexu. Nejjednoduššího to spravíme tak, že vnitřní jednotlivých tvarů zakážeme používat `| a |`. Tedy místo `*` musíme psát `[*|]` (strážka na začátku umožňuje známečně doplnit a pokud napíšeme `]` jako první, neukončí se tím množina, nýbrž už vložím znak `]`).

Nyní už to opravdň dělá, co má. Celý rozšiřací skript by mohl vypadat takto:

```
cat adresa | \
| while IFS=: read mail osloveni pohlavi co; do
t= $\$(echo $pohlavi | tr MF 23)$ 
<seblona sed -re "s/\\$osloveni/$osloveni/g" \
| sed -re "s/\\$co/$co/g" \
| sed -re 's/[[[:*]]\?([*])\?/[*|]\?/\$t/g
| mail -s "Pozvanka na soustředění" "$mail"
done
```

Kde `mail` je jedním z mnoha příkazů, které umožňují odeslat e-mail. K tomu samozřejmě potřebuje správně nastavení, které je nad rámec tohoto seriálu. Parametrem `-s` určujeme předmět.

sed pro pokročilé

Je na case naše sedové zaklamadlo rozklíčovat, `sed` ve skutečnosti dostane text a aplikuje na něj posloupnost příkazů. Tím mu můžeme předat buď jako parametr s použitím přepínače `-e` *příkaz*, nebo načíst ze souboru parametrem `-f` *soubor* (hodí se pro delší a složitější příkazy, v souboru nemusíme řešit shelové escapování). Jednotlivé příkazy se oddělují středníkem nebo koncem řádku. Též můžeme více příkazů zadat neklikanásobným použitím příkazů `-e`.

Přepínač `-r` zapíná rozšířenou syntaxi regexů, podobně jako u `grep`-u `-E`. Doporučujeme používat v podstatě vždy.

Všechny příkazy mají jednotnakerý název. My jsme dosud používali příkaz `s` (substituce) sloužící k nahrazení. Ten má tvar `s/regex/nahrada/modifikatory`. Mistro lomítka můžeme použít jakýkoli jiný znak (kromě písmen). Můžeme se tak vyhnout problémům s regexy obsahujícími lomítka. Oblíbenými volbami jsou např. `|` či `#` pro svou vizuální výraznost. Modifikátor `g` (globál) znamená nahrazení všech výskytů na řádku (jinak by se na každém řádku nahradil jen první). Dalším zajímavým modifikátorem `i` (ignoruj velkost písmen).

`sed` zpracovává vstup po řádcích a na každém z nich zvlášť provede všechny příkazy v pořadí, ve kterém jsou zapsány. Ve skutečnosti se každý řádek načte do řešcové proměnné, které se říká *pattern space*, na ní se provádějí jednotlivé příkazy a po jejich skončení je výsledná hodnota *pattern space* vypisována na výstup, není-li `sed` spuštěn s parametrem `-n` („nevypisovat“).

Před většinou příkazů lze napsat takzvanou *adresu* a tím určit, že se budou provádět jen na některém řádku či řádcích. Adresou může být například:

- `číslo` – Tého adresy vyhovuje právě tolikátý řádek vstupu, počítáno od jedničky.
- `$` – Poslední řádek.

¹⁴ <http://ksp.mff.cuni.cz/viz/kucharka/grafy>

- `/regex/` – Řádek obsahující výskyt regexu. Chceme-li použít místo lomítka jiný oddělovač, musíme na začátek napsat backslash, např. `\/dev/null#` adresuje řádky obsahující řetězec `\/dev/null`.
- `adresai` – Negace adresy; vyhovují řádky nevyhovující přírodní adrese.

Nyní už má smysl vysvětlit si některé další příkazy, které bez adresování nejsou příliš užitečné:

- `p` – Vypíše aktuální obsah pattern space na výstup. Nečastěji se používá spolu s parametrem `-n` pro selektivní vypisování řádků. Například `sed -nre /regex/ p` je to samé, jako `grep -E 'regex'`. Kombinaci `s/.*$/text/` a `p` můžeme na výstup vypsat libovolný text.
- `d` – „smazá“ řádek (zabrání jeho vypisování a skóci na další). Tedy `sed -re /regex/ d` funguje jako `grep -E`.
- `y/znakj/znakj/` – nahradí znaky z jednoho seznamu za odpovídající znaky z druhého, stejně jako `tr`, včetně stejného zápisu výčítí a rozsahů.

Příkazů můžeme dát místo jedné adresy také dvojici adres oddělených čárkami. Tím zadáváme rozsah – příkaz se provádí na všech řádcích od první adresy po druhou. Přesnější říčeno kdykoli se narazí na řádek vyhovující první adrese, příkaz se začne provádět a když se narazí na řádek vyhovující druhé adrese, zase se provádět přestane. Rozsah vždy zahrnuje oba krajní řádky (vyhovující přisluhujícím adresám). Takto se může příkaz provést i pro několik bloků řádek, pokud každý z nich začíná řádkem vyhovujícím první adrese a končí řádkem vyhovujícím druhé.

Některé formáty (například e-mailové zprávy) mají takovou strukturu, že obsahují nejprve hlavičky (s informacemi jako odesílatel, předmet atp.), potom prádný řádek, a teprve za tím tělo (text zprávy). Pokud bychom chtěli blok hlaviček odstranit, můžeme použít příkaz:

```
sed -re '1,/^$/ d'
```

Ukázky, jak pomocí `sed-u` nahradit `grep`, nebyly samozřejmě, `sed` totiž oproti `grep-u` má jednu velkou výhodou, totiž připravit `-i` (inplace). Jiz v prvním díle jsme se bavili o tom, že nemůžeme z jednoho souboru zároveň načíst a zároveň do něj zapisovat (`grep` slovo `<soubor >soubor` nemělo to, co byste chtěli). `sed -i` tohle umí zařídit. Jen mu dáte jako parametr (tedy náhle šlehněte přesměrování) název souboru, a on z něj náhle vstup a do toho samého souboru uloží svůj výstup. Interně to dělá tak, že vytvoří dočasný soubor, do kterého vstup zapisuje, a po svém skončení s ním nahradí (pomocí ekvivalentu příkazů `mv`) aktuální přírodní soubor. Tedy `sed -i -re ... soubor` je ekvivalentní posloupnosti příkazů:

```
sed -re '...' <soubor >soubor . tmp
mv soubor.tmp soubor
```

Toto je velmi častý mixový řádek, který je dobré si zapamatovat. Hodí se nejen když chceme zapisovat do stejného souboru, ze kterého čteme, ale víceméně kdykoli nahrazujeme či vytváříme nějaký soubor. Tím, že data zapisujeme do nějakého dočasného souboru, který kromě nás nikdo jiný nepoužívá, a teprve když je „hotový“, jej přeusme na cílové místo, se nemůže stát, že se nějaká jiná aplikace pokusí číst soubor v průběhu vytváření, když ještě neobsahuje

smysluplná data. Také pokud můžeme standardo zajistit (başlovným operátorem `&&`), že pokud náš příkaz modifikující nějaký soubor selže, `mv` se neprovede a zůstane zachovaná přírodní verze.

sed pro sílence

Se `sed-em` se dáji dělat i větší sílencosti. Jeho příkazy tvoří vlastně jednoduchý programovací jazyk.⁶ Kromě pattern space máte k dispozici ještě druhou(!) stringovou proměnnou, které se říká `hold` space. Pomocí příkazů `a` a `H` můžete aktuální obsah pattern space zapsat do hold space, resp. připojit na jeho konec. `g` a `G` dělají to samé opakujícím směrem. Při připojování je nový obsah od přírodního oddělen znakem konce řádku (`\n`). Příkazem `x` lze prohodit obsah pattern a hold space.

Pokud bychom si chtěli řídit načítání řádek nějak přesněji, než že pro každý řádek je náš program spuštěn znovu od začátku, můžeme. K tomu slouží příkazy `n`, který do pattern space dává řádek (přivodí zahodí) a `N`, který připojí další řádek na konec pattern space (oddělený `\n`, podobně jako `g`). To vše se děje stále v rámci jedné iterace našeho první ještě neaktivity řádek do pattern space a spustí náš program znovu od začátku. Tyto příkazy nám umožňují zpracovávat jednotlivé řádky jen nezavíme, ale dělat i nějaké složitější úpravy napří řádky.

Proměně už máme, ale správný programovací jazyk ještě potřebuje nějaké řídicí konstrukce. `sed` nabízí navští (`:` *název*), skoky (`b` *název*) a podmíněné skoky (`t` *název*). Podmíněný skok se provede, pokud od posledního podmíněného na tomto řádku vstupu došlo k alespoň jednomu úspěšnému nahrazení příkazem `s` (existuje i verze `T`, která má podmínku invertovanou).

Pokud chceme testovat, zda uspěl jeden konkrétní nahrazení příkaz, musíme si nejdřív připravit dřívější úspěchy na daném vstupním řádku „vyresetovat“ prostřednictvím příkazu `t`:

```
t reset; s/regex/nahradil/; t cil;
```

Nechceme-li nic nahrazovat a rádi bychom jen skákali podle toho, zda pattern space vyhovuje nějakému regexu, můžeme použít příkaz `b` podmíněný adresou: `/regex/ b`.

Všechny verze skoku při vymezení argumentu skácou za poslední příkaz. `q` a `Q` ukončí celý `sed` s vypisáním aktuálního pattern space, resp. bez něj. Umí nastavit návratovou hodnotu. Dále existují příkazy `x` a `w` pro čtení/zápis z/do pomocných souborů. Příkazy `[qQ]` jsou rozšířením GNU `sed-u` (nejběžněji se vyskytující verze) a nemají být dostupné v jiných verzích.

Zkusme si například napsat skript, který spojuje příkazy rozdělené na několik řádek pomocí zpětných lomítek do jednoho řádku:

```
sed -re ' : loop; s\\$//; T; N; s\\n//; t loop'
```

| | |
|------------------------|-------------------------|
| <i>Ukázkový vstup:</i> | <i>Ukázkový výstup:</i> |
| první | první |
| dr\\ | druhá |
| uh\\ | |
| y | |

Už byste měli zvládnout si rozmyslet, jak funguje.

⁶ Je Turingovsky úplný, ale asi v podobném smyslu, jako Brainfuck. Dá se najít implementace Turingova stroje v `sed-u`. Nebo Sokoban. Zagooglete si.

nicí, která nám v praxi může stačit, pokud podle ní třeba chceme věst náročný cyklistický závod.

O aproximaci je toho v literatuře napsáno mnoho zajímavého, pokud byste si o nich chtěli přečíst více v češtině, zkuste třeba kapitolu připravovaných skriptů z přednášky ADS na Matfyzu.¹²

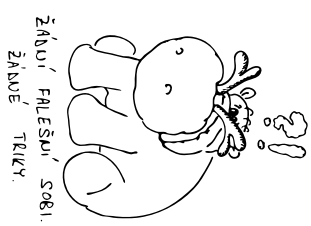
Více o třídě `NP` i o dalších aspektech složitosti můžete najít na stejné adrese, nebo zkuste vynikající anglicky psanou knihu *Algorithms* od profesora exotických jmen Dasgupta, Papadimitriou a Vazirani.

Jak už jsme zmínili, existují i problémy, které jsou mimo `P` i `NP`, a dokonce existuje spousta různých dalších tříd problémů. Je jich celá zoologická zahrada – pětky souborů můžete najít na stránkách Univerzity ve Waterloo.¹³

Seznam NP-úplných problémů

Seďte-li nad zátím nevyřešenou úlohou, kterou stále nemůžete rozluštnout, je možné, že budete `NP-úplná`. Abyste mohli mluvit `NP-úplnými` úlohami převádět, tak je dobré znát jich aspoň hrušku, podle toho, je-li problém `NP-úplný`, rovnováží, a tak dále.

V následujícím seznamu najdete několik úloh, které jsou zaručeně `NP-úplné`. Převodě se nám sem už sice nevedly, ale většinu z nich (už-i všechny) zvládnete vymyslet sami – zkuste to!



Název problému: Hamiltonovská kružnice

Vstup: Neorientovaný graf

Problém: Existuje v zadaném grafu kružnice procházející všemi vrcholů právě jednou?

Název problému: Hamiltonovská cesta.

Vstup: Neorientovaný graf, dva speciální vrcholy x a y .

Problém: Existuje cesta z x do y (posloupnost vrcholů, ve které se žádné dva neopakují), která prochází každým vrcholem právě jednou?

Název problému: Splnitelnost

Vstup: Logická formule. Tla tvoří proměnné a logické spojky negace \neg , konjunkce \wedge a disjunkce \vee . Například

$$(x \wedge \neg y) \vee z.$$

Problém: Můžeme proměnným přiřadit hodnoty 0 nebo 1 tak, že výsledná vyhodnocená formule má hodnotu 1?

Název problému: Součet podmnožiny

Vstup: Seznam nezáporných celých čísel, speciální číslo k .

Problém: Existuje podmnožina čísel, jejíž součet je přesně k ?

Název problému: Batoch

Vstup: Seznam dvojice nezáporných čísel, kde dvojice označuje hodnotu a váhu předmětu. Přirozené číslo b – možnost batohu, přirozené číslo k .

Problém: Umíme vložit do batohu předměty o hodnotě alespoň k , aniž bychom přišli přes limit váhy b ?

Název problému: Dva loupčiči

Vstup: Seznam nezáporných celých čísel.

Problém: Existuje rozdělení seznamu na dvě hromádky tak, že každé číslo bude v právě jedné hromádce a v každé hromádce bude stejný součet čísel?

Název problému: Klikla

Vstup: Neorientovaný graf, číslo k .

Problém: Existuje v grafu úplný podgraf o velikosti k , tedy k vrcholů takových, že mezi každými dvěma z nich vede hrana?

Název problému: Nezávislá množina

Vstup: Neorientovaný graf, číslo k .

Problém: Existuje v grafu prázdný podgraf o velikosti k , tedy k vrcholů, že žádné dva z nich nejsou spojeny hranou?

Název problému: Trojbarvený graf

Vstup: Neorientovaný graf

Problém: Lze vrcholy tohoto grafu obarvit třemi barvami tak, že každá hrana sousedí s vrcholy dvou různých barev?

Název problému: Rozpracování roviny

Vstup: Seznam bodů v rovině, kde každý má navíc přířazenou jednu z b barev, číslo k .

Problém: Umíme rozdělit rovinu pomocí k přímek tak, že v každé oblasti jsou jen body té samé barvy?

Název problému: 3D pátorování

Vstup: Seznam množin, žen a zvířátek, následovaný seznamem kompatibilních trojic tvaru {muž, žena, zvířátko}. Tyto trojice říkají, která trojice muž, žena a zvířátko by se dohromady snesla.

Problém: Můžeme všechny muže, ženy a zvířátka z prvního seznamu rozdělit do trojic tak, že každá trojice je kompatibilní a každá bytost je právě v jedné trojici?

Kuchařka sepsala

Martin Böhm & Jitka Smetička

¹² <http://mj.ucw.vyuaka/ads/49-prewody.pdf>
¹³ <https://complexityzoo.uwaterloo.ca/>

Ale je to správně, ba co víc, Cookova věta¹¹ říká, že existuje alespoň jeden takový problém. (Samotná deklace NP-třídního problému nezaručuje, že takový problém vůbec existuje.)

Ukazuje se však, že není sám, jsou jich stovky. Dokazovat existenci dalších NP-třídních problémů je však o dost jedli než dokázat Cookovu větu! Stačí totiž jen najít následující dva kroky:

- Dokázat, že problém je v NP – najít certifikát a polynomiální algoritmus, co jím využívá.
- Převést zadání libovolného NP-třídního problému na zadání našeho problému tak, že náš algoritmus vlastně vyřeší onen NP-třídní problém.

To postačí, protože pak libovolný jiný problém v NP nejprve převedeme na zvolený NP-třídní problém a pak postupně naní vymyšlený převod. Zveřejněn dvou polynomiálních algoritmů (převodů) je opět polynomiální algoritmus, takže podmínka převoditelnosti je splněna.

Ukažeme si důkaz NP-třídnosti jednoho problému na příkladu, pokud nám uvěříte, že již probíraný problém *Hamiltonovská kružnice* je NP-třídní. Neprve zadefiniujme jiný problém:

Název problému: Hamiltonovská cesta.

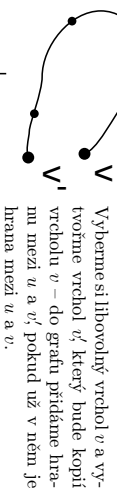
Vstup: Neorientovaný graf, dva speciální vrcholy x a y .

Problém: Existuje cesta z x do y (posloupnost vrcholů, ve které se žádné dva neopakují), která prochází každým vrcholem právě jednou?

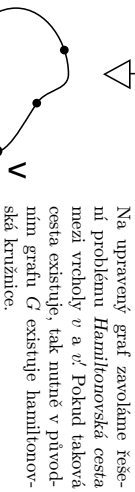
Certifikát: Posloupnost vrcholů tvořící správnou cestu.

Řešení v NP: Projdeme cestu z certifikátu a ověříme, že vrcholy jdou za sebou, je jich správný počet a žádné jsme neopakovali.

Důkaz NP-třídnosti: Převedeme předchozí problém (hamiltonovskou kružnici) na hledání hamiltonovské cesty. Uvažme graf G , ve kterém chceme najít hamiltonovskou kružnici.



Vyberme si libovolný vrchol v a vytvořme vrchol v' , který bude kopii vrcholu v – do grafu přidáme hranu mezi u a v' , pokud už v něm je hrna mezi u a v .

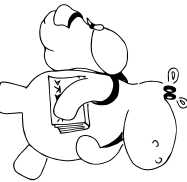


Na upravený graf zavoláme řešení problému *Hamiltonovská cesta* mezi vrcholy v a v' . Pokud taková cesta existuje, tak nutně v původním grafu G existuje hamiltonovská kružnice.

Pseudopolynomiální algoritmy

Znáte problém batolů? Jeho varianty jsou oblíbené na programovacích soutěžích. Zadat se může třeba takto: máme na vstupu seznam N dvojic kladných přirozených čísel, kde každá dvojice označuje váhu na nějakého předmětu. Nakonec dostaneme na vstupu ještě číslo B , které udává nosnost našeho batolu.

Otázka zní: jaký je nejmenší možný náklad, který přesto nepřesahuje váhový limit batolu?



Možná víte, že tůlha jde řešit dynamickým programováním – vytvoříme si pole *podbatoh*[i] a i do B , kde *podbatoh*[i] je maximální hodnota, kterou bych si odnesl v batolu o nosnosti i . Postupně od první větší do poslední pak projdeme celé pole *podbatoh*[i] – „pravá dolová“ od B do 1 a zjistíme, jestli je výhodnější do batolu vložit novou věc a volně místo doplnit starými (optimální volně místo pro předchozí věci máme napočítané), nebo si nechat jen ty staré. Tuto hodnotu pak zapíšeme jako aktuální pro váhu i na místo *podbatoh*[i].

Po N průchodech tohoto pole dostaneme řešení pro všechny větší dohromady na poli *podbatoh*[B]. Celková složitost je $O(NB)$, to je polynóm, algoritmus je tedy polynomiální.

Světě div se, toto řešení je ve skutečnosti exponenciální. Kde jsme v řešení udělali chybu? Nikde – naše složitost závisela na B , ovšem když se podíváme do vstupních dat, tak pokud jsou zapsána v binárním (nebo ternárním a vyšším) tvaru, zápis čísla B byl veliký $O(\log_2 B)$, ale naše složitost závisela na $B = 2^{O(\log_2 B)}$, tedy exponenciálně vzhledem k velikosti vstupu.

Problém batolů, respektive jeho rozhodovací verze, je dokonce NP-třídní problém.

Algoritmům, které řeší nějakou tůlhu a jsou polynomiální vůči *hodnotě čísel* na vstupu, ale exponenciální ve *velikosti zápisu* těchto čísel, říkáme *pseudopolynomiální algoritmy*. Některé další NP-třídní problémy mají pseudopolynomiální řešení (jako například *Dva soupeřící níže*), ale dá se dokázat, že na jiné problémy pseudopolynomiální algoritmus neexistuje (pokud $P \neq NP$).

Mimododem: pokud bychom na vstupu zapisovali čísla v nabitím zápisu (tedy místo každého čísla by bylo třeba tolik hvězdiček, jakou hodnotu představuje), každý pseudopolynomiální problém by ležel v P .

Poznámky na závěr

Otázku „je třída P rovna NP^P “, se již snažilo rozloupsnout mnoho matematiků a informaticků. Tato teorie přinesla spoustu zajímavých výsledků, například už se podařilo dokázat, že některými technikami tuto domněnku nelze nikdy dokázat, ani vyvrátit.

Kdyby platilo $P = NP$, pak by mnoho lidí zajásalo – mnoho přirozených problémů, které nastávají v reálném životě, by najednou bylo řešitelných rychle. Navíc by kračlo dosavadní sířování a bylo by možné najít rychle důkaz ke každému pravidelnému tvrzení výrokové logiky.

Tato rovnost by se dala hypoteticky ukázat velice snadno – stačilo by najít jeden polynomiální algoritmus pro libovolný NP-třídní problém! Většina informaticků studujících složitost se však domnívá, že se třídy nerovnájí.

To ale neznamená, že si to nemáme zkusit dokázat! Naopak, bojovat s NP-třídními problémy je užitečné i v reálném světě – mnohdy jde vymyslet třeba dobrá aproximace řešení. Například nemáme hamiltonovskou kružnici v polynomiálním čase, ale nalezneme nějakou relativně dlouhou kruž-

Další nástroje

Posledním významným nástrojem, který jsme nezmínili, je awk. To by nejspíš vyvedlo na samostatný díl. Slouží primárně k složitější práci s tabulkovými soubory. Na rozdíl od cut-u umí používat složitější než jednoznakové oddělovače, např. libovolnou posloupnost bílých znaků (ta je u awk dokonce výchozím oddělovačem). To se hodí při práci se soubory, jako je */etc/fstab*. Program v awk je podobný jako v sed-u posloupnost příkazů, která se spouští pro každé notifikaci řádků veji. Jazyk awk má daleko blíže k plnohodnotnému programovacímu jazyku než sed: má pojmenované proměnné, asociativní pole a další vycvičovací.

Uvedeme jen několik málo příkladů použití, měly by být zčásti samovyvěřující, zčásti interpretovatelné s pomocí manuálové stránky:

- Vyskání sloupečku odděleného obecnou posloupností bílých znaků z tabulky:

```
awk 'print $2' /etc/fstab
```
- Nalezení uživatele s daným ID:

```
awk -F: '{ $3 == 0 } { print $1 }' /etc/passwd
```
- Seřtání všech (číslených) řádků v souboru:

```
awk ' { sum += $1 } END { print sum }'
```

Postavení awk na pili cesty mezi jednohodnotou utilitkou a plnohodnotným programovacím jazykem z něj činí trocha zvláštní nástroj. Někdy je lepší použít místo něj cut či sed, jindy naopak opravdový programovací jazyk. Na zpracování textu nejlpeše poslouží Perl, který má bohatou podporu pro regxy a spoustu syntaktických zkraterek, jež v něm umožňují většinu jednohodnotých věcí napsat podobně krátce jako v jednořádkových nástrojích.

Jakou malou ochutnávku Perlu vám ukážeme skript, který z HTML dokumentu vypíše titulký všech hypertextových odkazů:

```
perl -0777 -ne 'for (m{<a.*?>.*?</a>}g) {
    s/<.*?//g; s/(-\s+|\s+$/g;
    print "$_<br>" if $_;
}'
```

Úkoly

V řešení se nebojte používat pomocné soubory, kde je to na místě. Klidně pro jednohodnotost s pevnými jmény. Můžete používat vše, co jsme se naučili, a další utilitky podobného režimu, můžete používat awk. Perl ani jiné „velké“ programovací jazyky nepoužívejte.

Úkol 1 [2b]: Ve vstupu souboru máte seznam jmen (srdý počet, jedno na řádek). Napište skript, který z nich vytvoří náhodně dvojice a vypíše je na výstup ve formátu *první osoba: druhá osoba*.

| | |
|------------------------|-------------------------|
| <i>Účázkový vstup:</i> | <i>Účázkový výstup:</i> |
| A | C:A |
| B | B:D |
| C | |
| D | |

Úkol 2 [4b]: Napište skript řešící tůlhu 27-Z3-2.⁷ Ve vstupu souboru dostanete slovník (jedno slovo na řádek), na výstup vypíšete nejdelší slovo, které má ve slovníku i svou vezu napsanou pozpátku.

| | |
|------------------------|-------------------------|
| <i>Účázkový vstup:</i> | <i>Účázkový výstup:</i> |
| recup | recup |
| ves | |
| vrabec | |
| pucek | |
| sev | |

Řešení s pomocí bashových cyklů je nudné, pro plný počet bodů to zkusete bez nich. Mohlo by se vám hodit awk a jeho funkce *length*.

Úkol 3 [4b]: V nějakém adresáři máte stáženou spoustu dlouhého oblibného seriálů (z legálních zdrojů, pochopitelně). A jak už to tak u legálních zdrojů drodí, soubory jsou pojmenované naprosto neuspořádaně. Jedně, čím si můžete být jisti, je, že název obsahuje číslo série a číslo epizody v tomto pořadí, mezi nimi je alespoň jeden nečíselný znak a pro jednohodnotost název žádné jiné číslo neobsahuje.

Napište skript, který stáhne z Wikipedie (či jiného příhodného zdroje) názvy dlouhých a všechny soubory přejmenuje tak, aby v jednotném formátu obsahoval číslo série a číslo a název epizody. Chcete-li, můžete předpokládat, že všechny soubory jsou ze stejné série.

Možná se vám snáz než HTML bude parsovat zdrojový wikitekst, který získáte příjopením `?action=raw` na konec URL článku.⁸

Úkol 4 [4b]: Vylepšete příklad vypisující všechny identifikatory v Cětkovém programu tak, aby ignoroval obsah řetězců a komentářů, připadně zkládání klíčová slova. Pro plný počet bodů by měl zvládnout jednořádkové (`// ...`) i víceřádkové (`* ... *`) komentáře a neměl by se nechat zmást escapováním uvozovkami a zpětnými lomítky v řetězcích. Ale určitě má smysl poslat i jednořádkové či částčné řešení. Oskumnosti jako komentář uvnitř řetězce (nebo naopak) ošetřovat nemusíte, pokud vytvozené nechte. Předpokládejte samozřejmě, že program je syntakticky správný.

| | |
|-------------------------------------------------------------------------------------|--------------------------|
| <i>Účázkový vstup:</i> | <i>Účázkový výstup:</i> |
| <code>/* print a greeting with quoted name */ printf("Hello, \"%s!\", name);</code> | <code>name printf</code> |

Pokud si s Cětkem nerozumíte, můžete si vybrat nějaký jiný srovnatelně složitý (tedy třeba by měl ideálně mít vícetáčekové komentáře či řetězce) programovací jazyk – třeba Python nebo Pascal.

Z cvičných důvodů zkusíte nenačítat celý vstupní soubor do paměti najednou. Ano, bylo by to jednodušší a pro praktické účely možná nejlepší řešení, ale tolik se tolo na něm nenačítá.

Filip Stědronský

¹¹ http://en.wikipedia.org/wiki/Cook%E2%80%9333Levin_theorem

⁷ <http://ksp.mf.cuni.cz/viz/27-Z3-2>

⁸ http://en.wikipedia.org/wiki/The_Big_Bang_Theory_season_4?action=raw

