

2 Hľadanie reťazcov



Vytvoril/a Peter Koscelanský [Administrator]
Posledná aktualizácia: Oct 19, 2021

Vytvorte konzolovú aplikáciu, ktorá dostane tri pozičné argumenty

- cestu k textovému súboru (ak súbor neexistuje alebo sa z neho nedá čítať skončíte s nenulovým error kódom),
- neprázdny reťazec znakov `x` (akceptuje maximálnu dĺžku ostro menšiu 256, pre väčšie reťazce treba skončiť s nenulovým error kódom),
- číslo `N` označujúce vzdialenosť (tu akceptujeme všetky `uint32_t` okrem nuly, inak nenulový error kód).

Program následne na štandardný výstup vypíše v riadkoch pozície, kde sa nachádza výskyt aký chceme. Prvý znak v súbore má pozíciu 0, rovnako aj prvý riadok je 0). Whitespace sa ráta do dĺžky súboru a konce riadkov sú vždy jeden znak (teda s tým nemusíte nič robiť, to za vás vyrieši C++). Súbory budú obsahovať iba znaky s hodnotou $0 < c \leq 127$, ak by bol v súbore aj iný znak, tak správanie je nešpecifikované.

Hľadáme výskyt reťazca `x`, ktoré majú vo svojom okolí (definovnom pomocou `+-N`) ďalší výskyt reťazca `x`.

Ak máme súbor s obsahom:

```
1 aaabaaabbbbbbaaa
```

`X = aaa` a `N = 4`, tak výsledkom budú pozície 0, 4. Lebo prvý výskyt má vo vzdialenosti 4 iný výskyt reťazca, konkrétne na pozícii 4 a $|4 - 0| \leq 4$ druhý výskyt tam tiež zarátame, lebo pred ním je vo vzdialenosti 4 ten náš prvý výskyt. Posledný výskyt nepoužijeme, lebo nemá nikde v okolí ešte jeden výskyt.

Keďže v príklade vyššie je všetko na jednom riadku, výsledok bude teda

```
1 0 0
2 0 4
```

Samotná podmienka blízkosti neberie do úvahy riadky, koniec riadku je vlastne jeden znak. Riadky a pozície v nich sa berú do úvahy iba vo výpise.

Ďalšie príklady

`X = aaa`, `N = 2`

```
1 aaaaabbbbb
2 bbaaaaaa
```

```
1 0 0
2 0 1
3 0 2
4 1 2
5 1 3
6 1 4
```

V príklade hore sa výskyt prelínajú, ale vôbec nevadí, je tam splnená podmienka blízkosti.

`X = aaabaaa`, `N = 2`

```
1 aaabaaabaaa
```

Tu bude výstupom nič, síce sa oba stringy prelínajú, ale začiatok prvého je 0 a začiatok druhého je 4 a teda ich rozdiel je väčší ako 2.

Pre vstupný súbor input.txt a `X = bbb, N = 20` je výsledok

```
1 4 41
2 4 45
3 6 29
4 6 30
```

Súbor môže byť obrovský, pokojne aj niekoľko desiatok GB, takže neukladajte celý jeho obsah do pamäte (pozor na `std::getline`). Ako vždy sa hodnotí to aby boli všetky chybové stavy ošetrené (nenulový error kód). Tentokrát sa bonusový body získava za rýchlosť. Čiže aby to aj nejak bežalo. Prvých desať implementácií dostane bonusový bod.

 Páči sa mi to Buďte prvý(-á), komu sa to páči

Žiadne označenia

 27 page comments



Peter Koscelanský [Administrator] Oct 19, 2021

Na cvičení bola otázka, či nový riadok sa rata ako znak a teda vlastne "resetuje" ten string čo hľadáme. Ano nový riadok je normálny znak a `ifstream::read` vám ho vráti.

Odpovedať • Páči sa mi to



Anonymný Oct 19, 2021

Dobry den, mam otazku k argumentom, na potvrdenie, ze som pochopil.

1) Takto zadane argumenty **su spravne**:

```
1 ./main input.txt aaa 6
```

subor = input.txt

hľadany reťazec = aaa

vzdialenosť = 6

2) Takto zadane argumenty **su tiež spravne**

```
1 ./main aaa input.txt 6
```

subor = aaa

hľadany reťazec = input.txt

vzdialenosť = 6

3) menej ako 3 alebo viac argumentov a ak tretí argument nie je `uint32_t` je **nenulový error**, napr:

```
1 ./main input.txt aaa 6 bbb
2 ./main input.txt aaa
3 ./main input.txt aaa bbb || ./main input.txt aaa -6 || ./main input.txt aaa 0
```

Dakujem.

Odpovedať • Páči sa mi to



Peter Koscelanský [Administrator] Oct 19, 2021

1. ano

2. nie, prvý musí ísť subor
3. ano všetky sú, že treba nenulový error kód

Odpovedať • Páči sa mi to



Anonymný Oct 19, 2021

```
subor = aaa  
hľadany reťazec = input.txt  
vzdialenosť = 6
```

no v tej 2) som to tak napísal, že bude otvárať subor "aaa" a ako reťazec zoberie "input.txt"

Odpovedať • Páči sa mi to



Peter Koscelanský [Administrator] Oct 20, 2021

Aha, ano máš pravdu, neviem som si to. Je to tak ako píšeš. Prvý je subor, potom reťazec a na záver vzdialenosť.

Odpovedať • Páči sa mi to



Anonymný Oct 19, 2021

Ak som načítal zo suboru string "abc\\naabc" a hľadám reťazec "abc" s $N = 4$, tak v tomto prípade `\n` sa počíta ako znak pri výpočte vzdialenosti a preto na výpis nič nedostanem, keďže `|5-0| > 4`?

Odpovedať • Páči sa mi to



Peter Koscelanský [Administrator] Oct 20, 2021

Ano znak nový riadok (`\n`) sa počíta do znakov.

Odpovedať • Páči sa mi to



Anonymný Oct 22, 2021

Aká je maximálna dĺžka riadka? Teda ak by som našiel zhodu v riadku na pozícii za veľkosťou `uint32_t`, tak by som nevedel vypísať asi číslo, že na akej som pozícii 😊.

Odpovedať • Páči sa mi to



Peter Koscelanský [Administrator] Oct 26, 2021

tak ja odporúčam použiť aspoň `int64_t`, to už isto nepresiahneme.

Odpovedať • Páči sa mi to



Anonymný Oct 26, 2021

ok, vďaka 😊

Odpovedať • Páči sa mi to



Anonymný Oct 30, 2021

Ja by som sa k tomuto tiež chcela pýtať - ja to načítavam do `size_t`, prišlo mi to ako najvhodnejšie riešenie. Avšak keď sa to spustí na 32 bitovom systéme, tak to bude málo :(. Môžem ratovať s tým, že také niečo sa nestane, alebo tam teda mám podávať fixne tých 64? :D

Odpovedať • Páči sa mi to



Peter Koscelanský [Administrator] Oct 30, 2021

Nie je to ideálne ale asi to na túto úlohu bude stačiť.

Odpovedať • Páči sa mi to

Anonymný Oct 22, 2021



Može hladany substring obsahovat \n?

Odpovedať • Páči sa mi to



Anonymný Oct 22, 2021

Súbory budú obsahovať iba znaky s hodnotou $0 < c \leq 127$, ak by bol v súbore aj iný znak, tak správanie je nešpecifikované.

\n je znak ako kazdy iny, c++ v tom nerobi rozdiel, preto nevidim problem v tom, aby sa mohol vyhľadavat. Ak si to chces vyskusať, tak daj argument do uvozoviek.

Odpovedať • Páči sa mi to



Anonymný Oct 22, 2021

No ja sa snazim o implementáciu pomocou KMP kvôli rýchlosti a efektívnosti. Tazko sa to vysvetľuje v komentári, ale bolo by dosť komplikované prepocítavať pozíciu začiatku substringu (teda stĺpec aj riadok) - musel by som si ukladať počet znakov v každom riadku, následne odpocítavať znaky z daného stĺpca aby som sa dostal k jeho pozícii a teda by to nebolo veľmi pamätovo efektívne. Preto by som rad vedel s istotou či toto mám riešiť alebo nie.

Odpovedať • Páči sa mi to



Peter Koscelanský [Administrator] Oct 26, 2021

Ano moze byt aj \n.

Odpovedať • Páči sa mi to



Anonymný Oct 28, 2021

Ako máme rozlišovať medzi tým či substring obsahuje \n ako znaky \ a n a \n ako newline?

Odpovedať • Páči sa mi to



Anonymný Oct 28, 2021

Myslím teraz pri načítavaní z argumentov, keďže tam nevieme načítať \n ako jeden symbol

Odpovedať • Páči sa mi to



Peter Koscelanský [Administrator] Oct 29, 2021

Ako? Nerozumiem, vy máte použiť to čo vám príde v `argv` a tam už nič neinterpretujte. Ono akým spôsobom vám tam ten parameter príde je jedno, vy iba zoberte ten parameter a jeho súčasťou môže byť aj \n teda new line, ale nie ako dva znaky \ a potom n, ale ako jeden znak, ktorý reprezentuje LINE FEED, teda '\n'.

Odpovedať • Páči sa mi to



Anonymný Oct 22, 2021

1.

Mozeme používať všetky knižnice na prácu so stringami? Lebo ak by som nechcel ísť na rýchlosť, tak s pomocou `std::find` je riešenie veľmi jednoduché. A myslím si, že nebol zamer tohto zadania, aby bolo ľahké. Alebo tu skor ide o to, aby sme sa zamerali hlavne na ošetrenie všetkých možných chybových a nedefinovaných stavov?

2. Čo sa týka počtu načítaných znakov na jeden raz, tak koľko by to cca malo byť? Je nejaké odporúčane maximum?

3. Keď pracujem na Linuxe, nemusím sa báť, že mi na testoch pustíte CRLF súbor?

Odpovedať • Páči sa mi to



Peter Koscelanský [Administrator] Oct 26, 2021

1. Zo štandardnej knižnice môžete používať všetko až na niektoré viacklakové konštrukcie (tie nepojdu vybudovať, takže to si rýchlo všimnete).

2. To je ťažko povedať, o tom si na internete diskutuje, ideálne asi v radoch tisícok a nie viac ako megabajt.

3. Nie, resp. ak ano tak \r je pre teba normálny znak.

Odpovedať • Páči sa mi to



Anonymný Oct 28, 2021

Zdravím, je potrebné nejako ošetrovať obsah parametru X?

Odpovedať • Páči sa mi to



Peter Koscelanský [Administrator] Oct 28, 2021

Iba veľkosť, budú tam len znaky ktoré sú ASCII <= 127, ale ak tam bude niečo iné, tak je správanie nedefinované, teda urobíte hocico, okrem nedefinovaného správania.

Odpovedať • Páči sa mi to



Anonymný Oct 29, 2021

Vie niekto ako poslať \n ako vstupný argument vo Visual Studiu?

Odpovedať • Páči sa mi to



Anonymný Oct 29, 2021 [↗](#)

Ja som to testoval tak, že som ho natvrdo vložil do kódu

Odpovedať • Páči sa mi to



Peter Koscelanský [Administrator] Oct 29, 2021

+1 to je super prístup. My používame `fork` a `exec` na spustenie procesu a tam sa dá dať všetko.

Odpovedať • Páči sa mi to



Anonymný Oct 29, 2021

500 IQ. Dakujem, fakt mi to nenapadlo :DD.

Odpovedať • Páči sa mi to