

Zadanie č. 2

Vo zvolenom (podľa seba) programovacom jazyku naprogramujte program, ktorý simuluje spracovanie regulárneho výrazu zadaného pomocou elementárnych regulárnych výrazov a operácií medzi nimi (de facto Thompsonova konštrukcia). Váš program bude pracovať nasledovným spôsobom:

1. Prvým argumentom programu pri spustení bude textový súbor, v ktorom sú uložené elementárne regulárne výrazy a operácie medzi nimi. Tieto spolu tvoria nejaký výsledný regulárny výraz. Formát súboru je popísaný na ďalšej strane.
2. Druhým argumentom programu bude textový súbor, v ktorom sú zapísané reťazce, každý na novom riadku. Počet týchto reťazcov je uvedený v prvom riadku súboru. Váš program postupne pre každý reťazec zistí, či patrí/nepatrí do jazyka daného regulárnym výrazom z prvého súboru a vypíše na obrazovku príslušný reťazec a to, či patrí/nepatrí do daného jazyka.

Vstupy a výstupy

Elementárne regulárne výrazy a príslušné operácie sa načítajú z textového súboru predstavujúceho **prvý argument** programu pri spustení.

Druhým argumentom programu pri spustení je textový súbor obsahujúci reťazce, ktorých príslušnosť do jazyka testujeme, pričom na prvom riadku súboru je uvedený počet reťazcov a na ďalších riadkoch sú jednotlivé reťazce, každý na samostatnom riadku. Prázdný reťazec ε je reprezentovaný prázdny riadkom. V neprázdnych reťazcoch sú použité len znaky $a, \dots, z, 0, \dots, 9$.

Ak by ste teda napríklad naprogramovali Váš program v jazyku Python a Váš kód by bol v súbore `zadanie2.py`, program by sa spúšťal:

```
python zadanie2.py regex.txt retazce.txt
```

kde `regex.txt` je príklad mena súboru, v ktorom sa nachádzajú elementárne regulárne výrazy s operáciami a `retazce.txt` je príklad mena súboru, v ktorom sa nachádzajú reťazce, ktoré chceme otestovať.

Zadanie môžete vypracovať v ľubovoľnom programovacom jazyku, Python je tu uvedený len ako príklad.

Implementácia algoritmov

Požaduje sa vaša samostatná práca, t.j. výsledný program musí byť výsledkom vašej samostatnej práce. Kopírovanie zdrojových kódov z internetu, prípadne od iných študentov bude hodnotené ako plagiátorstvo a v zmysle platného študijného poriadku hodnotené známku FX. Taktiež použité algoritmy musia byť výsledkom vašej vlastnej implementácie, t.j. ak napríklad používate konverziu NKA na DKA, musí ísť o Vašu vlastnú prácu.

Deadline zadania

Zdrojové kódy Vášho riešenia odovzdajte do AIS-u do príslušného miesta odovzdania. Deadline je **24. marec 2024, 23:59:59**.

FORMÁT PRVÉHO VSTUPNÉHO SÚBORU

Riadky sú implicitne číslované od 1, t.j. prvý riadok v súbore má číslo 1, druhý 2, atď. Na každom riadku môže byť:

- Prázdny riadok - reprezentuje regulárny výraz popisujúci prázdny reťazec ε
- Jeden symbol - reprezentuje regulárny výraz popisujúci príslušný symbol. Môžete predpokladať, že týmto symbolom je vždy alebo malé písmeno anglickej abecedy a-z, alebo číslica 0-9.
- Trojica U, i, j , kde U je písmeno veľké U (zo slova Union), i, j sú poradové čísla niektorých predchádzajúcich riadkov. Táto trojica predstavuje regulárny výraz, ktorý popisuje zjednotenie jazykov regulárnych výrazov z riadkov i a j
- Trojica C, i, j , kde C je písmeno veľké C (zo slova Concatenation), i, j sú poradové čísla niektorých predchádzajúcich riadkov. Táto trojica predstavuje regulárny výraz, ktorý popisuje zreťazenie jazykov regulárnych výrazov z riadkov i a j v poradí ij .
- Dvojica I, i , kde I je písmeno veľké I (zo slova Iteration), i je poradové číslo niektorého predchádzajúceho riadku. Táto dvojica predstavuje regulárny výraz, ktorý popisuje iteráciu jazyka regulárneho výrazu z riadku i
- Vstup bude vždy korektný, t.j. operácie I, C, U sa neodkazujú na neexistujúce riadky alebo na riadky s väčším poradovým číslom.

Príklad č. 1: Prvým argumentom by bol súbor `regex1.txt` s nasledovným obsahom:

```
a
I,1
C,1,2
```

Druhým argumentom by bol súbor `retazce1.txt` s nasledovným obsahom:

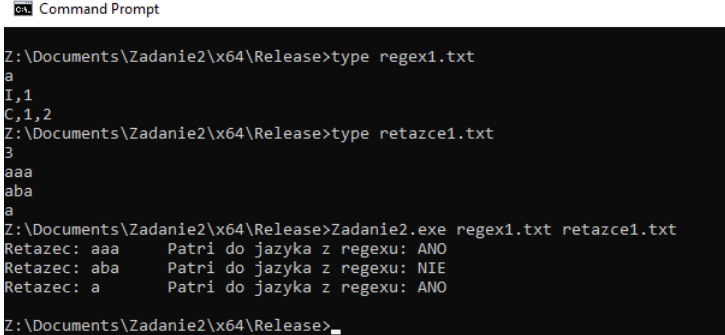
```
3
aaa
aba
a
```

K jednotlivým riadkom v prvom argumente teda prináležia nasledovné regulárne výrazy:

1. a
2. a^*
3. aa^*

Výsledný regulárny výraz (ten z posledného riadka) prvého súboru predstavuje regex aa^* , voči ktorému budeme porovnávať reťazce z druhého vstupného súboru. Následne sa z druhého vstupného súboru načítajú jednotlivé reťazce a program vypíše, či patria alebo nepatria do jazyka daného regulárnym výrazom z prvého súboru.

Príklad činnosti programu (s formátovaním výstupu sa nemusíte trápiť, dôležité je, aby bol vždy uvedený reťazec a to, či patrí/nepatrí do jazyka)):



```
Command Prompt
Z:\Documents\Zadanie2\x64\Release>type regex1.txt
a
I,1
C,1,2
Z:\Documents\Zadanie2\x64\Release>type retazce1.txt
3
aaa
aba
a
Z:\Documents\Zadanie2\x64\Release>Zadanie2.exe regex1.txt retazce1.txt
Retazec: aaa    Patri do jazyka z regexu: ANO
Retazec: aba    Patri do jazyka z regexu: NIE
Retazec: a      Patri do jazyka z regexu: ANO
Z:\Documents\Zadanie2\x64\Release>
```

Príklad č. 2: Prvým argumentom by bol súbor `regex2.txt` s nasledovným obsahom (prvý riadok je prázdny riadok, t.j. je tam regulárny výraz ε):

```
a
b
C,2,3
U,1,4
```

Druhým argumentom by bol súbor `retazce2.txt` s nasledovným obsahom (druhý reťazec je prázdny reťazec ε reprezentovaný prázdny riadkom):

```
3
aaa
```

```
ab
```

K jednotlivým riadkom v prvom argumente teda prináležia nasledovné regulárne výrazy:

1. ε
2. a
3. b
4. ab
5. $\varepsilon \mid (ab)$

Výstup programu (ak sa testuje prázdny reťazec, uveďte to napríklad tak, ako na obrázku):

```
Z:\Documents\Zadanie2\x64\Release>type regex2.txt
a
b
C,2,3
U,1,4
Z:\Documents\Zadanie2\x64\Release>type retazce2.txt
3
aaa

ab
Z:\Documents\Zadanie2\x64\Release>Zadanie2.exe regex2.txt retazce2.txt
Retazec: aaa          Patri do jazyka z regexu: NIE
Retazec: (prazdny)    Patri do jazyka z regexu: ANO
Retazec: ab           Patri do jazyka z regexu: ANO
Z:\Documents\Zadanie2\x64\Release>_
```

Príklad č. 3: Prvým argumentom by bol súbor `regex3.txt` s nasledovným obsahom:

```
a
b
U,1,2
I,3
```

Druhým argumentom by bol súbor `retazce3.txt` s nasledovným obsahom (prvý reťazec je prázdny reťazec ε reprezentovaný prázdny riadkom):

```
5

a
b
aa
ababab
```

K jednotlivým riadkom v prvom argumente teda prináležia nasledovné regulárne výrazy:

1. a
2. b
3. $a \mid b$
4. $(a \mid b)^*$

Výsledný regulárny výraz: $(a \mid b)^*$. Výstup by vyzeral napríklad nasledovne:

```
Z:\Documents\Zadanie2\x64\Release>type regex3.txt
a
b
U,1,2
I,3
Z:\Documents\Zadanie2\x64\Release>type retazce3.txt
5

a
b
aa
ababab
Z:\Documents\Zadanie2\x64\Release>Zadanie2.exe regex3.txt retazce3.txt
Retazec: (prazdny)      Patri do jazyka z regexu: ANO
Retazec: a              Patri do jazyka z regexu: ANO
Retazec: b              Patri do jazyka z regexu: ANO
Retazec: aa             Patri do jazyka z regexu: ANO
Retazec: ababab         Patri do jazyka z regexu: ANO
Z:\Documents\Zadanie2\x64\Release>_
```

Príklad č. 4: Prvým argumentom by bol súbor `regex4.txt` s nasledovným obsahom (tretí riadok je prázdny riadok, t.j. je tam regex ε):

```
a
b

C,2,1
U,4,3
I,5
c
C,6,7
```

Druhým argumentom by bol súbor `retazce4.txt` s nasledovným obsahom:

```
3
c
babac
bacba
```

K jednotlivým riadkom v prvom argumente teda prináležia nasledovné regulárne výrazy:

1. a
2. b
3. ε
4. ba
5. $\varepsilon \mid ba$
6. $(\varepsilon \mid ba)^*$
7. c
8. $(\varepsilon \mid ba)^*c$

Výsledný regulárny výraz: $(\varepsilon \mid ba)^*c$. Program by pre reťazce `c`, `babac` vypísal, že patria do jazyka daného regexom, pre reťazec `bacba` že nepatrí do jazyka daného regexom.

Príklad č. 5: Prvým argumentom by bol súbor `regex5.txt` s nasledovným obsahom (šiesty riadok je prázdny, t.j. je tam regex ε):

```
a
c
b
C,1,2
C,4,3

I,5
U,7,6
a
I,9
b
C,10,11
C,8,12
```

Druhým argumentom by bol súbor `retazce5.txt` s nasledovným obsahom:

```
3
acbacbb
acbaaab
acbaabb
```

K jednotlivým riadkom v prvom argumente teda prináležia nasledovné regulárne výrazy:

1. a
2. c
3. b
4. ac
5. acb
6. ε
7. $(acb)^*$
8. $(acb)^* \mid \varepsilon$
9. a
10. a^*

11. b

12. a^*b

13. $((acb)^* \mid \varepsilon)a^*b$

Výsledný regulárny výraz $((acb)^* \mid \varepsilon)a^*b$. Pre reťazce **acb**acbb, **acbaa**ab by program vypísal, že patria do jazyka daného regexom, pre reťazec **acba**abb by program vypísal, že nepatrí do jazyka daného regexom.

Changelog:
7.3.2024 Zverejnená prvá verzia zadania.