Slovenská technická univerzita v Bratislave Fakulta informatiky a informačných technológií

Syntaktický analyzátor Simple URL

Autori: Michal Zajic, Tímea Nagy

Predmet: Softvérové jazyky

Cvičenie: Streda 8:00

Cvičiaci: doc. Ing Michal Kompan, PhD.

Obsah

1.	Ana	lýza a návrh	1
	1.1.	Príklady viet z jazyka	1
	1.1.1.	Korektné príklady	1
	1.1.2.	Nekorektné príklady	1
	1.2.	Prepis gramatiky z BNF do gramatických pravidiel s alternatívami	2
	1.2.1.	Gramatika v BNF	2
	1.2.2.	Pravidlá s alternatívami	2
	1.3.	Transformácia gramatiky do LL(1)	3
	1.4.	Overenie LL(1)	7
	1.4.1.	First ₁	7
	1.4.2.	Follow ₁	7
	1.5.	Tabul'ka prechodov	9
2.	Imp	lementácia	11
	2.1.	Používateľské rozhranie	11
	2.2.	Testy	13
	2.3.	Tabuľka a pravidlá	14
	2.4.	Postup implementácie.	14

1. Analýza a návrh

Podiel práce:

- Timea Nagy 50%
- Michal Zajic 50%

1.1. Príklady viet z jazyka

Pre analyzátor sme vytvorili testy a v nich sme snažili pokryť všetky typy URL adries, ktoré patria do jazyka a aj niekoľko chybných. V tejto časti sme vybrali niektoré z nich.

1.1.1. Korektné príklady

Príklad 1: http://hostname:port/path?search

Napr. http://Youtube.com.foo:48484/Watch/Video?123456+foo+bar

Príklad 2: ftp://user:password@hostname:port/path

Napr. ftp://foo:fooBAR123456@foo:123456/foo/bar/foo/

Príklad 3: telnet://user:password@hostname:port

Napr. telnet://foo:123456@Foo:123456

Príklad 4: mailto::xalphas@xalphas.xalphas

Napr. mailto::foo@foo.BAR.foo.123

1.1.2. Nekorektné príklady

Príklad 1: nesprávne poradie search a path

Napr. http://foo.com?foo/bar

Príklad 2: nesprávna ftp adresa

Napr. fpt://a@b/

Príklad 3: nesprávny port (port môže byť len číselný)

Napr. telnet://foo@foo:foo

Príklad 4: chýbajúci xalphas

Napr. mailto::

1.2. Prepis gramatiky z BNF do gramatických pravidiel s alternatívami

1.2.1. Gramatika v BNF

```
url := httpaddress | ftpaddress | telnetaddress | mailtoaddress
httpaddress := ,,http://" hostport [ / path ] [ ? search ]
ftpaddress := ,,ftp://" login / path
telnetaddress := ,,telnet://" login
mailtoaddress := ,,mailto::" xalphas @ hostname
login := [ user [ : password ] @ ] hostport
hostport := hostname [ : port ]
hostname := xalphas { . xalphas }
port := digits
path := segment { / segment }
search := xalphas { + xalphas }
user := xalphas
password := xalphas
segment := { xalpha }
xalphas := xalpha { xalpha }
xalpha := alpha | digit
digits := digit {digit}
alpha := ,,A" | .. | ,,Z" | ,,a" | .. | ,,z"
digit := ,0" | .. | ,9"
```

1.2.2. Pravidlá s alternatívami

Zadefinovanie neterminálov:

A = url
B = httpaddress
C = ftpaddress
D = telnetaddress
E = mailtoaddress
F = login

```
G = hostport
H = hostname
I = port
J = path
K = search
L = user
M = password
N = segment
O = xalphas
P = xalpha
R = digits
S = alpha
T = digit
```

Prepis gramatiky v BNF pomocou neterminálov definovaných vyššie:

```
A \rightarrow B \mid C \mid D \mid E
B \rightarrow \text{,http://"} G \mid \text{,http://"} G \mid J \mid \text{,http://"} G \mid J \mid K \mid \text{,http://"} G \mid K
C \rightarrow ,ftp://"F/J
D \rightarrow ,,telnet://" F
E \rightarrow,,mailto::" O @ H
F \rightarrow L : M @ G | L @ G | G
G \rightarrow H \mid H : I
H \rightarrow OH
\mbox{H}\mbox{^{\smallfrown}} \rightarrow . O \mbox{H}\mbox{^{\backprime}} \mid \epsilon
I \rightarrow R
J \rightarrow N J
J \rightarrow / N J \mid \epsilon
K \rightarrow O K
K^{\hat{}} \rightarrow + O K^{\hat{}} \mid \epsilon
L \rightarrow O
M \rightarrow O
N \rightarrow PN \mid \epsilon
O \rightarrow PO
O' \rightarrow PO' \mid \epsilon
P \rightarrow S \mid T
R \rightarrow T R
R^{\cdot} \rightarrow T R^{\cdot} | \epsilon
S \rightarrow [A-Za-z]
T \rightarrow [0-9]
```

Červenou farbou sú vyznačené miesta, kde gramatika nesĺpňa podmienky LL(1). Pri vyznačených miestach nastáva First-First konflikt.

1.3. Transformácia gramatiky do LL(1)

Na odstránenie First-First kofliktov, ktoré sme opísali vyššie sme využili faktorizáciu:

$$A \rightarrow B \mid C \mid D \mid E$$

```
B \rightarrow ,,http://" G B
B \rightarrow / J B \mid ? K \mid \epsilon
B^{\sim} \rightarrow ? K \mid \epsilon
C \rightarrow ,ftp://" F / J
D \rightarrow , telnet://" F
E \rightarrow,,mailto::" O @ H
F \rightarrow F G
F \rightarrow L F \mid \epsilon
F`` \to @ | : M @
G \rightarrow H G
G \rightarrow : I \mid \epsilon
H \rightarrow O H
H^{\rightarrow} . O H^{\rightarrow} | \epsilon
I \rightarrow R
J \rightarrow N J
J \rightarrow / N J \mid \epsilon
K \rightarrow O K
K^{\hat{}} \rightarrow + O K^{\hat{}} \mid \epsilon
L \rightarrow O
M \rightarrow O
N \rightarrow PN \mid \epsilon
O \rightarrow PO
O' \rightarrow PO' \mid \epsilon
P \rightarrow S \mid T
R \rightarrow T R
R' \rightarrow T R' | \epsilon
S \rightarrow [A-Za-z]
T \rightarrow [0-9]
```

Červenou farbou je vyznačené miesto, kde gramatika nesĺpňa podmienky LL(1). Pri vyznačenom mieste nastáva First-First konflikt kvôli číslam a veľkým a malým písmenám. Výpočet First pre vybrané neterminálu pre znázornenie First-First konfliktu:

```
\begin{split} & First(F) = First(F^{`}) \setminus \epsilon \ U \ First(G) = \{ [0\text{-}9], \ [A\text{-}Za\text{-}z] \} \ U \ \{ \ [0\text{-}9], \ [A\text{-}Za\text{-}z] \} \\ & First(F^{`}) = First(L) \ U \ \{ \epsilon \} = \{ [0\text{-}9], \ [A\text{-}Za\text{-}z] \} \ U \ \{ \epsilon \} = \{ [0\text{-}9], \ [A\text{-}Za\text{-}z], \epsilon \} \\ & First(G) = First(H) = \{ [0\text{-}9], \ [A\text{-}Za\text{-}z] \} \\ & First(L) = First(O) = \{ [0\text{-}9], \ [A\text{-}Za\text{-}z] \} \\ & First(H) = First(O) = \{ [0\text{-}9], \ [A\text{-}Za\text{-}z] \} \\ & First(O) = First(S) \setminus \epsilon \ U \ First(T) = \{ [0\text{-}9], \ [A\text{-}Za\text{-}z] \} \\ & First(T) = \{ [0\text{-}9] \} \\ & First(S) = \{ [A\text{-}Za\text{-}z] \} \end{split}
```

Po viacerých pokusoch odstránenia First-First konfliktu, po lepšej analýze daného konfliktu a po konzultácií s cvičiacim sa prišlo na to, že hore uvedenú gramatiku nie je možné tranformovať na typ LL(1) prostriedkami obsiahnutými v predmete.

Problém nastáva pri riadkoch:

```
login := [ user [ : password ] @ ] hostport
hostport := hostname [ : port ]
```

Ako príklad si zoberieme vstup, ktorý obsahuje "usera" a vstup, ktorý neobsahuje "usera", ale len "hostname. "User" aj "hostname" sa môžu skladať z čísel a z malých a z veľkých písmen. Gramatika nešpecifikuje spôsob, ako zistiť, že aktuálne spracovávaný znak patrí pre "usera" alebo pre "hostname".

Na riešenie tohto problému sme zaviedli pravidlo, že sa "user" musí nachádzať v URL adrese práve raz.

Problémové riadky po zavedení pravidla vyzerajú nasledovne:

```
login := user [ : password ] @ hostport
hostport := hostname [ : port ]
```

Prepis zmenenej gramatiky v BNF pomocou neterminálov definovaných vyššie (zmeny oproti predošlej gramatike sú vyznačené farebne):

```
A \rightarrow B \mid C \mid D \mid E
B \rightarrow ,http://" G B
B \rightarrow /JB = |?K|\epsilon
B^{\sim} \rightarrow ? K \mid \varepsilon
C \rightarrow ,ftp://" F / J
D \rightarrow ,,telnet://" F
E \rightarrow "mailto::" O @ H
F \rightarrow L F @ G
F \rightarrow : M \mid \epsilon
G \rightarrow H G
G \rightarrow : I \mid \epsilon
H \rightarrow OH
H^{\hat{}} \rightarrow . O H^{\hat{}} \mid \epsilon
I \rightarrow R
J \rightarrow N J
J \rightarrow / N J \mid \varepsilon
K \rightarrow O K
K^{\hat{}} \rightarrow + O K^{\hat{}} \mid \epsilon
L \rightarrow O
M \rightarrow O
N \rightarrow PN \mid \epsilon
O \rightarrow PO
O' \rightarrow PO' \mid \epsilon
P \rightarrow S \mid T
R \rightarrow T R
R' \rightarrow T R' | \epsilon
S \rightarrow [A-Za-z]
T \rightarrow [0-9]
```

Tejto gramatike už neboli nájdené žiadne miesta, ktoré by nespĺňali podmienky LL(1).

Rozhodli sme sa však gramatiku ďalej upraviť, zjednodušiť. Ako prvé sme odstránili zbytočné neterminály s jedným pravidlom, ktoré obsahuje jeden neterminál. To sú neterminály I, L a M.

Všetky výskyty týchto neterminálov sme nahrali ich pravidlom. Gramatika po tejto zmene vyzerá nasledovne (zmeny oproti predošlej gramatike sú vyznačené farebne):

```
A \rightarrow B \mid C \mid D \mid E
B \rightarrow ,http://" G B
\overrightarrow{B} \, \stackrel{\text{\tiny ``}}{\rightarrow} \, / \, \overrightarrow{J} \, \overrightarrow{B} \, \stackrel{\text{\tiny ``}}{\mid} \, ? \, K \mid \epsilon
B^{\sim} \rightarrow ? K \mid \varepsilon
C \rightarrow ,ftp://" F / J
D \rightarrow ,,telnet://" F
E \rightarrow "mailto::" O @ H
F \rightarrow O F @ G
F \rightarrow : O \mid \varepsilon
G \rightarrow H G
G \rightarrow : R \mid \varepsilon
H \rightarrow O H
H \rightarrow .OH \mid \epsilon
J \rightarrow N J
J \rightarrow / N J \mid \epsilon
K \rightarrow O K
K^{\hat{}} \rightarrow + O K^{\hat{}} \mid \epsilon
N \rightarrow PN \mid \epsilon
O \rightarrow PO
O' \rightarrow PO' \mid \varepsilon
P \rightarrow S \mid T
R \rightarrow T R
R \rightarrow T R \mid \epsilon
S \rightarrow [A-Za-z]
T \rightarrow [0-9]
```

Následne sme gramatiku prepísali na krajší tvar, kde sme zmenili len pomenovanie neterminálov tak, aby išli v abecednom poradí a neobsahovali apostrof:

```
A \rightarrow ,,http://" F B | ,,ftp://" D / J | ,,telnet://" D | ,,mailto::" O @ H
B \rightarrow /JC \mid ?L \mid \epsilon
C \rightarrow ? L \mid \epsilon
D \rightarrow O E @ F
E \rightarrow : O \mid \epsilon
F \to H \; G
G \rightarrow : R \mid \epsilon
H \rightarrow OI
I \rightarrow . OI \mid \epsilon
J \rightarrow N K
K \rightarrow / N K \mid \epsilon
L \rightarrow OM
M \rightarrow + OM \mid \epsilon
N \rightarrow Q N \mid \epsilon
O \rightarrow Q P
P \to Q \; P \; | \; \epsilon
Q \rightarrow U \mid T
R \rightarrow T S
```

```
S \rightarrow T S \mid \epsilon
T \rightarrow [0-9]
U \rightarrow [A-Za-z]
    1.4.
             Overenie LL(1)
     1.4.1. First<sub>1</sub>
First(U) = \{[A-Za-z]\}\
First(T) = \{[0-9]\}\
First(S) = First(T) U \{\epsilon\} = \{[0-9]\} U \{\epsilon\} = \{[0-9], \epsilon\}
First(R) = First(T) = \{[0-9]\}
First(Q) = First(T) U First(U) = \{[0-9]\} U \{[A-Za-z]\} = \{[0-9], [A-Za-z]\}
First(P) = First(Q) U \{\epsilon\} = \{[0-9], [A-Za-z]\} U \{\epsilon\} = \{[0-9], [A-Za-z], \epsilon\}
First(O) = First(Q) = \{[0-9], [A-Za-z]\}
First(N) = First(Q) U \{\epsilon\} = \{[0-9], [A-Za-z]\} U \{\epsilon\} = \{[0-9], [A-Za-z], \epsilon\}
First(M) = \{+, \epsilon\}
First(L) = First(O) = \{[0-9], [A-Za-z]\}
First(K) = \{/, \epsilon\}
First(J) = First(N) \ \epsilon U First(K) = {[0-9], [A-Za-z]} U {/, \epsilon} = {[0-9], [A-Za-z], /, \epsilon}
First(I) = \{., \epsilon\}
First(H) = First(O) = \{ [0-9], [A-Za-z] \}
First(G) = \{:, \epsilon\}
First(F) = First(H) = \{[0-9], [A-Za-z]\}
First(E) = \{:, \epsilon\}
First(D) = First(O) = \{ [0-9], [A-Za-z] \}
First(C) = \{?, \epsilon\}
First(B) = \{/, ?, \epsilon\}
First(A) = {,,http://", ,,ftp://", ,,telnet://", ,,mailto::"}
    1.4.2. Follow<sub>1</sub>
Follow(A) = \{\$\}
Follow(B) = Follow(A) = \{\$\}
Follow(C) = Follow(B) = \{\$\}
Follow(D) = Follow(A) U \{/\} = \{\$\}\ U \{/\} = \{\$, /\}
Follow(E) = \{ @ \}
Follow(F) = First(B) \ \epsilon U Follow(A) U Follow(D) = {/, ?} U {$} U {$} = {/, ?, $}
Follow(G) = Follow(F) = \{/, ?, \$\}
Follow(H) = Follow(A) U First(G) \ \epsilon U Follow(F) = {$} U {:} U {/, ?, $} = {$, :, /, ?}
Follow(I) = Follow(H) U Follow(I) = \{\$, :, /, ?\}
Follow(J) = Follow(A) U First(C) \setminus \varepsilon U Follow(B) = {$} U {?} U {$} = {$, ?}
Follow(K) = Follow(J) U Follow(K) = \{\$, ?\}
Follow(L) = Follow(B) U Follow(C) = \{\$\} U \{\$\} = \{\$\}
```

 $Follow(M) = Follow(L) U Follow(M) = \{\$\}$

```
Follow(N) = First(K) \ \epsilon U Follow(J) U Follow(K) U Follow(N) = \{/\} U \{\$, ?\} U \{\$, ?\}
= \{/, \$, ?\}
Follow(O) = \{ @ \} U First(E) \setminus \varepsilon U Follow(D) U Follow(E) U First(I) \setminus \varepsilon U Follow(H)
U Follow(I) U First(M) \setminus \varepsilon U Follow (L) U Follow (M)
= \{ @ \} U \{ : \} U \{ \$, / \} U \{ @ \} U \{ . \} U \{ \$, :, /, ? \} U \{ \$, :, /, ? \} U \{ + \} U \{ \$ \} U \{ \$ \} = \{ @, :, \$, \}
/, ., ?, +}
Follow(P) = Follow(O) U Follow(P) = \{ @, :, \$, /, ., ?, + \}
Follow(Q) = First(N) \setminus \varepsilon U Follow(N) U First(P) \setminus \varepsilon U Follow(O) U Follow(P)
= \{[0-9], [A-Za-z]\} \cup \{/, \$, ?\} \cup \{[0-9], [A-Za-z]\} \cup \{@, :, \$, /, ., ?, +\} \cup \{@, :, \$, /, ., ?, +\}
=\{[0-9], [A-Za-z], /, \$, ?, @, :, ., +\}
Follow(R) = Follow(G) = \{/, ?, \$\}
Follow(S) = Follow(R) U Follow(S) = \{/, ?, \$\}
Follow(T) = Follow(Q) U First(S) \setminus \varepsilon U Follow(R) U Follow(S) = {[0-9], [A-Za-z], /, $, ?, @,
:, ., +} U {[0-9]} U {/, ?, $} U {/, ?, $}
=\{[0-9], [A-Za-z], /, \$, ?, @, :, ., +\}
Follow(U) = Follow(Q) = \{[0.9], [A-Za-z], /, \$, ?, @, :, ., +\}
```

Po analýze množín FIRST₁ a FOLLOW₁ sme sa presvedčili, že gramatika je LL(1).

1.5. Tabul'ka prechodov

Pravidlo A 1: http://FB 2: ftp://D/J 3: telnet://D	Pravidlo K 20 : /NK 21 : epsilon
4 : mailto::O@H	Pravidlo L 22 : OM
Pravidlo B	D
5 : /JC 6 : ?L	Pravidlo M 23:+OM
7 : epsilon	24 : epsilon
Provide C	Duarri dla M
Pravidlo C 8: ?L	Pravidlo N 25 : QN
9 : epsilon	26 : epsilon
Pravidlo D	Pravidlo O
10 : OE@F	27 : QP
	_
Pravidlo E	Pravidlo P
11::0	28 : QP
12 : epsilon	29 : epsilon
Pravidlo F	Pravidlo Q
13 : HG	30 : U
D 11 C	31 : T
Pravidlo G 14::R	Pravidlo R
15 : epsilon	32 : TS
Pravidlo H	32.15
16 : OI	Pravidlo S
	33 : TS
Pravidlo I	34 : epsilon
17:.OI	
18 : epsilon	Pravidlo T
Pravidlo J	35 : [0-9]
19: NK	Pravidlo U
17.111	36 : [A-Za-z]
	23.[11242]

	[A-Za-z]	[0-9]	+	/	•	:	@	http://	ftp://	telnet://	mailto::	?	\$
A	error	error	error	error	error	error	error	1	2	3	4	error	error
В	error	error	error	5	error	error	error	error	error	error	error	6	7
C	error	error	error	error	error	error	error	error	error	error	error	8	9
D	10	10	error	error	error	error	error	error	error	error	error	error	error
E	error	error	error	error	error	11	12	error	error	error	error	error	error
F	13	13	error	error	error	error	error	error	error	error	error	error	error
G	error	error	error	15	error	14	error	error	error	error	error	15	15
H	16	16	error	error	error	error	error	error	error	error	error	error	error
I	error	error	error	18	17	18	error	error	error	error	error	18	18
J	19	19	error	19	error	error	error	error	error	error	error	19	19
K	error	error	error	20	error	error	error	error	error	error	error	21	21
L	22	22	error	error	error	error	error	error	error	error	error	error	error
M	error	error	23	error	error	error	error	error	error	error	error	error	24
N	25	25	error	26	error	error	error	error	error	error	error	26	26
0	27	27	error	error	error	error	error	error	error	error	error	error	error
P	28	28	29	29	29	29	29	error	error	error	error	29	29
Q	30	31	error	error	error	error	error	error	error	error	error	error	error
R	error	32	error	error	error	error	error	error	error	error	error	error	error
S	error	33	error	34	error	error	error	error	error	error	error	34	34
T	error	35	error	error	error	error	error	error	error	error	error	error	error
U	36	error	error	error	error	error	error						

2. Implementácia

Podiel práce:

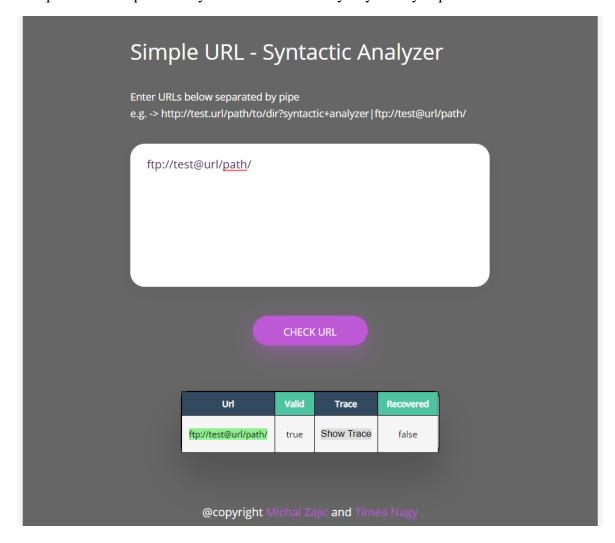
- Timea Nagy 50%
- Michal Zajic 50%

Implementovali sme tabuľkou riadený syntaktický analyzátor v jazyku JavaScript spolu s používateľským rozhraním a testami. Analyzátor rozpoznáva vety v jazyku simple URL.

2.1. Používateľské rozhranie

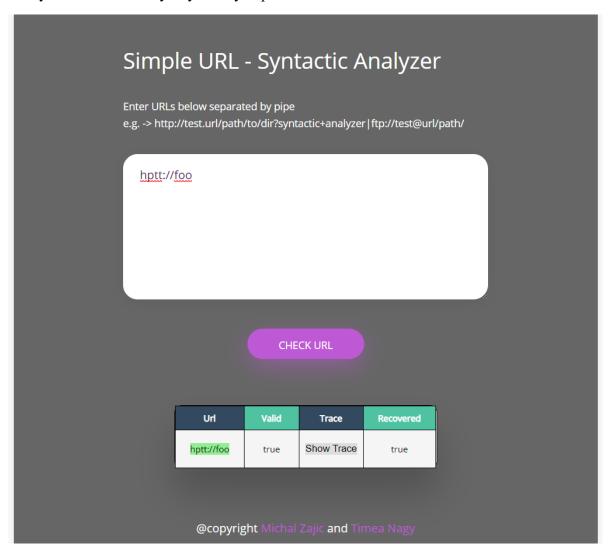
Analyzátor je možné spustiť spolu s používateľ ským rozhraním, ktorého základ je "index.html" súbor. Používateľ môže zadať ľubovoľný počet URL adresies oddelených pomocou znaku pipe "|". Po zadaní URL adresy analyzátor vypíše či je adresa platná alebo nie a či sa využilo zotavenie z chýb alebo nie. Ponúka aj možnosť zobrazenia krokov, ako bol vstup spracovaný.

Pri správnom vstupe a bez využitia zotavenia z chýb vyzerá výstup nasledovne:



Rules	Url
A	ftp://test@url/path/\$
ftp://D/J	ftp://test@url/path/\$
OE@F/J	test@url/path/\$
QPE@F/J	test@url/path/\$
UPE@F/J	test@url/path/\$
[A-Za-z]PE@F/J	test@url/path/\$
QPE@F/J	est@url/path/\$
UPE@F/J	est@url/path/\$

Pri využití zotavenia chýb vyzerá výstup nasledovne:



Rules	Url
A	hptt://foo\$
http://FB	http://foo\$
НСВ	foo\$
OIGB	foo\$
QPIGB	foo\$
UPIGB	foo\$
[A-Za-z]PIGB	foo\$
QPIGB	00\$

2.2. Testy

Pre analyzátor sme napísali testy tak, aby sme pokryli všetky možné URL adresy, ktoré sú platné a niekoľko neplatných adries. Nachádzajú sa v súbore "analyzer.test.js". Testy sme rozdelili na štyri časti podľa štyroch pravidiel, ktoré sa nachádzajú pri štartovacom neterminály – URL adresy začínajúce na "http://", "ftp://", "telnet://" a "mailto::". V rámci všetkých týchto typov sme zadefinovali platné URL adresy, adresy, pri ktorých sa využilo zotavenie z chýb a chybné URL adresy. Spustiť ich je možné pomocou príkazu "npm test".

K URL adresám máme informáciu či je adresa platná alebo nie a či bolo využité zotavenie z chýb alebo nie.

Ukážka platných URL adries:

```
["http://a", true, false],
  ["http://1", true, false],
  ["http://Z", true, false],
  ["http://5", true, false],
  ["http://foobarfoo", true, false],
  ["http://1234567", true, false],
  ["http://F00", true, false],
  ["http://foobarfoo123", true, false],
```

Ukážka URL adries, pri ktorých sa využilo zotavenie:

```
["tenet://foo:123456@foo:123456", true, true], // incorrect telnet address ["telnet://foo@foo:foo", true, true], // incorrect port
```

Ukážka neplatných URL adries:

```
["mailto::", false, false], // missing xalphas
["mailto::a@", false, false], // missing hostname
["mailto::a", false, false], // missing hostname and @ character
["mailto::a@b.", false, false], // incorrect hostname
["mailto::a.foo.bar@b", false, false], // incorrect xalphas
```

2.3. Tabuľka a pravidlá

Reprezentácia prechodovej tabuľky a pravidiel sa nachádza v súbore "table_and_rules.js". Má formu *dictionary*. Pre každý neterminál si pamätáme terminály, ktoré vieme vygenerovať pomocou tohto neterminálu a k terminálom pravidlo, ktorým dostaneme musíme pokračovať, aby sme dostali daný terminál. Ak neterminál môže vygenerovať hocijaké písmená alebo čísla, v *dictionary* máme pre tieto terminály špeciálne označenie. Na základe tohto označenia analyzátor na začiatku vygeneruje *key-value* hodnoty do *dictionary* pre každé písmeno alebo číslo. Kód pre generovanie týchto hodnôt sa nachádza v súbore "table_transform.js".

Ukážka reprezentácia tabuľky pre neterminál S:

```
"S" : {
    "number" : 33,
    "/" : 34,
    "?" : 34,
    "$" : 34
}
```

Ukážka reprezentácie pravidiel:

```
1: "http://FB",
2: "ftp://D/J",
3: "telnet://D",
4: "mailto::O@H",
5: "/JC",
6: "?L",
7: '',
8: "?L",
```

2.4. Postup implementácie

Samotný proces spracovávania vstupu sa nachádza v súbore "simple_url_syntactic_analyzer .js". Algoritmus začína predspracovaním vstupu, v rámci ktorého pridáme na koniec vstupu znak "\$".

Po predspracovaní inicializujeme zásobník, t. j. pridáme štartovací neterminál.

Samotný algoritmus odvodenia prebieha ako rekurzia. Má nasledovný priebeh:

- 1. Vyberieme prvý znak zo zásobníka, ktorý môže byť terminál alebo neterminál a vyberieme aj prvý znak z nespracovanej časti URL adresy.
- 2. Ak sú tieto znaky rovnaké a nejedná sa o neterminál, tak odstraňujeme tento znak zo zásobníka a aj z URL adresy. Pri štartovacom pravidle odstraňujeme celý terminál napr. "http://".
- 3. Ak sme rekurziu zavolali s regulárným výrazom napr. [A-Za-z] alebo [0-9], tak skontrolujeme či platí na prvý znak z nespracovanej časti URL adresy. Ak áno, odstránime daný znak.
- 4. Skontrolujeme či máme epsilonové pravidlo. Ak áno, zavoláme rekurziu znova so zvyšnými pravidlami.
- 5. Znova vyberieme prvý znak zo zásobníka a vyberieme aj prvý znak z nespracovanej časti URL adresy.
- 6. Ak narazíme na znak "\$", teda spracovali sme celú URL adresu a nemáme žiadne ďalšie pravidlá v zásobníku, adresa je validná a rekurziu skončíme.
- 7. Ak pre aktuálny neterminál nemáme žiadne pravidlo, s ktorým by sme sa dostali k terminálu na začiatku URL adresy, adresa nie je validná a končíme rekurziu.
- 8. Aplikujeme pravidlo (uložíme do zásobníka) pre aktuálny neterminál na základe terminálu, ktorý sa nachádza na začiatku nespracovanej časti URL adresy.
- 9. Ak pravidlo je regulárny výraz, musíme ju poslať do rekurzie.
- 10. Zavoláme rekurziu a celý proces začíname od znova.

2.5. Zotavenie z chýb

Zotavenie chýb nastane, ak sa používateľ pomýli na začiatku URL adresy alebo pri porte.

URL adresa môže začínať len so štyrmi rôznymi reťazcami znakov podľa typu URL adresy - "http://", "ftp://", "telnet://" a "mailto::". Ak sa používateľ pomýli v tejto časti URL adresy, tak sa pomocou edit distance algoritmu zistí, na ktorú z tých štyroch typov sa URL adresa podobá a na základe toho zmeníme začiatok URL adresy.

Port sa môže skladať len z čísel. Chceme však uznať aj URL adresy, v ktorých sa používateľ pomýlil pri porte a namiesto čísel tam napísal aj písmená. Pridali sme nové pravidlá do tabuľky:

$$R \to U S$$
$$S \to U S$$

Kde R a S sú terminály, ktoré generujú port a U generuje veľké a malé písmená.