

Problems with simplistic parser in regex to E-NFA

by Recursive Matcher

initial ideas:
(wrong)

regex	E-NFA
a	$\rightarrow (1) \xrightarrow{a} (2) \rightarrow$
a^*	$\rightarrow (1) \xrightarrow{a} (2) \rightarrow$ $\quad \quad \quad \downarrow \epsilon \quad \quad \quad \downarrow \epsilon$
$a + b$	$\rightarrow (1) \xrightarrow{a} (2) \rightarrow$ $\quad \quad \quad \downarrow \epsilon \quad \quad \quad \downarrow \epsilon$ $\quad \quad \quad (3) \xrightarrow{b} (4) \rightarrow$

pros: simple,
easily stackable
(no changing
previous decisions),
parses from left
to right.

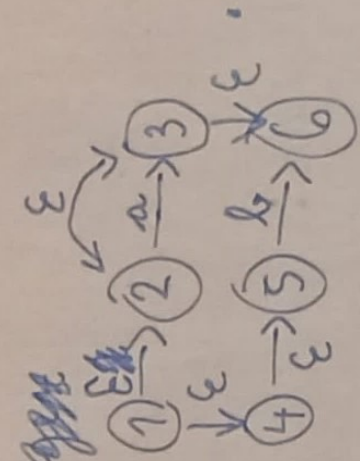
PROBLEMS:

1st noticed problem	regex	old - solution	new solution
①	$(a^*) + b$	$(1) \xrightarrow{a} (2) \rightarrow$ $\quad \quad \quad \downarrow \epsilon \quad \quad \quad \downarrow \epsilon$ $\quad \quad \quad (3) \xrightarrow{b} (4) \rightarrow$	accepts extra words. Example: "a b"

It can reach "b" after backtracking on "a" which isn't intended in the regex. Therefore, I came with a solution:

regex	E-NFA
a	$\rightarrow (1) \xrightarrow{a} (2) \rightarrow (3) \rightarrow$

graph for the $(a^*) + b$ is:



Now it works on this input. Still, it fails on the next (on pg. 2).

Evaluating grammar simplification

(Pascariu Matei)

regex	to ϵ -NFA	problem
$a + b^*$		Similarly, it accepts "a b".

(Ugly) solution: change "a" to $\epsilon \rightarrow 1 \xrightarrow{\epsilon} 2 \xrightarrow{a} 3 \xrightarrow{\epsilon} 4$.

Now, the problems seem all fixed. NOPE!

regex	to ϵ -NFA	problem:
$a + (a + b)^*$		accepts "ca"

solution: choose a different method.

Moral: be very careful when taking a shortcut/simplified solution.