Compiling a Regular Expression

RE	\rightarrow	ALT \$
ALT	\rightarrow	SEQ ALTLIST
ALTLIST	\rightarrow	pipe SEQ ALTLIST
		λ
SEQ	\rightarrow	ATOM SEQLIST
		λ
SEQLIST	\rightarrow	ATOM SEQLIST
		λ
ATOM	\rightarrow	NUCLEUS ATOMMOD
ATOMMOD	\rightarrow	kleene
		plus
		λ
NUCLEUS	\rightarrow	open ALT close
		char CHARRNG
	1	dot
CHARRNG	\rightarrow	dash char
		λ

To the left is an LL(1) compatible grammar for simple regular expressions, **sufficient for lexing** (scanning) simple program source. Notice the **order of operations** within the grammar:

 The lowest precedence operator is ALTternation; it is also wrapped by parenthesis

 $NUCLEUS \rightarrow openALT \ close$

which makes the form highest precedence.

- ALT s are made up of one or more
 SEQuences, separated by the *pipe* symbol,
- and SEQ s are a list of **ATOM**s

Compiling a Regular Expression

RE	\rightarrow	ALT \$
ALT	\rightarrow	SEQ ALTLIST
ALTLIST	\rightarrow	pipe SEQ ALTLIST
		λ
SEQ	\rightarrow	ATOM SEQLIST
		λ
SEQLIST	\rightarrow	ATOM SEQLIST
		λ
ATOM	\rightarrow	NUCLEUS ATOMMOD
ATOMMOD	\rightarrow	kleene
		plus
		λ
NUCLEUS	\rightarrow	open ALT close
		char CHARRNG
	- İ	dot
CHARRNG	\rightarrow	dash char
		λ
	,	

To the left is an LL(1) compatible grammar for simple regular expressions, **sufficient for lexing** (scanning) simple program source. Notice the **order of operations** within the grammar:

- An ATOM is made up of its NUCLEUS and possibly an operator: kleene (*, any number of) or plus (one or more).
- A NUCLEUS may be either a parenthetical group ((a|bc)*), a character range (just t or Q-V), or the any-character dot (.)

The Concrete ("Raw") Parse Tree

The raw (concrete) parse tree result shows there are no syntax errors in the source regular expression. But compared to its simplified abstract syntax tree it has several downsides:

- 1. Large (as in a much larger data structure),
- Most of the internal nodes are not needed to represent the semantic meaning of the original regex,
- 3. Did we mention big?





An Abstract Syntax Tree (AST)

We would like an efficient process for translating raw parse trees to ASTs, ideally we could do this **during the parse**.

The **simplified AST** is < 25% the size of the original version, and this is for a very small input.



Stack RE	
	dash
	D
	dot
	g
	plus
	\$

Operation: begin PARSE TREE



Operation: char predicts rule 2 $ALT \rightarrow SEQALTLIST$



PARSE TREE







PARSE TREE



Operation: char predicts rule 9 $ATOM \rightarrow NUCLEUS ATOMMOD$



<u>_*</u>

PARSE TREE



Operation: char predicts rule 14 $NUCLEUS \rightarrow char CHARRNG$



ALTLIST

*

\$

Operation: token char match (char)

STACK	QUEUE	Parse Tree
CHARRNG	dash	Root
*>	D	RE
ATOMMOD	dot	ALT
*	g	SEQ
SEQLIST	plus	ATOM
*>	\$	NUCLEUS
ALTLIST		Â
*		

\$

Operation: dash predicts rule 16 $CHARRNG \rightarrow dash char$

STACK	QUEUE
dash	dash
char	D
*	dot
*	g
ATOMMOD	plus
*	\$
SEQLIST	









Operation: token dash match (dash)

STACK	QUEUE	Parse Tree
char	D	Root
*>	dot	RE
*>	g	ALT
ATOMMOD	plus	SEQ
*	\$	ATOM
SEQLIST		NUCLEUS
*		A
ALTLIST		dash

< *`

\$

Operation: token char match (char)

STACK	QUEUE	PARSE TREE
*>	dot	Root
*>	g	RE
ATOMMOD	plus	ALT
*>	\$	SEQ
SEQLIST		ATOM
*>		NUCLEUS
ALTLIST		A CHARRNG
<*>		dash D

\$ /*`



Operation: end of CHARRNG production

D



Operation: end of NUCLEUS production — before rule 14 SDT procedure



Operation: end of NUCLEUS production — after rule 14 SDT procedure



Operation: end of NUCLEUS production — after rule 14 SDT procedure



Notice that our tree now contains node types not from the grammar (range), That's OK, it's your tree, you can do with it as you please :) Input RE: A-D.g+

Operation: end of NUCLEUS production

STACK	QUEUE	Parse Tree
ATOMMOD	dot	Root
*	g	RE
SEQLIST	plus	ALT
*	\$	SEQ
ALTLIST		Атом
*		NUCLEUS
\$		range
*		





Operation: end of ATOMMOD production



Operation: end of ATOM production — before rule 9 SDT procedure



if ($\mathit{node.ATOMMOD.child}$ is λ) then ($\mathit{replace}\ \mathit{node}\ with\ \mathit{node.children}[0].child$ in parse tree T return

Operation: end of ATOM production — after rule 9 SDT procedure



if ($\mathit{node.ATOMMOD.child}$ is λ) then ($\mathit{replace}\ \mathit{node}\ with\ \mathit{node.children}[0].child$ in parse tree T return

*

Operation: end of ATOM production

STACK	QUEUE	Parse Tree
SEQLIST	dot	Root
*>	g	RE
ALTLIST	plus	ALT
*>	\$	SEQ
\$		range
*>		



Operation: dot predicts rule 9 $ATOM \rightarrow NUCLEUS ATOMMOD$



ALTLIST

*

\$



Operation: dot predicts rule 15 NUCLEUS \rightarrow dot



*

\$





Operation: end of NUCLEUS production





Operation: char predicts rule 12 $\mathit{ATOMMOD} \rightarrow \lambda$

ATOMMOD





Input RE: A-D.g+

ATOMMOD

λ

Operation: end of ATOMMOD production





Operation: end of ATOM production — before rule 9 SDT procedure



if ($\mathit{node.ATOMMOD.child}$ is λ) then ($\mathit{replace}\ \mathit{node}\ with\ \mathit{node.children}[0].child$ in parse tree T return

Operation: end of ATOM production — after rule 9 SDT procedure



if ($\mathit{node.ATOMMOD.child}$ is λ) then ($\mathit{replace}\ \mathit{node}\ with\ \mathit{node.children}[0].child$ in parse tree T return

·*

Operation: end of ATOM production






Operation: char predicts rule 9 $ATOM \rightarrow NUCLEUS ATOMMOD$

STACK	QUEUE
NUCLEUS	g
ATOMMOD	plus
*	\$
SEQLIST	
*	
*	
*	
ALTLIST	

*

\$



PARSE TREE



Operation: char predicts rule 14 $NUCLEUS \rightarrow char CHARRNG$

STACK	QUEUE
char	g
CHARRNG	plus
<*>	\$
ATOMMOD	
*>	
SEQLIST	
*>	
*	

`*

ALTLIST



Operation: token char match (char)









Operation: end of CHARRNG production



Operation: end of NUCLEUS production — before rule 14 SDT procedure



Operation: end of NUCLEUS production — after rule 14 SDT procedure

SEOLIST

ATOM

NUCLEUS

g



Operation: end of NUCLEUS production



Root RE ALT SEQ SEQLIST range D dot SEQLIST ATOM NUCLEUS g



Operation: token plus match (plus)



۲×.

Operation: end of ATOMMOD production



Operation: end of ATOM production — before rule 9 SDT procedure



Operation: end of ATOM production — after rule 9 SDT procedure

SEOLIST

plus

q



Operation: end of ATOM production









Whoa! There was an SDT procedure in there... What was its logic?



Whoa! There was an SDT procedure in there... What was its logic? remove parent's right most child Input RE: A-D.g+



There is an implementation nuance here, are the SDT procedures expected to manage the notion of **the current node** of the parsing engine? In this particular case, the procedure could make the (new) right most node of *node.parent* the current node (**plus**); alternatively a **current node stack** could be used instead of a singular notion of the current node.



Yet **another approach** to avoiding *current node corruption* is to permit SDT procedures to manipulate **only their own and their descendent's structures**. This "don't talk back to your parent" approach works particularly well for SDT procedures in LR parses. Avoid mixing and matching your SDT philosophies within one grammar and implementation.

Operation: end of SEQLIST production — before rule 7 SDT procedure



Operation: end of SEQLIST production — after rule 7 SDT procedure





Input RE: A-D.g+

g

Operation: end of SEQLIST production — before rule 7 SDT procedure



Operation: end of SEQLIST production — after rule 7 SDT procedure

QUEUE PARSE TREE STACK Root RE ALTLIST ALT SEQ \$ SEOLIST range * dot plus **procedure** SEQLIST (node, parse tree T) α let $parent \leftarrow parent of node in T$ if (node.child is λ) then (trim last child of *parent.children* away) else if (parent is SEQLIST) then (let myChildren = node.children trim last child of *parent.children* away append parent.children with myChildren

Operation: end of SEQLIST production PARSE TREE STACK QUEUE * \$ Root ALTLIST RE ALT * \$ SEQ <u>^</u> SEQLIST range dot D plus

Input RE: A-D.g+

g

Operation: end of SEQ production — before rule 5 SDT procedure



Operation: end of SEQ production — after rule 5 SDT procedure





```
procedure SEQ( node, parse tree T )
    if ( node.children.SEQLIST exists ) then (
        replace node.children.SEQLIST with node.children.SEQLIST's children
)
    if ( |node.children| = 1 ) then (
        replace node with node.child in parse tree T
)
```







PARSE TREE Root RE ALT ALTLIST SEQ plus dot λ D g Α

Operation: end of ALTLIST production — before rule 4 SDT procedure







What will the SDT procedure for $\ ALTLIST \rightarrow \lambda$ production do? Input RE: **A-D.g+**

Operation: end of ALTLIST production — after rule 4 SDT procedure



PARSE TREE



Remove its parent's right most child.

Here is another instance where the *current node* notion in the parsing algorithm can fail if not coded carefully (because the current node was just lobbed off the tree). Input RE: **A-D.g+**





What should $ALT \rightarrow SEQALTLIST$ production do with only one child?
Operation: end of ALTLIST production





What should $ALT \rightarrow SEQALTLIST$ production do with only one child? Replace itself with its child. Input RE: A-D.g+

STACK QUEUE





Input RE: A-D.g+





PARSE TREE



And what about the $RE \rightarrow ALT$ \$ production?

Input RE: A-D.g+





PARSE TREE



And what about the $RE \rightarrow ALT$ \$ production? Replace itself with its left most child. Beware, the farther up the parse tree you progress, the less likely that a node's children will match a production rule's RHS! Input RE: **A**-**D**.**g**+

Operation: end of RE production

STACK QUEUE

PARSE TREE



Voilà, a proper RE expression tree. Your LGA tonight will ask you to complete the pseudo code for the other production rules in our RE grammar.