

# STRUCTURE

**in text : regexs and grammars**

# ARE THESE TEXTS "STRUCTURED"?

SQL schema definition or query:

```
DELETE FROM DomesticStudentsFor2022
WHERE mark = 'E';
```

Java statement:

```
while ( A[k] != x ) { k++; }
```

XML documents:

```
<html><head><title>My Web
Page</title></head>
<body><p> Thanks for viewing
</p></body></html>
```

Whakatauki:

I orea te tuatara ka puta ki  
waho

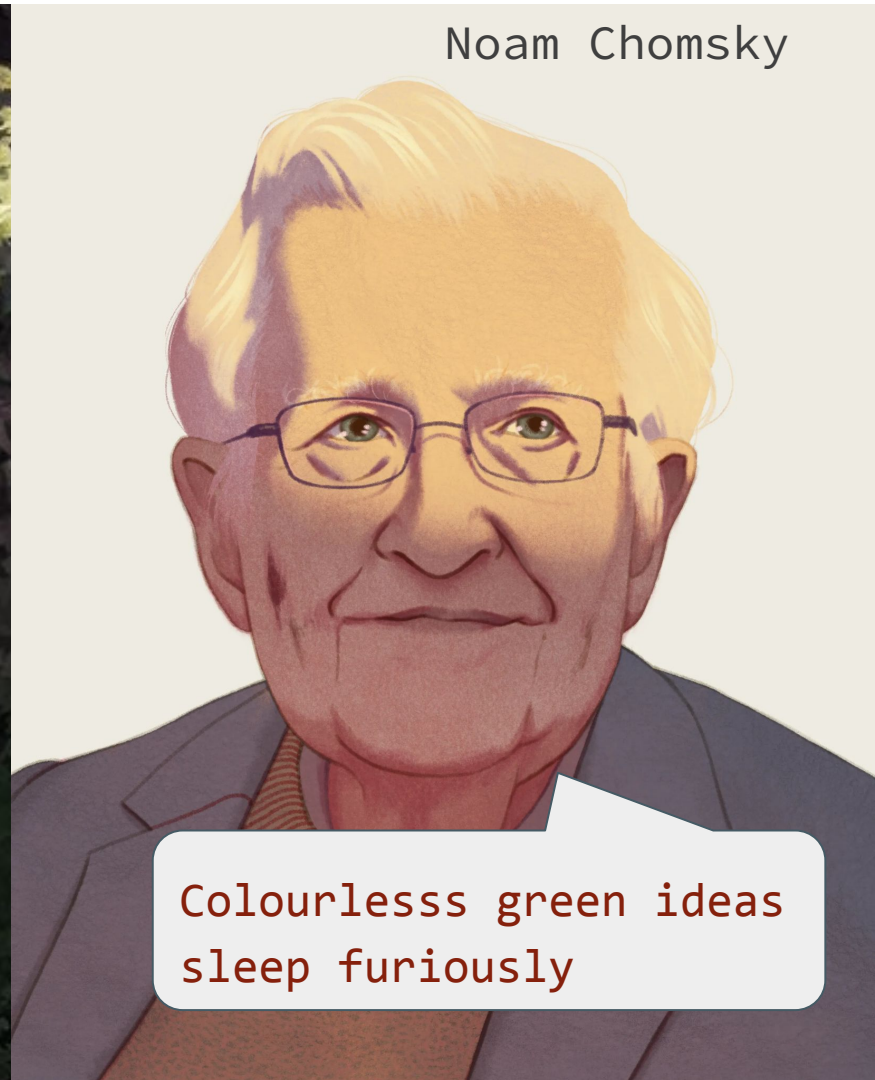
*(a problem is solved by  
continuing to find solutions)*

Poem:

I've hunted near,  
I've hunted far  
I even looked inside my car.  
I've lost my glasses,  
I'm in need,  
To have them now so I can read.

*Anne Scott, 2014*

Noam Chomsky



Colourless green ideas  
sleep furiously

# HOW TO DEFINE SYNTAX?

- Specify a finite set of acceptable sentences..?
- Instead: look for a finite recipe that describes all acceptable sentences

So a grammar is a **finite description** of a possibly infinite set of **acceptable sentences**



like a filter..?

Today we will look at the simplest grammars, or rather an equivalent representation we call regular expressions (“regex”)

# REGULAR EXPRESSIONS

there's the whole set of everything that's "legal", e.g.

**legal:** a, ab, aaa, aab, aaaaaaaaaabbbbb, b, bbb, ...

**not legal:** ba, aba, abba, bbbbbbaaaaaa, ...

A very condensed way of saying this is a **regular expression:**

**a\*b\***

➡ note the 'regex' is finite even though the list might not be.

# REGEX IN JAVA

Here's the start of the relevant [java docs](#):

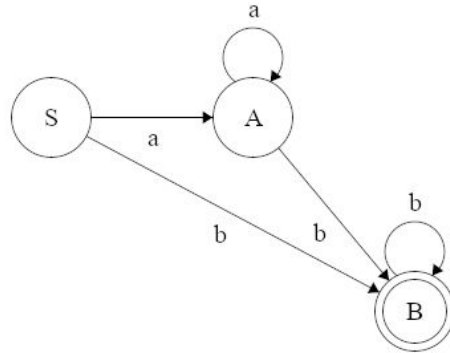
•	Any character
<code>\d</code>	A digit: <code>[0-9]</code>
<code>^</code>	“not something”, e.g. <code>[^0-9]</code> means <i>not</i> a digit
<code>\s</code>	A whitespace character: <code>[ \t\n\x0B\f\r]</code>
<code>\S</code>	A non-whitespace character: <code>[^\s]</code>
<code>XY</code>	X followed by Y
<code>X Y</code>	Either X or Y
<code>(X)</code>	X, as a capturing group
<code>X*</code>	X, zero or more times
<code>X+</code>	X, one or more times
<code>X{n}</code>	X, exactly n times

helpful:

<https://regexone.com/references/java>

# REGEX $\rightarrow$ AUTOMATA

the regex is just a code: javac turns it into a finite state machine



Wikipedia:

“A **finite-state machine (FSM)** or **finite-state automaton** is a mathematical model of computation.”

Constructing this machine isn't always trivial.

# JAVA DOES IT LIKE THIS

```
boolean b = Pattern.matches("a*b", "aaaaab");
```

combines the **generation** of the FSM (for the first arg) and its **use** on the particular string (second arg).

But because it's expensive to build that FSM, if you want to match lots of strings it's better to do it like this (from the docs):

```
Pattern p = Pattern.compile("a*b");
```

```
Matcher m = p.matcher("aaaaab");
```

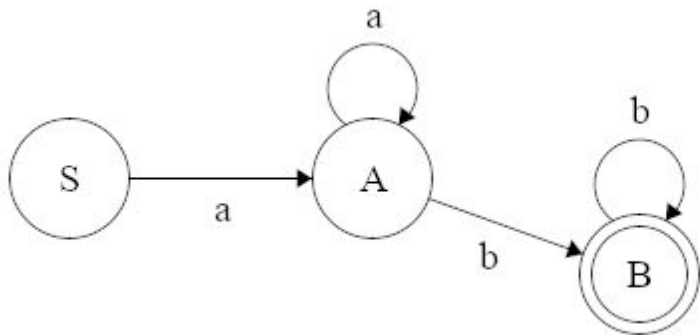
```
boolean b = m.matches();
```



## WAYS OF REPRESENTING A REGULAR LANGUAGE #2:

## FINITE STATE MACHINE

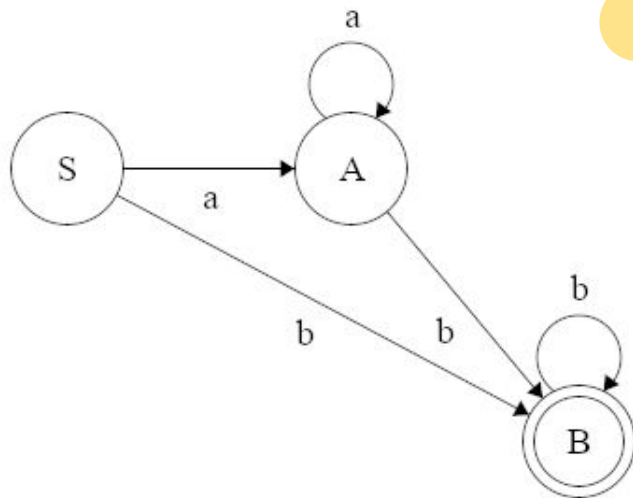
- $a+b+$
- $aa^*bb^*$



double circle  
means  
acceptable

- $a^*b^+$

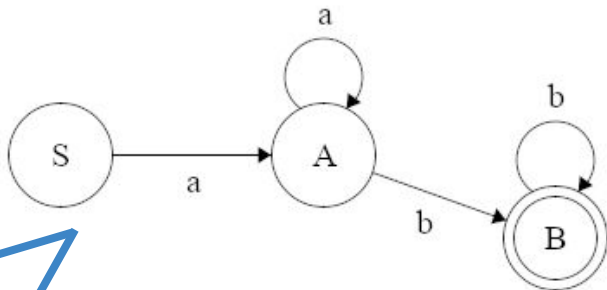
optional: we  
could have just  
made A the start  
state...



## GRAMMAR

Rules that spell out what's possible, eg:

- $S \rightarrow aA$
- $A \rightarrow aA \mid bB$
- $B \rightarrow bB \mid \text{end}$



(this is same as the regex:  $a^+b^+$ )

regex's correspond to just the simplest grammars - the "regular" ones

regular grammars  $\Leftrightarrow$  regex  
(describe the same set of sequences)

each rule component is either

- a "terminal" (a,b...)
- OR
- a terminal and a non-terminal (A,B,...)



# EXAMPLE OF A "CONTEXT FREE" GRAMMAR

HTMLFILE	□	"<html>"	[ HEAD ]	BODY	"</html>"			
HEAD	□	"<head>"	TITLE		"</head>"			
TITLE	□	"<title>"	TEXT		"</title>"			
BODY	□	"<body>"	[ BODYTAG ]*		"</body>"			
BODYTAG	□	H1TAG		PTAG		OLTAG		ULTAG
H1TAG	□	"<h1>"	TEXT		"</h1>"			
PTAG	□	"<p>"	TEXT		"</p>"			
OLTAG	□	"<ol>"	[ LITAG ]+		"</ol>"			
ULTAG	□	"<ul>"	[ LITAG ]+		"</ul>"			
LITAG	□	"<li>"	TEXT		"</li>"			
TEXT	□	<i>sequence of characters</i>						
		<i>other than &lt; and &gt;</i>						

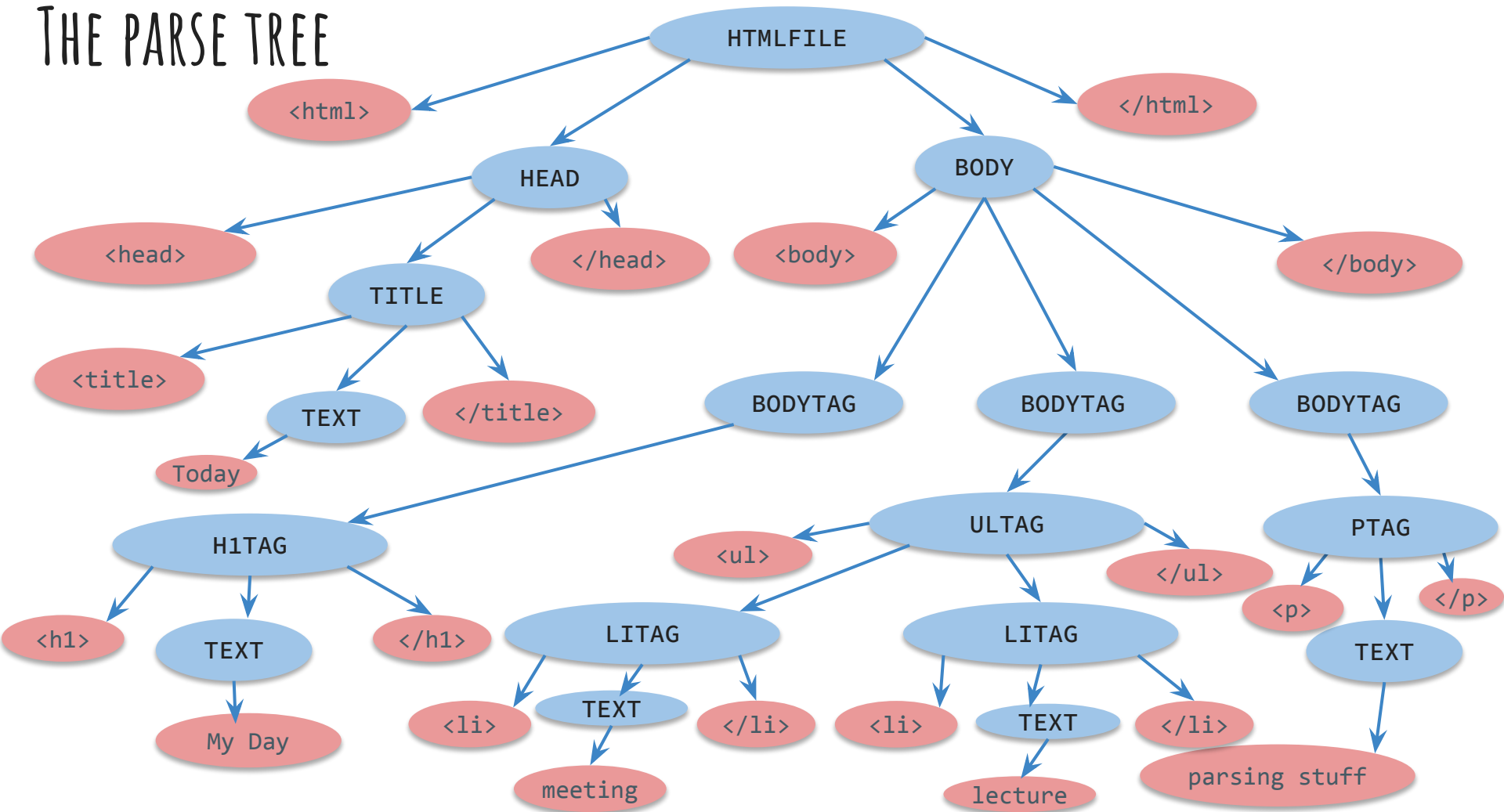
**more** than a  
"regular" grammar

regular grammars  $\Leftrightarrow$  regex  
(describe the same set of  
sequences)

each rule component is  
either

- a "terminal" (a,b,...)
- OR
- a terminal and a  
non-terminal  
(A,B,...)

# THE PARSE TREE



# CONTEXT-FREE GRAMMARS

we've gone beyond regex ( == **regular** grammars) to what are called **context free** grammars (CFG).

Regular grammars can't do "nesting", but CFG can.

e.g. can you write a regex for a language that includes

(x) ((x)) (((x))) (((((x)))) ...

?

Specific cases are fine, but the general matching between number of brackets isn't. This CFG can do it though:

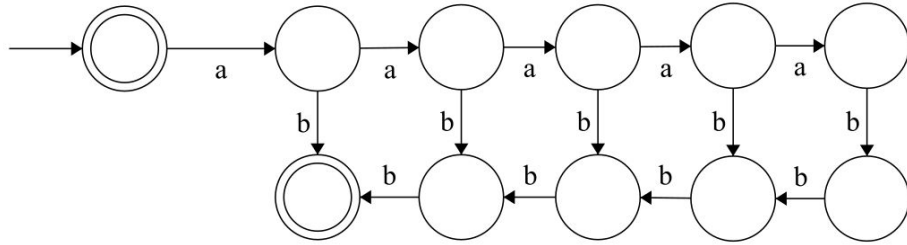
$$\text{EXPR} ::= \text{"x"} \quad | \quad \text{"(" EXPR ")"}$$

# HOW ABOUT THIS?

can you write a regex for a language that includes (only)

ab aabb aaabbb aaaabbbb ...

?



regex / FSM /  
regular grammar

attempts (fail)

Intuitively, you'd need some kind of a “**stack**”, right?

This Context Free Grammar can do it:

EXPR  $\rightarrow$  “a” “b” | “a” EXPR “b”

# THERE ARE LOTS OF GREAT WEB RESOURCES

just a couple of fun examples if you're interested:

- <https://regexone.com/>

(note in Java you have to use two \ 's )

- and ...

<http://web.mit.edu/6.005/www/fa16/classes/17-regex-grammars/>

# WHAT CODE SHALL WE WRITE TOMORROW?

We will read in a file that contains a grammar, such as

**START** -> **NP VP**

**NP** -> the **NOUN** | a **NOUN**

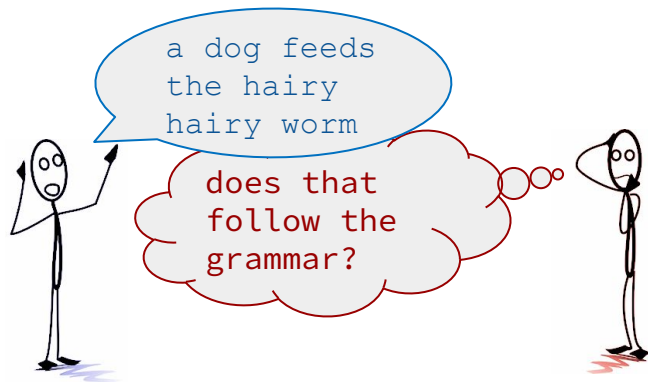
**VP** -> **VERB NP**

**NOUN** -> **THING** | **ADJ NOUN**

**THING** -> dog | cat | worm | ferret | door knob

**VERB** -> eats | squashes | feeds | loves

**ADJ** -> hairy | smooth



And then we will generate lots of strings from this grammar, like this:

the cat loves a dog

the hairy cat loves the door knob

the smooth hairy smooth hairy cat squashes the smooth ferret

the smooth ferret eats a hairy dog

a smooth dog eats the hairy door knob