

$r ::=$ constant string
 match this string &
 nothing else

Regular expressions
 are a
Domain specific
Language
 (DSL)

| $r_1 + r_2$
 match this or match that

| $r_1 \cdot r_2$
 match the prefix against this
 & suffix against that

| r^* (Kleene-*)
 match this again & again & again

| Bottom \perp \emptyset | Intersection $r_1 \cap r_2$
 match nothing at all. | Complement \overline{r}

Problem 1: Match sentences where the word
 "John" appears
 at least twice.

Problem 0: "John" appears at least once

string _____ pattern _____

"John" matches Constant "John"

"John went to ~~market~~ C "John"
the market"

P0: $[a-z]^* \cdot \text{"John"} \cdot [a-z]^*$

P1: $[a-z]^* \text{"John"} [a-z]^* \cdot \text{"John"} [a-z]^*$

$\cup [a-z]^* \text{John John } [a-z]^*$

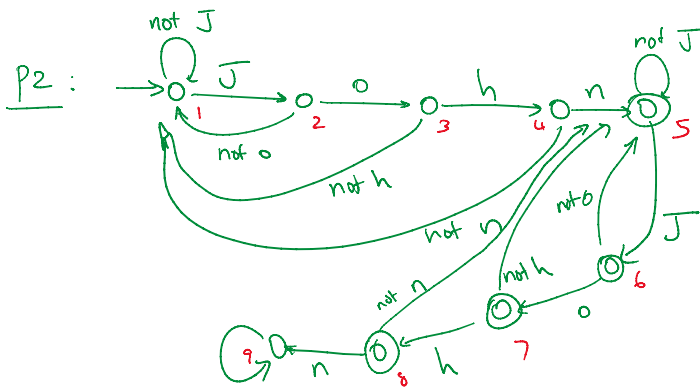
P2: Contains "John" exactly once.

P2': Does not contain John at all

$\neg ([a-z]^* \text{John } [a-z]^*) \cdot \text{John}$

$\rightarrow ([a-z]^* \text{John } [a-z]^*)$

John



Curr state	Curr symbol	Next state
1	J	2
1	not J	1
2	o	3
2	not o	1
⋮	⋮	⋮

$q = 1$
for each char a
in the string:
 $q := \text{table}[q, a]$

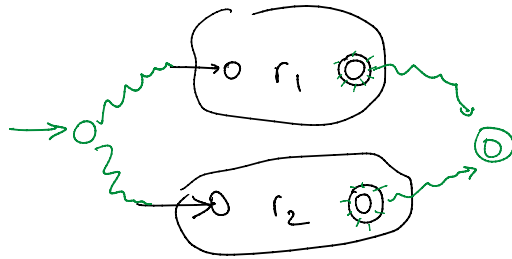
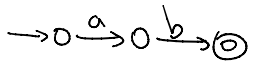
2 | not 0 | 1
 ⋮ | ⋮ | ⋮

1 | union {q, a}

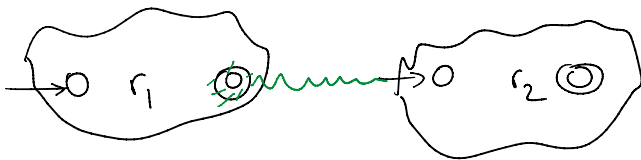
Games to match regular expressions

① Constant "ab"

② $r_1 + r_2$



③ Concatenation $r_1 \cdot r_2$



$a^* \cdot b^*$ "Sequence of a-s followed by a sequence of b-s"

$\underbrace{\quad}_{M_1} \quad \underbrace{\quad}_{M_2}$

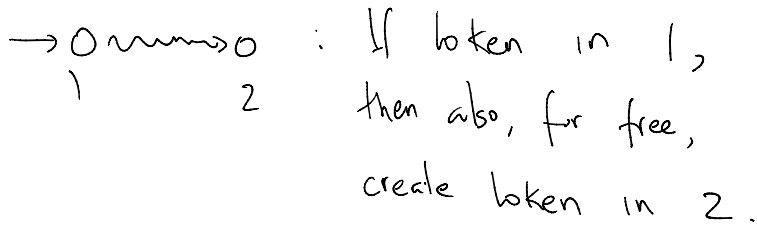
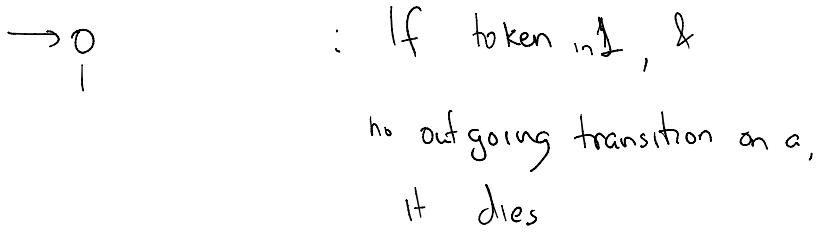
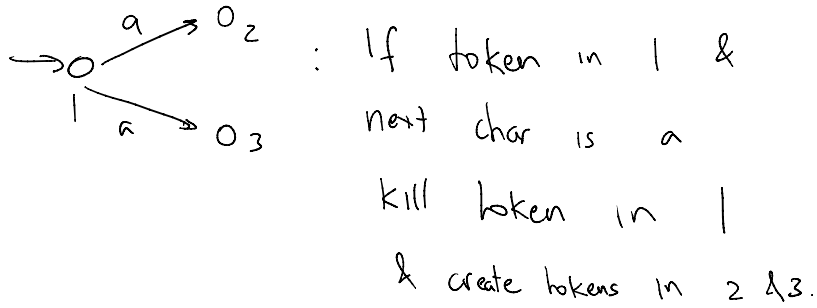
a a a | a a b b b



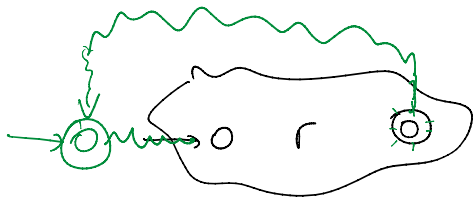
$\rightarrow 0 \xrightarrow{a} 0$
 1 2

: If token in 1 & next char is a,

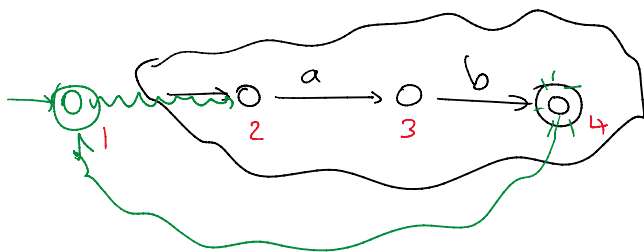
move token to 2



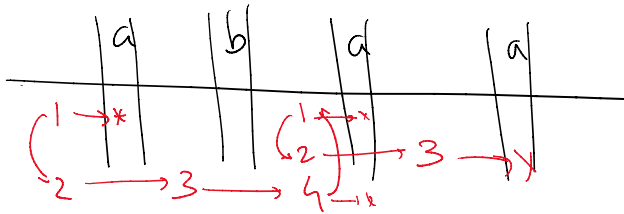
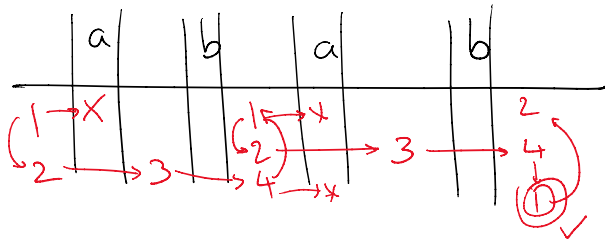
④ Kleene $*$, r^*



Ex: Pattern $(ab)^*$



|a| |b| |a| |b|



Machines with states & arrows : automata

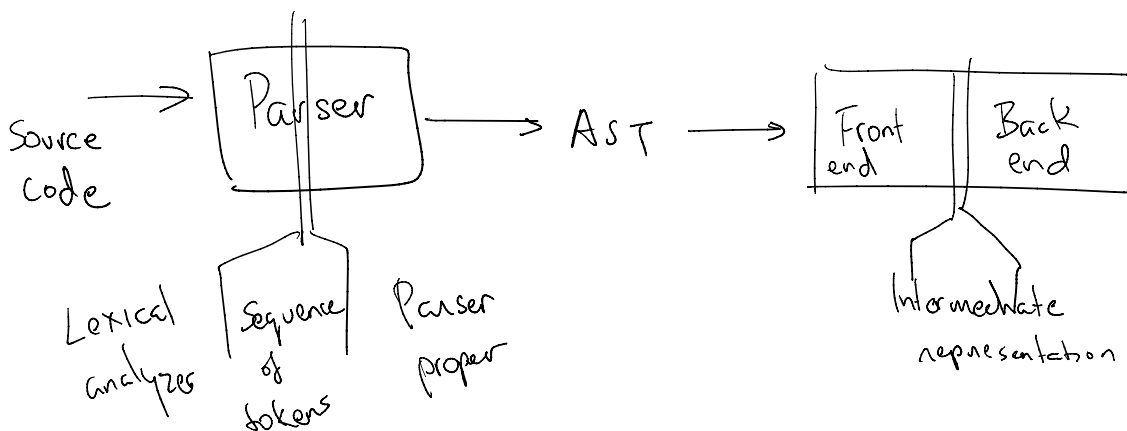
findely many states : finite automata

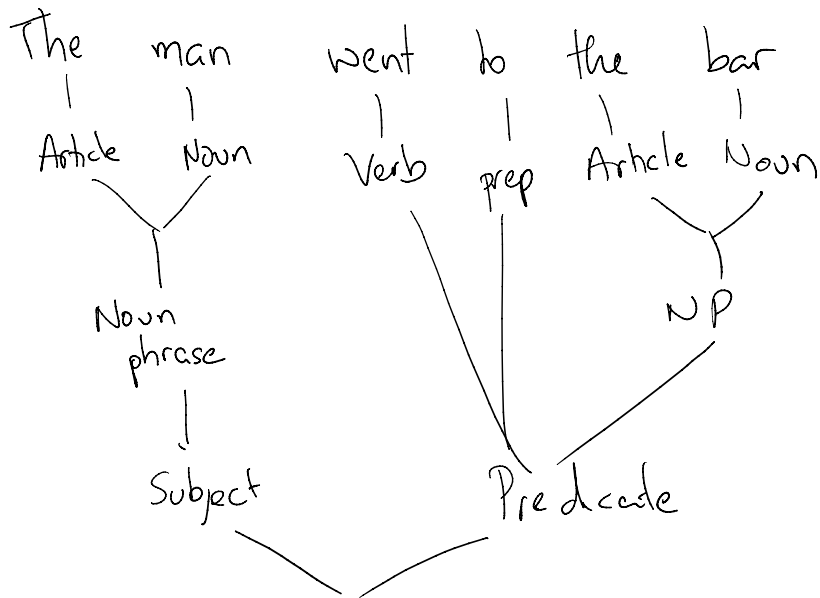
multiple tokens
simultaneously in
play

: non-deterministic

Subset
Construction

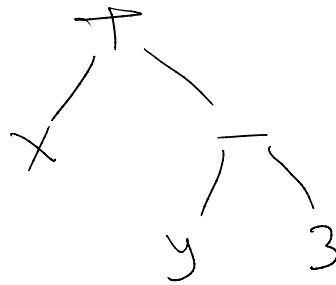
deterministic





" $x + (y - 3)$ "

→



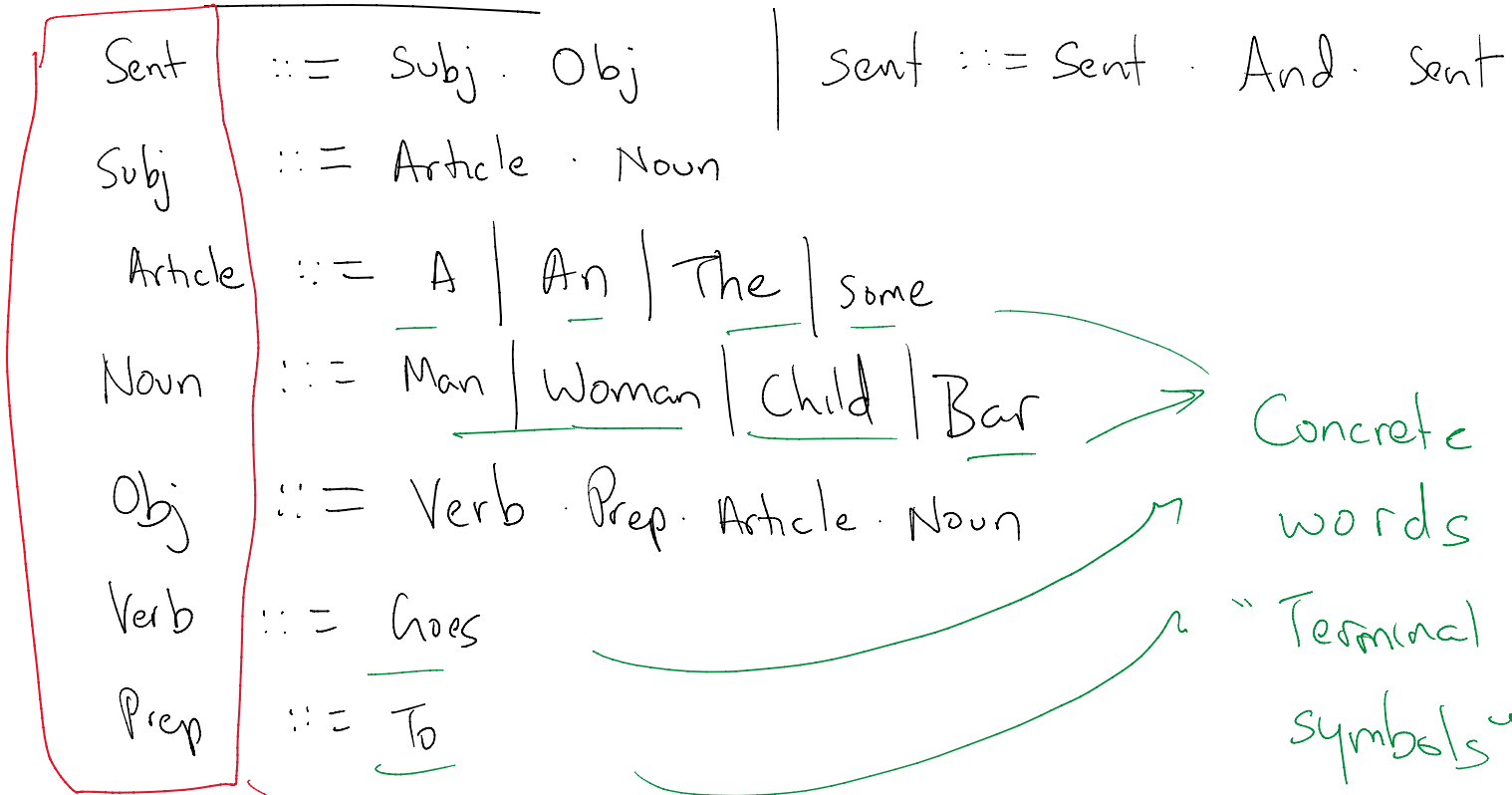
Problem: Given a sequence of tokens, w

& a context-free grammar, G

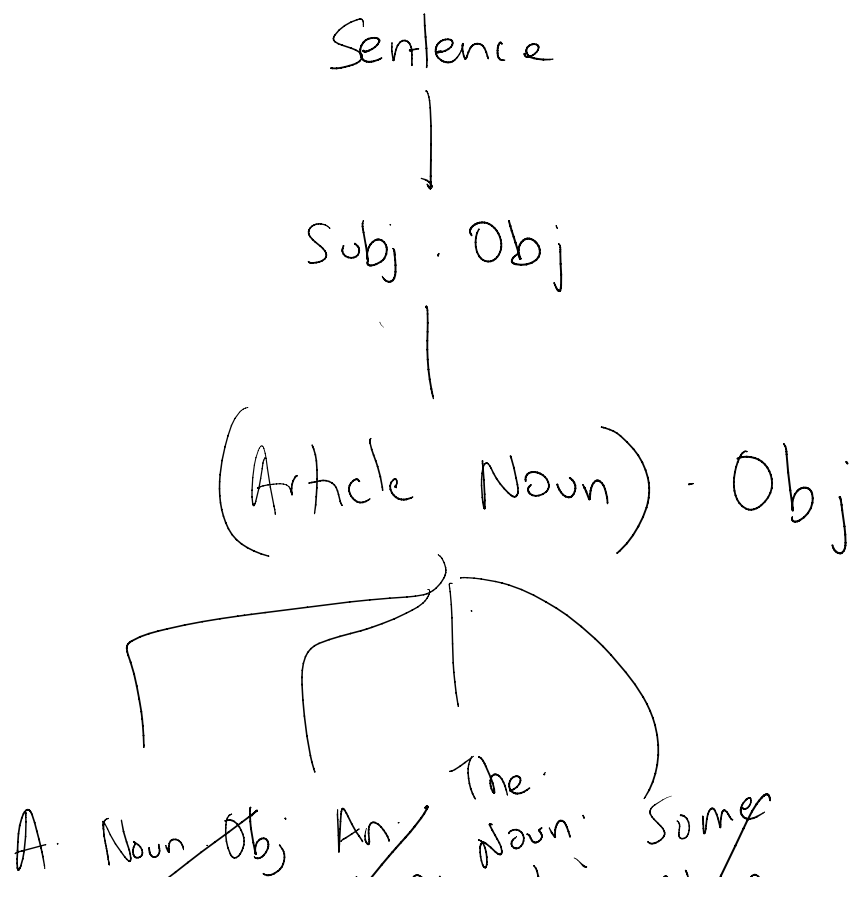
(a) Find if w is a legal production according to G .

(b) If it is, construct the parse tree.

G



The Man Goes To The Bar → Non-terminal symbols



A. ~~Noun Obj~~ ~~An. noun. Obj~~ ~~the noun. Obj~~ ~~Some noun. Obj~~

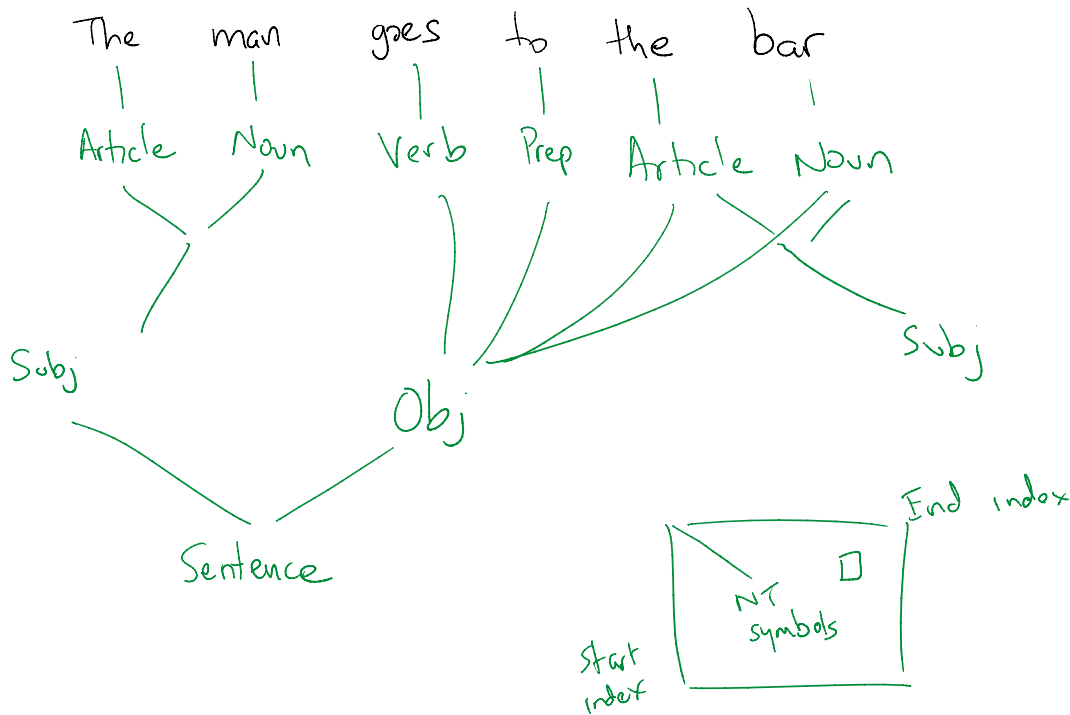
SLD decomposition

Algorithm using Dynamic Programming

For each substring w

for each non-terminal symbol S

in memoization table, note if w matches S .



$O(n^3 |G|)$ time

Cocke - Young - Kasami (CYK) parser

LL(k) , LR(D) , LALR (*),

...