

CONTEXT-FREE

LANGUAGES

Recall: A context-free grammar is a 4-tuple  $G = (NT, T, R, S)$ ,  
where

$NT$ : finite set of non-terminal symbols

$T$ : finite set of terminal symbols

$R$ : finite set of production rules, each rule of the form

$$X ::= \alpha$$

$X \in NT$        $\alpha \in (NT \cup T)^*$

only a single  $X$

$S$ : start symbol,  $S \in NT$

$$L(G) = \{ \omega \mid S ::= \omega \text{ via an application of a finite sequence of rules in } R \}$$

Any language  $\mathcal{L}$  s.t.  $\mathcal{L} = \mathcal{L}(G)$  for some CFG  $G$  is called a **context-free language**.

$\mathcal{L}_{ab} = \{ \omega \mid \omega \text{ has an equal number of 'a's and 'b's} \}$  is context-free as is the language of balanced parentheses.

We said that the main application of CFGs is in parsing. We decide whether or not a string is well-formed by checking if there is some sequence of rules which generates it.

Consider  $L_{ab} = \{w \mid w \text{ has an equal number of 'a's and 'b's}\}$

$L_{ab}$  is generated by the grammar

$$S ::= \epsilon \mid aSb \mid bSa \mid SS$$

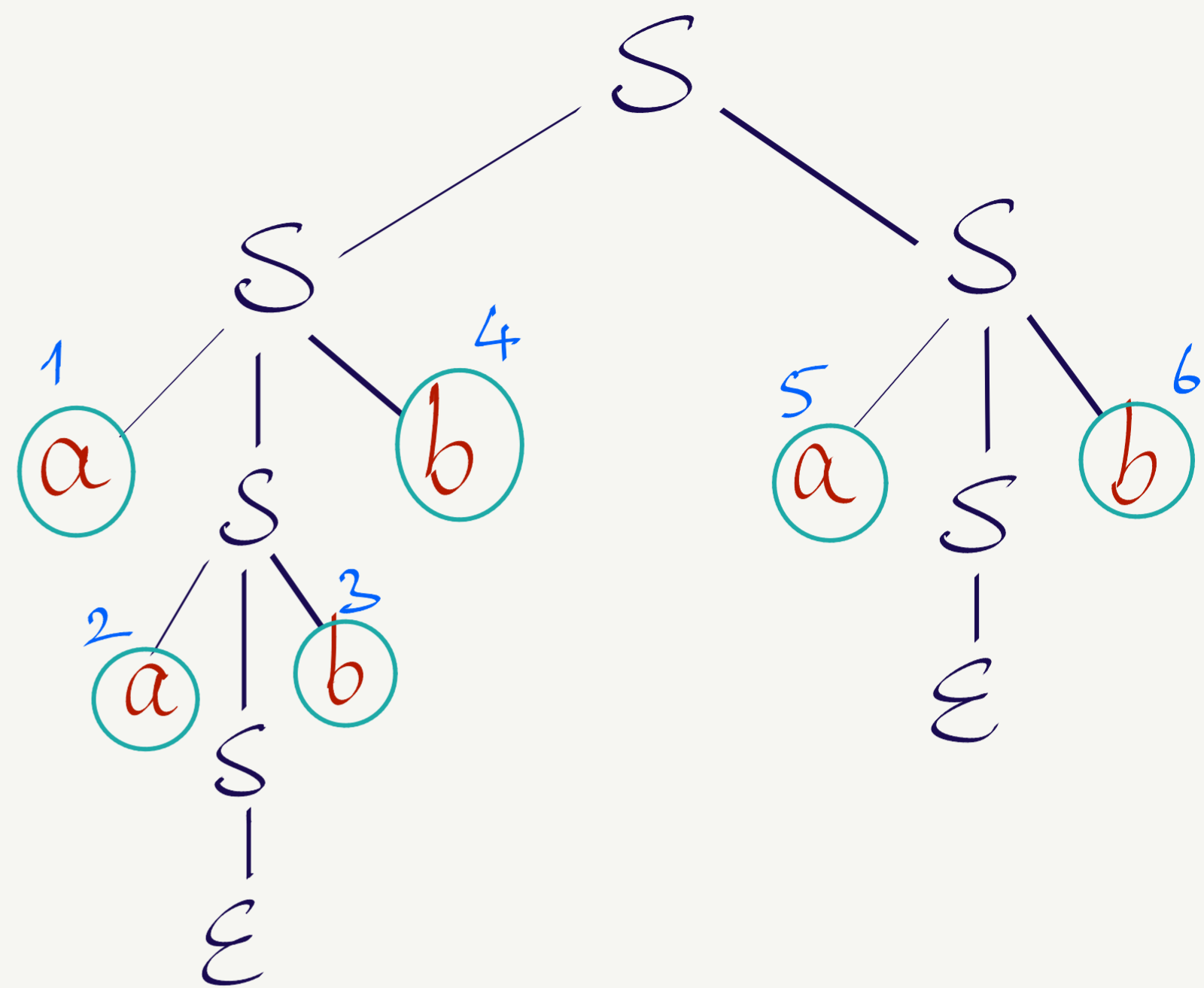
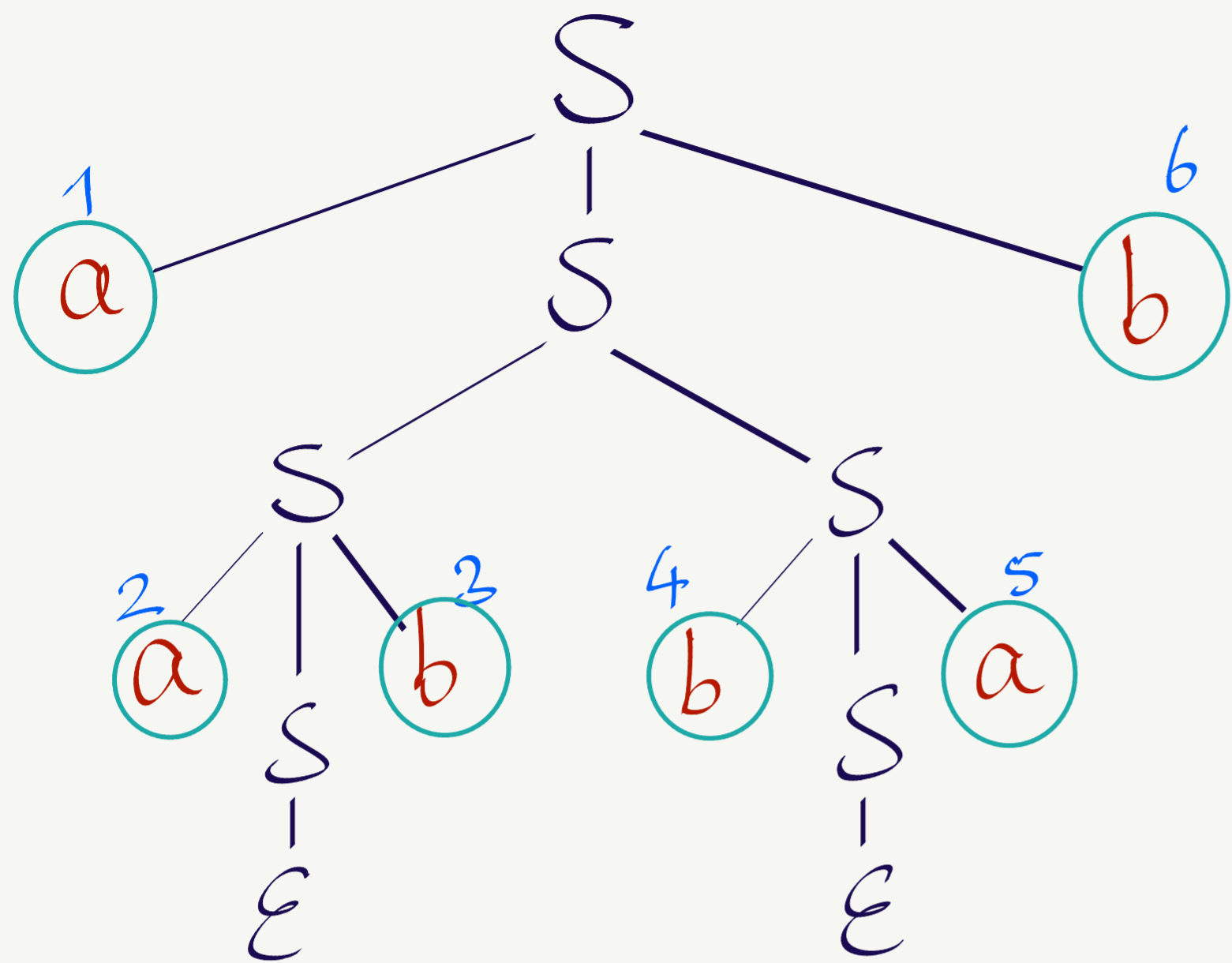
Consider  $aabbab \in L_{ab}$ . How might this grammar generate it?

Have to guess the rule which was applied to get this string,  
and continue recursively!

So suppose we got  $aabbab$  by using  $S ::= aSb$ .

Now I need to check if the grammar can generate  $abba$  etc.

The easiest way to keep track of this is via **parse trees**.



A grammar which can generate multiple parse trees for a string is called *ambiguous* (otherwise, *unambiguous*)

The grammar  $S ::= \epsilon \mid aSbS \mid bSaS$  is also ambiguous

To remove ambiguity, one must somehow ensure that matches are unique.

$$S ::= \epsilon \mid aBS \mid \underline{b}AS$$

match the first 'b' against this 'a', then the rest

$$B ::= b \mid aBB$$

needs to match two 'b's against two already-read 'a's

$$A ::= a \mid bAA$$

needs to match two 'a's against two already-read 'b's

Proving that a grammar is unambiguous can be difficult!

Can sometimes do induction on strings in the language, but not always!

Much like we provided a machine model for regexes via DFAs/NFAs, we would like a machine model for CFGs as well.

We said that DFAs cannot count, so there was no DFA for  $L_{ab}$ , because recognizing  $L_{ab}$ , intuitively, required the machine to

- count # 'a's
- count # 'b's
- check that these numbers were equal.

What is a small extension we can do to a DFA/NFA so it can recognize  $L_{ab}$ ?

Suppose we add a counter  $ctr$  which can increment by one, decrement by one, and check whether it is zero.

Initially:  $ctr = 0$

If you see an 'a', increment  $ctr$

If you see a 'b', decrement  $ctr$  (cannot do this if  $ctr = 0$ )

If this machine has an accepting run on a string  $w$ ,  $w$  has as many 'a's as 'b's.

Exercise: Try to formalize this as a new kind of finite automaton.

Each state needs to track the value of the counter

Cannot decrement at a state if counter is zero.



For a language like  $L_{pal} = \{w \cdot \text{rev}(w) \mid w \in \Sigma^*\}$ ,

it is not clear how to use a single counter to help recognize it.

But what we want is that if we somehow guess the end of  $w$ ,  
the letter we read *last* in  $w$  should also be  
the letter we read *first* in whatever follows,

and that this inside-out matching continues till the end.

A *stack* could help us keep track of this!