

CS21

Decidability and Tractability

Lecture 5

January 17, 2018

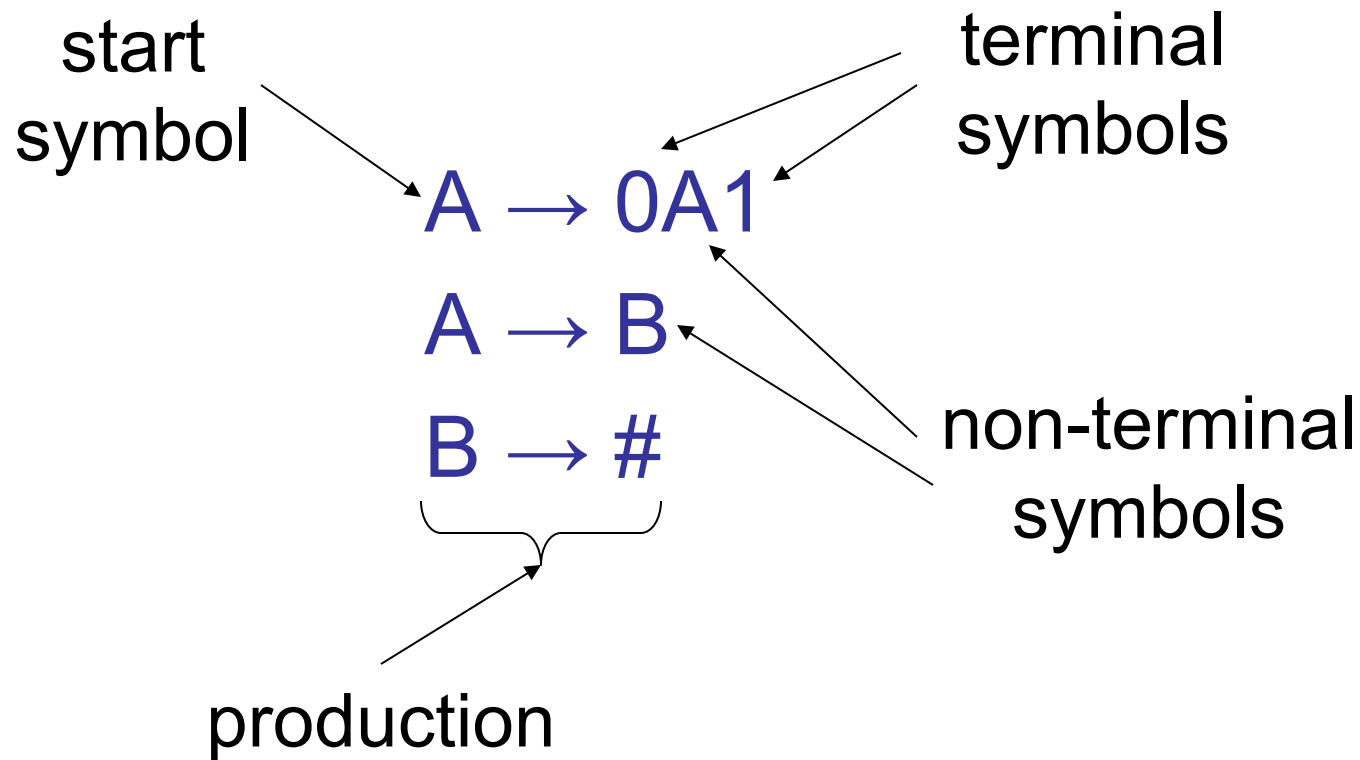
Outline

- Context-Free Grammars and Languages
- equivalence of NPDAs and CFGs
- non context-free languages

Context-free grammars and languages

- languages recognized by a (N)FA are exactly the languages described by regular expressions, and they are called the regular languages
- languages recognized by a NPDA are exactly the languages described by context-free grammars, and they are called the context-free languages

Context-Free Grammars



Context-Free Grammars

- generate strings by repeated replacement of **non-terminals** with **string of terminals and non-terminals**
 - write down start symbol (non-terminal)
 - replace a non-terminal with the right-hand-side of a rule that has that non-terminal as its left-hand-side.
 - repeat above until no more non-terminals

Context-Free Grammars

Example:

$A \Rightarrow 0A1 \Rightarrow 00A11 \Rightarrow$
 $000A111 \Rightarrow 000B111 \Rightarrow$
 $000\#111$

$A \rightarrow 0A1$

$A \rightarrow B$

$B \rightarrow \#$

- a **derivation** of the string 000#111
- set of all strings generated in this way is the **language of the grammar** $L(G)$
- called a **Context-Free Language**

Context-Free Grammars

- Natural languages (e.g. English) structure:

<sentence> → <noun-phrase><verb-phrase>

<noun-phrase> → <cpx-noun> | <cpx-noun><prep-phrase>

<verb-phrase> → <cpx-verb> | <cpx-verb><prep-phrase>

<prep-phrase> → <prep><cpx-noun>

<cpx-noun> → <article><noun>

<cpx-verb> → <verb> | <verb><noun-phrase>

<article> → a | the

<noun> → dog | cat | flower

<verb> → eats | sees

<prep> → with

shorthand for
multiple rules
with same lhs

Generate a string in
the language of this
grammar.

Context-Free Grammars

- CFGs don't capture natural languages completely
- computer languages often **defined** by CFG
 - hierarchical structure
 - slightly different notation often used “Backus-Naur form”
 - see next slide for example

Example CFG

$\langle \text{stmt} \rangle \rightarrow \langle \text{if-stmt} \rangle \mid \langle \text{while-stmt} \rangle \mid \langle \text{begin-stmt} \rangle$
 $\mid \langle \text{asgn-stmt} \rangle$
 $\langle \text{if-stmt} \rangle \rightarrow \text{IF } \langle \text{bool-expr} \rangle \text{ THEN } \langle \text{stmt} \rangle \text{ ELSE } \langle \text{stmt} \rangle$
 $\langle \text{while-stmt} \rangle \rightarrow \text{WHILE } \langle \text{bool-expr} \rangle \text{ DO } \langle \text{stmt} \rangle$
 $\langle \text{begin-stmt} \rangle \rightarrow \text{BEGIN } \langle \text{stmt-list} \rangle \text{ END}$
 $\langle \text{stmt-list} \rangle \rightarrow \langle \text{stmt} \rangle \mid \langle \text{stmt} \rangle ; \langle \text{stmt-list} \rangle$
 $\langle \text{asgn-stmt} \rangle \rightarrow \langle \text{var} \rangle := \langle \text{arith-expr} \rangle$
 $\langle \text{bool-expr} \rangle \rightarrow \langle \text{arith-expr} \rangle \langle \text{compare-op} \rangle \langle \text{arith-expr} \rangle$
 $\langle \text{compare-op} \rangle \rightarrow < \mid > \mid \leq \mid \geq \mid =$
 $\langle \text{arith-expr} \rangle \rightarrow \langle \text{var} \rangle \mid \langle \text{const} \rangle$
 $\mid (\langle \text{arith-expr} \rangle \langle \text{arith-op} \rangle \langle \text{arith-expr} \rangle)$
 $\langle \text{arith-op} \rangle \rightarrow + \mid - \mid * \mid /$
 $\langle \text{const} \rangle \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$
 $\langle \text{var} \rangle \rightarrow a \mid b \mid c \mid \dots \mid x \mid y \mid z$

CFG formal definition

- A **context-free grammar** is a 4-tuple

$$(V, \Sigma, R, S)$$

where

- V is a finite set called the **non-terminals**
- Σ is a finite set (disjoint from V) called the **terminals**
- R is a finite set of **productions** where each production is a non-terminal and a string of terminals and non-terminals.
- $S \in V$ is the **start variable** (or start non-terminal)

CFG formal definition

- u, v, w are strings of non-terminals and terminals, and $A \rightarrow w$ is a production:

“ uAv yields uwv ” notation: $uAv \Rightarrow uwv$

also: “yields in 1 step” notation: $uAv \Rightarrow^1 uwv$

- in general:

“yields in k steps” notation: $u \Rightarrow^k v$

– meaning: there exists strings u_1, u_2, \dots, u_{k-1} for which $u \Rightarrow u_1 \Rightarrow u_2 \Rightarrow \dots \Rightarrow u_{k-1} \Rightarrow v$

CFG formal definition

- notation: $u \Rightarrow^* v$
 - meaning: $\exists k \geq 0$ and strings u_1, \dots, u_{k-1} for which $u \Rightarrow u_1 \Rightarrow u_2 \Rightarrow \dots \Rightarrow u_{k-1} \Rightarrow v$
- if $u =$ start symbol, this is a **derivation of v**
- The **language of G** , denoted $L(G)$ is:

$$\{w \in \Sigma^* : S \Rightarrow^* w\}$$

CFG example

- Balanced parentheses:
 - $()$
 - $((()((()())))$
- a string w in $\Sigma^* = \{ (,) \}^*$ is balanced iff:
 - # “(”s equals # “)”s, and
 - for any prefix of w , # “(”s \geq # “)”s

Exercise: design a CFG for balanced parentheses.

CFG example

- Arithmetic expressions over $\{+, *, (,), a\}$

– $(a + a) * a$

– $a * a + a + a + a + a$

- A CFG generating this language:

$\langle \text{expr} \rangle \rightarrow \langle \text{expr} \rangle * \langle \text{expr} \rangle$

$\langle \text{expr} \rangle \rightarrow \langle \text{expr} \rangle + \langle \text{expr} \rangle$

$\langle \text{expr} \rangle \rightarrow (\langle \text{expr} \rangle) \mid a$

CFG example

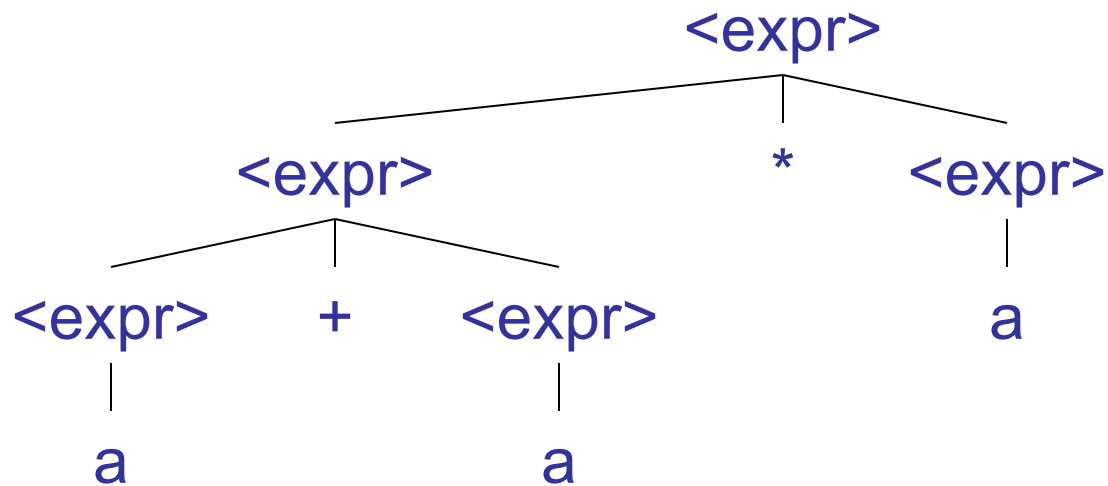
$$\begin{aligned} \langle \text{expr} \rangle &\rightarrow \langle \text{expr} \rangle * \langle \text{expr} \rangle \\ \langle \text{expr} \rangle &\rightarrow \langle \text{expr} \rangle + \langle \text{expr} \rangle \\ \langle \text{expr} \rangle &\rightarrow (\langle \text{expr} \rangle) \mid a \end{aligned}$$

- A derivation of the string: **a+a*a**

$$\begin{aligned} \langle \text{expr} \rangle &\Rightarrow \langle \text{expr} \rangle * \langle \text{expr} \rangle \\ &\Rightarrow \langle \text{expr} \rangle + \langle \text{expr} \rangle * \langle \text{expr} \rangle \\ &\Rightarrow a + \langle \text{expr} \rangle * \langle \text{expr} \rangle \\ &\Rightarrow a + a * \langle \text{expr} \rangle \\ &\Rightarrow a + a * a \end{aligned}$$

Parse Trees

- Easier way to picture derivation: **parse tree**



- grammar encodes grouping information; this is captured in the parse tree.

CFGs and parse trees

```
<expr> → <expr> * <expr>  
<expr> → <expr> + <expr>  
<expr> → (<expr>) | a
```

- Is this a good grammar for arithmetic expressions?
 - can group wrong way (+ precedence over *)

Solution to problem

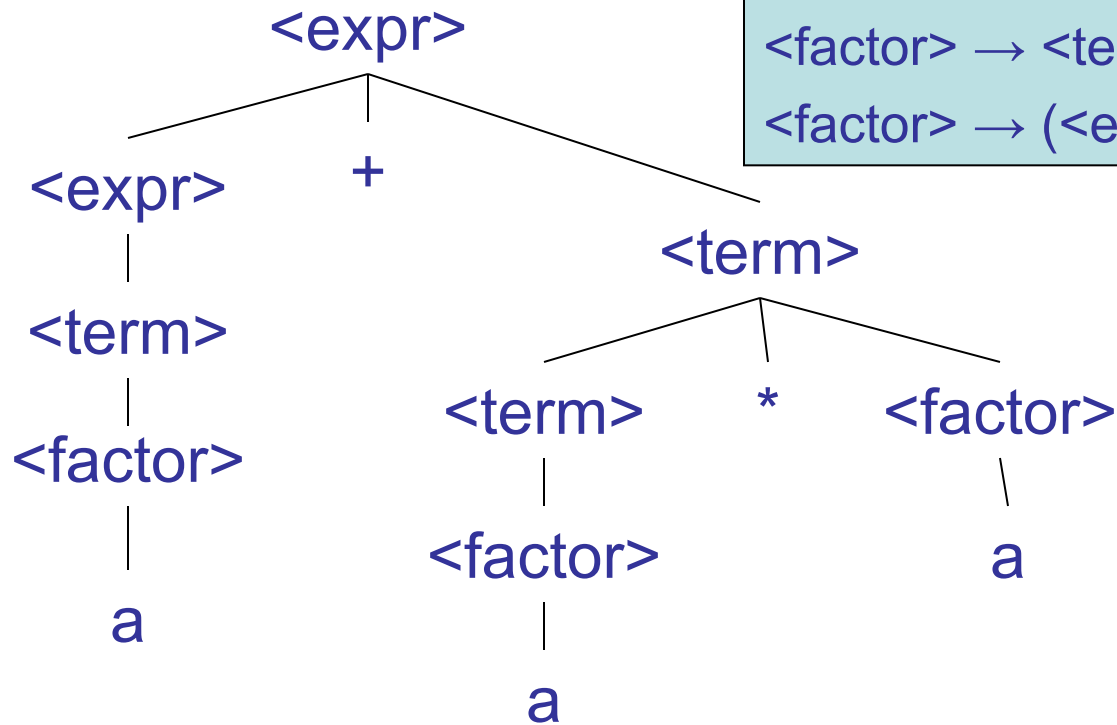
```
<expr> → <expr> + <term> | <term>  
<term> → <term> * <factor> | <factor>  
<factor> → <term> * <factor>  
<factor> → (<expr>) | a
```

- forces correct precedence in parse tree grouping
 - within parentheses, * cannot occur as ancestor of + in the parse tree.

Parse Trees

- parse tree for $a + a * a$ in new grammar:

$\langle \text{expr} \rangle \rightarrow \langle \text{expr} \rangle + \langle \text{term} \rangle \mid \langle \text{term} \rangle$
 $\langle \text{term} \rangle \rightarrow \langle \text{term} \rangle * \langle \text{factor} \rangle \mid \langle \text{factor} \rangle$
 $\langle \text{factor} \rangle \rightarrow \langle \text{term} \rangle * \langle \text{factor} \rangle$
 $\langle \text{factor} \rangle \rightarrow (\langle \text{expr} \rangle) \mid a$



Some facts about CFLs

- CFLs are closed under
 - union (proof?)
 - concatenation (proof?)
 - star (proof?)
- Every regular language is a CFL
 - proof?

NPDA, CFG equivalence

Theorem: a language L is recognized by a NPDA *iff* L is described by a CFG.

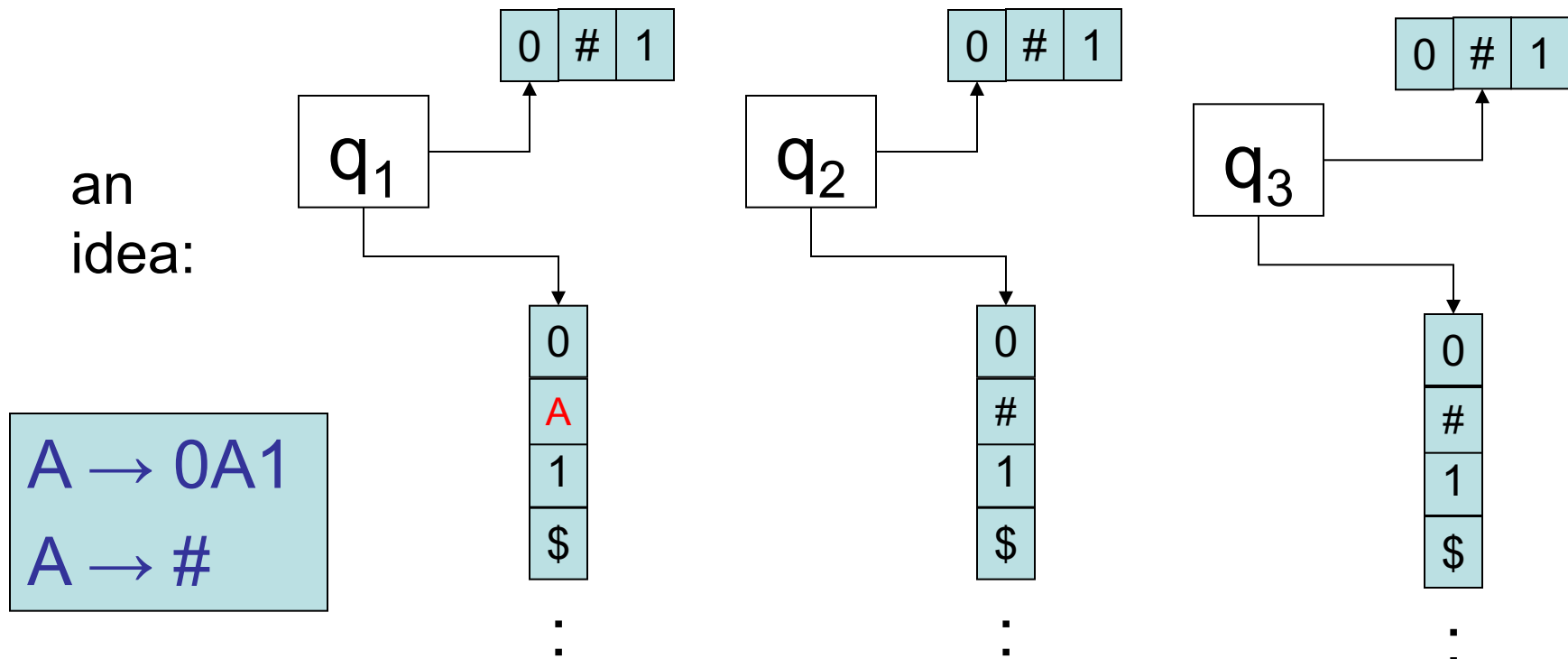
Must prove *two* directions:

(\Rightarrow) L is recognized by a NPDA *implies* L is described by a CFG.

(\Leftarrow) L is described by a CFG *implies* L is recognized by a NPDA.

NPDA, CFG equivalence

Proof of (\Leftarrow): L is described by a CFG
implies L is recognized by a NPDA.



NPDA, CFG equivalence

1. we'd like to non-deterministically guess the derivation, forming it on the stack
2. then scan the input, popping matching symbol off the stack at each step
3. accept if we get to the bottom of the stack at the end of the input.

what is wrong with this approach?