# Compiler Construction

## Lecture 3: Context-free grammars

Jeremy Yallop

jeremy.yallop@cl.cam.ac.uk

Lent 2023

Parsing
● ○

CFGs

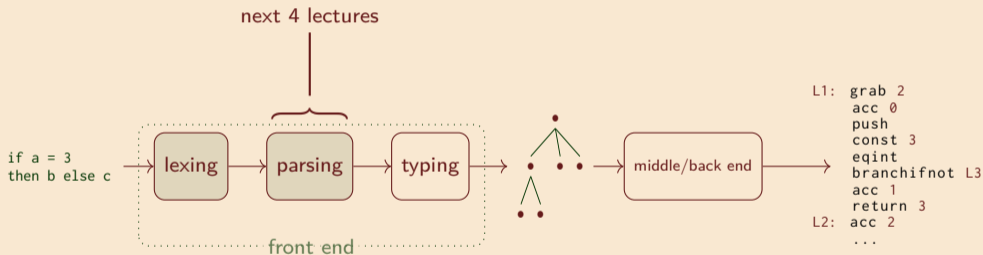Derivations

PDAs

Ambiguity

Top-down &
bottom-up

next 4 lectures

```
if a = 3
then b else c
```

lexing → parsing → typing

front end

middle/back end

```
L1: grab 2
    acc 0
    push
    const 3
    eqint
    branchifnot L3
    acc 1
    return 3
L2: acc 2
    ...
```

**Parsing**

CFGs

Derivations

PDAs

Ambiguity

Top-down &
bottom-up

A small fragment of the C standard:

---

**6.7 Declarations**
**Syntax**

*declaration:*
  *declaration-specifiers init-declarator-list$_{opt}$ ;*
  *static-assert-declaration*

*declaration-specifiers:*
  *storage-class-specifier declaration-specifiers$_{opt}$*
  *type-specifier declaration-specifiers$_{opt}$*
  *type-qualifier declaration-specifiers$_{opt}$*
  *function-specifier declaration-specifiers$_{opt}$*
  *alignment-specifier declaration-specifiers$_{opt}$*

*init-declarator-list:*
  *init-declarator*
  *init-declarator-list , init-declarator*

*init-declarator:*
  *declarator*
  *declarator = initializer*

---

Today's Q: how can we turn this declarative specification into a program?

# Context-free grammars

nonterminals $N$       productions $P \subseteq N \times (N \cup T)*$

$$M = \langle N, T, P, S \rangle$$

terminals $T$       start symbol $S \in N$

Each $\langle A, \alpha \rangle \in P$ is written as $A \rightarrow \alpha$

$$G_1 = \langle N_1, T_1, P_1, E \rangle$$

*where*
$$N_1 = \{E\}$$
$$T_1 = \{+, *, (, ), \mathsf{id}\}$$

$$P_1 = E \rightarrow \begin{array}{l} E + E \\ \mid E * E \\ \mid (E) \\ \mid \mathsf{id} \end{array}$$

NB: $P_1$ definition is shorthand for

$$P_1 = \{\langle E, E + E \rangle, \langle E, E * E \rangle, \langle E, (E) \rangle, \langle E, \mathsf{id} \rangle\}$$

# Derivations

Notation conventions:

$$\alpha, \beta, \gamma \ldots \in (N \cup T)*$$
$$A, B, C, \ldots \in N$$

Given: $\alpha A \beta$ and a production $A \to \gamma$ a derivation step is written as

$$\alpha AB \Rightarrow \alpha \gamma \beta$$

$\Rightarrow^+$ means one or more derivation steps

$\Rightarrow^*$ means zero or more derivation steps.

A **leftmost** derivation

$$
\begin{aligned}
E &\Rightarrow E*E \\
&\Rightarrow (E)*E \\
&\Rightarrow (E+E)*E \\
&\Rightarrow (x+E)*E \\
&\Rightarrow (x+y)*E \\
&\Rightarrow (x+y)*(E) \\
&\Rightarrow (x+y)*(E+E) \\
&\Rightarrow (x+y)*(z+E) \\
&\Rightarrow (x+y)*(z+x)
\end{aligned}
$$

A **rightmost** derivation

$$
\begin{aligned}
E &\Rightarrow E*E \\
&\Rightarrow E*(E) \\
&\Rightarrow E*(E+E) \\
&\Rightarrow E*(E+x) \\
&\Rightarrow E*(z+x) \\
&\Rightarrow (E)*(z+x) \\
&\Rightarrow (E+E)*(z+x) \\
&\Rightarrow (E+y)*(z+x) \\
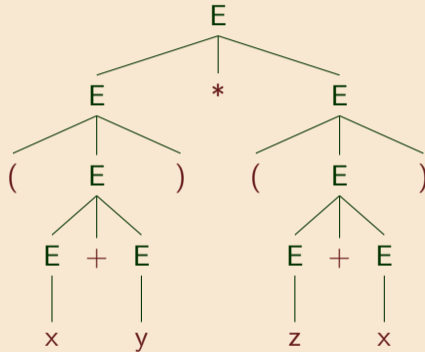&\Rightarrow (x+y)*(z+x)
\end{aligned}
$$

Parsing

CFGs

**Derivations**
●●●○○

PDAs

Ambiguity

Top-down &
bottom-up

```
                          E
                ┌─────────┼─────────┐
                E         *         E
           ┌────┼────┐         ┌────┼────┐
          (    E    )         (    E    )
             ┌──┼──┐             ┌──┼──┐
             E  +  E             E  +  E
             │     │             │     │
             x     y             z     x
```

The derivation tree for (x+y) * (z+x). All derivations of this expression will
produce the same derivation tree.

Parsing
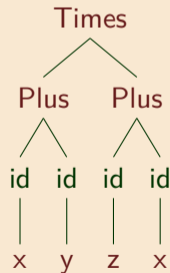
CFGs

**Derivations**
● ● ● ● ○

PDAs

Ambiguity

Top-down &
bottom-up

(Terminology: $=$ **parse** tree
$=$ **derivation** tree
$=$ **concrete syntax** tree)

An **abstract syntax** tree contains only the information needed to generate an intermediate representation

$L(G)$: the language generated by $G$

$$L(G) \;=\; \{w \in T* \mid S \Rightarrow^+ w\}$$

For example, if $G$ has productions

$$S \rightarrow aSb \mid \epsilon$$

then

$$L(G) \;=\; \{a^n b^n \mid n \geq 0\}$$
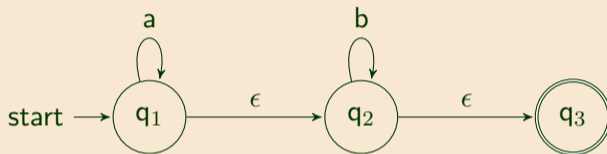
So CFGs can capture more than regular languages.

# Pushdown automata

Regular languages are accepted by finite automata:

$a^*b^*$



Context-free languages are accepted by pushdown automata, finite automata augmented with stacks.

$a^n b^n$

stack symbols $\Gamma$

states $Q$

start state $q_0 \in Q$

$$M = \langle Q, \Sigma, \Gamma, \delta, q_0, Z \rangle$$

alphabet $T$

initial stack symbol $z \in \Gamma S$

transitions $\delta$:
$\forall q \in Q,$
$\quad a \in (\Sigma \cup \{\epsilon\}),$
$\quad\quad X \in \Gamma,$
$\quad\quad\quad \delta(q, a, X) \subseteq Q \times \Gamma^*$

**Pushdown Automata (PDAs)**

Parsing

CFGs

Derivations

PDAs
● ○ ○

Ambiguity

Top-down &
bottom-up

$\langle q', \beta \rangle \in \delta(q, a, X)$ means:

When the machine is $\begin{cases} \text{in state } q, \text{ and} \\ \text{reading } a \text{ and} \\ \text{with } X \text{ on top of the stack,} \end{cases}$

it can $\begin{cases} \text{move to state } q' \text{ and} \\ \text{replace } X \text{ with } \beta. \end{cases}$

i.e. it *pops* $X$ from the
stack and *pushes* $\beta$.

For $q \in Q, w \in \Sigma^*, \alpha \in \Gamma^*$, $\langle q, w, \alpha \rangle$ is called an **instantaneous description** (ID).
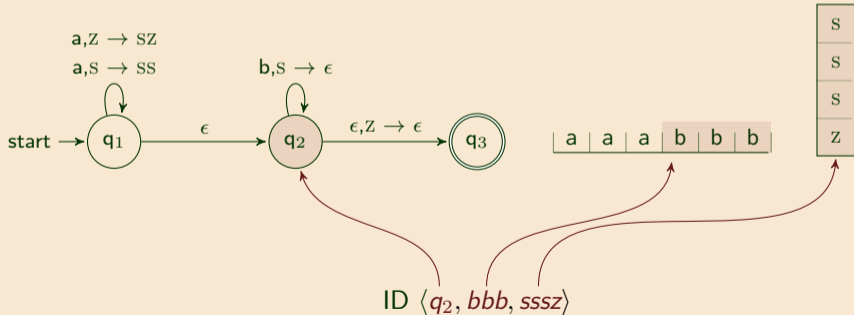
It denotes the PDA looking at the first symbol of $w$
with $\alpha$ on the stack



ID $\langle q_2, bbb, sssz \rangle$

For $\langle q', \beta \rangle \in \delta(q, a, X), a \in \Sigma$, define the relation $\rightarrow$ on IDs as

$$\langle q, aw, X\alpha \rangle \rightarrow \langle q', w, \beta\alpha \rangle$$

and for $\langle q', \beta \rangle \in \delta(q, \epsilon, X)$ as

$$\langle q, w, X\alpha \rangle \rightarrow \langle q', w, \beta\alpha \rangle$$

Then the **language accepted by** $M$, $L(M)$, is:

$$L(M) = \{ w \in \Sigma* \mid \exists q \in Q, \langle q_0, w, Z \rangle \rightarrow^+ \langle q, \epsilon, \epsilon \rangle \}$$
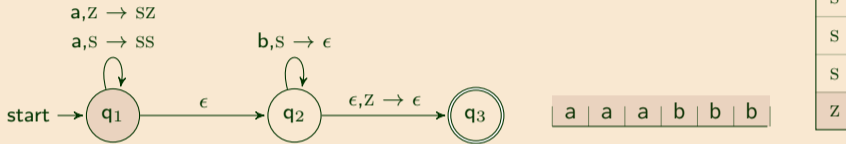
;



$\langle q_1, \text{aaabbb}, \text{Z} \rangle$

;



a,$Z \to SZ$
a,$S \to SS$

b,$S \to \epsilon$

start $\to$ $q_1$ $\xrightarrow{\epsilon}$ $q_2$ $\xrightarrow{\epsilon,Z \to \epsilon}$ $q_3$

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

| S |
| S |
| S |
| Z |

$\langle q_1, \textit{aaabbb}, Z \rangle$
$\langle q_1, \textit{aabbb}, SZ \rangle$

;



$\langle q_1, \textit{aaabbb}, \text{z} \rangle$
$\langle q_1, \textit{aabbb}, \text{sz} \rangle$
$\langle q_1, \textit{abbb}, \text{ssz} \rangle$

Parsing

CFGs

Derivations

PDAs
● ● ●

Ambiguity

Top-down &
bottom-up

;



a,$z \to sz$
a,$s \to ss$          b,$s \to \epsilon$

start → $q_1$ —$\epsilon$→ $q_2$ —$\epsilon$,$z \to \epsilon$→ $q_3$

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

| s |
|---|
| s |
| s |
| z |

$\langle q_1, \textit{aaabbb}, z \rangle$
$\langle q_1, \textit{aabbb}, sz \rangle$
$\langle q_1, \textit{abbb}, ssz \rangle$
$\langle q_2, \textit{bbb}, sssz \rangle$

;



$\langle q_1, \textit{aaabbb}, \textsc{z} \rangle$
$\langle q_1, \textit{aabbb}, \textsc{sz} \rangle$
$\langle q_1, \textit{abbb}, \textsc{ssz} \rangle$
$\langle q_2, \textit{bbb}, \textsc{sssz} \rangle$
$\langle q_2, \textit{bb}, \textsc{ssz} \rangle$

;



$\langle q_1, \textit{aaabbb}, \text{Z} \rangle$
$\langle q_1, \textit{aabbb}, \text{SZ} \rangle$
$\langle q_1, \textit{abbb}, \text{SSZ} \rangle$
$\langle q_2, \textit{bbb}, \text{SSSZ} \rangle$
$\langle q_2, \textit{bb}, \text{SSZ} \rangle$
$\langle q_2, \textit{b}, \text{SZ} \rangle$

;



$\langle q_1, \textit{aaabbb}, \textsc{z} \rangle$
$\langle q_1, \textit{aabbb}, \textsc{sz} \rangle$
$\langle q_1, \textit{abbb}, \textsc{ssz} \rangle$
$\langle q_2, \textit{bbb}, \textsc{sssz} \rangle$
$\langle q_2, \textit{bb}, \textsc{ssz} \rangle$
$\langle q_2, \textit{b}, \textsc{sz} \rangle$
$\langle q_2, \epsilon, \textsc{z} \rangle$

;



| s |
|---|
| s |
| s |
| z |

a,$\text{z} \rightarrow \text{sz}$
a,$\text{s} \rightarrow \text{ss}$      b,$\text{s} \rightarrow \epsilon$

start $\rightarrow$ $q_1$ $\xrightarrow{\epsilon}$ $q_2$ $\xrightarrow{\epsilon,\text{z} \rightarrow \epsilon}$ $q_3$

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

$\langle q_1, aaabbb, \text{z} \rangle$
$\langle q_1, aabbb, \text{sz} \rangle$
$\langle q_1, abbb, \text{ssz} \rangle$
$\langle q_2, bbb, \text{sssz} \rangle$
$\langle q_2, bb, \text{ssz} \rangle$
$\langle q_2, b, \text{sz} \rangle$
$\langle q_2, \epsilon, \text{z} \rangle$
$\langle q_3, \epsilon, \epsilon \rangle$

PDA and CFG facts:

For every CFG $G$
   there is a PDA $M$
      such that $L(G) = L(M)$

For every PDA $M$
   there is a CFG $G$
      such that $L(G) = L(M)$

Is the parsing problem solved? Given a CFG $G$ we can construct the PDA $M$.

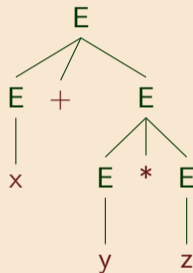No! For programming languages we want $M$ to be **deterministic**

# Ambiguity

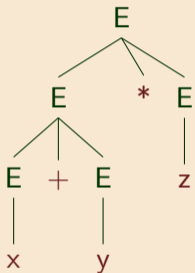Both derivation trees correspond to x + y * z.

But (x + y) * z is not the same as x + (y * z).

**Ambiguity** causes problems going from program texts to derivation trees.

We can often modify the grammar to eliminate ambiguity.

$$G_2 \;=\; \langle N_2, T_1, P_2, E \rangle$$

*where*

$$P_2 \;=\; \begin{array}{rcll} E & \rightarrow & E + T \mid T & \text{(expressions)} \\ T & \rightarrow & T * F \mid F & \text{(terms)} \\ F & \rightarrow & (E) \mid id & \text{(factors)} \end{array}$$

(Can you prove that $L(G_1) = L(G_2)$?)

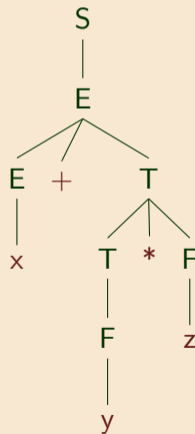The modified grammar eliminates ambiguity. The following is now the unique derivation tree for x + y * z:

Parsing

CFGs

Derivations

PDAs

**Ambiguity**
● ● ● ●

Top-down &
bottom-up

1. Some context-free languages are inherently ambiguous — every CFG for them is ambiguous. For example

$$
\begin{aligned}
L \ = \ & \{a^n b^n c^m d^m \mid m \geq 1, n \geq 1\} \\
\cup \ & \{a^n b^m c^m d^n \mid m \geq 1, n \geq 1\}
\end{aligned}
$$

2. Checking for ambiguity in an arbitrary CFG is not decidable.
3. Given two grammars $G1$ and $G2$, checking $L(G1) = L(G2)$ is not decidable.

(See Hopcroft & Ullman, "Introduction to Automata Theory, Languages, and Computation")

# Top-down & bottom-up

**Top-down**: attempts a left-most derivation. We'll look at two techniques:

Recursive          Predictive
descent            parsing
(hand coded)   (table driven)

**Bottom-up**: attempts a right-most derivation backwards. We'll look at two techniques:

SLR(1)                 LR(1)
(Simple LR(1))

Bottom-up techniques are strictly more powerful (can parse more grammars)

```
type token =
  ADD | MUL | LPAREN | RPAREN | IDENT of string

let rec
    e toks = e' (t toks)
and e' = function
  | ADD :: toks → e' (t toks)
  | toks           → toks (* ε *)
and t toks = t' (f toks)
and t' = function
  | MUL :: toks → t' (f toks)
  | toks           → toks (* ε *)
and f = function
  | LPAREN :: toks →
    (match e toks with
    | RPAREN :: toks → toks
    | _              → failwith "RPAREN")
  | IDENT _ :: toks → toks
  | _              → failwith "F"
```

$$
\begin{aligned}
E &\rightarrow T E' \\
E' &\rightarrow + T E' \mid \epsilon \\
T &\rightarrow F T' \\
T' &\rightarrow * F T' \mid \epsilon \\
F &\rightarrow ( E ) \mid id
\end{aligned}
$$

Parse corresponds to a left-most derivation constructed in a top-down manner

Recursive descent parsing is not suitable for $G_2$.

Left-recursion $E \rightarrow E + T$ will lead to an infinite loop:

```
let rec
    e toks = match e toks (* loop! *) with
             | PLUS :: toks → ...
```

$$
\begin{array}{rcl}
E & \rightarrow & E + T \mid T \\
T & \rightarrow & T * F \mid F \\
F & \rightarrow & (E) \mid \text{id}
\end{array}
$$

$$G_2 = \langle N_2, T_1, P_2, E \rangle \qquad\qquad G_3 = \langle N_3, T_1, P_3, E \rangle$$

*where*                                    *where*

$\cdots$                                           $\cdots$

$$
\begin{aligned}
P_2 = \quad
\begin{array}{rcl}
E &\rightarrow& E + T \mid T \\
T &\rightarrow& T * F \mid F \\
F &\rightarrow& (E) \mid id
\end{array}
\qquad
P_3 = \quad
\begin{array}{rcl}
E &\rightarrow& T\,E' \\
E' &\rightarrow& +\,T\,E' \mid \epsilon \\
T &\rightarrow& F\,T' \\
T' &\rightarrow& *\,F\,T' \mid \epsilon \\
F &\rightarrow& (\,E\,) \mid id
\end{array}
\end{aligned}
$$

*(Can you prove that $L(G_2) = L(G_3)$?)*

```
let rec
    e toks = e' (t toks)
and e' = function
  | ADD :: toks → e' (t toks)
  | toks         → toks (* ε *)
and t toks = t' (f toks)
and t' = function
  | MUL :: toks → t' (f toks)
  | toks         → toks (* ε *)
and f = function
  | LPAREN :: toks →
    (match e toks with
    | RPAREN :: toks → toks
    | _  → failwith "RPAREN")
  | IDENT _ :: toks → toks
  | _              → failwith "F"
```

Parsing x + y * z, i.e.

```
[IDENT "x";
 PLUS;
 IDENT "y";
 TIMES;
 IDENT "z"]
```

Evaluation trace:

          e toks

⇝    e' (t toks)

⇝    e' (t' (f toks))

⇝    ...

# Next time: LL parsing