# COS320: Compiling Techniques

Zak Kincaid

March 7, 2019
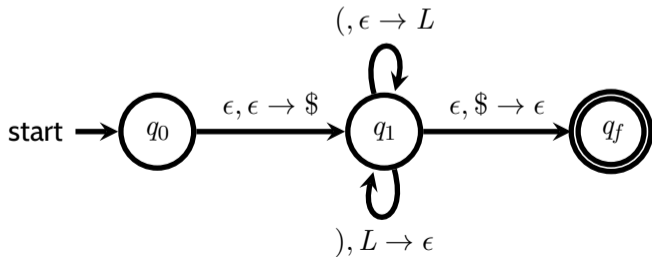
*Parsing II: LL parsing*

# Recall: Context-free grammars

- A *context-free grammar* $G = (N, \Sigma, R, S)$ consists of:
  - $N$: a finite set of *non-terminal symbols*
  - $\Sigma$: a finite alphabet (or set of *terminal symbols*)
  - $R \subseteq N \times (N \cup \Sigma)^*$ a finite set of *rules* or *productions*
  - $S \in N$: the starting non-terminal.
- A *derivation* consists of a finite sequence of words $\gamma_1, ..., \gamma_n \in (N \cup \Sigma)^*$ such that $\gamma_1 = S$ and for each $i$, $\gamma_{i+1}$ is obtained from $\gamma_i$ by replacing a non-terminal symbol with the right-hand-side of one of its rules
- The set of all strings $w \in \Sigma^*$ such that $G$ has a derivation of $w$ is the *language* of $G$, written $\mathcal{L}(G)$.
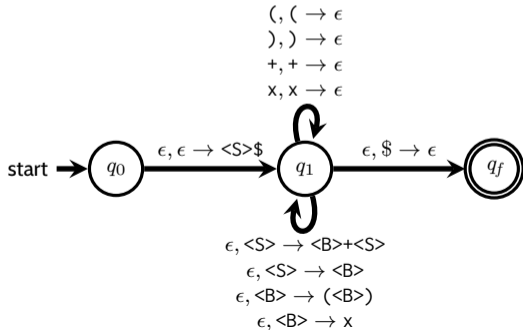
# Parsing

- Context-free grammars are *generative*: easy to find strings that belongs to $\mathcal{L}(G)$, not so easy determine whether a *given* string belongs to $\mathcal{L}(G)$
- *Pushdown automata* (PDA) are a kind of automata that recognize context-free languages
- Pushdown automaton recognizing <S> ::= <S><S> | (<S>) | $\epsilon$:
  - *Stack alphabet*: $\$$ marks bottom of the stack, $L$ marks unbalanced left paren

# Top-down parsing

- Stack represents intermediate state of a derivation, minus the consumed part of the input string.
- Start with $S$ on the stack
- Any time top of the stack is a non-terminal $A$, non-deterministically choose a rule $A ::= \gamma \in R$. Pop $A$ off the stack, and push $\gamma$
- If the top of the stack is a terminal $a$, consume $a$ from the input string and pop $a$ off the stack
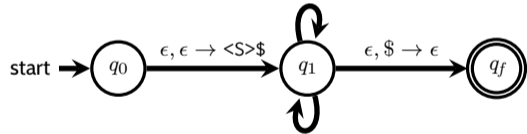- Accept when stack is empty

$$\langle S \rangle ::= \langle B \rangle + \langle S \rangle \mid \langle B \rangle$$
$$\langle B \rangle ::= (\langle S \rangle) \mid x$$



$(, ( \to \epsilon$
$), ) \to \epsilon$
$+, + \to \epsilon$
$x, x \to \epsilon$

start $\to q_0$  $\epsilon, \epsilon \to \langle S \rangle \$ \quad q_1 \quad \epsilon, \$ \to \epsilon \quad q_f$

$\epsilon, \langle S \rangle \to \langle B \rangle + \langle S \rangle$
$\epsilon, \langle S \rangle \to \langle B \rangle$
$\epsilon, \langle B \rangle \to (\langle S \rangle)$
$\epsilon, \langle B \rangle \to x$

<S> ::= <B>+<S> | <B>

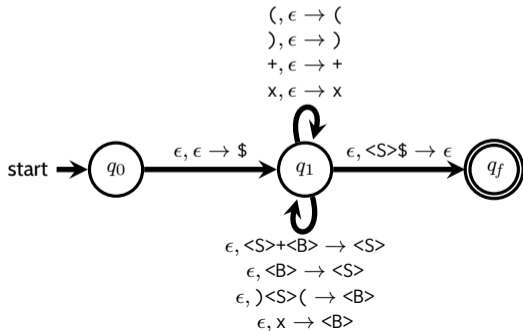<B> ::= (<S>) | x

$(, ( \rightarrow \epsilon$
$), ) \rightarrow \epsilon$
$+, + \rightarrow \epsilon$
$x, x \rightarrow \epsilon$

start $\rightarrow$ $q_0$ $\quad \epsilon, \epsilon \rightarrow$ <S>\$ $\quad q_1$ $\quad \epsilon, \$ \rightarrow \epsilon$ $\quad q_f$

$\epsilon,$ <S> $\rightarrow$ <B>+<S>
$\epsilon,$ <S> $\rightarrow$ <B>
$\epsilon,$ <B> $\rightarrow$ (<B>)
$\epsilon,$ <B> $\rightarrow$ x

| State | Stack | Input |
|-------|-------|-------|
| $q_0$ | $\epsilon$ | (x+x)+x |
| $q_1$ | <S>\$ | (x+x)+x |
| $q_1$ | <B>+<S>\$ | (x+x)+x |
| $q_1$ | (<S>)+<S>\$ | (x+x)+x |
| $q_1$ | <S>)+<S>\$ | x+x)+x |
| $q_1$ | <B>+<S>)+<S>\$ | x+x)+x |
| $q_1$ | x+<S>)+<S>\$ | x+x)+x |
| $q_1$ | +<S>)+<S>\$ | +x)+x |
| $q_1$ | <S>)+<S>\$ | x)+x |
| $q_1$ | <B>)+<S>\$ | x)+x |
| $q_1$ | x)+<S>\$ | x)+x |
| $q_1$ | )+<S>\$ | )+x |
| $q_1$ | +<S>\$ | +x |
| $q_1$ | <S>\$ | x |
| $q_1$ | <B>\$ | x |
| $q_1$ | x\$ | x |
| $q_1$ | \$ | $\epsilon$ |
| $q_f$ | $\epsilon$ | $\epsilon$ |

# Bottom-up parsing

- Stack holds a word in $(N \cup \Sigma)^*$ from which it is possible to derive the part of the input string that has been consumed
- At any time, may read a letter from input string and push it on top of the stack
- At any time, may non-deterministically choose a rule $A ::= \gamma_1...\gamma_n$ and apply it in reverse: pop $\gamma_n...\gamma_1$ off the top of the stack, and push $A$.
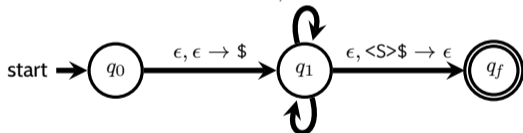- Accept when stack just contains start non-terminal

$$<S> ::= <B>+<S> \mid <B>$$
$$<B> ::= (<S>) \mid x$$

<S> ::= <B>+<S> | <B>

<B> ::= (<S>) | x

$(, \epsilon \rightarrow ($
$), \epsilon \rightarrow )$
$+, \epsilon \rightarrow +$
$x, \epsilon \rightarrow x$

start → $q_0$   $\epsilon, \epsilon \rightarrow \$$   $q_1$   $\epsilon, <S>\$ \rightarrow \epsilon$   $q_f$

$\epsilon, <S>+<B> \rightarrow <S>$
$\epsilon, <B> \rightarrow <S>$
$\epsilon, )<S>( \rightarrow <B>$
$\epsilon, x \rightarrow <B>$

| State | Stack | Input |
|---|---|---|
| $q_0$ | $\epsilon$ | (x+x)+x |
| $q_1$ | $ | (x+x)+x |
| $q_1$ | ($ | x+x)+x |
| $q_1$ | x($ | +x)+x |
| $q_1$ | <B>($ | +x)+x |
| $q_1$ | +<B>($ | x)+x |
| $q_1$ | x+<B>($ | )+x |
| $q_1$ | <B>+<B>($ | )+x |
| $q_1$ | <S>+<B>($ | )+x |
| $q_1$ | <S>($ | )+x |
| $q_1$ | )<S>($ | +x |
| $q_1$ | <B>$ | +x |
| $q_1$ | +<B>$ | x |
| $q_1$ | x+<B>$ | $\epsilon$ |
| $q_1$ | <B>+<B>$ | $\epsilon$ |
| $q_1$ | <S>+<B>$ | $\epsilon$ |
| $q_1$ | <S>$ | $\epsilon$ |
| $q_f$ | $\epsilon$ | $\epsilon$ |

# Parsing overview

- Basic problem with both top-down and bottom-up construction: *non-determinism*
  - Non-deterministic search is inefficient
    - E.g., consider <S> ::= <S>a | <S>b | $\epsilon$. Top-down parser must "guess" the entire input string at the beginning (breadth-first backtracking search takes exponential time in length of input string, depth-first does not terminate).
  - Algorithms for parsing any context free grammar in cubic[1] time, based on dynamic programming (Earley, and Cocke-Younger-Kasami).

---

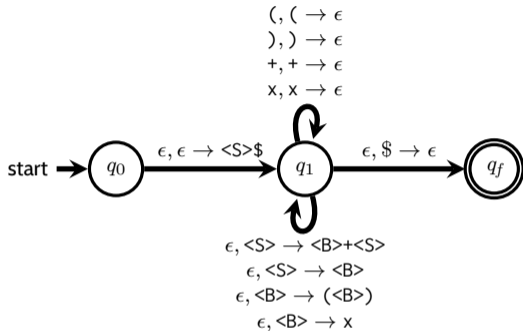[1]Also sub-cubic galactic algorithms

# Parsing overview

- Basic problem with both top-down and bottom-up construction: *non-determinism*
  - Non-deterministic search is inefficient
    - E.g., consider <S> ::= <S>a | <S>b | $\epsilon$. Top-down parser must "guess" the entire input string at the beginning (breadth-first backtracking search takes exponential time in length of input string, depth-first does not terminate).
  - Algorithms for parsing any context free grammar in cubic[1] time, based on dynamic programming (Earley, and Cocke-Younger-Kasami).
- Parser generators use these same ideas, but restricted to cases where we can eliminate non-determinism.
- Possible for both top-down and bottom-up style
  - **Today:** *LL* (*L*eft-to-right, *L*eftmost derivation) parsers: top-down
    - Easy to understand & write by hand
  - **Next week:** *LR* (*L*eft-to-right, *R*ightmost derivation) parsers: bottom-up
    - More general, (variations) implemented in parser generators

---

[1]Also sub-cubic galactic algorithms

# LL parsing



$$<S> ::= <B>+<S> \mid <B>$$
$$<B> ::= (<S>) \mid x$$

Automaton transitions:

On $q_1$ self-loop (top):
$$(, ( \rightarrow \epsilon$$
$$), ) \rightarrow \epsilon$$
$$+, + \rightarrow \epsilon$$
$$x, x \rightarrow \epsilon$$

$$start \rightarrow q_0 \quad \epsilon, \epsilon \rightarrow <S>\$ \quad q_1 \quad \epsilon, \$ \rightarrow \epsilon \quad q_f$$

On $q_1$ self-loop (bottom):
$$\epsilon, <S> \rightarrow <B>+<S>$$
$$\epsilon, <S> \rightarrow <B>$$
$$\epsilon, <B> \rightarrow (<B>)$$
$$\epsilon, <B> \rightarrow x$$

- "Any time top of the stack is a non-terminal $A$, non-deterministically choose a production $A ::= \gamma \in R$. Pop $A$ off the stack, and push $\gamma$"
    - Key problem: need to deterministically choose which production to use
    - Solution: Look at the next input symbol, but don't consume it (*lookahead*)
        - This is $LL(1)$ parsing. $LL(k)$ allows $k$ lookahead tokens

- We say that a grammar is $LL(k)$ if we look ahead $k$ symbols in a top-down parser, we know which rule we should apply.
  - Let $G = (N, \Sigma, R, S)$ be a grammar. $G$ is $LL(k)$ iff: for any $S \Rightarrow^* \alpha A\beta$, for any word $w \in \Sigma^k$, if there is some $A ::= \gamma \in R$ such that $\gamma\beta \Rightarrow^* w\beta'$ (for some $\beta'$), then $\gamma$ is unique.
- Not every context-free language has an $LL(k)$ grammar.
  - $\{a^i b^j : i = j \vee 2i = j\}$ is not $LL(k)$ for any $k$
- Which of the following are $LL(1)$ grammars?
  - <S> ::= a<S> | b<S> | $\epsilon$

  - <S> ::= <S>a | <S>b | $\epsilon$
  - <S> ::= <B>+<S> | <B>
    <B> ::= (<S>) | x

- We say that a grammar is $LL(k)$ if we look ahead $k$ symbols in a top-down parser, we know which rule we should apply.
    - Let $G = (N, \Sigma, R, S)$ be a grammar. $G$ is $LL(k)$ iff: for any $S \Rightarrow^* \alpha A \beta$, for any word $w \in \Sigma^k$, if there is some $A ::= \gamma \in R$ such that $\gamma \beta \Rightarrow^* w \beta'$ (for some $\beta'$), then $\gamma$ is unique.
- Not every context-free language has an $LL(k)$ grammar.
    - $\{a^i b^j : i = j \lor 2i = j\}$ is not $LL(k)$ for any $k$
- Which of the following are $LL(1)$ grammars?
    - <S> ::= a<S> | b<S> | $\epsilon$
      More generally, any grammar that results from our DFA→CFG conversion
    - <S> ::= <S>a | <S>b | $\epsilon$
    - <S> ::= <B>+<S> | <B>
      <B> ::= (<S>) | x

# Left-factoring

- The grammar

$$<S> ::= <B>+<S> \mid <B>$$
$$<B> ::= (<S>) \mid x$$

  is not LL(1): ( lookahead can't distinguish the two <S> rules
- However, there is an LL(1) grammar for the language

# Left-factoring

- The grammar

$$<S> ::= <B>+<S> \mid <B>$$
$$<B> ::= (<S>) \mid x$$

  is not LL(1): ( lookahead can't distinguish the two <S> rules
- However, there is an LL(1) grammar for the language

$$<S> ::= <B><R>$$
$$<R> ::= +<S> \mid \epsilon$$
$$<B> ::= (<S>) \mid x$$

- General strategy: factor out rules with common prefixes ("left factoring")

# Eliminating left recursion

- A grammar is **left-recursive** if there is a non-terminal $A$ such that $A \Rightarrow^+ A\gamma$ (for some $\gamma$)
- Left-recursive grammars are not $LL(k)$ for any $k$
- Consider:

$$<S> ::= <S>+<B> \mid <B>$$
$$<B> ::= (<S>) \mid x$$

Can remove left recursion as follows:

$$<S> ::= <B><S'>$$
$$<S'> ::= +<B><S'> \mid \epsilon$$
$$<B> ::= (<S>) \mid x$$

(Recognizes the same language, but parse trees are different!)

# Mechanical construction of LL(1) parsers

- Fix a grammar $G = (N, \Sigma, R, S)$
- For any word $\gamma \in (N \cup \Sigma)^*$, define **first**$(\gamma) = \{a \in \Sigma : \gamma \Rightarrow^* aw\}$
- For any word $\gamma \in (N \cup \Sigma)^*$, say that $\gamma$ is **nullable** if $\gamma \Rightarrow^* \epsilon$
- For any non-terminal $A$, define **follow**$(A) = \{a \in \Sigma : \exists \gamma, \gamma'. S \Rightarrow \gamma A a \gamma'\}$
- Transition table for $G$ can be computed using **first**, **follow**, and **nullable**:
  1. For each non-terminal $A$ and letter $a$, initialize $\Gamma(A, a)$ to $\emptyset$
  2. For each rule $A ::= \gamma$
     - Add $\gamma$ to $\Gamma(A, a)$ for each $a \in$ **first**$(\gamma)$
     - If $\gamma$ is nullable, add $\gamma$ to $\Gamma(A, a)$ for each $a \in$ **follow**$(A)$

# Mechanical construction of LL(1) parsers

- Fix a grammar $G = (N, \Sigma, R, S)$
- For any word $\gamma \in (N \cup \Sigma)^*$, define **first**$(\gamma) = \{a \in \Sigma : \gamma \Rightarrow^* aw\}$
- For any word $\gamma \in (N \cup \Sigma)^*$, say that $\gamma$ is **nullable** if $\gamma \Rightarrow^* \epsilon$
- For any non-terminal $A$, define **follow**$(A) = \{a \in \Sigma : \exists \gamma, \gamma'. S \Rightarrow \gamma A a \gamma'\}$
- Transition table for $G$ can be computed using **first**, **follow**, and **nullable**:
  1. For each non-terminal $A$ and letter $a$, initialize $\Gamma(A, a)$ to $\emptyset$
  2. For each rule $A ::= \gamma$
     - Add $\gamma$ to $\Gamma(A, a)$ for each $a \in$ **first**$(\gamma)$
     - If $\gamma$ is nullable, add $\gamma$ to $\Gamma(A, a)$ for each $a \in$ **follow**$(A)$
- $G$ is $LL(1)$ iff $\Gamma(A, a)$ is empty or singleton for all $A$ and $a$

# Mechanical construction of LL(1) parsers

- Fix a grammar $G = (N, \Sigma, R, S)$
- For any word $\gamma \in (N \cup \Sigma)^*$, define **first**$(\gamma) = \{a \in \Sigma : \gamma \Rightarrow^* aw\}$
- For any word $\gamma \in (N \cup \Sigma)^*$, say that $\gamma$ is **nullable** if $\gamma \Rightarrow^* \epsilon$
- For any non-terminal $A$, define **follow**$(A) = \{a \in \Sigma : \exists \gamma, \gamma'.S \Rightarrow \gamma A a \gamma'\}$
- Transition table for $G$ can be computed using **first**, **follow**, and **nullable**:
  - ❶ For each non-terminal $A$ and letter $a$, initialize $\Gamma(A, a)$ to $\emptyset$
  - ❷ For each rule $A ::= \gamma$
    - Add $\gamma$ to $\Gamma(A, a)$ for each $a \in$ **first**$(\gamma)$
    - If $\gamma$ is nullable, add $\gamma$ to $\Gamma(A, a)$ for each $a \in$ **follow**$(A)$
- $G$ is $LL(1)$ iff $\Gamma(A, a)$ is empty or singleton for all $A$ and $a$
- Operation of the parser on a word $w$:
  - Start with stack <S>
  - While $w$ not empty
    - If top of the stack is a terminal $a$ and $w = aw'$, pop and set $w = w'$
    - If top of the stack is a non-terminal $A$ and $w = aw'$, pop and push (singleton) $\Gamma(A, w)$ (or reject of $\Gamma(A, w)$ is empty)
  - Accept if stack is empty; reject otherwise.

# Computing nullable

- **nullable** is the *smallest set* of non-terminals such that if there is some $A ::= \gamma_1...\gamma_n \in R$ with $\gamma_1, ..., \gamma_n \in$ **nullable** implies $A \in$ **nullable**
  - Fixpoint computation:
    - $\text{nullable}_0 = \emptyset$
    - $\text{nullable}_{i+1} = \{A : \exists \gamma_1, ..., \gamma_n \in \text{nullable}_i. A ::= \gamma_1...\gamma_n \in R\}$
    - $\text{nullable} = \bigcup_{i=0}^{\infty} \text{nullable}_i$

  > nullable $\leftarrow \emptyset$;
  > changed $\leftarrow$ true;
  > **while** *changed* **do**
  >     changed $\leftarrow$ false;
  >     **for** $A := \gamma_1...\gamma_n \in R$ **do**
  >         **if** $A \notin nullable \land \gamma_1, ..., \gamma_n \in nullable$ **then**
  >             $nullable \leftarrow nullable \cup \{A\}$;
  >             changed $\leftarrow$ true;

- Fixpoint computations appear everywhere!
  - Later we will see how they are used in dataflow analysis

# Computing first and follow

- **first** is the *smallest function*[2] such that
  - For each $a \in \Sigma$, $\mathbf{first}(a) = \{a\}$
  - For each $A ::= \gamma_1...\gamma_i...\gamma_n \in R$, with $\gamma_1, ..., \gamma_{i-1}$ nullable, $\mathbf{first}(A) \supseteq \mathbf{first}(\gamma_i)$
- **follow** is the *smallest function* such that
  - For each $A ::= \gamma_1...\gamma_i...\gamma_n \in R$, with $\gamma_{i+1}, ..., \gamma_n$ nullable, $\mathbf{follow}(\gamma_i) \supseteq \mathbf{follow}(A)$
  - For each $A ::= \gamma_1...\gamma_i...\gamma_j...\gamma_n \in R$, with $\gamma_{i+1}, ..., \gamma_{j-1}$ nullable, $\mathbf{follow}(\gamma_i) \supseteq \mathbf{first}(A)$
- Both can be computed using a fixpoint algorithm, like **nullable**

---

[2]Pointwise order: $f \leq g$ if for all $x, f(x) \leq g(x)$