# COMP2022: Formal Languages and Logic
## 2017, Semester 1, Week 5

Joseph Godbehere

Adapted from slides by A/Prof Kalina Yacef

April 4, 2017

THE UNIVERSITY OF
SYDNEY

## ANNOUNCEMENTS

Mid-semester survey

- ▶ On eLearning, until Friday
- ▶ Help us improve the course for you
- ▶ Anonymous

Quizzes!

- ▶ Content: always similar to questions from the prior tutorial
- ▶ How to do well: revise the tutorial content

## OUTLINE

- Revision

- Push Down Automata

- LL(k) Table-Descent Parsing

## REVISION

A grammar $G$ is sets of variables, terminals, production rules, and a start variable.

The production rules in *context-free grammars* always have the form $A \Rightarrow \alpha$, where $A$ is a variable and $\alpha$ is a string of variables and/or terminals (or $\varepsilon$)

$w$ can be derived by $G$ if there exists a sequence of substitutions such that $S \Rightarrow^+ w$
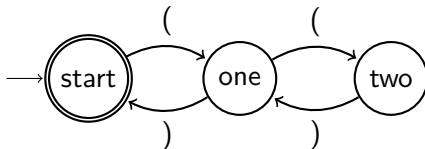
- ▶ Leftmost derivation: substitute the leftmost *variable* in $w$
- ▶ Rightmost derivation: substitute the rightmost *variable* in $w$

Language of $G$ is the set of strings which can be derived from $S$
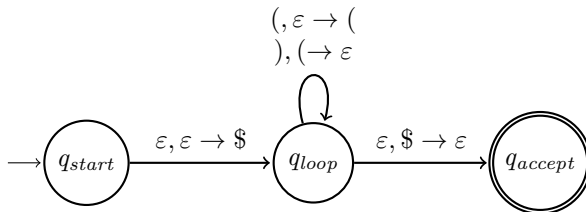
Ambiguous strings have multiple parse trees (equivalently, multiple left derivations, or multiple right derivations.)

LIMITS OF DFA/NFA

Recall this NFA for balanced parentheses:



It accepts strings of balanced parentheses, nested up to two deep.
What if we wanted a depth of 3?
What if we want any level of nesting?
If we added a *stack* to the automata, we could accept any level of
nesting

# PARSING $((()))$



| Current state | Unscanned input | Stack contents |
|:---:|:---:|:---:|
| $q_{start}$ | $((()))$ | $\emptyset$ |
| $q_{loop}$ | $((()))$ | $\$$ |
| $q_{loop}$ | $(())$ | $($\$$ |
| $q_{loop}$ | $()))$ | $(($\$$ |
| $q_{loop}$ | $)))$ | $((($\$$ |
| $q_{loop}$ | $))$ | $(($\$$ |
| $q_{loop}$ | $)$ | $($\$$ |
| $q_{loop}$ | $\varepsilon$ | $\$$ |
| $q_{accept}$ | $\varepsilon$ | $\emptyset$ |

# PUSH DOWN AUTOMATA

Finite automata *with a stack*.

Revision:
A *stack* is a structure with a LIFO property (Last In, First Out)

Stack operations:

- ▶ PEEK: read the top symbol without removing it
- ▶ POP: read the top symbol and remove it
- ▶ PUSH: write a symbol onto the top of the stack
- ▶ NOP: perform no operation on the stack

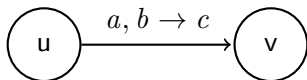## PUSH DOWN AUTOMATA

Finite automata *with a stack*.

A PDA is a machine with the ability to:

- ▶ Read the letters of input string (like FA)
- ▶ Make state changes (like FA)
- ▶ Perform stack operations (new!)

In addition to storing terminals and variables, the stack will accept a special end of string symbol, which we will denote $.

## PDA TRANSITIONS

PDA Transitions include stack operations

$$\left(\; u \;\right) \xrightarrow{\;a,\, b \to c\;} \left(\; v \;\right)$$

$a,\, b \to c$ means:

- ▶ Read (and remove) symbol $a$ from the front of the input
- ▶ Pop symbol $b$ from the top of the stack
- ▶ Push $c$ onto the stack

We follow a transition if it is possible to perform these operations

## PDA TRANSITIONS

We use $\varepsilon$ to denote "no operation". For example:

- $a, \varepsilon \rightarrow c$
    - read (and remove) $a$ from the front of the input
    - do not pop anything from the stack
    - push $c$ onto the stack

## PDA TRANSITIONS

We use $\varepsilon$ to denote "no operation". For example:

- $a, \varepsilon \to c$
    - read (and remove) $a$ from the front of the input
    - do not pop anything from the stack
    - push $c$ onto the stack

- $\varepsilon, \varepsilon \to c$
    - do not read anything from the input
    - do not pop anything from the stack
    - push $c$ onto the stack

## PDA TRANSITIONS

We use $\varepsilon$ to denote "no operation". For example:

▶ $a, \varepsilon \to c$
  ▶ read (and remove) $a$ from the front of the input
  ▶ do not pop anything from the stack
  ▶ push $c$ onto the stack

▶ $\varepsilon, \varepsilon \to c$
  ▶ do not read anything from the input
  ▶ do not pop anything from the stack
  ▶ push $c$ onto the stack

▶ $\varepsilon, b \to \varepsilon$
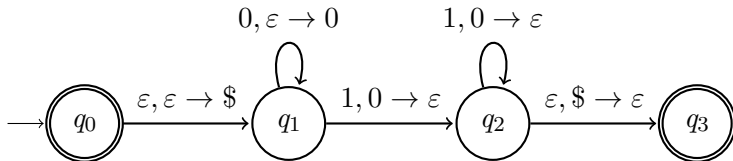  ▶ do not read anything from the input
  ▶ pop $b$ from the top of the stack
  ▶ do not push anything onto the stack

# EXAMPLE: PDA $M_2$ FOR $\{0^n 1^n \mid n \geq 0\}$



1. Initialise the stack with $\$$
2. As each 0 is scanned, push it onto the stack
3. As each 1 is scanned, pop a 0 from the stack
4. If all the 0's have been popped, $\$$ will be on top the stack. If so, pop it and transition to $q_3$
5. If we have read all the input, accept the string

## PDA: FORMAL DEFINITION

A pushdown automata is a 6-tuple $(Q, \Sigma, \Gamma, \delta, q_0, F)$ where:

1. $Q$ is a finite set called the states
2. $\Sigma$ is a finite set called the input alphabet
3. $\Gamma$ is a finite set called the stack alphabet
4. $\delta : Q \times \Sigma_\varepsilon \times \Gamma_\varepsilon \to \mathcal{P}(Q \times \Gamma_\varepsilon)$ is the transition function
5. $q_0 \in Q$ is the start state
6. $F \subseteq Q$ is the set of accept states

## PDA: FORMAL DESCRIPTION OF $M_2$

1. $Q = \{q_0, q_1, q_2, q_3\}$
2. $\Sigma = \{0, 1\}$
3. $\Gamma = \{0, \$\}$
4. $\delta : Q \times \Sigma_\varepsilon \times \Gamma_\varepsilon \to \mathcal{P}(Q \times \Gamma_\varepsilon)$ is

|              | $\varepsilon$        | 0                    | $\$$                 |
|--------------|----------------------|----------------------|----------------------|
| $q_0, \varepsilon$ | $\{(q_1, \$)\}$ |                      |                      |
| $q_1, 0$     | $\{(q_1, 0)\}$       |                      |                      |
| $q_1, 1$     |                      | $\{(q_2, \varepsilon)\}$ |                  |
| $q_2, 1$     |                      | $\{(q_2, \varepsilon)\}$ |                  |
| $q_2, \varepsilon$ |                |                      | $\{(q_3, \varepsilon)\}$ |

5. $q_0$ is the start state
6. $F = \{q_0, q_3\}$

14/52

LANGUAGE RECOGNISED BY A PDA

A string is *accepted* by a PDA if there is some path from the start
state to a final state, which consumes all the letters of the string.
Otherwise the string is *rejected*.

The language *recognised* by a PDA is the set of strings which it
accepts.

## LANGUAGE RECOGNISED BY A PDA

Formal definition:

Let $M$ be a PDA and $w = w_1 w_2 ... w_n$ be a string in $\Sigma^\star$. $M$ accepts $w$ if and only if there exists a string path:

$(r_0, s_0) \rightarrow_{a_1} (r_1, s_1) \rightarrow_{a_2} ... \rightarrow_{a_n} (r_n, s_n)$ such that:

1. $r_0 = q_0$ and $s_0 = \varepsilon$
2. $\delta(r_i, a_{i+1}, b) = (r_{i+1}, c)$ for $i = 0, ..., n-1$ where $s_i = bt$ and $s_{i+1} = ct$ with $t \in \Gamma^\star$
3. $r_n \in F$

## Language recognised by a PDA

Formal definition:

Let $M$ be a PDA and $w = w_1 w_2 ... w_n$ be a string in $\Sigma^\star$. $M$ accepts $w$ if and only if there exists a string path:

$(r_0, s_0) \to_{a_1} (r_1, s_1) \to_{a_2} ... \to_{a_n} (r_n, s_n)$ such that:

1. $r_0 = q_0$ and $s_0 = \varepsilon$
2. $\delta(r_i, a_{i+1}, b) = (r_{i+1}, c)$ for $i = 0, ..., n-1$ where $s_i = bt$ and $s_{i+1} = ct$ with $t \in \Gamma^\star$
3. $r_n \in F$

▶ This is slightly oversimplified! Actually we just need there to be some sequence of transitions which only consume one input symbol (i.e. we can also follow as many $\delta(\varepsilon, x, y)$ transitions as needed.)

## EQUIVALENT FORMS OF ACCEPTANCE

Final-state acceptance:
A string is accepted if there is some path from the start state to a final state which consumes all of the letters of the string. This is the most common definition, and the one we used earlier.

Empty-stack acceptance:
A string is accepted if the input string has been consumed and the stack is empty (no requirement about the machine being in any particular state.)

## DETERMINISM

A PDA is *deterministic* if for all states there is at most one possible move, given the unscanned input and stack contents.

More formally, suppose that:

▶ The PDA is in state $q$

▶ The next input symbol is $a$

▶ The symbol on top of the stack is $b$

Then we can transition to any pair (next state, symbol to push) in:

$$\delta(q, a, b) \cup \delta(q, \varepsilon, b) \cup \delta(q, a, \varepsilon) \cup \delta(q, \varepsilon, \varepsilon)$$

The PDA is *non-deterministic* if this set contains more than one pair.

DETERMINISM

DFA and NFA have the same recognition power

Non-deterministic PDA are *more* powerful than deterministic PDA.

- ▶ Some context-free languages can be recognised by a non-deterministic PDA but *not* by a deterministic one.
  - ▶ Example: Language of palindromes (why?)

## CONTEXT-FREE GRAMMARS AND PDA

Theorem:
Context-free languages are exactly the languages accepted by PDAs

Proof:

1. For any PDA, there is a CFG generating the same language
   ▸ We will not show this construction (Sipser 2.20 – 2.31)
   ▸ The method usually produces very large grammars

2. For any CFG, there is a PDA recognising the same language
   ▸ We will explore this algorithm
   ▸ The PDA non-deterministically performs a leftmost derivation

## CFG → PDA ALGORITHM

The PDA will non-deterministically attempt a leftmost derivation.

Push the empty stack marker $, then the start symbol on the stack.

Repeat the following steps:

# CFG $\rightarrow$ PDA ALGORITHM

The PDA will non-deterministically attempt a leftmost derivation.

Push the empty stack marker \$, then the start symbol on the stack.

Repeat the following steps:

▶ If the top of the stack is a variable $A$, then non-deterministically choose one of the rules for $A$. Substitute $A$ on the stack with the string produced by the rule.

## CFG → PDA ALGORITHM

The PDA will non-deterministically attempt a leftmost derivation.

Push the empty stack marker \$, then the start symbol on the stack.

Repeat the following steps:

- If the top of the stack is a variable $A$, then non-deterministically choose one of the rules for $A$. Substitute $A$ on the stack with the string produced by the rule.
- If the top of the stack is a terminal symbol $a$, then read the next symbol from the input. If they match, pop the stack. Otherwise, reject this branch of the non-determinism.

## CFG → PDA ALGORITHM

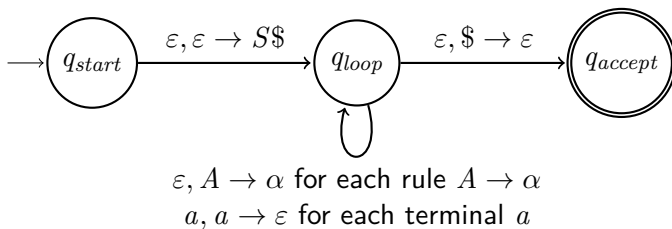The PDA will non-deterministically attempt a leftmost derivation.

Push the empty stack marker \$, then the start symbol on the stack.

Repeat the following steps:

- ▶ If the top of the stack is a variable $A$, then non-deterministically choose one of the rules for $A$. Substitute $A$ on the stack with the string produced by the rule.
- ▶ If the top of the stack is a terminal symbol $a$, then read the next symbol from the input. If they match, pop the stack. Otherwise, reject this branch of the non-determinism.
- ▶ If the top of the stack is \$, enter the accept state. This will accept the input if it has all been read.
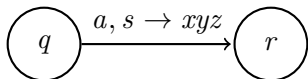
## CFG → PDA CONSTRUCTION METHOD

First build a "reduced PDA" as follows:



$\varepsilon, A \to \alpha$ for each rule $A \to \alpha$
$a, a \to \varepsilon$ for each terminal $a$

# CFG → PDA CONSTRUCTION METHOD

Now expand the transitions in the "reduced PDA" to a normal one, such that only one symbol is push or popped to the stack in each transitions:

$$q \xrightarrow{\quad a, s \to xyz \quad} r$$

becomes:

$$q \xrightarrow{\quad a, s \to z \quad} \bigcirc \xrightarrow{\quad \varepsilon, \varepsilon \to y \quad} \bigcirc \xrightarrow{\quad \varepsilon, \varepsilon \to x \quad} r$$

i.e. we push the string of symbols $xyz$ onto the stack in reverse order, with each transition after the first one reading nothing from the input or the stack.

# CFG $\rightarrow$ PDA example for $B \rightarrow (B)B \mid \varepsilon$

First build a "reduced PDA" as follows:

# CFG $\rightarrow$ PDA EXAMPLE FOR $B \rightarrow (B)B \mid \varepsilon$

Expand $\varepsilon, \varepsilon \rightarrow S\$$



$$\varepsilon, B \rightarrow \varepsilon$$
$$\varepsilon, B \rightarrow (B)B$$
$$(, ( \rightarrow \varepsilon$$
$$), ) \rightarrow \varepsilon$$

# CFG $\to$ PDA example for $B \to (B)B \mid \varepsilon$

Expand $\varepsilon, B \to (B)B$



$\varepsilon, B \to \varepsilon$
$(, ( \to \varepsilon$
$), ) \to \varepsilon$

$\to q_{start} \quad \varepsilon, \varepsilon \to \$ \quad \bigcirc \quad \varepsilon, \varepsilon \to S \quad q_{loop} \quad \varepsilon, \$ \to \varepsilon \quad q_{accept}$

$\varepsilon, B \to B$

$\varepsilon, \varepsilon \to ($

$\varepsilon, \varepsilon \to )$

$\varepsilon, \varepsilon \to B$

Grammars are fundamental for constructing parsers.

Parsers should be efficient, so determinism is important.

**Top-down parsing:**

**Bottom-up parsing:**

Grammars are fundamental for constructing parsers.

Parsers should be efficient, so determinism is important.

**Top-down parsing:** A parse tree is constructed from the root, proceeding downward towards the leaves

- ▶ Constructs the derivation by starting with the grammar's start symbol, then working towards the string. i.e. $S \Rightarrow^+ w$
- ▶ LL(k) parsing
    - ▶ **L**eft-to-right, **L**eftmost derivation

**Bottom-up parsing:**

Grammars are fundamental for constructing parsers.

Parsers should be efficient, so determinism is important.

**Top-down parsing:** A parse tree is constructed from the root, proceeding downward towards the leaves

- ▶ Constructs the derivation by starting with the grammar's start symbol, then working towards the string. i.e. $S \Rightarrow^+ w$
- ▶ LL(k) parsing
  - ▶ **L**eft-to-right, **L**eftmost derivation

**Bottom-up parsing:** The parse tree is constructed starting from the leaves, working up towards the root

- ▶ Constructs the derivation by starting with the string, and working backwards to the start symbol. i.e. $w \Leftarrow ... \Leftarrow S$
- ▶ LR(k) parsing
  - ▶ **L**eft-to-right, **R**ightmost derivation in reverse

# FROM CFG → PDA → TABLE DRIVEN PARSERS

The interesting part of the PDA is the stack and the $Q_{loop}$ state

The stack contains right-hand side of the production rules, in an order which allows a left-most derivation.

This can be programmed as a descent table-driven parser

- ► Input: current token and variable on top of the stack
- ► Output: which rule to use

Non-determinism in the table is a problem: how do we choose which rule to use?

We will see how to construct a deterministic table whenever possible

- ► Not all CFG have an equivalent *deterministic* PDA

## LOOKAHEAD SYMBOLS

Suppose we have a grammar with two rules for $S$:

$$S \rightarrow XY \mid YZ \qquad \text{(other rules not shown)}$$

How do we choose which rule should replace $S$?

Suppose $XY$ can only derive strings which start with $a$, and $YZ$ can only derive strings which start with $b$ or $c$.

Then if we look ahead one symbol, we know which rule to select:

- to derive $abc$ we must choose $S \Rightarrow XY$
- to derive $cab$ we must choose $S \Rightarrow YZ$

# LL(K) PARSING

LL(k) grammar

- ▶ L: Left to right scanning of input
- ▶ L: Leftmost derivation
- ▶ k: uses $k$ look ahead symbols

Deterministic derivation by looking ahead $k$ symbols

Using less lookahead symbols is usually more efficient

# LL(1) GRAMMAR: EXAMPLE

$L = \{a^n bc^n \mid n \geq 0\}$

$$S \rightarrow aSc \mid b$$

This grammar is LL(1), because the right side of each production rule lead to strings beginning with different letters

Each step of the derivation can be deterministically determined by examining the current symbol (1 lookahead symbol)

- ► If the remaining input starts with $a$, use $S \rightarrow aSc$
- ► If the remaining input starts with $b$, use $S \rightarrow b$
- ► (If it starts with anything else, no derivation exists)

# LL(1) GRAMMAR: EXAMPLE

$$L = \{ a^n b c^n \mid n \geq 0 \}$$

$$S \to aSc \mid b$$

Derivation of $aabcc$

$$
\begin{aligned}
S &\Rightarrow aSc & & \text{Use } S \to aSc \text{ because } aabcc \text{ begins with } a \\
&\Rightarrow aaScc & & \text{Use } S \to aSc \text{ because } abcc \text{ begins with } a \\
&\Rightarrow aabcc & & \text{Use } S \to b \text{ because } bcc \text{ begins with } b
\end{aligned}
$$

# LL(2) GRAMMAR: EXAMPLE

$L = \{a^m b^n c \mid n \geq 0\}$

$$S \rightarrow AB$$
$$A \rightarrow aA \mid a$$
$$B \rightarrow bB \mid c$$

Each step of the derivation can be deterministically determined by examining the current symbol and the next one (2 lookahead symbols). e.g. if we need to replace a variable $A$ and:

- ▶ If the remaining input starts with $aa$, use $A \rightarrow aA$
- ▶ If the remaining input starts with $ab$ or $ac$, use $A \rightarrow a$
- ▶ (If it starts with anything else, no derivation exists)

# LL(2) GRAMMAR: EXAMPLE

$$S \rightarrow AB$$
$$A \rightarrow aA \mid a$$
$$B \rightarrow bB \mid c$$

Derivation of $aabbc$

| $S \Rightarrow AB$ | No other choice |
|---|---|
| $\Rightarrow aAB$ | Use $A \rightarrow aA$ because $aabbc$ begins with $aa$ |
| $\Rightarrow aaB$ | Use $A \rightarrow a$ because $abbc$ begins with $ab$ |
| $\Rightarrow aabB$ | Use $B \rightarrow bB$ because $bbc$ begins with $b$ |
| $\Rightarrow aabbB$ | Use $B \rightarrow bB$ because $bc$ begins with $b$ |
| $\Rightarrow aabbc$ | Use $B \rightarrow c$ because $c$ begins with $c$ |

# NON-LL(K) GRAMMAR: EXAMPLE

$$S \rightarrow aS \mid T$$
$$T \rightarrow aBb \mid \varepsilon$$

Not LL(1): next symbol $a$ is not enough to determine which production to use ($S \rightarrow aS$ and $S \rightarrow T$ can both generate strings starting with $a$)

Not LL(2): the input $aa$ is not enough either

Not LL(k): we need to know how many $b$'s there are. For any $k$ we can choose $n > k$ such that $a^n b^n \in L(G)$ but we would need to lookahead $2n > k$ symbols to decide which rule to use first.

# TABLE-DRIVEN LL(1) PARSING

In a PDA: the stack contains the right hand side of the rules for a leftmost derivation

In a *descent table-driven parser*: Given the current input and the variable on top of the stack, the table specifies which production rule to use.

# DESCENT TABLE-DRIVEN LL(1) PARSER

```
loop
    T = symbol on top of the stack
    c = current input symbol
    if T == c == $ then accept
    else if T is a terminal or T == $ then
        if T == c then pop T and consume c
        else error
    else if P[T,c] = α is defined then
        pop T and push α onto the stack
                //(in reverse order)
    else error
endloop
```

## Rules starting with non-terminals

$$S \to A \mid B$$
$$A \to aA \mid \varepsilon$$
$$B \to bB \mid c$$

It's clear that the productions for $A$ and $B$ are LL(1), but how do we choose which rule to use when deriving from $S$?

We look at the possible symbols which can start strings derived after $S \to A$, or after $S \to B$.

The set of symbols which could start any string derived from $\alpha$ is called the FIRST set of $\alpha$

# FIRST AND FOLLOW SETS

In order to fill in the entries of the table-drived parser, we need to compute some FIRST and FOLLOW sets.

$FIRST(\alpha)$ is the set of all terminals which could start strings derived from $\alpha$. We will need to calculate these for every production of $G$ (i.e. the *right hand side* of each rule).

$FOLLOW(V)$ is the set of all terminals which which could follow the variable $V$ at any stage of the derivation. Needed whenever $V$ can derive $\varepsilon$.

## TABLE CONSTRUCTION: FIRST SETS

If $\alpha$ is a string, then $FIRST(\alpha)$ is the set of terminals which can begin strings derived from $\alpha$. Iff $\alpha \Rightarrow^+ \varepsilon$ then $\varepsilon \in FIRST(\alpha)$.

Construction rules: Where $a$ is a terminal, and $\alpha$ is some string of terminals and variables:

## TABLE CONSTRUCTION: FIRST SETS

If $\alpha$ is a string, then $FIRST(\alpha)$ is the set of terminals which can begin strings derived from $\alpha$. Iff $\alpha \Rightarrow^+ \varepsilon$ then $\varepsilon \in FIRST(\alpha)$.

Construction rules: Where $a$ is a terminal, and $\alpha$ is some string of terminals and variables:

1. $FIRST(\varepsilon) = \{\varepsilon\}$ (by definition)

## TABLE CONSTRUCTION: FIRST SETS

If $\alpha$ is a string, then $FIRST(\alpha)$ is the set of terminals which can begin strings derived from $\alpha$. Iff $\alpha \Rightarrow^+ \varepsilon$ then $\varepsilon \in FIRST(\alpha)$.

Construction rules: Where $a$ is a terminal, and $\alpha$ is some string of terminals and variables:

1. $FIRST(\varepsilon) = \{\varepsilon\}$ (by definition)
2. $FIRST(a\alpha) = FIRST(a) = \{a\}$

## TABLE CONSTRUCTION: FIRST SETS

If $\alpha$ is a string, then $FIRST(\alpha)$ is the set of terminals which can begin strings derived from $\alpha$. Iff $\alpha \Rightarrow^+ \varepsilon$ then $\varepsilon \in FIRST(\alpha)$.

Construction rules: Where $a$ is a terminal, and $\alpha$ is some string of terminals and variables:

1. $FIRST(\varepsilon) = \{\varepsilon\}$ (by definition)
2. $FIRST(a\alpha) = FIRST(a) = \{a\}$
3. If $A \to \alpha_1 \mid ... \mid \alpha_n$ then
   $FIRST(A) = FIRST(\alpha_1) \cup ... \cup FIRST(\alpha_n)$

## TABLE CONSTRUCTION: FIRST SETS

If $\alpha$ is a string, then $FIRST(\alpha)$ is the set of terminals which can begin strings derived from $\alpha$. Iff $\alpha \Rightarrow^+ \varepsilon$ then $\varepsilon \in FIRST(\alpha)$.

Construction rules: Where $a$ is a terminal, and $\alpha$ is some string of terminals and variables:

1. $FIRST(\varepsilon) = \{\varepsilon\}$ (by definition)

2. $FIRST(a\alpha) = FIRST(a) = \{a\}$

3. If $A \to \alpha_1 \mid ... \mid \alpha_n$ then
   $FIRST(A) = FIRST(\alpha_1) \cup ... \cup FIRST(\alpha_n)$

4. If $\alpha \neq \varepsilon$ then
   - If $\varepsilon \notin FIRST(A)$ then $FIRST(A\alpha) = FIRST(A)$
   - If $\varepsilon \in FIRST(A)$ then
     $FIRST(A\alpha) = FIRST(A) \setminus \{\varepsilon\} \cup FIRST(\alpha)$

## Examples calculating FIRST sets

$$B \to (B)B \mid \varepsilon$$

$FIRST(``(B)B") =$

## Examples calculating FIRST sets

$$B \rightarrow (B)B \mid \varepsilon$$

$$FIRST(\text{``}(B)B\text{''}) = \{(\}$$            rule 2
$$FIRST(B) =$$

## Examples calculating FIRST sets

$$B \rightarrow (B)B \mid \varepsilon$$

$$FIRST(\text{``}(B)B\text{''}) = \{(\} \qquad\qquad \text{rule 2}$$
$$FIRST(B) = FIRST(\text{``}(B)B\text{''}) \cup FIRST(\varepsilon) \quad \text{rule 3}$$
$$= $$

## EXAMPLES CALCULATING FIRST SETS

$$B \rightarrow (B)B \mid \varepsilon$$

$$
\begin{aligned}
FIRST(\text{``}(B)B\text{''}) &= \{( \} & \text{rule 2} \\
FIRST(B) &= FIRST(\text{``}(B)B\text{''}) \cup FIRST(\varepsilon) & \text{rule 3} \\
&= FIRST(\text{``}(B)B\text{''}) \cup \{\varepsilon\} & \text{rule 1} \\
&=
\end{aligned}
$$

## EXAMPLES CALCULATING FIRST SETS

$$B \rightarrow (B)B \mid \varepsilon$$

$$
\begin{aligned}
FIRST(\text{``}(B)B\text{''}) &= \{(\} & \text{rule 2} \\
FIRST(B) &= FIRST(\text{``}(B)B\text{''}) \cup FIRST(\varepsilon) & \text{rule 3} \\
&= FIRST(\text{``}(B)B\text{''}) \cup \{\varepsilon\} & \text{rule 1} \\
&= \{(\} \cup \{\varepsilon\} \\
&= \{(, \varepsilon\}
\end{aligned}
$$

## EXAMPLES CALCULATING FIRST SETS

$$S \to BC \mid a$$
$$B \to bB \mid \varepsilon$$
$$C \to cC \mid \varepsilon$$

We want to calculate $FIRST(BC)$. We will need to calulate some other sets first to help us.

$FIRST(\varepsilon) =$

# EXAMPLES CALCULATING FIRST SETS

$$S \to BC \mid a$$
$$B \to bB \mid \varepsilon$$
$$C \to cC \mid \varepsilon$$

We want to calculate $FIRST(BC)$. We will need to calulate some other sets first to help us.

$FIRST(\varepsilon) = \{\varepsilon\}$            rule 1

$FIRST(B) =$

## EXAMPLES CALCULATING FIRST SETS

$$S \rightarrow BC \mid a$$
$$B \rightarrow bB \mid \varepsilon$$
$$C \rightarrow cC \mid \varepsilon$$

We want to calculate $FIRST(BC)$. We will need to calulate some other sets first to help us.

$$FIRST(\varepsilon) = \{\varepsilon\} \qquad \text{rule 1}$$
$$FIRST(B) = FIRST(bB) \cup FIRST(\varepsilon) \qquad \text{rule 3}$$
$$=$$

## Examples calculating FIRST sets

$$S \to BC \mid a$$
$$B \to bB \mid \varepsilon$$
$$C \to cC \mid \varepsilon$$

We want to calculate $FIRST(BC)$. We will need to calulate some other sets first to help us.

$$
\begin{aligned}
FIRST(\varepsilon) &= \{\varepsilon\} && \text{rule 1} \\
FIRST(B) &= FIRST(bB) \cup FIRST(\varepsilon) && \text{rule 3} \\
&= \{b, \varepsilon\} && \text{rules 1 and 2} \\
FIRST(C) &=
\end{aligned}
$$

## EXAMPLES CALCULATING FIRST SETS

$$S \to BC \mid a$$
$$B \to bB \mid \varepsilon$$
$$C \to cC \mid \varepsilon$$

We want to calculate $FIRST(BC)$. We will need to calulate some other sets first to help us.

$$
\begin{aligned}
FIRST(\varepsilon) &= \{\varepsilon\} && \text{rule 1} \\
FIRST(B) &= FIRST(bB) \cup FIRST(\varepsilon) && \text{rule 3} \\
&= \{b, \varepsilon\} && \text{rules 1 and 2} \\
FIRST(C) &= \{c, \varepsilon\} && \text{similarly} \\
FIRST(BC) &=
\end{aligned}
$$

## EXAMPLES CALCULATING FIRST SETS

$$S \to BC \mid a$$
$$B \to bB \mid \varepsilon$$
$$C \to cC \mid \varepsilon$$

We want to calculate $FIRST(BC)$. We will need to calulate some other sets first to help us.

$$
\begin{aligned}
FIRST(\varepsilon) &= \{\varepsilon\} & \text{rule 1} \\
FIRST(B) &= FIRST(bB) \cup FIRST(\varepsilon) & \text{rule 3} \\
&= \{b, \varepsilon\} & \text{rules 1 and 2} \\
FIRST(C) &= \{c, \varepsilon\} & \text{similarly} \\
FIRST(BC) &= FIRST(B) \setminus \{\varepsilon\} \cup FIRST(C) & \text{rule 4b} \\
&=
\end{aligned}
$$

## Examples calculating FIRST sets

$$S \rightarrow BC \mid a$$
$$B \rightarrow bB \mid \varepsilon$$
$$C \rightarrow cC \mid \varepsilon$$

We want to calculate $FIRST(BC)$. We will need to calulate some other sets first to help us.

$$
\begin{aligned}
FIRST(\varepsilon) &= \{\varepsilon\} && \text{rule 1} \\
FIRST(B) &= FIRST(bB) \cup FIRST(\varepsilon) && \text{rule 3} \\
&= \{b, \varepsilon\} && \text{rules 1 and 2} \\
FIRST(C) &= \{c, \varepsilon\} && \text{similarly} \\
FIRST(BC) &= FIRST(B) \setminus \{\varepsilon\} \cup FIRST(C) && \text{rule 4b} \\
&= \{b\} \cup \{c, \varepsilon\} \\
&= \{b, c, \varepsilon\}
\end{aligned}
$$

# WHEN RULES CAN YIELD $\varepsilon$

Consider the grammar

$$S \rightarrow AB$$
$$A \rightarrow aA \mid \varepsilon$$
$$B \rightarrow bB \mid c$$

$B$ can start with $b$ or $\varepsilon$, i.e. $FIRST(A) = \{b, \varepsilon\}$

Suppose we are parsing the string $bc$, and so far we have derived $S \Rightarrow AB$. How does the parser know which production to use next?

## WHEN RULES CAN YIELD $\varepsilon$

Consider the grammar

$$S \to AB$$
$$A \to aA \mid \varepsilon$$
$$B \to bB \mid c$$

$B$ can start with $b$ or $\varepsilon$, i.e. $FIRST(A) = \{b, \varepsilon\}$

Suppose we are parsing the string $bc$, and so far we have derived $S \Rightarrow AB$. How does the parser know which production to use next?

$b$ is not in $FIRST(aA)$ or $FIRST(\varepsilon)$, but we can clearly see that using $A \to \varepsilon$ will give us $AB \Rightarrow B \Rightarrow bB \Rightarrow bc$

## WHEN RULES CAN YIELD $\varepsilon$

Consider the grammar

$$S \to AB$$
$$A \to aA \mid \varepsilon$$
$$B \to bB \mid c$$

$B$ can start with $b$ or $\varepsilon$, i.e. $FIRST(A) = \{b, \varepsilon\}$

Suppose we are parsing the string $bc$, and so far we have derived $S \Rightarrow AB$. How does the parser know which production to use next?

$b$ is not in $FIRST(aA)$ or $FIRST(\varepsilon)$, but we can clearly see that using $A \to \varepsilon$ will give us $AB \Rightarrow B \Rightarrow bB \Rightarrow bc$

When a variable can derive $\varepsilon$, we need to look at the terminal symbols which can begin strings which could FOLLOW that variable in an derivation. These are called the FOLLOW sets.

## TABLE CONSTRUCTION: FOLLOW SETS

Definition:

If $A$ is a variable, then $FOLLOW(A)$ is the set of terminals which can be derived from any string which can appear immediately to the right of $A$ in some stage of the derivation.

Construction rules:

Where $\alpha, \beta$ are strings of symbols, and $S, X, Y$ are variables:

## TABLE CONSTRUCTION: FOLLOW SETS

Definition:

If $A$ is a variable, then $FOLLOW(A)$ is the set of terminals which can be derived from any string which can appear immediately to the right of $A$ in some stage of the derivation.

Construction rules:

Where $\alpha, \beta$ are strings of symbols, and $S, X, Y$ are variables:

1. If $S$ is the start symbol, then $\$ \in FOLLOW(S)$

   (i.e. the start symbol can be followed by the end of the string)

# TABLE CONSTRUCTION: FOLLOW SETS

Definition:

If $A$ is a variable, then $FOLLOW(A)$ is the set of terminals which can be derived from any string which can appear immediately to the right of $A$ in some stage of the derivation.

Construction rules:

Where $\alpha, \beta$ are strings of symbols, and $S, X, Y$ are variables:

1. If $S$ is the start symbol, then $\$ \in FOLLOW(S)$

   (i.e. the start symbol can be followed by the end of the string)

2. If $X \rightarrow \alpha Y$ then $FOLLOW(X) \subset FOLLOW(Y)$

   (i.e. anything that can follow $Y$ can follow $X$)

## TABLE CONSTRUCTION: FOLLOW SETS

Definition:

If $A$ is a variable, then $FOLLOW(A)$ is the set of terminals which can be derived from any string which can appear immediately to the right of $A$ in some stage of the derivation.

Construction rules:

Where $\alpha, \beta$ are strings of symbols, and $S, X, Y$ are variables:

1. If $S$ is the start symbol, then $\$ \in FOLLOW(S)$

   (i.e. the start symbol can be followed by the end of the string)

2. If $X \rightarrow \alpha Y$ then $FOLLOW(X) \subset FOLLOW(Y)$

   (i.e. anything that can follow $Y$ can follow $X$)

3. If $X \rightarrow \alpha Y \beta$ then $FIRST(\beta) \setminus \{\varepsilon\} \subset FOLLOW(Y)$

   (i.e. any terminal which can start $\beta$ can follow $Y$)

# TABLE CONSTRUCTION: FOLLOW SETS

Definition:

If $A$ is a variable, then $FOLLOW(A)$ is the set of terminals which can be derived from any string which can appear immediately to the right of $A$ in some stage of the derivation.

Construction rules:

Where $\alpha, \beta$ are strings of symbols, and $S, X, Y$ are variables:

1. If $S$ is the start symbol, then $\$ \in FOLLOW(S)$

   (i.e. the start symbol can be followed by the end of the string)

2. If $X \to \alpha Y$ then $FOLLOW(X) \subset FOLLOW(Y)$

   (i.e. anything that can follow $Y$ can follow $X$)

3. If $X \to \alpha Y \beta$ then $FIRST(\beta) \setminus \{\varepsilon\} \subset FOLLOW(Y)$

   (i.e. any terminal which can start $\beta$ can follow $Y$)

4. If $X \to \alpha Y \beta$, $\varepsilon \in FIRST(\beta)$ then
   $FOLLOW(X) \subset FOLLOW(Y)$

   (i.e. if $X$ can derive a string ending in $Y$, anything that follows $X$ can follow $Y$)

## EXAMPLES CALCULATING FOLLOW SETS

$$S \to BC \mid a$$
$$B \to bB \mid \varepsilon$$
$$C \to cC \mid \varepsilon$$

$FOLLOW(S) =$

## EXAMPLES CALCULATING FOLLOW SETS

$$S \rightarrow BC \mid a$$
$$B \rightarrow bB \mid \varepsilon$$
$$C \rightarrow cC \mid \varepsilon$$

$FOLLOW(S) = \{\$\}$          only rule 1 applied

$FOLLOW(C) =$

## EXAMPLES CALCULATING FOLLOW SETS

$$S \rightarrow BC \mid a$$
$$B \rightarrow bB \mid \varepsilon$$
$$C \rightarrow cC \mid \varepsilon$$

$FOLLOW(S) = \{\$\}$      only rule 1 applied

$FOLLOW(C) = FOLLOW(S)$      only rule 4 applied

$\qquad\qquad\quad =$

## EXAMPLES CALCULATING FOLLOW SETS

$$S \to BC \mid a$$
$$B \to bB \mid \varepsilon$$
$$C \to cC \mid \varepsilon$$

$FOLLOW(S) = \{\$\}$        only rule 1 applied

$FOLLOW(C) = FOLLOW(S)$        only rule 4 applied

$\qquad\qquad = \{\$\}$

$FOLLOW(B) =$

## EXAMPLES CALCULATING FOLLOW SETS

$$S \rightarrow BC \mid a$$
$$B \rightarrow bB \mid \varepsilon$$
$$C \rightarrow cC \mid \varepsilon$$

$FOLLOW(S) = \{\$\}$      only rule 1 applied

$FOLLOW(C) = FOLLOW(S)$      only rule 4 applied

$\qquad\qquad = \{\$\}$

$FOLLOW(B) = FIRST(C) \setminus \{\varepsilon\}$      rule 3

## EXAMPLES CALCULATING FOLLOW SETS

$$S \to BC \mid a$$
$$B \to bB \mid \varepsilon$$
$$C \to cC \mid \varepsilon$$

$$FOLLOW(S) = \{\$\} \qquad \text{only rule 1 applied}$$
$$FOLLOW(C) = FOLLOW(S) \qquad \text{only rule 4 applied}$$
$$= \{\$\}$$
$$FOLLOW(B) = FIRST(C) \setminus \{\varepsilon\} \qquad \text{rule 3}$$
$$\cup \, FOLLOW(C) \qquad \text{rule 4}$$
$$=$$

## Examples calculating FOLLOW sets

$$S \rightarrow BC \mid a$$
$$B \rightarrow bB \mid \varepsilon$$
$$C \rightarrow cC \mid \varepsilon$$

$$FOLLOW(S) = \{\$\} \qquad \text{only rule 1 applied}$$
$$FOLLOW(C) = FOLLOW(S) \qquad \text{only rule 4 applied}$$
$$= \{\$\}$$
$$FOLLOW(B) = FIRST(C) \setminus \{\varepsilon\} \qquad \text{rule 3}$$
$$\cup FOLLOW(C) \qquad \text{rule 4}$$
$$= \{c\} \cup \{\$\}$$
$$=$$

## EXAMPLES CALCULATING FOLLOW SETS

$$S \rightarrow BC \mid a$$
$$B \rightarrow bB \mid \varepsilon$$
$$C \rightarrow cC \mid \varepsilon$$

$$
\begin{aligned}
FOLLOW(S) &= \{\$\} && \text{only rule 1 applied} \\
FOLLOW(C) &= FOLLOW(S) && \text{only rule 4 applied} \\
&= \{\$\} \\
FOLLOW(B) &= FIRST(C) \setminus \{\varepsilon\} && \text{rule 3} \\
&\quad \cup FOLLOW(C) && \text{rule 4} \\
&= \{c\} \cup \{\$\} \\
&= \{c, \$\}
\end{aligned}
$$

## Examples calculating FOLLOW sets

$$B \rightarrow (B)B \mid \varepsilon$$

$FOLLOW(B) =$

# EXAMPLES CALCULATING FOLLOW SETS

$$B \rightarrow (B)B \mid \varepsilon$$

$$
\begin{aligned}
FOLLOW(B) &= \{\$\} & \text{rule 1} \\
&\cup FIRST\big({}^{\text{``}})B{}^{\text{''}}\big) \setminus \{\varepsilon\} & \text{rule 3} \\
&=
\end{aligned}
$$

## EXAMPLES CALCULATING FOLLOW SETS

$$B \rightarrow (B)B \mid \varepsilon$$

$$
\begin{aligned}
FOLLOW(B) &= \{\$\} && \text{rule 1}\\
&\quad \cup FIRST\big(``)B"\big) \setminus \{\varepsilon\} && \text{rule 3}\\
&= \{\$\} \cup \{)\}\\
&= \{\$, )\}
\end{aligned}
$$

CONSTRUCTING THE PARSE TABLE

Rows: one for each variable of the grammar

Columns: one for each terminal of the grammar, and for the end of string marker \$

Steps to fill the table $T$:

1. If there is a rule $R \rightarrow \alpha$ with $a \in FIRST(\alpha)$ then put $\alpha$ in $T[R, a]$

2. If there is a rule $R \rightarrow \alpha$ with $\varepsilon \in FIRST(\alpha)$ and $a \in FOLLOW(R)$, then put $\alpha$ in $T[R, a]$

## EXAMPLE

| $F \rightarrow \alpha$ | $FIRST(\alpha)$ | $FOLLOW(F)$ if $\varepsilon \in FIRST(\alpha)$ |
|---|---|---|
| $S \rightarrow BC$ | | |
| $S \rightarrow a$ | | |
| $B \rightarrow bB$ | | |
| $B \rightarrow \varepsilon$ | | |
| $C \rightarrow cC$ | | |
| $C \rightarrow \varepsilon$ | | |

Parse table:

|   | $a$ | $b$ | $c$ | \$ |
|---|---|---|---|---|
| S | | | | |
| B | | | | |
| C | | | | |

48/52

## EXAMPLE

| $F \to \alpha$ | $FIRST(\alpha)$ | $FOLLOW(F)$ if $\varepsilon \in FIRST(\alpha)$ |
|---|---|---|
| $S \to BC$ | $FIRST(BC) = \{b, c, \varepsilon\}$ | |
| $S \to a$ | $FIRST(a) = \{a\}$ | |
| $B \to bB$ | $FIRST(bB) = \{b\}$ | |
| $B \to \varepsilon$ | $FIRST(\varepsilon) = \{\varepsilon\}$ | |
| $C \to cC$ | $FIRST(cC) = \{c\}$ | |
| $C \to \varepsilon$ | $FIRST(\varepsilon) = \{\varepsilon\}$ | |

Parse table:

|   | $a$ | $b$ | $c$ | \$ |
|---|---|---|---|---|
| S |   |   |   |   |
| B |   |   |   |   |
| C |   |   |   |   |

48/52

## EXAMPLE

| $F \to \alpha$ | $FIRST(\alpha)$ | $FOLLOW(F)$ if $\varepsilon \in FIRST(\alpha)$ |
|---|---|---|
| $S \to BC$ | $FIRST(BC) = \{b, c, \varepsilon\}$ | $FOLLOW(S) = \{\$\}$ |
| $S \to a$ | $FIRST(a) = \{a\}$ | |
| $B \to bB$ | $FIRST(bB) = \{b\}$ | |
| $B \to \varepsilon$ | $FIRST(\varepsilon) = \{\varepsilon\}$ | $FOLLOW(B) = \{c, \$\}$ |
| $C \to cC$ | $FIRST(cC) = \{c\}$ | |
| $C \to \varepsilon$ | $FIRST(\varepsilon) = \{\varepsilon\}$ | $FOLLOW(C) = \{\$\}$ |

Parse table:

|   | $a$ | $b$ | $c$ | $\$$ |
|---|---|---|---|---|
| S |   |   |   |   |
| B |   |   |   |   |
| C |   |   |   |   |

48/52

## EXAMPLE

| $F \to \alpha$ | $FIRST(\alpha)$ | $FOLLOW(F)$ if $\varepsilon \in FIRST(\alpha)$ |
|---|---|---|
| $S \to BC$ | $FIRST(BC) = \{b, c, \varepsilon\}$ | $FOLLOW(S) = \{\$\}$ |
| $S \to a$ | $FIRST(a) = \{a\}$ | |
| $B \to bB$ | $FIRST(bB) = \{b\}$ | |
| $B \to \varepsilon$ | $FIRST(\varepsilon) = \{\varepsilon\}$ | $FOLLOW(B) = \{c, \$\}$ |
| $C \to cC$ | $FIRST(cC) = \{c\}$ | |
| $C \to \varepsilon$ | $FIRST(\varepsilon) = \{\varepsilon\}$ | $FOLLOW(C) = \{\$\}$ |

Parse table:

| | $a$ | $b$ | $c$ | $\$$ |
|---|---|---|---|---|
| S | $a$ | $BC$ | $BC$ | $BC$ |
| B | | $bB$ | $\varepsilon$ | $\varepsilon$ |
| C | | | $cC$ | $\varepsilon$ |

## PARSING WITH A PARSING TABLE

1. Append the end of input marker $ to the input and push $ to the stack

2. Put the *start variable* on the stack and scan the *first token*

3. Repeat the following:
    3.1 If *the top of the stack is a variable symbol* $V$, and the current token is $a$, then pop $V$ and push the string from the table entry $(V, a)$. If the entry was empty, reject the input.
    3.2 else *if the top of the stack is a terminal symbol* $t$, compare $t$ to $a$. If they match, pop the stack and scan the next token. Otherwise, reject the input.
    3.3 else *if the top of the stack and the token are both $*, then accept the input (the stack is empty and we have used all the input.)
    3.4 else reject the input (the stack is empty but there is unread input.)

# PARSING *bcc*

| Remaining input | Stack |
|---|---|

|   | $a$ | $b$ | $c$ | \$ |
|---|---|---|---|---|
| S | $a$ | $BC$ | $BC$ | $BC$ |
| B |   | $bB$ | $\varepsilon$ | $\varepsilon$ |
| C |   |   | $cC$ | $\varepsilon$ |

## PARSING $bcc$

| Remaining input | Stack |
|---|---|
| $bcc\$$ | $S\$$ |
| $bcc\$$ | $BC\$$ |
| $bcc\$$ | $bBC\$$ |
| $cc\$$ | $BC\$$ |
| $cc\$$ | $C\$$ |
| $cc\$$ | $cC\$$ |
| $c\$$ | $C\$$ |
| $c\$$ | $cC\$$ |
| $\$$ | $C\$$ |
| $\$$ | $\$$ |
| ACCEPTED | |

|  | $a$ | $b$ | $c$ | $\$$ |
|---|---|---|---|---|
| S | $a$ | $BC$ | $BC$ | $BC$ |
| B | | $bB$ | $\varepsilon$ | $\varepsilon$ |
| C | | | $cC$ | $\varepsilon$ |

# PARSING *bcbc*

| Remaining input | Stack |
|---|---|

|   | $a$ | $b$ | $c$ | $\$$ |
|---|---|---|---|---|
| S | $a$ | $BC$ | $BC$ | $BC$ |
| B |   | $bB$ | $\varepsilon$ | $\varepsilon$ |
| C |   |   | $cC$ | $\varepsilon$ |

## PARSING *bcbc*

| Remaining input | Stack |
|---|---|
| *bcbc*$ | *S*$ |
| *bcbc*$ | *BC*$ |
| *bcbc*$ | *bBC*$ |
| *cbc*$ | *BC*$ |
| *cbc*$ | *C*$ |
| *cbc*$ | *cC*$ |
| *bc*$ | *C*$ |
| REJECTED | |

|   | *a* | *b* | *c* | $ |
|---|---|---|---|---|
| S | *a* | *BC* | *BC* | *BC* |
| B |   | *bB* | *ε* | *ε* |
| C |   |   | *cC* | *ε* |

Push-down Automata

- ▶ "NFA with a stack"
- ▶ CFG to PDA construction method
- ▶ Recognises the set of CFL
- ▶ Non-deterministic PDA are *more powerful* than D-PDA

Parsing

- ▶ LL(k) parsers look ahead up to k symbols
- ▶ Not all CFG are LL(k)
- ▶ $FIRST(\alpha)$: set of terminals (or $\varepsilon$) which start strings derived from $\alpha$
- ▶ $FOLLOW(X)$: the set of terminals (or \$) which could start strings following $X$ in a derivation
- ▶ How to build a parse table for an LL(1) CFG
- ▶ How to parse a string using an LL(1) parse table