

Compilation

0368-3133 (Semester A, 2013/14)

Lecture 4: Syntax Analysis
(Top-Down Parsing)

Modern Compiler Design: Chapter 2.2

Noam Rinetzky

Slides credit: Roman Manevich, Mooly Sagiv, Jeff Ullman, Eran Yahav

1

Admin

- Next week: Trubowicz 101 (Law school)
- Mobiles ...

2

What is a Compiler?

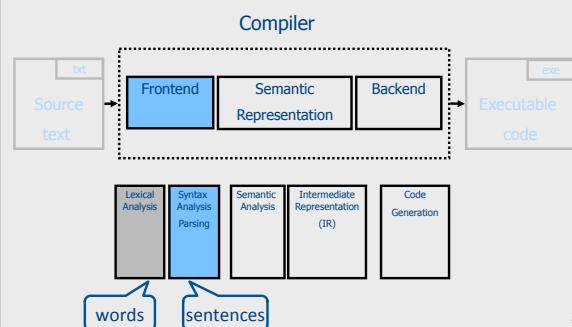
“A compiler is a **computer program** that **transforms** source code written in a programming language (**source language**) into another language (**target language**).

The most common reason for wanting to transform source code is to create an **executable program**.”

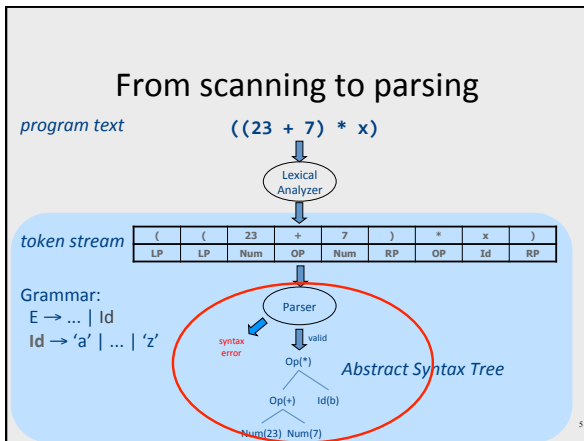
--Wikipedia

3

Conceptual Structure of a Compiler



4



- ### Context Free Grammars
- $G = (V, T, P, S)$
- V – non terminals (syntactic variables)
 - T – terminals (tokens)
 - P – derivation rules
 - Each rule of the form $V \rightarrow (T \cup V)^*$
 - S – start symbol
- 6

CFG terminology

$S \rightarrow S ; S$

$S \rightarrow \text{id} := E$

$E \rightarrow \text{id}$

$E \rightarrow \text{num}$

$E \rightarrow E + E$

Symbols:

Terminals (tokens): ; := () id num print

Non-terminals: S E L

Start non-terminal: S

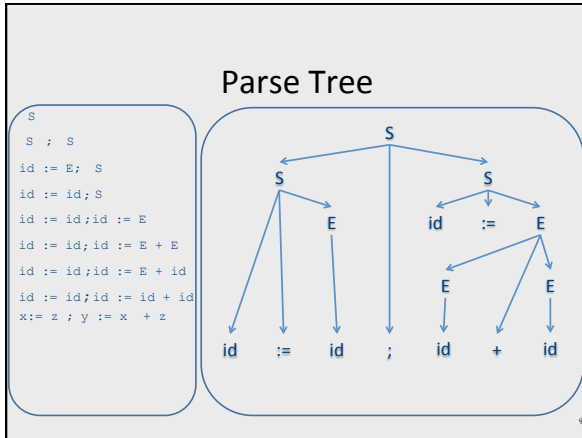
Convention: the non-terminal appearing in the first derivation rule

Grammar productions (rules)

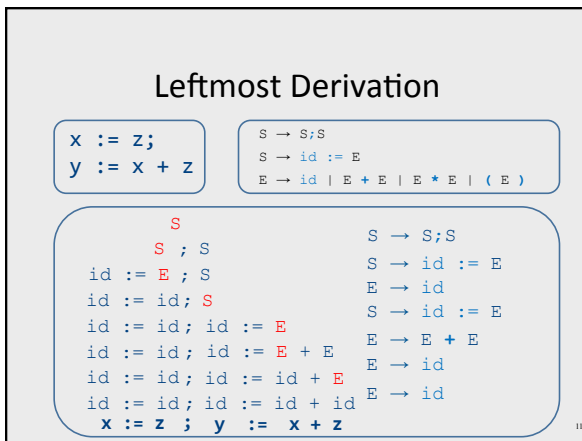
$N \rightarrow \mu$

7

- ### CFG terminology
- **Derivation** - a sequence of replacements of non-terminals using the derivation rules
 - **Language** - the set of strings of terminals derivable from the start symbol
 - **Sentential form** - the result of a partial derivation in which there may be non-terminals
- 8



- ### Leftmost/rightmost Derivation
- Leftmost derivation
 - always expand leftmost non-terminal
 - Rightmost derivation
 - Always expand rightmost non-terminal
 - Allows us to describe derivation by listing the sequence of rules
 - always know what a rule is applied to



- ### Broad kinds of parsers
- Parsers for **arbitrary** grammars
 - Earley's method, CYK method
 - Usually, not used in practice (though might change)
 - **Top-Down** parsers
 - Construct parse tree in a top-down matter
 - Find the leftmost derivation
 - Linear **Bottom-Up** parsers
 - Construct parse tree in a bottom-up manner
 - Find the rightmost derivation in a reverse order

Intuition: Top-Down Parsing

- Begin with start symbol
- “Guess” the productions
- Check if parse tree yields user's program

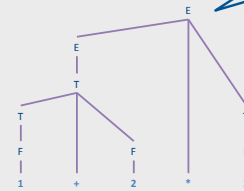
13

Intuition: Top-Down parsing

Unambiguous grammar

$E \rightarrow E * T$
 $E \rightarrow T$
 $T \rightarrow T + F$
 $T \rightarrow F$
 $F \rightarrow \text{id}$
 $F \rightarrow \text{num}$
 $F \rightarrow (E)$

Recall: Non standard precedence ...



14

Intuition: Top-Down parsing

Unambiguous grammar

$E \rightarrow E * T$
 $E \rightarrow T$
 $T \rightarrow T + F$
 $T \rightarrow F$
 $F \rightarrow \text{id}$
 $F \rightarrow \text{num}$
 $F \rightarrow (E)$

We need this rule to get the *

E

1 + 2 * 3

15

Intuition: Top-Down parsing

Unambiguous grammar

$E \rightarrow E * T$
 $E \rightarrow T$
 $T \rightarrow T + F$
 $T \rightarrow F$
 $F \rightarrow \text{id}$
 $F \rightarrow \text{num}$
 $F \rightarrow (E)$



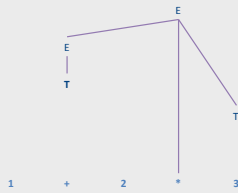
1 + 2 * 3

16

Intuition: Top-Down parsing

Unambiguous grammar

$E \rightarrow E * T$
 $E \rightarrow T$
 $T \rightarrow T + F$
 $T \rightarrow F$
 $F \rightarrow \text{id}$
 $F \rightarrow \text{num}$
 $F \rightarrow (E)$

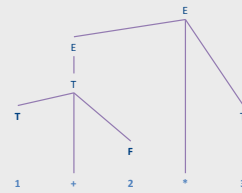


17

Intuition: Top-Down parsing

Unambiguous grammar

$E \rightarrow E * T$
 $E \rightarrow T$
 $T \rightarrow T + F$
 $T \rightarrow F$
 $F \rightarrow \text{id}$
 $F \rightarrow \text{num}$
 $F \rightarrow (E)$

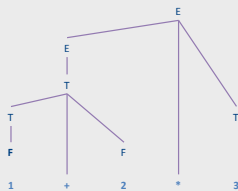


18

Intuition: Top-Down parsing

Unambiguous grammar

$E \rightarrow E * T$
 $E \rightarrow T$
 $T \rightarrow T + F$
 $T \rightarrow F$
 $F \rightarrow \text{id}$
 $F \rightarrow \text{num}$
 $F \rightarrow (E)$

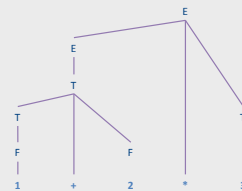


19

Intuition: Top-Down parsing

Unambiguous grammar

$E \rightarrow E * T$
 $E \rightarrow T$
 $T \rightarrow T + F$
 $T \rightarrow F$
 $F \rightarrow \text{id}$
 $F \rightarrow \text{num}$
 $F \rightarrow (E)$

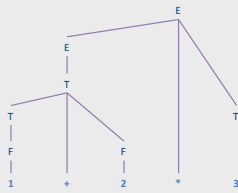


20

Intuition: Top-Down parsing

Unambiguous grammar

$E \rightarrow E * T$
 $E \rightarrow T$
 $T \rightarrow T + F$
 $T \rightarrow F$
 $F \rightarrow \text{id}$
 $F \rightarrow \text{num}$
 $F \rightarrow (E)$

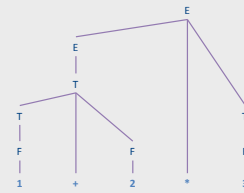


21

Intuition: Top-Down parsing

Unambiguous grammar

$E \rightarrow E * T$
 $E \rightarrow T$
 $T \rightarrow T + F$
 $T \rightarrow F$
 $F \rightarrow \text{id}$
 $F \rightarrow \text{num}$
 $F \rightarrow (E)$

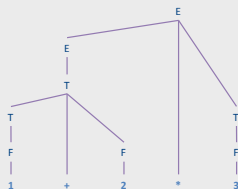


22

Intuition: Top-Down parsing

Unambiguous grammar

$E \rightarrow E * T$
 $E \rightarrow T$
 $T \rightarrow T + F$
 $T \rightarrow F$
 $F \rightarrow \text{id}$
 $F \rightarrow \text{num}$
 $F \rightarrow (E)$



23

Challenges in top-down parsing

- Top-down parsing begins with virtually no information
 - Begins with just the start symbol, which matches every program
- How can we know which productions to apply?

24

Recursive Descent

- Blind exhaustive search
 - Goes over all possible production rules
 - Read & parse prefixes of input
 - Backtracks if guesses wrong
- Implementation
 - Uses (possibly recursive) functions for every production rule
 - Backtracks → “rewind” input

25

Recursive descent

```
bool A() { // A → A1 | ... | An
    pos = recordCurrentPosition();
    for (i = 1; i ≤ n; i++) {
        if (Ai())
            return true;
        rewindCurrent(pos);
    }
    return false;
}

bool Ai() { // Ai = X1X2...Xk
    for (j=1; j ≤ k; j++)
        if (Xj is a terminal)
            if (Xj == current) match(current);
            else return false;
        else if (!Xj()) return false;
    return true;
}
```

current

↓

((23	+	7)	*	x)
LP	LP	Num	OP	Num	RP	OP	Id	RP

token stream

26

Example

- Grammar
 - $E \rightarrow T \mid T + E$
 - $T \rightarrow \text{int} \mid \text{int} * T \mid (E)$
- Input: (5)
- Token stream: LP int RP

27

Problem: Left Recursion

p. 127

$E \rightarrow E - \text{term} \mid \text{term}$

```
int E() {
    return E() && match(token('-')) && term();
}
```

- What happens with this procedure?
- Recursive descent parsers cannot handle left-recursive grammars

28

p. 130

Left recursion removal

$N \rightarrow N\alpha \mid \beta$

→

$N \rightarrow \beta N'$
 $N' \rightarrow \alpha N' \mid \epsilon$

G_1 G_2

- $L(G_1) = \beta, \beta\alpha, \beta\alpha\alpha, \beta\alpha\alpha\alpha, \dots$
- $L(G_2) = \text{same}$

Can be done algorithmically.
 Problem: grammar becomes mangled beyond recognition

- For our 3rd example:

$E \rightarrow E - \text{term} \mid \text{term}$

→

$E \rightarrow \text{term } TE \mid \text{term}$
 $TE \rightarrow - \text{term } TE \mid \epsilon$

29

Challenges in top-down parsing

- Top-down parsing begins with virtually no information
 - Begins with just the start symbol, which matches every program
- How can we know which productions to apply?
- **Wanted:** Top-Down parsing without backtracking

30

Predictive parsing

- Given a grammar G and a word w derive w using G
 - Apply production to leftmost nonterminal
 - Pick production rule based on next input token
- General grammar
 - More than one option for choosing the next production based on a token
- Restricted grammars (LL)
 - Know exactly which single rule to apply based on
 - Non terminal
 - Next (k) tokens (lookahead)

31

Boolean expressions example

$E \rightarrow LIT \mid (E \text{ OP } E) \mid \text{not } E$
 $LIT \rightarrow \text{true} \mid \text{false}$
 $OP \rightarrow \text{and} \mid \text{or} \mid \text{xor}$

not (not true or false)

32

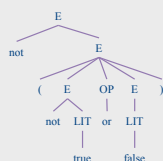
Boolean expressions example

$E \rightarrow \text{LIT} \mid (E \text{ OP } E) \mid \text{not } E$
 $\text{LIT} \rightarrow \text{true} \mid \text{false}$
 $\text{OP} \rightarrow \text{and} \mid \text{or} \mid \text{xor}$

production to apply known from next token

`not (not true or false)`

$E \Rightarrow$
 $\text{not } E \Rightarrow$
 $\text{not } (E \text{ OP } E) \Rightarrow$
 $\text{not } (\text{not } E \text{ OP } E) \Rightarrow$
 $\text{not } (\text{not } \text{LIT} \text{ OP } E) \Rightarrow$
 $\text{not } (\text{not } \text{true} \text{ OP } E) \Rightarrow$
 $\text{not } (\text{not } \text{true} \text{ or } E) \Rightarrow$
 $\text{not } (\text{not } \text{true} \text{ or } \text{LIT}) \Rightarrow$
 $\text{not } (\text{not } \text{true} \text{ or } \text{false})$



33

Problem: productions with common prefix

$\text{term} \rightarrow \text{ID} \mid \text{ID} [\text{expr}]$

- Cannot tell which rule to use based on lookahead (ID)

34

Solution: left factoring

- Rewrite the grammar to be in LL(1)

$\text{term} \rightarrow \text{ID} \mid \text{ID} [\text{expr}]$



$\text{term} \rightarrow \text{ID after_ID}$
 $\text{After_ID} \rightarrow [\text{expr}] \mid \epsilon$

Intuition: just like factoring $x*y + x*z$ into $x*(y+z)$

35

Left factoring – another example

$S \rightarrow \text{if } E \text{ then } S \text{ else } S$
 $\quad \mid \text{if } E \text{ then } S$
 $\quad \mid T$



$S \rightarrow \text{if } E \text{ then } S'$
 $\quad \mid T$
 $S' \rightarrow \text{else } S \mid \epsilon$

36

LL(k) Parsers

- Predictive parser
 - Can be generated automatically
 - Does not use recursion
 - Efficient
- In contrast, recursive descent
 - Manual construction
 - Recursive
 - Expensive

37

LL(k) parsing via pushdown automata and prediction table

- Pushdown automaton uses
 - Prediction stack
 - Input token stream
 - Transition table
 - nonterminals x tokens -> production alternative
 - Entry indexed by nonterminal N and token t contains the alternative of N that must be predicted when current input starts with t

38

Example transition table

- (1) $E \rightarrow LIT$
- (2) $E \rightarrow (E O P E)$
- (3) $E \rightarrow \text{not } E$
- (4) $LIT \rightarrow \text{true}$
- (5) $LIT \rightarrow \text{false}$
- (6) $OP \rightarrow \text{and}$
- (7) $OP \rightarrow \text{or}$
- (8) $OP \rightarrow \text{xor}$

Input tokens

	()	not	true	false	and	or	xor	\$
E	2		3	1	1				
LIT				4	5				
OP						6	7	8	

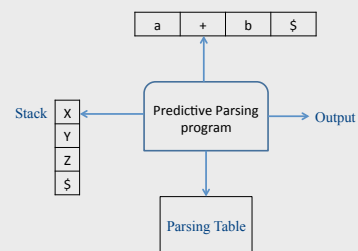
Nonterminals

Table

Which rule should be used

39

Non-Recursive Predictive Parser



40

LL(k) parsing via pushdown automata and prediction table

- Two possible moves
 - Prediction**
 - When top of stack is nonterminal N, pop N, lookup table[N,t]. If table[N,t] is not empty, push table[N,t] on prediction stack, otherwise – syntax error
 - Match**
 - When top of prediction stack is a terminal T, must be equal to next input token t. If (t = T), pop T and consume t. If (t ≠ T) syntax error
- Parsing terminates when prediction stack is empty
 - If input is empty at that point, success. Otherwise, syntax error

41

Running parser example

aacbb\$

A → aAb | c

Input suffix	Stack content	Move
aacbb\$	A\$	predict(A,a) = A → aAb
aacbb\$	aAb\$	match(a,a)
acbb\$	Ab\$	predict(A,a) = A → aAb
acbb\$	aAbb\$	match(a,a)
cbb\$	Abb\$	predict(A,c) = A → c
cbb\$	cbb\$	match(c,c)
bb\$	bb\$	match(b,b)
b\$	b\$	match(b,b)
\$	\$	match(\$,\$) – success

	a	b	c
A	A → aAb		A → c

42

FIRST sets

- $FIRST(\alpha) = \{ t \mid \alpha \rightarrow^* t \beta \} \cup \{ \epsilon \mid \alpha \rightarrow^* \epsilon \}$
 - FIRST(α) = all terminals that α can appear as first in some derivation for α
 - ϵ if α can be derived from X
- Example:
 - FIRST(LIT) = { true, false }
 - FIRST((E OP E)) = { '(' }
 - FIRST(not E) = { not }

43

Computing FIRST sets

- $FIRST(t) = \{ t \}$ // "t" non terminal
- $\epsilon \in FIRST(X)$ if
 - $X \rightarrow \epsilon$ or
 - $X \rightarrow A_1 \dots A_k$ and $\epsilon \in FIRST(A_i) \ i=1\dots k$
- $FIRST(\alpha) \subseteq FIRST(X)$ if
 - $X \rightarrow A_1 \dots A_k \ \alpha$ and $\epsilon \in FIRST(A_i) \ i=1\dots k$

44

FIRST sets computation example

STMT → if EXPR then STMT
 | while EXPR do STMT
 | EXPR ;
 EXPR → TERM
 | zero? TERM
 | not EXPR
 | ++ id
 | -- id
 TERM → id
 | constant

STMT	EXPR	TERM

45

1. Initialization

STMT → if EXPR then STMT
 | while EXPR do STMT
 | EXPR ;
 EXPR → TERM
 | zero? TERM
 | not EXPR
 | ++ id
 | -- id
 TERM → id
 | constant

STMT	EXPR	TERM
if while	zero? Not ++ --	id constant

46

2. Iterate 1

STMT → if EXPR then STMT
 | while EXPR do STMT
 | EXPR ;
 EXPR → TERM
 | zero? TERM
 | not EXPR
 | ++ id
 | -- id
 TERM → id
 | constant

STMT	EXPR	TERM
if while	zero? Not ++ --	id constant
zero? Not ++ --		

47

2. Iterate 2

STMT → if EXPR then STMT
 | while EXPR do STMT
 | EXPR ;
 EXPR → TERM -> id
 | zero? TERM
 | not EXPR
 | ++ id
 | -- id
 TERM → id
 | constant

STMT	EXPR	TERM
if while	zero? Not ++ --	id constant
zero? Not ++ --	id constant	

48

2. Iterate 3 – fixed-point

STMT \rightarrow if EXPR then STMT
 | while EXPR do STMT
 | EXPR ;
 EXPR \rightarrow TERM
 | zero? TERM
 | not EXPR
 | ++ id
 | -- id
 TERM \rightarrow id
 | constant

STMT	EXPR	TERM
if while	zero? Not ++ --	id constant
zero? Not ++ --	id constant	
id constant		

49

FOLLOW sets



- $\text{FOLLOW}(X) = \{ t \mid S \rightarrow^* \alpha X t \beta \}$
 - FOLLOW(X) = set of tokens that can immediately follow X in some sentential form
 - NOT related to what can be derived from X
- Intuition: $X \rightarrow A B$
 - $\text{FIRST}(B) \subseteq \text{FOLLOW}(A)$
 - $\text{FOLLOW}(B) \subseteq \text{FOLLOW}(X)$
 - If $B \rightarrow^* \epsilon$ then $\text{FOLLOW}(A) \subseteq \text{FOLLOW}(X)$

50

FOLLOW sets: Constraints

- $\$ \in \text{FOLLOW}(S)$
- $\text{FIRST}(\beta) - \{\epsilon\} \subseteq \text{FOLLOW}(X)$
 - For each $A \rightarrow \alpha X \beta$
- $\text{FOLLOW}(A) \subseteq \text{FOLLOW}(X)$
 - For each $A \rightarrow \alpha X \beta$ and $\epsilon \in \text{FIRST}(\beta)$

51

Example: FOLLOW sets

- $E \rightarrow TX \quad X \rightarrow + E \mid \epsilon$
- $T \rightarrow (E) \mid \text{int } Y \quad Y \rightarrow * T \mid \epsilon$

Terminal	+	(*)	int
FOLLOW	int, (int, (int, (), \$	*,), +, \$

Non-Term.	E	T	X	Y
FOLLOW), \$	+,), \$,\$)), \$

52

Prediction Table

- $A \rightarrow \alpha$
- $T[A,t] = \alpha$ if $t \in \text{FIRST}(\alpha)$
- $T[A,t] = \alpha$ if $\epsilon \in \text{FIRST}(\alpha)$ and $t \in \text{FOLLOW}(A)$
 - t can also be \$
- T is not well defined \rightarrow the grammar is not LL(1)

53

Problem: Non LL Grammars

```
S → A a b
A → a | ε
```

```
bool S() {
    return A() && match(token('a')) && match(token('b'));
}
```

```
bool A() {
    return match(token('a')) || true;
}
```

- What happens for input "ab"?
- What happens if you flip order of alternatives and try "aab"?

54

Problem: Non LL Grammars

```
S → A a b
A → a | ε
```

- $\text{FIRST}(S) = \{ a \}$ $\text{FOLLOW}(S) = \{ \$ \}$
- $\text{FIRST}(A) = \{ a \epsilon \}$ $\text{FOLLOW}(A) = \{ a \}$
- **FIRST/FOLLOW conflict**

55

Solution: substitution

```
S → A a b
A → a | ε
```

↓ Substitute A in S

```
S → a a b | a b
```

↓ Left factoring

```
S → a after_A
after_A → a b | b
```

56

LL(k) grammars

- A grammar is LL(k) if it can be derived via:
 - Top-down derivation
 - Scanning the input from left to right (L)
 - Producing the leftmost derivation (L)
 - With lookahead of k tokens (k)
 - T is well defined
- A language is said to be LL(k) when it has an LL(k) grammar

57

LL(k) grammars

- A grammar is not LL(k) if it is
 - Ambiguous
 - Left recursive
 - Not left factored
 - ...

58

Earley Parsing

- Invented by Earley [PhD. 1968]
- Handles arbitrary CFG
- Can handle ambiguous grammars
- Complexity $O(N^3)$ when $N = |\text{input}|$
- Uses dynamic programming
 - Compactly encodes ambiguity

59

Dynamic programming

- Break a problem P into subproblems $P_1 \dots P_k$
 - Solve P by combining solutions for $P_1 \dots P_k$
 - Memo-ize (store) solutions to subproblems instead of re-computation
- Bellman-Ford shortest path algorithm
 - $\text{Sol}(x,y,i) = \text{minimum of}$
 - $\text{Sol}(x,y,i-1)$
 - $\text{Sol}(t,y,i-1) + \text{weight}(x,t)$ for edges (x,t)

60

Earley Parsing

- Dynamic programming implementation of a recursive descent parser
 - $S[N+1]$ Sequence of sets of “Earley states”
 - $N = |\text{INPUT}|$
 - Earley state (item) s is a sentential form + aux info
 - $S[i]$ All parse tree that can be produced (by a RDP) after reading the first i tokens
 - $S[i+1]$ built using $S[0] \dots S[i]$

61

Earley States

- $s = \langle \text{constituent, back} \rangle$
 - constituent (dotted rule) for $A \rightarrow \alpha\beta$
 - $A \rightarrow \bullet\alpha\beta$ *predicated* constituents
 - $A \rightarrow \alpha\bullet\beta$ *in-progress* constituents
 - $A \rightarrow \alpha\beta\bullet$ *completed* constituents
 - back previous Early state in derivation

62

Earley Parser

```

Input = x[1...N]
S[0] = <E' → •E, 0>; S[1] = ... S[N] = {}
for i = 0 ... N do
  until S[i] does not change do
    foreach s ∈ S[i]
      if s = <A → ...•a..., b> and a=x[i+1] then
        S[i+1] = S[i+1] ∪ {<A → ...a•..., b> } // scan
      if s = <A → ... •X ..., b> and X → α then
        S[i] = S[i] ∪ {<X → •α, i > } // predict
      if s = <A → ... •, b> and <X → ...•A..., k> ∈ S[b] then
        S[i] = S[i] ∪ {<X → ...A•..., k > } // complete
    
```

63

Example

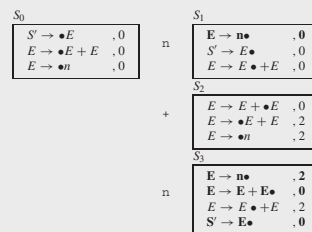


FIGURE 1. Earley sets for the grammar $E \rightarrow E + E \mid n$ and the input $n + n$. Items in bold are ones which correspond to the input's derivation.

64