

Syntax Analysis

Earley Algorithm

Prof. Licia Sbattella

aa 2007-08

Translated and adapted by L. Breveglieri

A GENERAL METHOD FOR SYNTAX ANALYSIS: EARLEY

The EARLEY method (or algorithm) allows to deal with any grammar, even with an ambiguous one, as it constructs all the possible derivations of the string to recognise.

Its time complexity is proportional to the third power of the string length, but is lower if the grammar is unambiguous, and even less if it is deterministic.

The Earley method originates from the LR(k) method, but it does not attempt to build a deterministic pushdown automaton; rather it exploits a data structure more complex than the stack is (actually it uses an array of sets), which represents efficiently many possible stacks, sharing some contents.

In the practice, the Earley algorithm simulates a non-deterministic pushdown automaton, but without falling into the exponential time complexity that indeterminism usually implies.

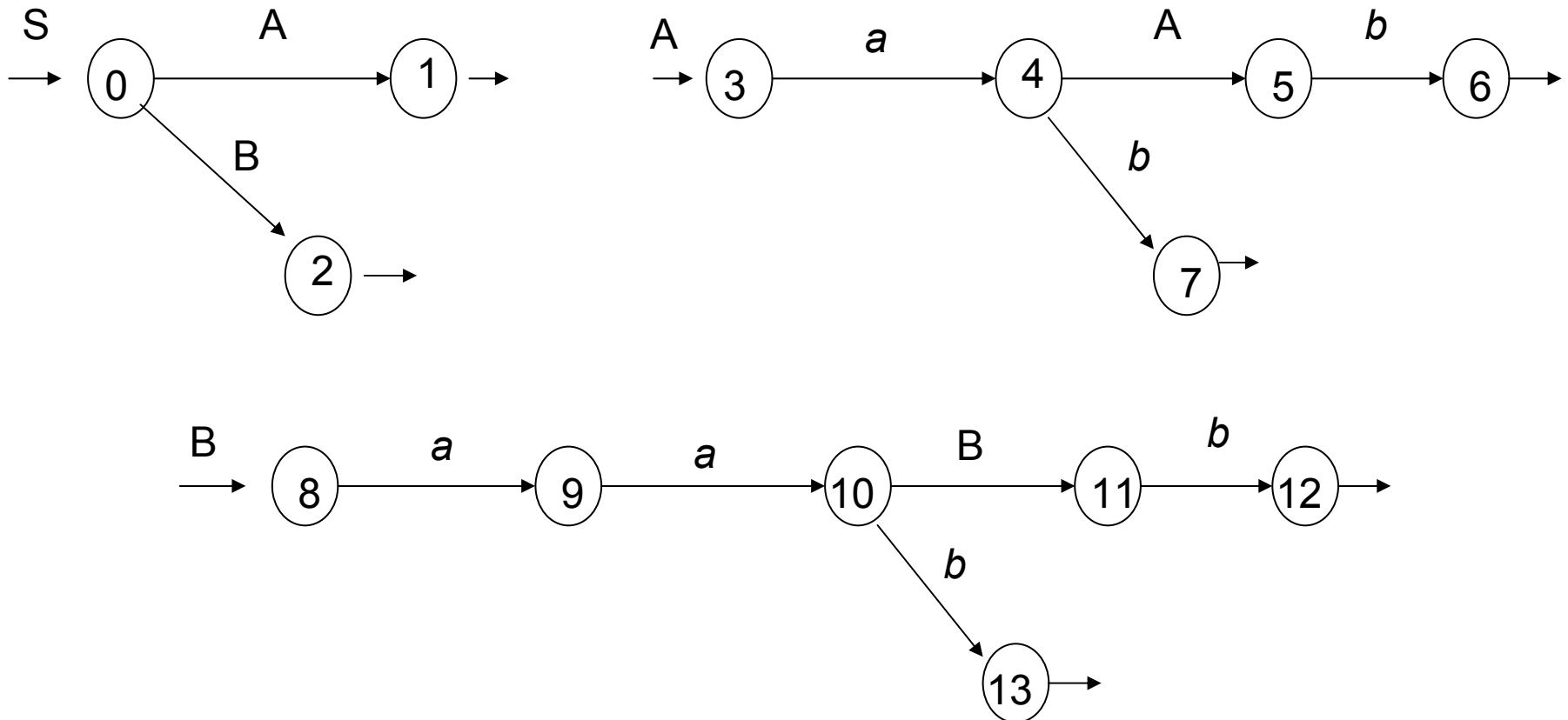
STARTING HYPOTHESIS: THE GRAMMAR DOES NOT CONTAIN NULL RULES

The Earley algorithm can work with or without lookahead, but the asymptotic time complexity is not influenced. For simplicity here lookahead is not considered.

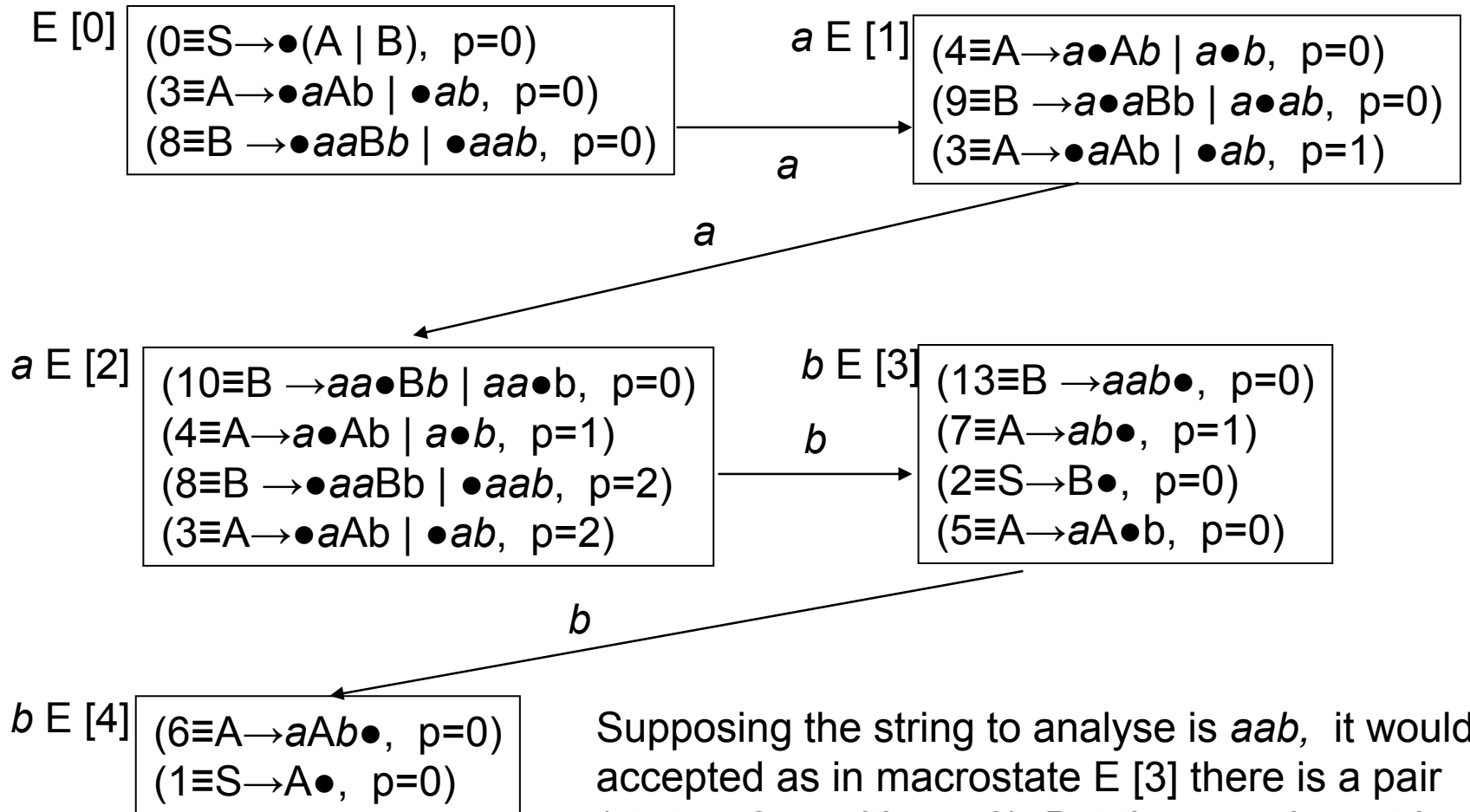
EXAMPLE - a non-LR(k) grammar

$$L = \{a^n b^n \mid n \geq 1\} \cup \{a^{2^n} b^n \mid n \geq 1\}$$

$$S \rightarrow A \mid B \quad A \rightarrow aAb \mid ab \quad B \rightarrow aaBb \mid aab$$



string to analyse: *aabb*



Supposing the string to analyse is *aab*, it would be accepted as in macrostate E [3] there is a pair (state = 2, position = 0). But the complete string is *aabb*, and is accepted in E [4] as such a macrostate contains the pair (state = 1, position = 0).

The Earley algorithm is a top-down analyzer that constructs simultaneously all the possible left derivations of the string. The algorithm reads the string $x_1 \dots x_n$ from left to right and, when it examines character x_i , produces pairs of two elements, as follows:

<machine state, position in the string> denoted as $(s = \dots, p = \dots)$ with $0 \leq p \leq i$

As already seen, the state can be alternatively denoted as marked rule $s \equiv A \rightarrow \alpha \bullet \beta$

Intuitively pair $(s \equiv A \rightarrow \alpha \bullet \beta, p)$ represents a statement and a target:

Statement: there has been found a substring

$x_{j+1 \dots i}$ with $0 \leq j < i$ that derives from α

Target: find all the positions $k > i$ such that

substring $x_{i+1 \dots k}$ derives from β

If the algorithm finds such a position k , it can state that from non-terminal A there derives substring $x_{j+1 \dots k}$

$$\begin{array}{c} * \\ \alpha \Rightarrow x_{j+1 \dots i} \\ * \\ \beta \Rightarrow x_{i+1 \dots k} \\ * \\ A \Rightarrow x_{j+1 \dots k} \end{array}$$

A pair $(q \equiv A \rightarrow \alpha \bullet, j)$, where state q is final (end-marked rule), is said to be completed.

ALGORITHM – Earley recognizer

The task of the Earley recognizer is that of finding the derivation of the entire string x , but the algorithm generally produces more than it is strictly requested: it also finds whether any prefix u of the string $x = uv$ belongs to the language itself.

The Earley algorithm builds a vector $E [0 \dots n]$ (n is the string length) dimensioned to the string length, the elements of which are pairs.

PASS 1: *Initialization* - set the targets to find all the prefixes of x that are derivable from the axiom S . The initial set is filled with pairs obtained from the axiom S and all the other sets are initially empty.

$$\begin{aligned} E[0] &:= \left\{ (q_\alpha, 0) \mid q_\alpha \text{ is the network initial state} \right\} \\ E[i] &:= 0, \quad \text{for } i = 1, \dots, n \\ i &:= 0 \end{aligned}$$

PASS 2: *Derivation Construction* - then the operations of prediction, completion and shift are applied in the natural order $0, 1, \dots, n$, in order to compute all the sets $E[i]$.

At step i the algorithm can add elements only to the current set $E[i]$ and to the consecutive one $E[i + 1]$.

If none of the operations has added new pairs to $E[i]$, the algorithm moves to consider set $E[i + 1]$.

The algorithm stops and rejects the string if $E[i + 1]$ is empty and $i < n$.

The algorithm terminates with success and accepts the string when the set $E[n]$ has been fully filled and contains (at least) one completed pair $(q \equiv S \rightarrow \alpha \bullet, 0)$ of the axiom (clearly q is the network initial state).

PREDICTION OPERATION (CLOSURE)

Every target in the set $E[i]$ can add new sub-targets to the set itself. Whenever a new sub-target is added, its position pointer is set to the current position i .

for every pair $(q \equiv A \rightarrow \alpha \bullet B \gamma, j)$ already in $E[i]$
add to the set $E[i]$ the new pair $(r \equiv B \rightarrow \bullet \dots, i)$
where r is the initial state of the machine expanding B

SHIFT OPERATION

Update the targets in the set $E[i + 1]$ depending on the next character in the input string. The position pointer of the updated targets is the same as that of the current target.

for every pair $(q \equiv A \rightarrow \alpha \bullet a \gamma, j)$ already in $E[i]$
if $x_{i+1} = a$ then
add to the set $E[i + 1]$ the new pair $(r \equiv A \rightarrow \alpha a \bullet \gamma, j)$

Each character x_{i+1} of the input string is examined by the algorithm only once.

COMPLETION OPERATION

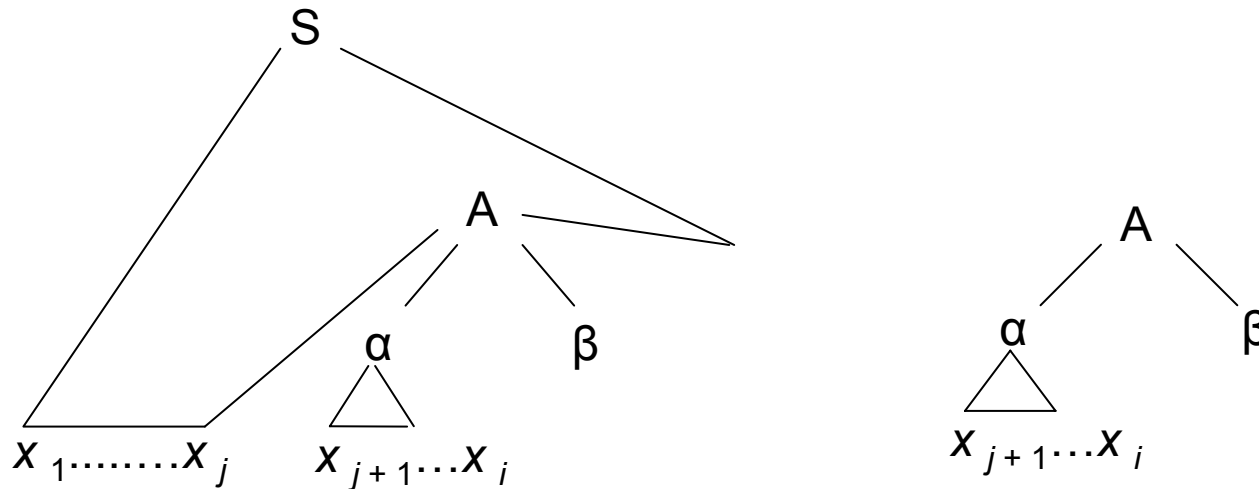
A completed pair ($q \equiv A \rightarrow \alpha \bullet, j$) states that the derivation of string $x_{j+1 \dots i}$ has been identified starting from the non-terminal A . It is necessary to update set $E [i]$ and add to it such a statement.

for every completed pair ($q \equiv A \rightarrow \alpha \bullet, j$) already in $E [i]$
for every pair ($r \equiv B \rightarrow \beta \bullet A \gamma, k$) already in $E [j]$
add to the set $E [i]$ the new pair ($s \equiv B \rightarrow \beta A \bullet \gamma, k$)

On completion, if pair ($r \equiv B \rightarrow \beta \bullet A \gamma, k$) $\in E [j]$ and the string γ is null, the new pair ($s \equiv B \rightarrow \beta A \bullet \gamma, k$), just added to the set $E [i]$, is itself completed, and therefore it is necessary to iterate the completion operation again.

It is easy to prove that the Earley algorithm accepts only strings belonging to the language $L(G)$: a new pair is added to a set only if the derivation that it states is possible in G . It is more complex to prove that every phrase of the language is recognised by the algorithm (proof omitted here).

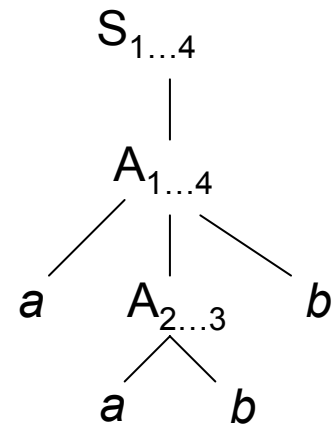
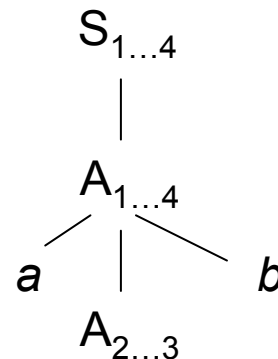
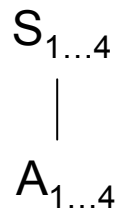
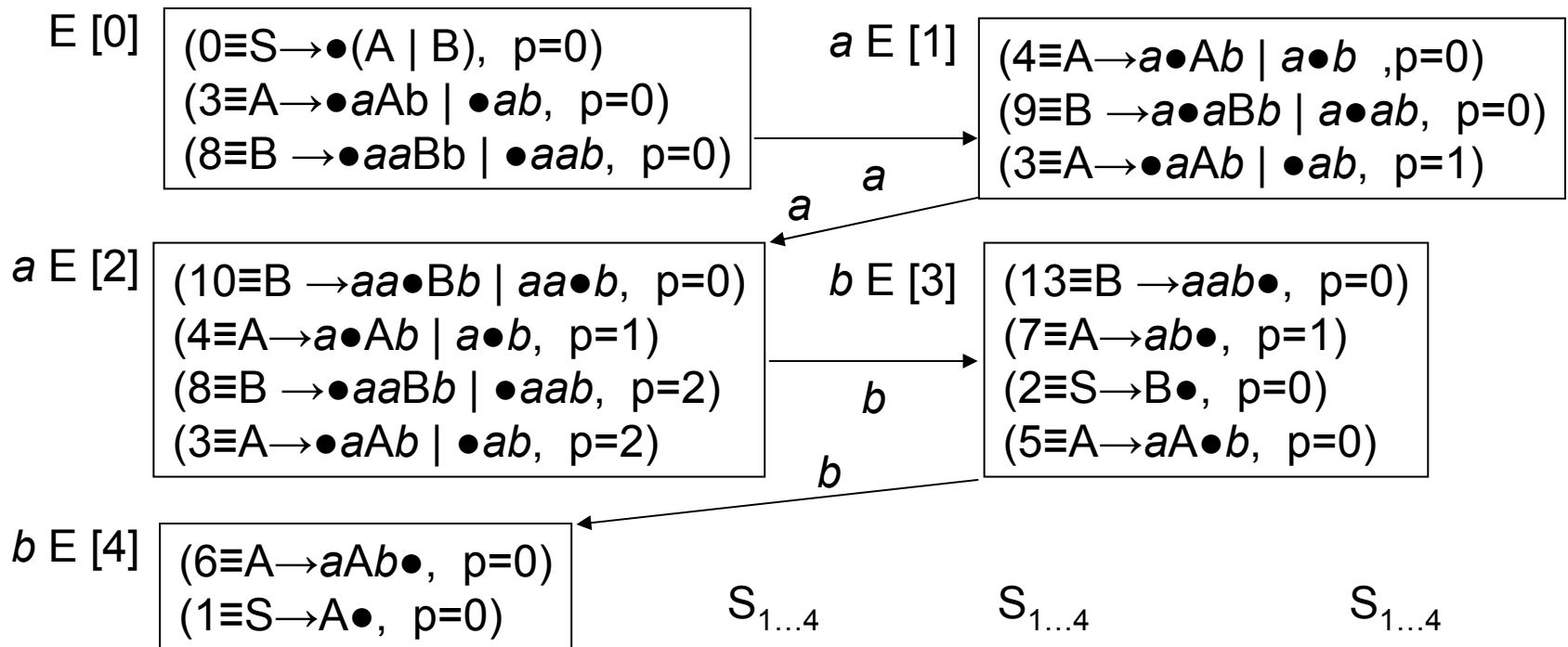
HOW TO RECONSTRUCT THE SYNTAX TREE PROGRESSIVELY



It is instructive to imagine the three operations as a progressive construction of the syntax tree: to each array set E there corresponds a set of trees.

In $E[i]$ there is a pair $(s \equiv A \rightarrow \alpha \bullet \beta, p = j)$ if and only if grammar G admits a syntax tree as shown above (left tree). The right tree shows the specific subtree associated with the pair.

At the end of the algorithm, the string is accepted if and only if one of the trees associated with the last set $E[n]$, is a complete syntax tree with axiomatic root.



If the grammar is ambiguous, analysing a phrase produces all the possible syntax trees, represented by sharing common factors.

EXAMPLE – parsing an ambiguous language – the grammar contains a two-sided recursive rule, hence is ambiguous

phrase $a + a + a$

$$S \rightarrow E \quad E \rightarrow E + E \quad E \rightarrow a$$

$$E [0] \begin{array}{|l} S \rightarrow \bullet E, 0 \\ E \rightarrow \bullet E + E \mid \bullet a, 0 \end{array}$$

$$a E [1] \begin{array}{|l} E \rightarrow a \bullet, 0 \\ S \rightarrow E \bullet, 0 \\ E \rightarrow E \bullet + E, 0 \end{array}$$

$$+ E [2] \begin{array}{|l} E \rightarrow E + \bullet E, 0 \\ E \rightarrow \bullet E + E \mid \bullet a, 2 \end{array}$$

$$a E [3] \begin{array}{|l} E \rightarrow a \bullet, 2 \\ E \rightarrow E + E \bullet, 0 \\ E \rightarrow E \bullet + E, 2 \\ S \rightarrow E \bullet, 0 \\ E \rightarrow E \bullet + E, 0 \end{array}$$

$$+ E [4] \begin{array}{|l} E \rightarrow E + \bullet E, 0 \\ E \rightarrow E + \bullet E, 2 \\ E \rightarrow \bullet E + E \mid \bullet a, 4 \end{array}$$

$$a E [5] \begin{array}{|l} E \rightarrow a \bullet, 4 \\ E \rightarrow E + E \bullet, 2 \\ E \rightarrow E \bullet + E, 4 \\ E \rightarrow E + E \bullet, 0 \\ E \rightarrow E \bullet + E, 2 \\ S \rightarrow E \bullet, 0 \\ E \rightarrow E \bullet + E, 0 \end{array}$$

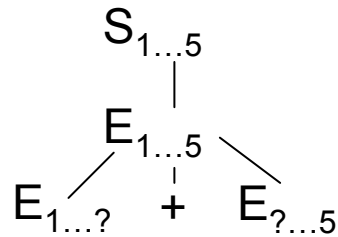
Pair

1st Tree

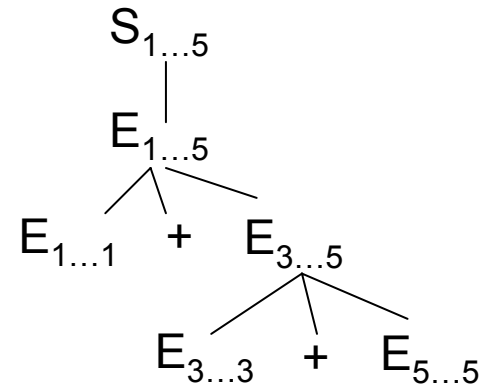
Pair

2nd Tree

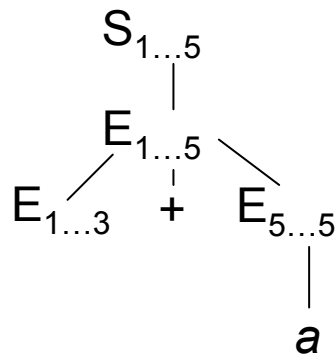
E [5]
S → E●, 0
E → E+E●, 0



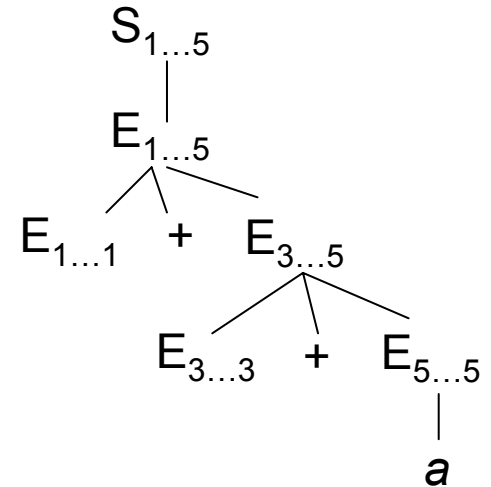
E [5]
S → E●, 0
E → E+E●, 0
E → E+E●, 2



E → a●, 4



E → a●, 4



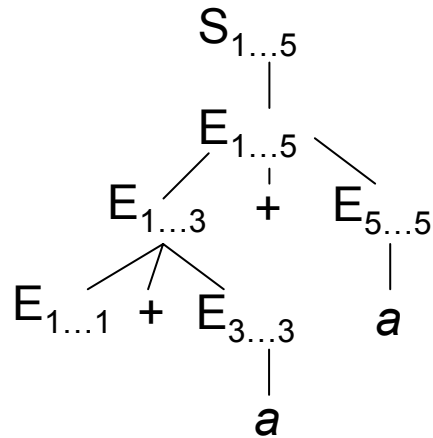
Pair

1st Tree

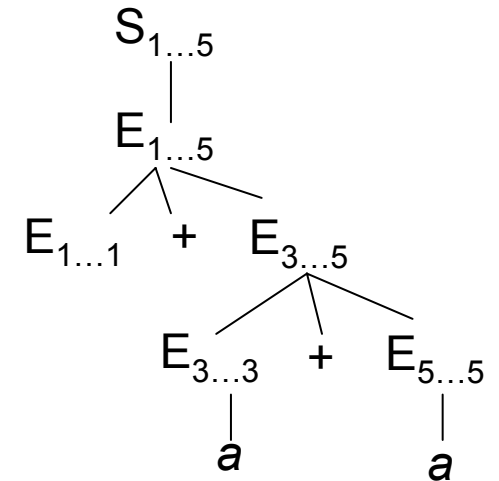
Pair

2nd Tree

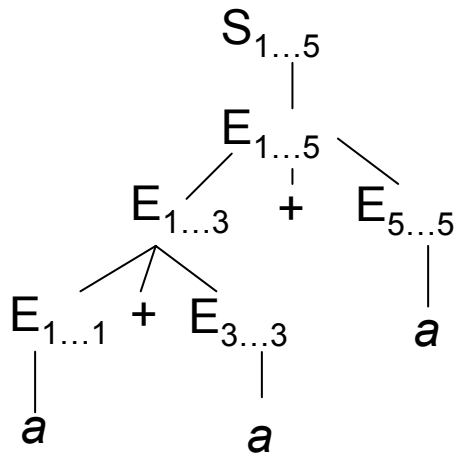
$E [3]$
$E \rightarrow E + E \bullet, 0$
$E \rightarrow a \bullet, 2$



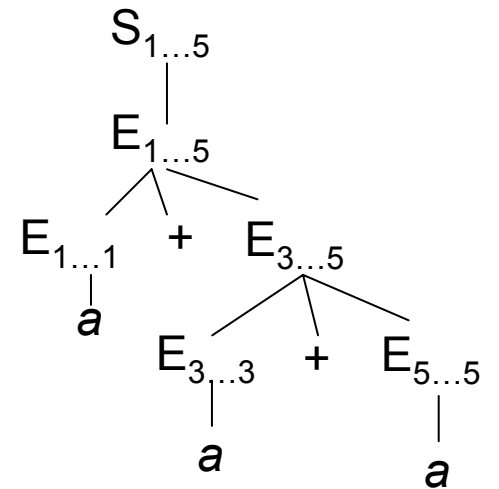
$E [3]$
$E \rightarrow a \bullet, 2$



$E [1]$
$E \rightarrow a \bullet, 0$



$E [1]$
$E \rightarrow a \bullet, 0$



HOW TO DEAL WITH NULL RULES

The Early algorithm, as it has been illustrated so far, builds two sets of pairs: $E[i]$, which holds the pairs produced by prediction and completion, and $E[i + 1]$, which holds pairs produced by shift.

In the presence of ε -rules, the completion operation should examine set $E[i]$ built so far, and then invoke prediction, which should in turn invoke completion, and so on until the two operations do not have anything left to add to the set. This method is correct, but inefficient.

The *Aycock* and *Horspool* method is more efficient and modifies the prediction operation as follows:

PREDICTION OPERATION (with ε -rules)

for every pair $(q \equiv A \rightarrow \alpha \bullet B \gamma, j)$ already in $E [i]$
 add to the set $E [i]$ the new pair $(r \equiv B \rightarrow \bullet \delta, i)$
 where r is the initial state of the machine expanding B
 if B is a nullable non-terminal
 add to the set $E [i]$ also the new pair $(s \equiv A \rightarrow \alpha B \bullet \gamma, j)$

The new fact is that the marker \bullet is moved to the right side of the nullable non-terminal symbol B , in agreement with the observation that the derivation might end up with replacing the non-terminal B by ε .

EXAMPLE – a grammar where all non-terminal symbols are nullable

$$S' \rightarrow S \quad S \rightarrow AAAA \quad A \rightarrow a \quad A \rightarrow E \quad E \rightarrow \varepsilon$$

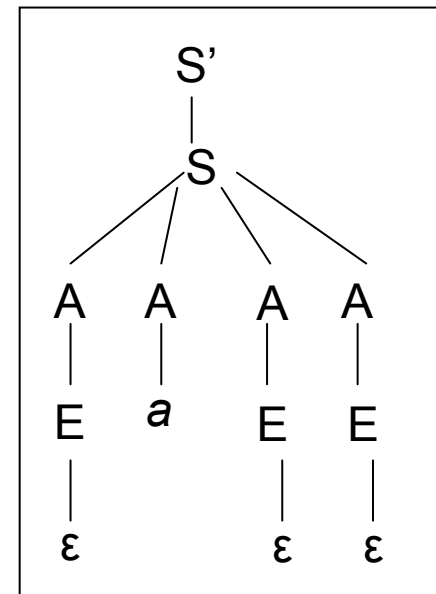
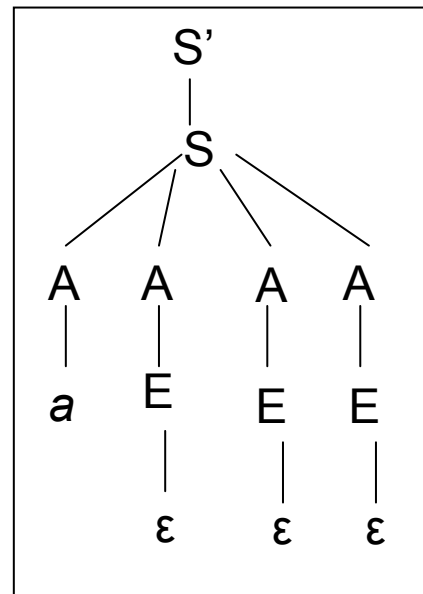
phrase to recognise: a

E [0]

$S' \rightarrow \bullet S$, 0
$S \rightarrow \bullet AAAA$, 0
$S' \rightarrow S \bullet$, 0
$A \rightarrow \bullet a$, 0
$A \rightarrow \bullet E$, 0
$S \rightarrow A \bullet AAA$, 0
$E \rightarrow \bullet$, 0
$A \rightarrow E \bullet$, 0
$S \rightarrow AA \bullet AA$, 0
$S \rightarrow AAA \bullet A$, 0
$S \rightarrow AAAA \bullet$, 0

a E [1]

$A \rightarrow a \bullet$, 0
$S \rightarrow A \bullet AAA$, 0
$S \rightarrow AA \bullet AA$, 0
$S \rightarrow AAA \bullet A$, 0
$S \rightarrow AAAA \bullet$, 0
$A \rightarrow \bullet a$, 1
$A \rightarrow \bullet E$, 1
$S' \rightarrow S \bullet$, 0
$E \rightarrow \bullet$, 1
$A \rightarrow E \bullet$, 1



and so on

IMPROVEMENTS OF THE EARLEY METHOD

1) HOW TO COMBINE EARLEY AND LOOKAHEAD

The pairs contained in the Earley sets can be extended by adding lookahead information as well. Lookahead sets are computed as in the LR(1) case.

If pairs are associated with lookahead sets, the Earley algorithm can avoid to add to a set a pair that is certainly doomed to failure. However, adding lookahead may sometimes cause the number of pairs to increase (as each pair contains more information). The effective advantage of adding lookahead to Earley is still a controversial matter.

The most efficient implementations of the Earley method do not use lookahead.

2) HOW TO WORK EARLEY OUT WITH EBNF GRAMMARS

The Earley algorithm can be modified so as to work correctly with EBNF rules (actually, representing a grammar by means of a machine network does so).

Bibliography

- S. Crespi Reghizzi, *Linguaggi Formali e Compilazione*, Pitagora Editrice Bologna, 2006
- Hopcroft, Ullman, *Formal Languages and their Relation to Automata*, Addison Wesley, 1969
- A. Salomaa – *Formal Languages*, Academic Press, 1973
- D. Mandrioli, C. Ghezzi – *Theoretical Foundations of Computer Science*, John Wiley & Sons, 1987
- L. Breveglieri, S. Crespi Reghizzi, *Linguaggi Formali e Compilatori: Temi d'Esame Risolti*, web site (eng + ita)