# Syntax Analysis
# of Non-Deterministic Grammars
# Earley (Tabular) Method

*Translated and adapted by L. Breveglieri*

# A GENERAL SYNTAX ANALYSIS METHOD – EARLEY (TABULAR)

The *Earley* method (also called *tabular* method) deals with any grammar type, even ambiguous, but we fully apply it only to non-deterministic grammars

It builds in parallel all the possible derivations of the string prefix scanned so far

Earley is similar to *ELR*, but it does not use the stack; instead, it uses a vector of sets, which efficiently represents stacks that have common parts

It resembles closely the implementation of the *ELR* deterministic *PDA* by means of a vectored stack

It simulates a non-deterministic pushdown analyzer, but with a polynomial time complexity
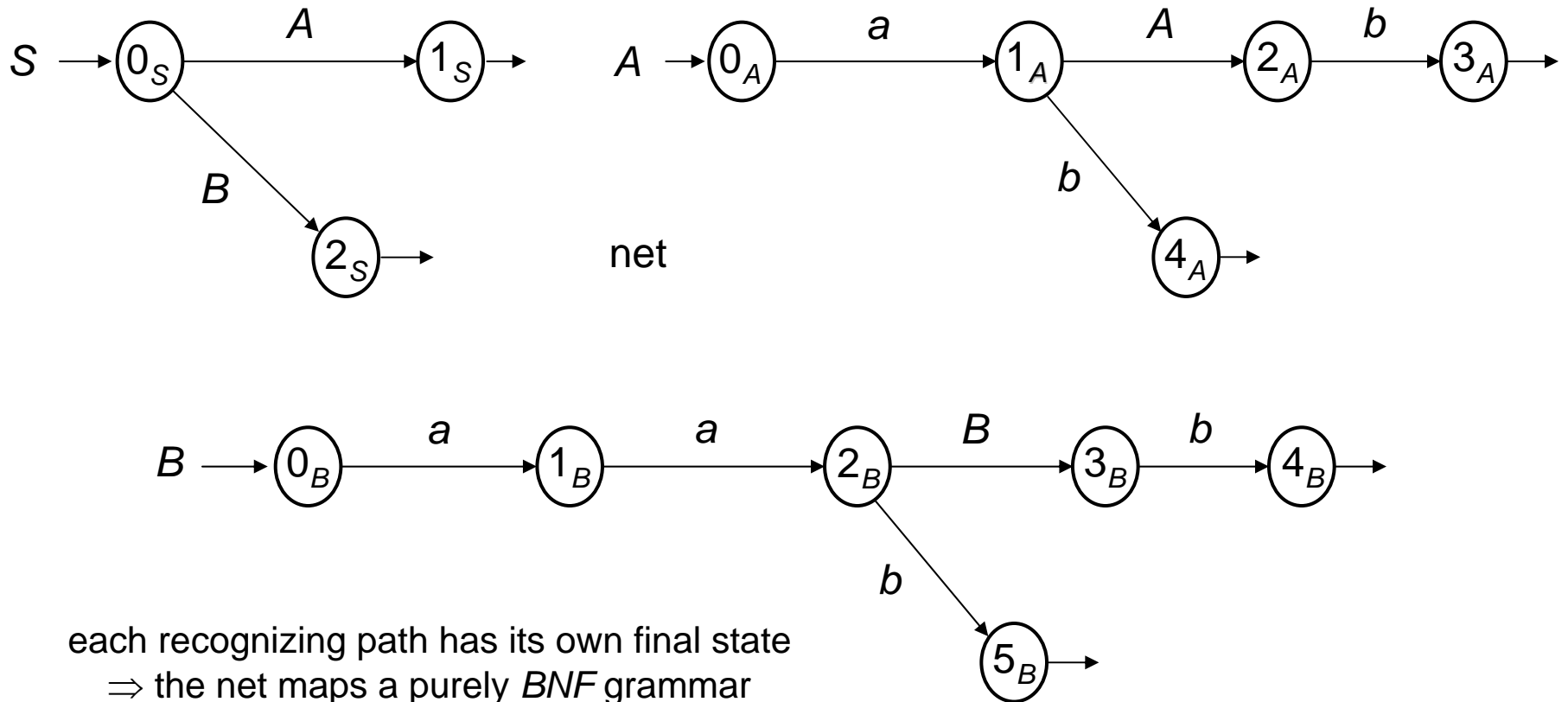
Notice that also the *ELR* method simulates parallel analysis threads, but only until reduction time

The Earley algorithm has variants without or with look-ahead, but here for simplicity look-ahead is not used
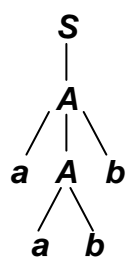
EXAMPLE – a grammar not *LR* (*k*)

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

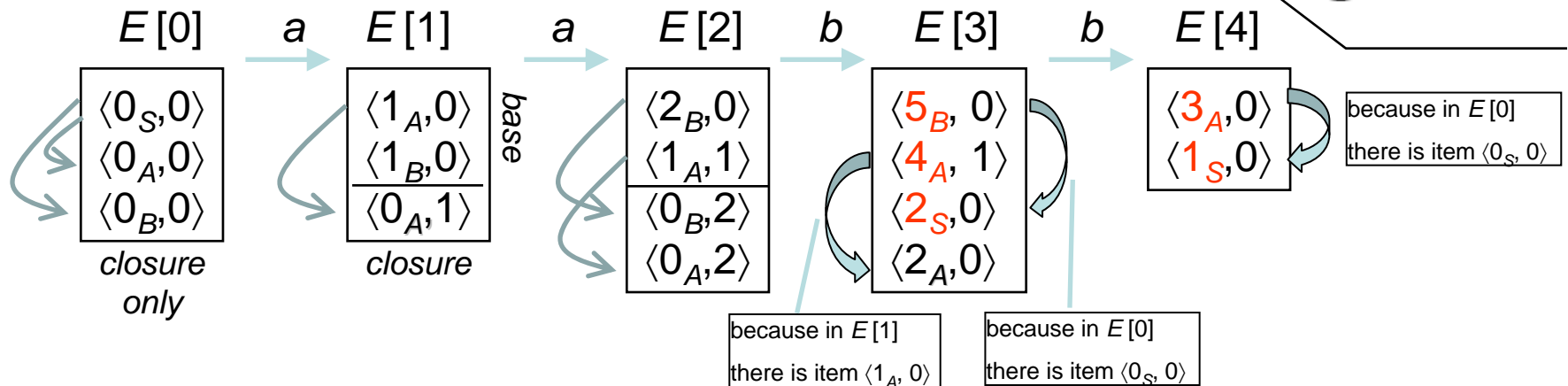$$S \to A \mid B \quad A \to aAb \mid ab \quad B \to aaBb \mid aab$$



net

each recognizing path has its own final state
  $\Rightarrow$ the net maps a purely *BNF* grammar

ANALYSIS – string *"a a b"*



net

syntax tree

$$S \rightarrow A \mid B$$

$E[0]$ — $a$ — $E[1]$ — $a$ — $E[2]$ — $b$ — $E[3]$ — $b$ — $E[4]$

$E[0]$:
$\langle 0_S, 0\rangle$
$\langle 0_A, 0\rangle$
$\langle 0_B, 0\rangle$
*closure only*

$E[1]$:
$\langle 1_A, 0\rangle$
$\langle 1_B, 0\rangle$
$\langle 0_A, 1\rangle$
*base / closure*

$E[2]$:
$\langle 2_B, 0\rangle$
$\langle 1_A, 1\rangle$
$\langle 0_B, 2\rangle$
$\langle 0_A, 2\rangle$

$E[3]$:
$\langle 5_B, 0\rangle$
$\langle 4_A, 1\rangle$
$\langle 2_S, 0\rangle$
$\langle 2_A, 0\rangle$

$E[4]$:
$\langle 3_A, 0\rangle$
$\langle 1_S, 0\rangle$

because in $E[0]$ there is item $\langle 0_S, 0\rangle$

because in $E[1]$ there is item $\langle 1_A, 0\rangle$

because in $E[0]$ there is item $\langle 0_S, 0\rangle$

$\xrightarrow{a}$ = terminal shift      $\searrow$ = closure      $\circlearrowright$ = non-terminal shift
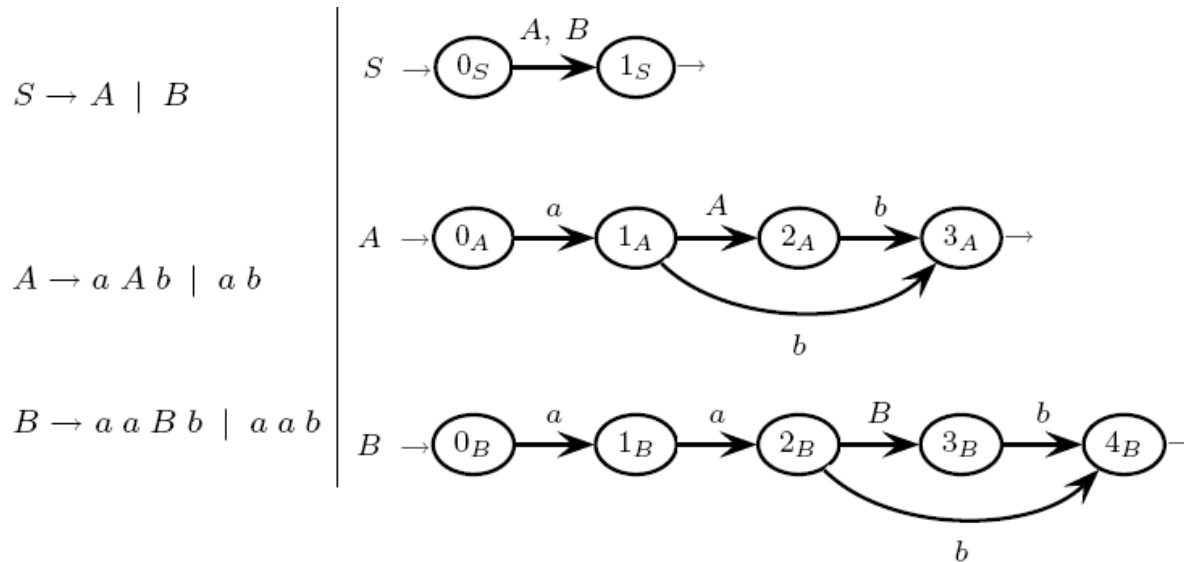
final states are in red

If the string were *"a a b"*, it would be accepted because $\langle 2_S, 0\rangle \in E[3]$ and $2_S$ is final for $M_S$

But it is the string *"a a b b"* that is accepted, because $\langle 1_S, 0\rangle \in E[4]$ and $1_S$ is final for $M_S$

Besides deciding if it holds $x \in L$, the algorithm decides also for all the prefixes of string $x$

$S \to A \mid B$

$A \to a\,A\,b \mid a\,b$

$B \to a\,a\,B\,b \mid a\,a\,b$



grammar still *BNF*
but the final states
are unified

Earley vector $E\,[0\ldots 4]$

| $E\,(0)$ | | $x_1$ | $E\,(1)$ | | $x_2$ | $E\,(2)$ | | $x_3$ | $E\,(3)$ | | $x_4$ | $E\,(4)$ | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $q$ | $j$ | $a$ | $q$ | $j$ | $a$ | $q$ | $j$ | $b$ | $q$ | $j$ | $b$ | $q$ | $j$ |
| $0_S$ | 0 | | $1_A$ | 0 | | $2_B$ | 0 | | $4_B$ | 0 | | $3_A$ | 0 |
| $0_A$ | 0 | | $1_B$ | 0 | | $1_A$ | 1 | | $3_A$ | 1 | | $1_S$ | 0 |
| $0_B$ | 0 | | $0_A$ | 1 | | $0_B$ | 2 | | $1_S$ | 0 | | | |
| | | | | | | $0_A$ | 2 | | $2_A$ | 0 | | | |

ALGORITHM – the Earley (tabular) recognizer

The algorithm builds a vector $E[0 \ldots n]$ of $n = |x|$ elements

Each element $E[i]$ is a set of items (or pairs) $\langle s, j \rangle$ like the items
in the $LR$ stack, but without look-ahead

The vector $E$ contains all the analysis threads in a compacted form,
i.e., the common parts are shared

We define the move types that are found in the steps of the algorithm

*terminal shift*        scanning: there are at most $n$ shifts, one per each character in $x$
*closure*              same as the closure in $LR$
*non-terminal shift*    same as the non-terminal shift in $LR$
*completion*          $=$ closure $+$ non-terminal shift, possibly repeated two or more times

call *i* the *current position*, the max index value *i* such that the set $E[i]$ is not empty

$TerminalShift(E, i)$      - - with index $1 \leq i \leq n$

- - loop that computes the *terminal shift* operation
- - for each preceding pair that shifts on terminal $x_i$

**for** $\left(\text{each pair } \langle p, j \rangle \in E[i-1] \text{ and } q \in Q \text{ s.t. } p \xrightarrow{x_i} q\right)$ **do**

    add pair $\langle q, j \rangle$ to element $E[i]$

**end for**

---

$Completion(E, i)$      - - with index $0 \leq i \leq n$

**do**

    - - loop that computes the *closure* operation
    - - for each pair that launches machine $M_X$

    **for** $\left(\text{each pair } \langle p, j \rangle \in E[i] \text{ and } X, q \in V, Q \text{ s.t. } p \xrightarrow{X} q\right)$ **do**

        add pair $\langle 0_X, i \rangle$ to element $E[i]$

    **end for**      *closure*

    - - nested loops that compute the *nonterminal shift* operation
    - - for each final pair that enables a shift on nonterminal $X$

    **for** $\left(\text{each pair } \langle f, j \rangle \in E[i] \text{ and } X \in V \text{ such that } f \in F_X\right)$ **do**

        - - for each pair that shifts on nonterminal $X$

        **for** $\left(\text{each pair } \langle p, l \rangle \in E[j] \text{ and } q \in Q \text{ s.t. } p \xrightarrow{X} q\right)$ **do**

            add pair $\langle q, l \rangle$ to element $E[i]$

        **end for**

    **end for**      *non-terminal shift*

**while** (some pair has been added)

---

- - analyze the terminal string $x$ for possible acceptance
- - define the Earley vector $E[0\ldots n]$ with $|x| = n \geq 0$

$E[0] := \{ \langle 0_S, 0 \rangle \}$      - - initialize the first elem. $E[0]$

**for** $i := 1$ **to** $n$ **do**      - - initialize all elem.s $E[1\ldots n]$

    $E[i] := \emptyset$

**end for**

$Completion(E, 0)$      - - complete the first elem. $E[0]$

$i := 1$

- - while the vector is not finished and the previous elem. is not empty

**while** $(i \leq n \wedge E[i-1] \neq \emptyset)$ **do**

    $TerminalShift(E, i)$      - - put into the current elem. $E[i]$

    $Completion(E, i)$      - - complete the current elem. $E[i]$

    $i++$

**end while**      *Earley*

every string character $x_{i+1}$ is examined
only once in the scanning
the machine input head is one-way
and there is not backtracking.

string *x* is accepted if and only if item $\langle f_S, 0 \rangle \in E[n]$      $f_S \in F_S$ is a final state for machine *S*

ext. grammar (EBNF)

machine network

$S \to a^+ \ (b \, B \, a)^* \ | \ A^+$



grammar is *EBNF*
$\Rightarrow$ its net has loops

$A \to a \, A \, b \ | \ a \, b$

$B \to b \, B \, a \ | \ c \ | \ \varepsilon$
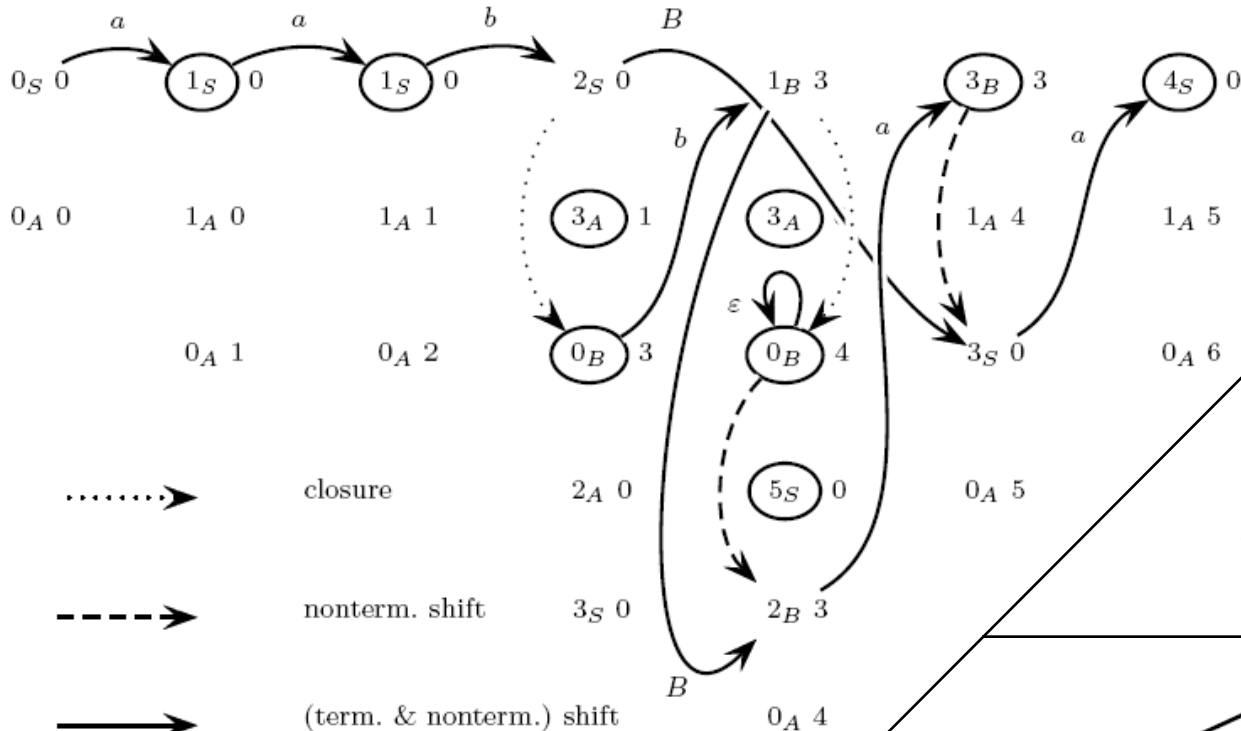
NON-DETERMINISM: if    - the adjacent substrings of *a* and *b* have the same length
    - there is not any character *c*

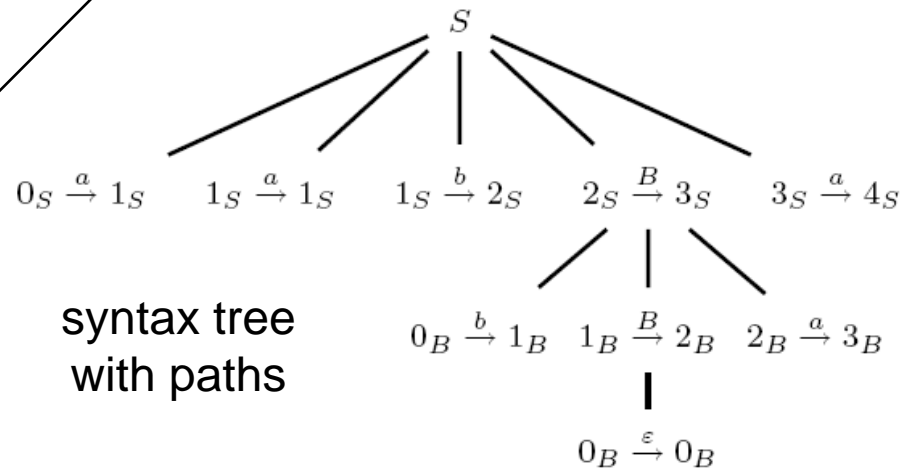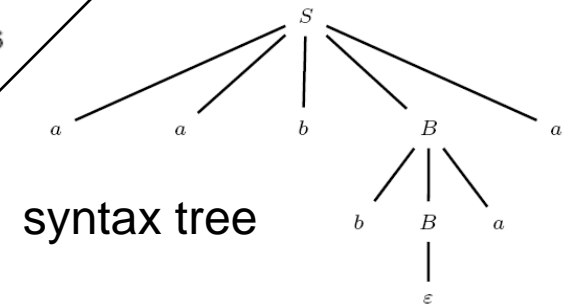Then the string deep structure can be determined only at the end of the scanning

syntax tree of string
"*a a b b a a*"

E [0]  a  E [1]  a  E [2]  b  E [3]  b  E [4]  a  E [5]  a  E [6]



ANALYSIS
includes the case
of a nullable
non-terminal

$0_S$ 0    $1_S$ 0    $1_S$ 0    $2_S$ 0    $1_B$ 3    $3_B$ 3    $4_S$ 0

$0_A$ 0    $1_A$ 0    $1_A$ 1    $3_A$ 1    $3_A$    $1_A$ 4    $1_A$ 5

$0_A$ 1    $0_A$ 2    $0_B$ 3    $0_B$ 4    $3_S$ 0    $0_A$ 6

........▸  closure       $2_A$ 0   $5_S$ 0    $0_A$ 5

- - - ▸  nonterm. shift   $3_S$ 0   $2_B$ 3

——▸  (term. & nonterm.) shift   $0_A$ 4

syntax tree

syntax tree
with paths

$0_S \xrightarrow{a} 1_S$   $1_S \xrightarrow{a} 1_S$   $1_S \xrightarrow{b} 2_S$   $2_S \xrightarrow{B} 3_S$   $3_S \xrightarrow{a} 4_S$

$0_B \xrightarrow{b} 1_B$   $1_B \xrightarrow{B} 2_B$   $2_B \xrightarrow{a} 3_B$

$0_B \xrightarrow{\varepsilon} 0_B$

the grammar (or the net)
is not ambiguous
$\Rightarrow$ only one tree can be identified

# WORST-CASE ASYMPTOTIC COMPLEXITY OF EARLEY

Let $n \geq 1$ be the length of the string to be analyzed

space complexity = size of the Earley vector

| | |
|---|---|
| = # of sets E [j] × max size of an element | |
| = $(n + 1) \times kn$ | $k$ is the number of net states |
| = $O(n^2)$ | is quadratic |

time compexity = # of terminal shifts + # of closures + # of nonterminal shifts

| | |
|---|---|
| = $O(n^2) + O(n^2) + O(n^2) \times kn$ | $k$ is the number of net states |
| = $O(n^3)$ | is cubic |

So in general Earley is polynomial, though it is not linear

Special case: $n = 0$ (null string), then both complexities are constant, as only closure is used and is limited by the number of states $k$, which is a constant; anyway, the complexity is still polynomial

In most practical cases with applicative grammars (or machine nets), Earley performs better than above and approaches a linear behaviour

# RECONSTRUCTION OF THE SYNTAX TREE
## (for non-ambiguous grammars only)

$T$he recursive procedure *BuildTree* (*BT*) builds the one tree starting from vector $E$

*BT* has four parameters:      nonterminal $X$, subtree root, with $X \in V$

final state of machine $M_X$

integer $i$, substring left position

integer $j$, substring right position

*BT* returns:      the (sub)tree of the derivation $X \Rightarrow^* x_{i+1}, \ldots x_j$
that has root $X$

$I$nitial call $BT (S, f_S, 0, n)$

$BuildTree\ (\ X,\ f,\ j,\ i\ )$

- - $X$ is a nonterminal, $f$ is a final state of $M_X$ and $0 \leq j \leq i \leq n$
- - return as parenthesized string the syntax tree rooted at node $X$

- - node $X$ will have a list $\mathcal{C}$ of terminal and nonterminal child nodes
- - either list $\mathcal{C}$ will remain empty or it will be filled from right to left

$\mathcal{C} := \varepsilon$                     - - set to $\varepsilon$ the list $\mathcal{C}$ of child nodes of $X$

$q := f$                     - - set to $f$ the state $q$ in machine $M_X$

$k := i$                     - - set to $i$ the index $k$ of vector $E$

In its recostruction loop, function *BT* performs the following operations:
- scanning the vector *E* backwards by iterating the while loop
- rebuilding terminal and non-terminal shift moves
- identifying the child nodes of the same parent node
- recursively calling itself if set *E*[ *i* ] contains a pair from non-terminal shift
$\Rightarrow$ the recursive call builds the subtree of the shifted non-terminal

-- walk back the sequence of term. & nonterm. shift oper.s in $M_X$

**while** ( $q \neq 0_X$ ) **do**          -- while current state $q$ is not initial

  -- try to backwards recover a terminal shift move $p \xrightarrow{x_k} q$, i.e.,
  -- check if node $X$ has terminal $x_k$ as its current child leaf

(a)   **if** $\begin{pmatrix} \exists\, h = k-1 \quad \exists\, p \in Q_X \quad \text{such that} \\ \langle\, p,\, j\, \rangle \in E\,[h] \;\wedge\; \text{net has } p \xrightarrow{x_k} q \end{pmatrix}$ **then**

    $\mathcal{C} := x_k \cdot \mathcal{C}$          -- concatenate leaf $x_k$ to list $\mathcal{C}$

  **end if**

terminal shift

  -- try to backwards recover a nonterm. shift oper. $p \xrightarrow{Y} q$, i.e.,
  -- check if node $X$ has nonterm. $Y$ as its current child node

(b)   **if** $\begin{pmatrix} \exists\, Y \in V \quad \exists\, e \in F_Y \quad \exists\, h\ \ j \leq h \leq k \leq i \quad \exists\, p \in Q_X \quad \text{s.t.} \\ \langle\, e,\, h\, \rangle \in E\,[k] \;\wedge\; \langle\, p,\, j\, \rangle \in E\,[h] \;\wedge\; \text{net has } p \xrightarrow{Y} q \end{pmatrix}$ **then**

    -- recursively build the subtree of the derivation:
    -- $Y \xRightarrow[G]{+} x_{h+1} \ldots x_k$ if $h < k$  or  $Y \xRightarrow[G]{+} \varepsilon$ if $h = k$
    -- and concatenate to list $\mathcal{C}$ the subtree of node $Y$

    $\mathcal{C} := BuildTree\ (\, Y,\, e,\, h,\, k\, ) \cdot \mathcal{C}$

  **end if**

non-terminal shift

  $q := p$          -- shift the current state $q$ back to $p$
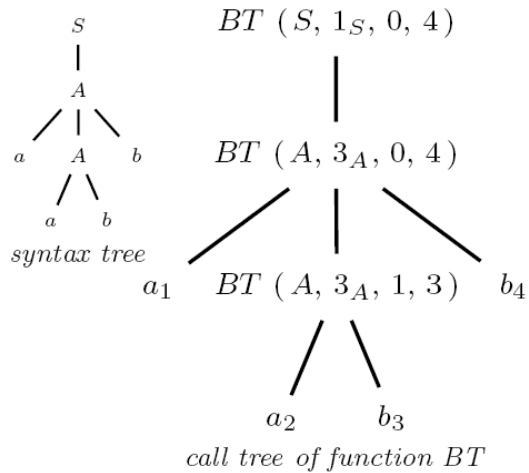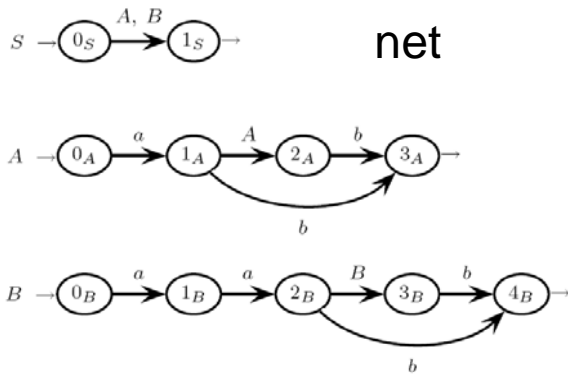
  $k := h$          -- drag the current index $k$ back to $h$

**end while**

**return** $(\, \mathcal{C}\, )_X$          -- return the tree rooted at node $X$   □

non-ambiguous grammar $\Rightarrow$
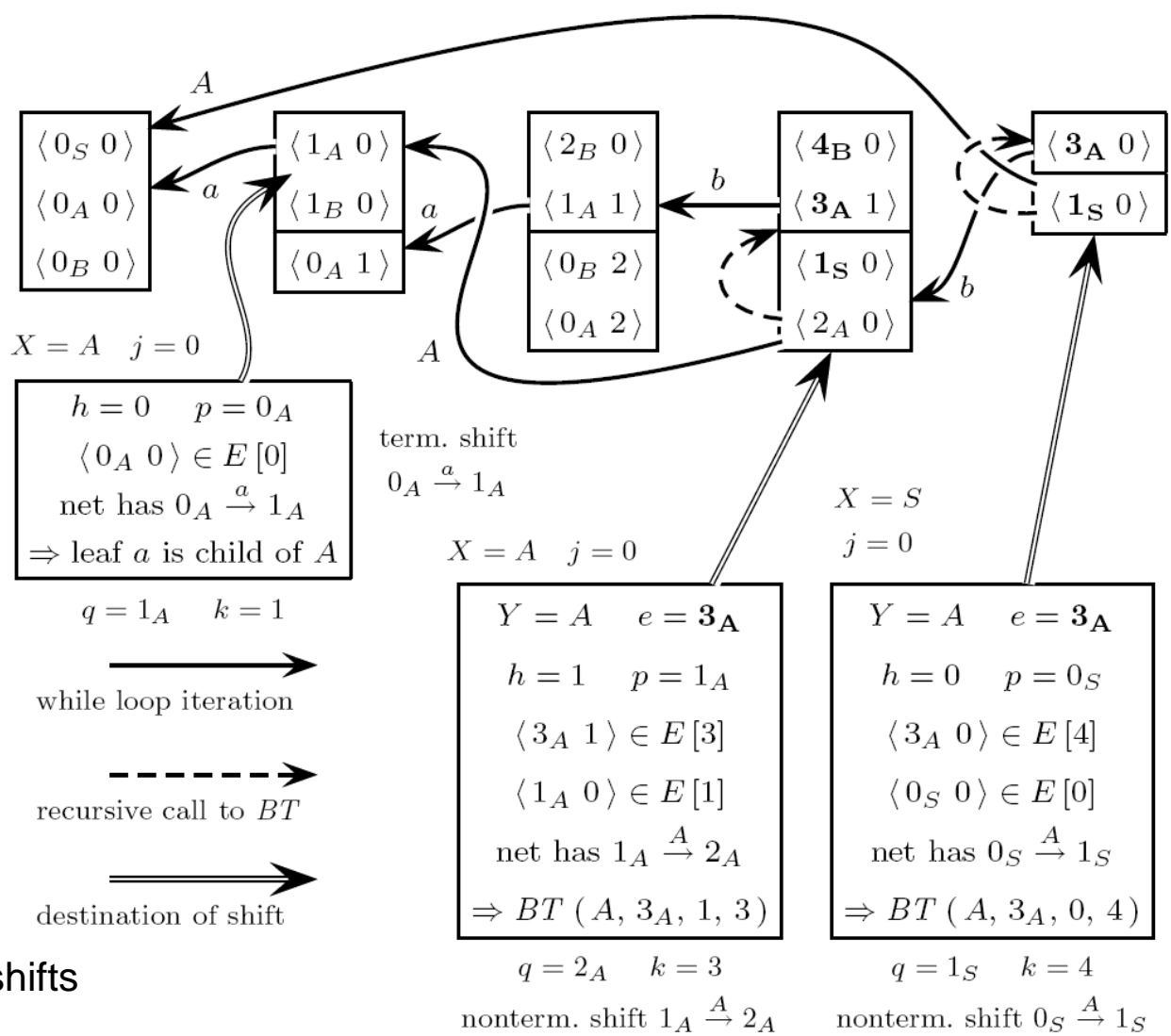deterministic reconstruction

$E\,[0]$     $a_1$     $E\,[1]$     $a_2$     $E\,[2]$     $b_3$     $E\,[3]$     $b_4$     $E\,[4]$

net



$X = A \quad j = 0$

$$h = 0 \quad p = 0_A$$
$$\langle 0_A\ 0 \rangle \in E\,[0]$$
$$\text{net has } 0_A \xrightarrow{a} 1_A$$
$$\Rightarrow \text{leaf } a \text{ is child of } A$$

term. shift
$$0_A \xrightarrow{a} 1_A$$

$X = A \quad j = 0$

$q = 1_A \quad k = 1$

→ while loop iteration

⇢ recursive call to $BT$

⇒ destination of shift

$X = S$
$j = 0$

$BT\,(S, 1_S, 0, 4)$

$BT\,(A, 3_A, 0, 4)$

*syntax tree*

$a_1 \quad BT\,(A, 3_A, 1, 3) \quad b_4$

$a_2 \qquad b_3$

*call tree of function BT*

$$Y = A \quad e = \mathbf{3_A}$$
$$h = 1 \quad p = 1_A$$
$$\langle 3_A\ 1 \rangle \in E\,[3]$$
$$\langle 1_A\ 0 \rangle \in E\,[1]$$
$$\text{net has } 1_A \xrightarrow{A} 2_A$$
$$\Rightarrow BT\,(A, 3_A, 1, 3)$$

$q = 2_A \quad k = 3$
nonterm. shift $1_A \xrightarrow{A} 2_A$

$$Y = A \quad e = \mathbf{3_A}$$
$$h = 0 \quad p = 0_S$$
$$\langle 3_A\ 0 \rangle \in E\,[4]$$
$$\langle 0_S\ 0 \rangle \in E\,[0]$$
$$\text{net has } 0_S \xrightarrow{A} 1_S$$
$$\Rightarrow BT\,(A, 3_A, 0, 4)$$

$q = 1_S \quad k = 4$
nonterm. shift $0_S \xrightarrow{A} 1_S$

there are three more terminal shifts
(not shown here)

## WORST-CASE COMPLEXITY OF BUILDTREE

Since the grammar is supposed to be not ambiguous, the tree size is linear, so the time complexity of *BT* is also cubic (like Earley)

Extending *BT* to ambiguous grammars is possible as well, but to efficiently represent all the trees a Sparse Tree Forest (*STF*) is necessary, which is a graph type more general than a tree (not considered here)


## OPTIMIZATION OF THE EARLEY METHOD BY USING LOOK-AHEAD

The items in the sets *E*[i] of the Earley vector *E* can be extended by including look-ahead, computed as it is done in the *ELR*(1) method

By siding a look-ahead to each item, the Earley algorithm avoids to put into the sets *E*[i] items that correspond to choices doomed to failure. This way however, for some grammars the number of items may increase rather than decrease. So the advantage of using look-ahead in Earley is controversial

At present the most efficient implementations of Earley do not use look-ahead