# Earley Parser

Christopher Millar and Ekaterina Volkova

Seminar für Sprachwissenschaft

Universität Tübingen

January 2007

In general, **breadth-first bottom-up** parsers are attractive since:

- they work on-line;
- can handle left-recursion;
- can be doctored to handle ε-rules.

Still the question remains:

**How to curb their needless activity?**

A method that will *restrict the fan-out* to reasonable proportions while still *retaining full generality* was developed by Earley .

# Earley Parser: Basic Concept

**Main problem**: the **spurious reductions** can never derive from the **start symbol**.

**Solution**: give a method to **restrict the reductions** only to those that derive from the **start symbol**.

The resulting parser takes **at most $n^3$** units of time for input of length n **rather than $C^n$**.

# Earley Parser: Definition

Earley's parser can also be described as a **breadth-first top-down** parser with **bottom-up recognition**, Still, we prefer to treat it as a **bottom-up method,** for it can handle left-recursion directly but needs special measures to handle ε-rules.

An Earley item is an **item** with an indication of the **position** of the symbol at which the recognition of the recognized part started.

$$E->E \bullet QF@3 \longleftarrow \textit{Position}$$

The sets of items contain exactly those items...

 a) of which the part before the dot has been recognized so far **...and...**

b) are useful in reaching the start symbol.

# Earley Parser: Methods

The Earley Parser uses methods called **Scanner, Completer** and **Predictor**.

- **Scanner** is like "shift".
- **Completer** is like "reduce".
- **Predictor** is unique to the Earley parser.

Scanner

Completer

Predictor

The Scanner, Completer and Predictor deal with **four sets of items** for each token in the input.

We'll refer to a token as **sigma@p** or as

$$\boldsymbol{\delta}_{\mathbf{p}}$$

# Earley Parser: The Four Sets

sigma@p is surrounded by four sets:

- **itemset@p-1**

- **completed@p**

- **active@p**

- **predicted@p**

itemset@p-1

completed@p

active@p

predicted@p

- **itemset@p-1 -** items available just before sigma@p;

- **completed@p -** items that have become completed after sigma@p;

- **active@p -** non-completed items after sigma@p:

- **predicted@p** - the set of newly predicted items.

## The Scanner :

looks at **sigma@p** -> goes through **itemset@p-1**

-> makes copies of all items that contain •**sigma**

-> changes them to **sigma** • -> adds them...

a) to the set **completed@p** if the item@p was completed ...or...

b) to the set **active@p** if the **item@p** is not yet completed

**Rules not containing ●sigma are discarded!**

**The Completer** inspects **completed@p**, which contains the completely recognized items and can now be *reduced*.

For each item of the form **R --> sigma@m** the Completer goes to **itemset@(m-1)**, and calls the Scanner; which goes to work on **R**.

The Scanner will make copies of all items in **itemset@(m-1)** featuring a •R, replace the •R by R• and store them in either **completed@p** or **active@p**. At this stage items could be added to the set **completed@p**.

**Eventually the Completer stops completing.**

(When it has <u>completely</u> <u>completed</u>

the set **completed@p** :)  )

# Earley Parser: The Predictor

The Predictor goes through the sets **active@p** (which was filled by the Scanner) and **predicted@p** (which is empty initially), and considers all non-terminals which have a • before them.

For each expected non-terminal **N** and each rule for that non-terminal **N --> P...,** the Predictor adds an item to the set **predicted@p**.

This may introduce new predicted non-terminals (for instance, P) to **predicted@p** which causes more work for the Predictor.

**Eventually the Predictor stops predicting.**

The sets **active@p** and **predicted@p** together form the new **itemset@p**. If the completed set for the last symbol in the input contains an item **S-->...●@1**. Then the input is recognized.

# Earley Parser: Example

Consider an example with the following grammar and the input: **a - a + a.**

S --> E

E --> EQF

E --> F

Q --> +

Q --> -

F --> a

There is one **Predictor, Scanner** and **Completer** stage for each symbol.

Parsing begins by **calling the Predictor** on the initial active set containing **S --> E@1** which generates **itemset@0.**

$$\text{act/pred}_0$$

$$
\begin{array}{ll}
S \rightarrow \bullet E & @1 \\
\hdashline
E \rightarrow \bullet EQF & @1 \\
E \rightarrow \bullet F & @1 \\
F \rightarrow \bullet a & @1 \\
\end{array}
$$

$$= \text{itemset}_0$$

The Predictor, reads **active@0, {S-> •E@1 }** and **predicted@0**, which is initially empty, and fills the set **predicted@0**.

{act.@0} U {pred.@0} = {itemset@0}

**completed₁**

$$F \to a \bullet @1$$
$$E \to F \bullet @1$$
$$S \to E \bullet @1$$

$a_1$

**act/pred₁**

$$E \to E \bullet QF @1$$
$$Q \to \bullet + \quad @2$$
$$Q \to \bullet - \quad @2$$

= itemset₁

After scanning **δ@1** the Completer completes some rules, and puts the other possible rules in **active@1**. Predictor makes predictions from those that are in the active set.

Continue as before until the input is consumed.

completed$_3$

| | |
|---|---|
| F->a• | @3 |
| E->EQF• | @1 |
| S->E• | @1 |

a$_3$

act/pred$_3$

| | |
|---|---|
| E->E•QF | @1 |
| Q->•+ | @4 |
| Q->•- | @4 |

= itemset$_3$

As you can see we already have few possibilities...

$$\text{completed}_4$$

$$\boxed{Q\text{->}+\bullet@4}$$

$$+_4$$

$$\text{act/pred}_4$$

$$\boxed{\begin{array}{ll} E\text{->}EQ\bullet F@1 \\ \hdashline F\text{->}\bullet a & @5 \end{array}}$$

$$= \text{itemset}_4$$

**completed$_5$**

| | |
|---|---|
| F->a• | @5 |
| E->EQF•@1 | |
| S->E• | @1 |

a$_5$

**active$_5$**

| |
|---|
| E->E•QF@1 |

**S --> E• @1** is in the set completed and the last input symbol has been read.

**Therefore the sentence is recognized!!!**
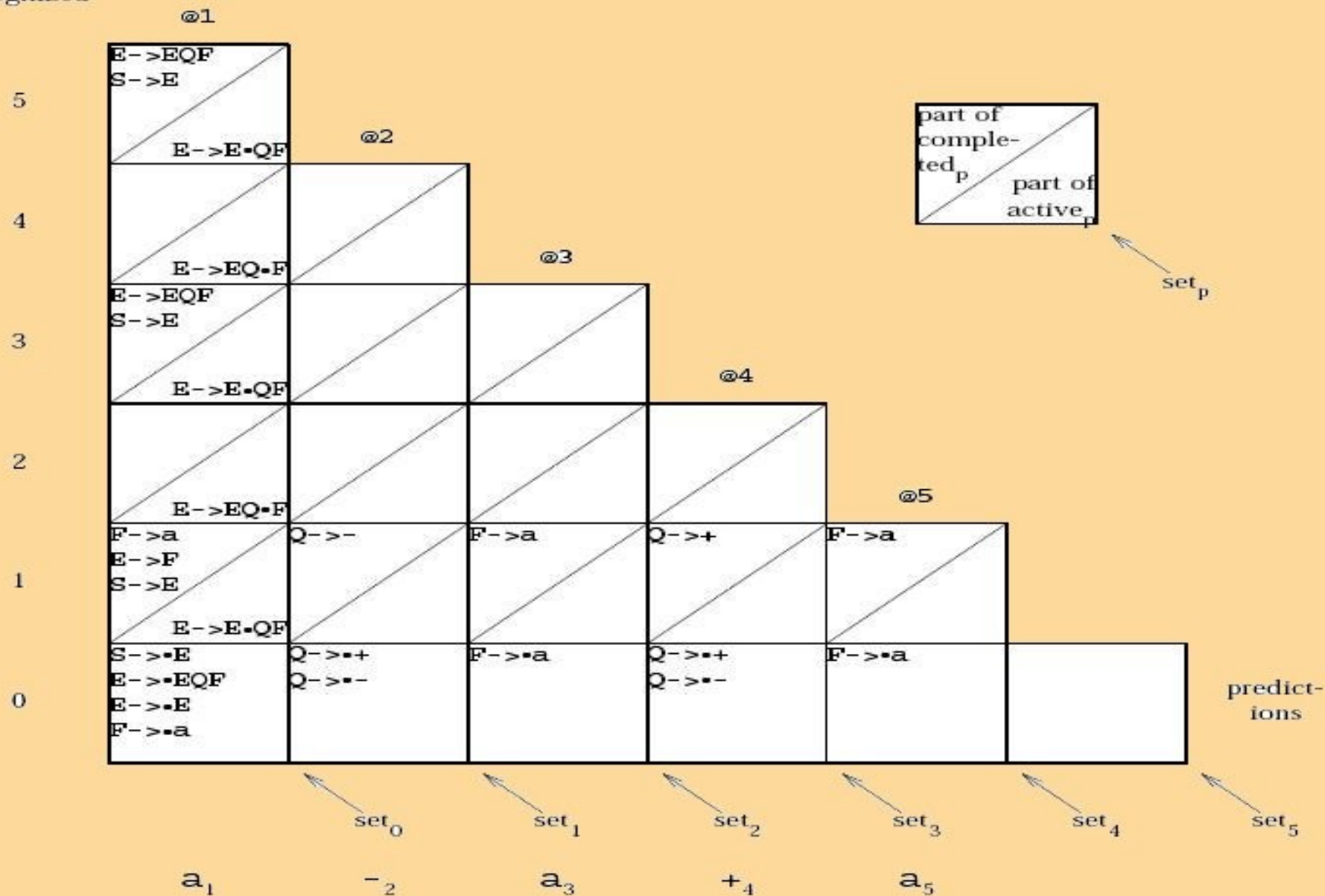
# Earley Parser: Comparison to CYK

Similarities:

- are Chart Parsers

- worst case memory requirements $O(n^2)$

- worst case time complexity $O(n^3)$

- use bottom-up recognition

- use a top-down parser to build trees

The Early Parser however eliminates rules which will not be useful as we go along, with non ambiguous grammars such as the example shown we get a worst time complexity of $O(n^2)$.
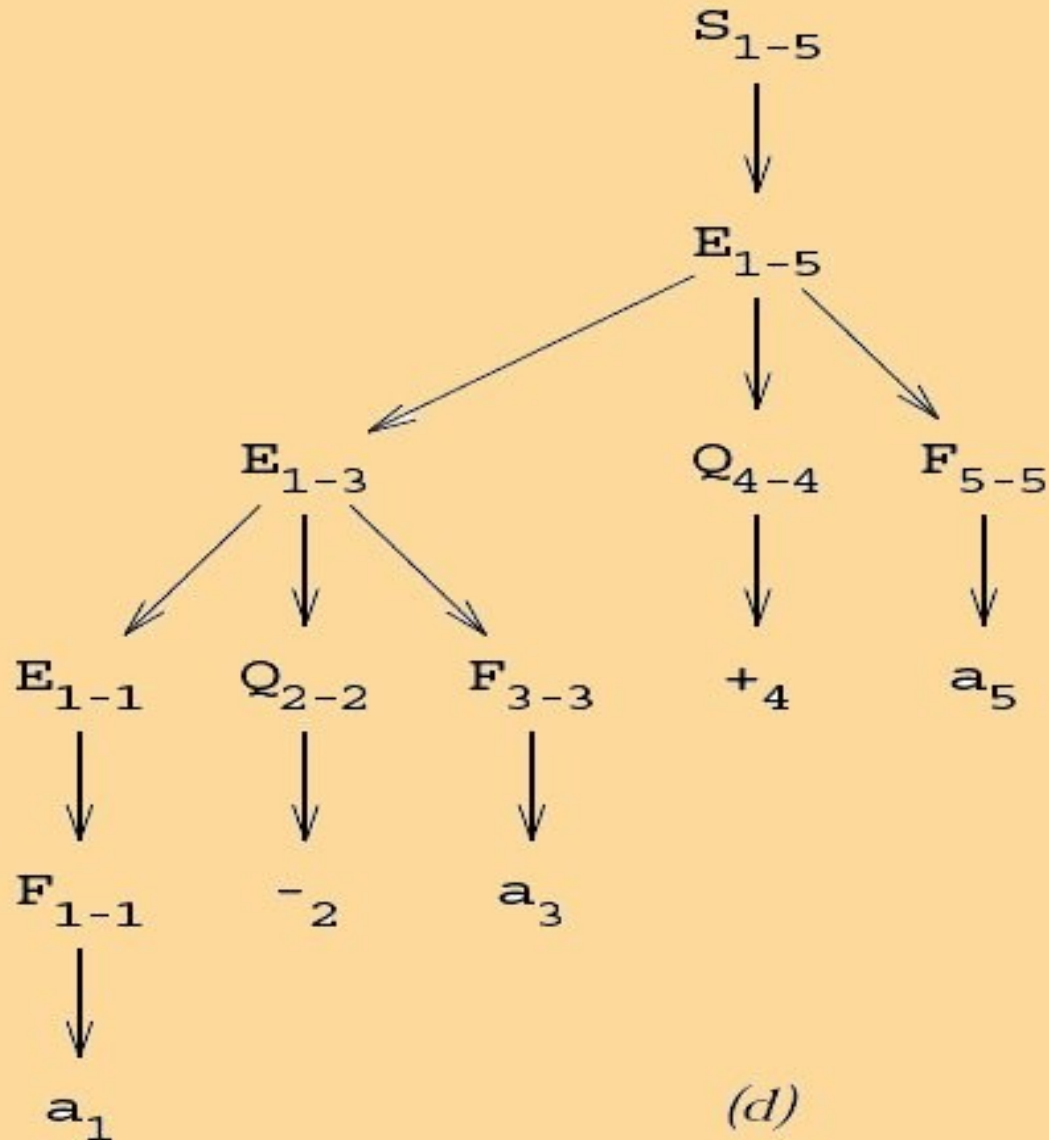
# Earley Parser: Recognition Chart

# Earley Parser: CYK Recognition Chart



length
recognized

|   | @1 | @2 | @3 | @4 | @5 |
|---|---|---|---|---|---|
| 5 | E S | | | | |
| 4 | | | | | |
| 3 | E S | | E S | | |
| 2 | | | | | |
| 1 | F E S | Q | F E S | Q | F E S |

|   | $a_1$ | $-_2$ | $a_3$ | $+_4$ | $a_5$ |

# Earley Parser: Parsing Tree



As with the CYK parser, a simple top-down Unger-type parser can be used to reconstruct all possible parse trees from a chart.

# Earley Parser: A Worse Example

We get worst case behaviour when we have to deal with ambiguous grammars like:

S --> SS

S --> x

completed₁

S->x• @1

act/pred₀

S->•SS@1
............................
S->•x @1

= itemset₀

x₁

act/pred₁

S->S•S@1
............................
S->•SS@2

S->•x @2

= itemset₁

**completed₂**

    S->x• @2
    S->SS•@1

**act/pred₂**

    S->S•S@2
    S->S•S@1
    ........................
    S->•SS@3
    S->•x @3

= itemset₂

x₂

x₃

**completed₃**

    S->x• @3
    S->SS•@2
    S->SS•@1

**active₃**

    S->S•S@3
    S->S•S@2
    S->S•S@1

The **active@p** and **predicted@p** sets keep growing untill the final symbol is read. When building a parse tree from the resulting chart we find two possible derivations, but if the input would be longer the the situation would be worse!

**The Earley parser doesn't like ε-rules!**

(Does **anybody** like them?)

Consider the following non-e-free grammar with the input a a / a.

S --> E          E --> EQF          E --> F


Q --> *

Q --> /

Q --> e

F --> a

After reading a1 we have a situation where every time the predictor predicts a •Q it must also predict a Q•

$\text{act/pred}_1$

E->E•QF@1

E->EQ•F@1

. . . . . . . . . . . . . . . . . .

Q->•x      @2

Q->•/      @2

F->•a      @2

$= \text{itemset}_1$

This can effect the behaviour of the Completer which is working on itemset@1.

In the end we can find a parse with this grammar.

What would happen to the itemset if we had a

rule Q --> QQ ?

An Early parser would resolve it but not without

inefficiency.

E --> E•QF    E --> EQ•F

Q --> •QQ    Q --> Q•Q        Q --> QQ•

Q --> *

Q --> /        ε-rules add significantly to the

F --> a        complexity time

Prediction Lookahead reduces the number of incorrect predictions made by the Predictor by considering next input symbol before adding items to **predicted@p**. It uses a set of FIRST terminal symbols, for each non terminal.

S -> A | AB | B     FIRST(S) = {p, q}

A -> C              FIRST(A) = {p}

B -> D              FIRST(B) = {q}

C -> p              FIRST(C) = {p}

D -> q              FIRST(D) = {q}

act/pred₀

| | |
|---|---|
| S'->•S | @1 |
| S->•A | @1 |
| S->•AB | @1 |
| S->•B | @1 |
| A->•C | @1 |
| B->•D | @1 |
| C->•p | @1 |
| D->•q | @1 |

= itemset₀

q₁

completed₁

| | |
|---|---|
| D->q• | @1 |
| B->D• | @1 |
| S->B• | @1 |
| S'->S• | @1 |

act/pred₁

= itemset₁

Without

lookahead

act/pred₀

S' -> •S @1
S -> •B @1
B -> •D @1
D -> •q @1

= itemset₀

q₁

completed₁

D -> q• @1
B -> D• @1
S -> B• @1
S' -> S• @1

act/pred₁

= itemset₁

With

lookahead

# Earley Parser: Conclusion

Earley Parser shows a very successful combination of strong sides of top-down and bottom-up methods, handles well left recursion and ε-rules, and, being armoured by lookahead, takes the optimal possible amount of memory.

**Earley rules!**