

Main page Contents Featured content Current events Random article Donate to Wkipedia Wkipedia store

Interaction

Help About Wikipedia Community portal Recent changes Contact page

Tools

What links here Related changes Upload file Special pages Permanent link Page information Wkidata item Cite this page

Print/export

Create a book Download as PDF Printable version

O

Languages

Български Čeština

Deutsch

Español فارسی

Français

日本語

★ Polski Português Русский

> Српски / srpski Українська

Æ Edit links

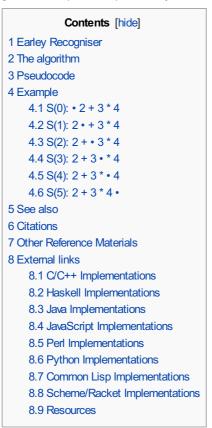
Article Talk Read Edit More ▼ Search Q

# Earley parser

From Wikipedia, the free encyclopedia

In computer science, the **Earley parser** is an algorithm for parsing strings that belong to a given context-free language, though (depending on the variant) it may suffer problems with certain nullable grammars.<sup>[1]</sup> The algorithm, named after its inventor, Jay Earley, is a chart parser that uses dynamic programming; it is mainly used for parsing in computational linguistics. It was first introduced in his dissertation<sup>[2]</sup> in 1968 (and later appeared in abbreviated, more legible form in a journal<sup>[3]</sup>).

Earley parsers are appealing because they can parse all context-free languages [discuss], unlike LR parsers and LL parsers, which are more typically used in compilers but which can only handle restricted classes of languages. The Earley parser executes in cubic time in the general case  $O(n^3)$ , where n is the length of the parsed string, quadratic time for unambiguous grammars  $O(n^2)$ , and linear time for almost all LR(k) grammars. It performs particularly well when the rules are written left-recursively.



## Earley Recogniser [edit]

The following algorithm describes the Earley recogniser. The recogniser can be easily modified to create a parse tree as it recognises, and in that way can be turned into a parser.

## The algorithm [edit]

In the following descriptions,  $\alpha$ ,  $\beta$ , and  $\gamma$  represent any string of terminals/nonterminals (including the empty string), X and Y represent single nonterminals, and *a* represents a terminal symbol.

Earley's algorithm is a top-down dynamic programming algorithm. In the following, we use Earley's dot notation: given a production  $X \to \alpha \beta$ , the notation  $X \to \alpha \bullet \beta$  represents a condition in which  $\alpha$  has already been parsed and  $\beta$  is expected.

Input position 0 is the position prior to input. Input position n is the position after accepting the nth token. (Informally, input positions can be thought of as locations at token boundaries.) For every input position, the parser generates a *state set*. Each state is a tuple  $(X \to \alpha \bullet \beta, i)$ , consisting of

- the production currently being matched  $(X \rightarrow \alpha \beta)$
- our current position in that production (represented by the dot)
- the position i in the input at which the matching of this production began: the origin position

(Earley's original algorithm included a look-ahead in the state; later research showed this to have little practical effect on the parsing efficiency, and it has subsequently been dropped from most implementations.)

The state set at input position k is called S(k). The parser is seeded with S(0) consisting of only the top-level rule. The parser then repeatedly executes three operations: *prediction*, *scanning*, and *completion*.

- Prediction: For every state in S(k) of the form (X → α Y β, j) (where j is the origin position as above), add (Y → γ, k) to S(k) for every production in the grammar with Y on the left-hand side (Y → γ).
- **Scanning**: If a is the next symbol in the input stream, for every state in S(k) of the form  $(X \to \alpha \bullet a \beta, j)$ , add  $(X \to \alpha \bullet \beta, j)$  to S(k+1).
- **Completion**: For every state in S(k) of the form  $(X \to \gamma \bullet, j)$ , find states in S(j) of the form  $(Y \to \alpha \bullet X \beta, i)$  and add  $(Y \to \alpha X \bullet \beta, i)$  to S(k).

It is important to note that duplicate states are not added to the state set, only new ones. These three operations are repeated until no new states can be added to the set. The set is generally implemented as a queue of states to process, with the operation to be performed depending on what kind of state it is.

## Pseudocode [edit]

Adapted from [4] by Daniel Jurafsky and James H. Martin

```
function EARLEY-PARSE(words, grammar)
     ENQUEUE((\gamma \rightarrow \bullet S, 0), chart[0])
     for i ← from 0 to LENGTH(words) do
          for each state in chart[i] do
                if INCOMPLETE?(state) then
                     \textbf{if} \ \text{NEXT-CAT(state)} \ \textbf{is} \ \textbf{a} \ \text{nonterminal} \ \textbf{then}
                                                                            // non-terminal
                          PREDICTOR(state, i, grammar)
                     else do
                          SCANNER(state, i)
                                                                           // terminal
                else do
                    COMPLETER (state, i)
           end
     end
     return chart
procedure PREDICTOR((A \rightarrow \alpha \cdot B, i), j, grammar)
     for each (B \rightarrow \gamma) in GRAMMAR-RULES-FOR(B, grammar) do
          ADD-TO-SET((B \rightarrow •\gamma, j), chart[j])
     end
procedure SCANNER((A \rightarrow \alpha \cdot B, i), j)
     if B ⊂ PARTS-OF-SPEECH(word[j]) then
          ADD-TO-SET((B \rightarrow word[j], j), chart[j + 1])
procedure COMPLETER((B \rightarrow Y^{\bullet}, j), k)
     for each (A \rightarrow \alpha \bullet B\beta, i) in chart[j] do
          ADD-TO-SET((A \rightarrow \alpha B \cdot \beta, i), chart[k])
     end
```

## Example [edit]

Consider the following simple grammar for arithmetic expressions:

```
<P> ::= <S>  # the start rule

<S> ::= <S> "+" <M> | <M>

<M> := <M> "*" <T> | <T>

<T> ::= "1" | "2" | "3" | "4"
```

With the input:

```
2 + 3 * 4
```

This is the sequence of state sets:

```
(state no.) Production (Origin) # Comment
```

### S(0): • 2 + 3 \* 4 [edit]

#### S(1): 2 • + 3 \* 4 [edit]

```
(1) T \rightarrow \text{number} \bullet (0) # scan from S(0)(6)

(2) M \rightarrow T \bullet (0) # complete from (1) and S(0)(5)

(3) M \rightarrow M \bullet * T (0) # complete from (2) and S(0)(4)

(4) S \rightarrow M \bullet (0) # complete from (2) and S(0)(3)

(5) S \rightarrow S \bullet + M (0) # complete from (4) and S(0)(2)

(6) P \rightarrow S \bullet (0) # complete from (4) and S(0)(1)
```

#### S(2): 2 + • 3 \* 4 [edit]

```
(1) S \to S + \bullet M (0) # scan from S(1) (5)

(2) M \to \bullet M * T (2) # predict from (1)

(3) M \to \bullet T (2) # predict from (1)

(4) T \to \bullet number (2) # predict from (3)
```

#### S(3): 2 + 3 • \* 4 [edit]

```
(1) T \rightarrow \text{number} \bullet (2) # scan from S(2)(4)

(2) M \rightarrow T \bullet (2) # complete from (1) and S(2)(3)

(3) M \rightarrow M \bullet * T (2) # complete from (2) and S(2)(2)

(4) S \rightarrow S + M \bullet (0) # complete from (2) and S(2)(1)

(5) S \rightarrow S \bullet + M (0) # complete from (4) and S(0)(2)

(6) P \rightarrow S \bullet (0) # complete from (4) and S(0)(1)
```

#### S(4): 2 + 3 \* • 4 [edit]

```
(1) M \rightarrow M * \bullet T (2) # scan from S(3)(3)
(2) T \rightarrow \bullet number (4) # predict from (1)
```

### S(5): 2 + 3 \* 4 • [edit]

```
(1) T \rightarrow \text{number} \bullet (4) # scan from S(4)(2)

(2) M \rightarrow M * T \bullet (2) # complete from (1) and S(4)(1)

(3) M \rightarrow M \bullet * T (2) # complete from (2) and S(2)(2)

(4) S \rightarrow S + M \bullet (0) # complete from (2) and S(2)(1)

(5) S \rightarrow S \bullet + M (0) # complete from (4) and S(0)(2)

(6) P \rightarrow S \bullet (0) # complete from (4) and S(0)(1)
```

The state (P  $\rightarrow$  S  $\bullet$ , 0) represents a completed parse. This state also appears in S(3) and S(1), which are

## See also [edit]

- · CYK algorithm
- · Context-free grammar
- Parsing Algorithms

## Citations [edit]

- 1. ^ Kegler, Jeffrey. "What is the Marpa algorithm?" & Retrieved 20 August 2013.
- 2. \* Earley, Jay (1968). An Efficient Context-Free Parsing Algorithm (PDF). Camegie-Mellon Dissertation.
- 3. ^ Earley, Jay (1970), "An efficient context-free parsing algorithm", Communications of the ACM 13 (2): 94–102, doi:10.1145/362007.362035 ₺
- A Jurafsky, D. (2009). Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition 

  P. Pearson Prentice Hall. ISBN 9780131873216.

## Other Reference Materials [edit]

- Aycock, John; Horspool, R. Nigel (2002). *Practical Earley Parsing. The Computer Journal* **45** (6). pp. 620–630. doi:10.1093/comjnl/45.6.620 ₭.
- Leo, Joop M. I. M. (1991), "A general context-free parsing algorithm running in linear time on every LR(*k*) grammar without using lookahead", *Theoretical Computer Science* **82** (1): 165–176, doi:10.1016/0304-3975(91)90180-A ₺, MR 1112117 ₺.
- Tomita, Masaru (1984). "LR parsers for natural languages". COLING. 10th International Conference on Computational Linguistics. pp. 354–357.

## External links [edit]

#### C/C++ Implementations [edit]

- 'Earley Parser' & An Earley parser C library.
- 'C Earley Parser' 

  An Earley parser C.

#### Haskell Implementations [edit]

'Earley' ☑ An Earley parser DSL in Haskell

#### Java Implementations [edit]

- Pep 

   A Java library that implements the Earley algorithm and provides charts and parse trees as parsing artifacts.
- [1] A Java implementation of Earley parser.

#### JavaScript Implementations [edit]

- 'JavaScript Moony Parser' 

  A type of Earley parser written in JavaScript.
- Nearley & An Earley parser that's beginning to integrate the improvements that Marpa adopted.

#### Perl Implementations [edit]

- Marpa::R2 ₺, a Perl module. Marpa ₺ is an Earley's algorithm that includes the improvements made by Joop Leo, and by Aycock and Horspool.
- Parse::Earley ☑ A Perl module that implements Jay Earley's original algorithm.

## Python Implementations [edit]

- Charty ☑ a Python implementation of an Earley parser.
- NLTK a Python toolkit that has an Earley parser.
- Spark 🗗 an Object Oriented "little language framework" for Python that implements an Earley parser.
- earley3.py A stand-alone implementation of the algorithm in less than 150 lines of code, including generation of the parsing-forest and samples.

#### Common Lisp Implementations [edit]

• CL-EARLEY-PARSER ☑ A Common Lisp library that implements an Earley parser.

## Scheme/Racket Implementations [edit]

• Charty-Racket ☑ A Scheme / Racket implementation of an Earley parser.

### Resources [edit]

 $\bullet$  The Accent compiler-compiler  $\ensuremath{\mathbb{Z}}$ 

Categories: Parsing algorithms | Dynamic programming

This page was last modified on 21 June 2015, at 13:25.

Text is available under the Creative Commons Attribution-ShareAlike License; additional terms may apply. By using this site, you agree to the Terms of Use and Privacy Policy. Wikipedia® is a registered trademark of the Wikimedia Foundation, Inc., a non-profit organization.

Privacy policy About Wikipedia Disclaimers Contact Wikipedia Developers Mobile view



