# Declarative Programming and Natural Language

Søren Jakob Løvborg

## Abstract

This paper analyzes benefits and challenges (together with possible solutions) of using natural language processing for data entry and computer programming.

The paper looks at data entry in existing declarative languages and the underlying relational model is analyzed, also covering the subject of semantic networks. Context-free grammars are described in the context of natural language parsing, and modified Earley parsing algorithm is described and implemented, dispensing with some unnecessary complexities of the original.

A number of challenges of natural language are described, from the difficulties of identifying and classifying lexemes, to the ambiguous constructions of everyday English.

The paper concludes that while natural language is too complex for computers to grok in the general case, natural language may be viable in highly domain-specific areas.

# Contents

# 1   Introduction

Ordinarily, computers have been programmed using imperative programming languages (such as assembler code, C or Java), characterized in that algorithms are explicitly spelled out.

Declarative programming (which include logic programming and functional programming) takes a higher-level, more mathematically stringent approach, in which only a goal is stated, and the computer is then tasked with finding an appropriate algorithm to reach this goal.

The fact that most programming remains imperative must be attributed to the fact that programmers are much better at designing algorithms than computers. Indeed, many declarative languages (such as Prolog, ML and Lisp) include imperative constructs, sacrificing purity for increased control over program execution.

Declarative programming is not without its advantages, however. Because it does away with the imperative notion of *state*, declarative programming is suitable for massive parallelization.

More interestingly, one may argue that declarative programming more closely matches the way humans represent knowledge in writing.

While imperative programming requires understanding of state and flow control, declarative programming only requires understanding of knowledge representation and syntax. This may give non-specialists insight into the operations of a program.

Another way to make the life easier for non-specialists is to use *natural language* (i.e. English) for programming. This is considerably more difficult to implement properly, as can be seen by the relatively few attempts (successful or otherwise) at doing this. However, its feasibility for domain-specific tasks is demonstrated by such diverse experiments as SHRDLU and the Inform 7 programming language.

## 1.1   Problem description

This paper analyzes benefits and challenges (together with possible solutions) of using natural language processing for data entry and computer programming. However, many of the problems are beyond the scope of this paper, and are indeed unlikely to be solved for many years.

The paper describes a number of existing declarative programming languages and natural language based languages in section 2.

The nature of words in the English language is discussed in section 3, while section 4 discuss the mathematics of language and grammar. Section 5 describes a modified Earley parser, and section 6 an actual implementation.

Representing knowledge using semantic networks is discussed in section 7.

# 2 Overview of existing solutions

Since the 1950s, much work has gone into designing languages for declarative programming, resulting in logic programming languages like Prolog and Gödel, and functional languages like Lisp, ML and Haskell.

The topic of natural programming languages has not received as much attention. Practical English-like languages have been crude, as well as few and far between. Notable attempts include AppleScript and Macromedia's Lingo (both inspired by Dan Winkler's 1987 HyperTalk language), SQL and COBOL, all of which are still in use in some form or other, for application scripting, database queries, and financial systems.

## 2.1 Mathematical notation

Although not strictly a *programming* language, mathematical notation is the most widespread formalized declarative language in use, and a language that has evolved over thousands of years.

Recall that a mathematical $n$-ary *relation* is an $n$ argument boolean-valued function (a *predicate*), defining how values relate to eachother. A simple example is the binary equals relation, defined

$$\text{equals}(a, b) \quad \Leftrightarrow \quad a = b$$

As an example, the following relations are introduced:

$\text{parent}(p, c)$        $p$ is the parent of $c$.
$\text{father}(p, c)$        $p$ is the father of $c$.
$\text{mother}(p, c)$      $p$ is the mother of $c$.
$\text{childof}(c, m, f)$    $c$'s mother is $m$, and $c$'s father is $f$.
$\text{sibling}(c_1, c_2)$     $c_1$ is the sibling of $c_2$.

The relations are defined in terms of each other as follows:

$$\text{parent}(p, c) \Leftarrow \text{father}(p, c)$$

$$\text{parent}(p, c) \Leftarrow \text{mother}(p, c)$$

$$\text{childof}(c, m, f) \Leftrightarrow \text{mother}(m, c) \wedge \text{father}(f, c)$$

$$\text{sibling}(c_1, c_2) \Leftarrow c_1 \neq c_2 \wedge \exists p : (\text{parent}(p, c_1) \wedge \text{parent}(p, c_2))$$

We can then make assertions about people to construct a family tree. An example drawn from norse mythology is given in figure 1.

$$\text{childof}(\text{odin}, \text{bestla}, \text{borr})$$
$$\text{childof}(\text{vili}, \text{bestla}, \text{borr})$$
$$\text{childof}(\text{ve}, \text{bestla}, \text{borr})$$



Figure 1: *A family tree in terms of* childof*-relations, and as a pedigree chart*

## 2.2 Prolog

Prolog is a logic programming language, with a pure declarative core and a number of imperative extensions.

In predicate logic, a *definite Horn clause* is an expression on the form

$$A \Leftarrow B_1 \wedge B_2 \wedge \cdots \wedge B_n$$

where $A$ and $B_i$ are all predicates of the form $P(p_1, \ldots, p_m)$.

A basic Prolog program consists of a number of Horn clauses, and the user can then ask the program to prove predicates. Central to Prolog is *negation as failure*, in which anything not provably true is assumed to be false.

Since Prolog is limited to Horn clauses, the rule

$$\text{childof}(c, m, f) \Leftrightarrow \text{mother}(m, c) \wedge \text{father}(f, c)$$

cannot be described. The best we can do is

$$\text{childof}(c, m, f) \Leftarrow \text{mother}(m, c) \wedge \text{father}(f, c)$$

This restriction greatly reduces the complexity of the Prolog theorem prover, but complicates programming.

```
parent(P,C) :- father(P,C).
parent(P,C) :- mother(P,C).
childof(C,M,F) :- mother(M,C), father(F,C).
sibling(C1,C2) :- parent(P,C1), parent(P,C2), \dif(C1, C2).
```

Due to the modified childof definition, the family tree cannot be specified using

```
childof(odin, bestla, borr).
childof(ve, bestla, borr).
childof(vili, bestla, borr).
```

since childof no longer implies the mother and father relation. Instead, the family tree must be specified more verbosely,

```
mother(bestla, odin).
mother(bestla, ve).
mother(bestla, vili).
father(borr, odin).
father(borr, ve).
father(borr, vili).
```

Once this is done, queries can be made against the data set:

```
| ?- childof(odin,bestla,borr).
true ? a
no

| ?- childof(C,M,F).
C = odin
F = borr
M = bestla ? a

C = ve
F = borr
M = bestla

C = vili
F = borr
M = bestla

yes
```

## 2.3 Relational database management systems and SQL

Codd's classic paper on relational database management revolutionized database construction by introducing the mathematical concept of relations to data management. [Codd70] Where earlier databases had used a simple table–record structure, and the merging of data from different tables would have had to be done in the application, this merger could now be performed in the database.

Codd's definition of *named relations* match what SQL refers to as *tables*: A named set of $n$-tuples satisfying the $n$-ary relation.

The father-relation could be created as a table with this command:

```
CREATE TABLE father (parent VARCHAR(50), father VARCHAR(50))
```

However, working with father and mother in this manner would be impractical; in SQL, it's easiest to take an "object-centrist" approach, in which each table row (or tuple) tells us as much as possible about a single object. The result is more akin to the childof relation, the names being the only difference:

```
CREATE TABLE Person (name VARCHAR(50), mother VARCHAR(50), father VARCHAR(50))
```

4

## 2.4   SHRDLU

Terry Winograd's SHRDLU from the late 1960s was a computer program that allowed the user to interact with a virtual world of three-dimensional shapes through natural language. At the time, its apparent intelligence was a tremendous breakthrough in artificial intelligence and natural language research, but its successful understanding of language was largely dependent on the narrow confines of its virtual world, and later research failed to make the program less *domain-specific* while retaining its (apparent) intelligence. Winograd himself declared it to be a dead-end in AI research.

Nevertheless, SHRDLU demonstrates a number of desirable qualities in a natural language processing program, as seen in these excerpts: [Winograd71]

*Basic sentences*
\> Pick up a big red block.
OK.

*Comparative adjectives ("taller") and pronouns ("it")*
\> Find a block which is taller than the one you are holding and put it into the box.
By "it", I assume you mean the block which is taller than the one I am holding. OK.

*Negation and counting*
\> How many blocks are not in the box?
Four of them.

*Disambiguation*
\> How many things are on top of green cubes?
I'm not sure what you mean by "On top of" in the phrase "on top of green cubes".
Do you mean:
1 – Directly on the surface
2 – Anywhere on top of?

*Memory of past events*
\> Had you touched any pyramid before you put the green one on the little cube?
Yes, the green one.

*Learning by definition*
\> A "steeple" is a stack which contains two green cubes and a pyramid.
I understand.
\> Are there any steeples now?
No.
\> Build one.
OK.

Greatly simplifying SHRDLU's task was the fact that it had perfect knowledge of the state of its world, although not of the *rules* – it had to *learn* that triangles wouldn't stack, for instance.

## 2.5  Inform 7

Inform 7 is a language and authoring system designed largely for writers, not programmers. A notable feature is that the programming language is based on natural language.

Thus, we can define the father, mother and childof relations as such:

**Chapter 1 - Definitions**

Fatherhood relates one man (called the father) to various people.
Motherhood relates one woman (called the mother) to various people.

The verb to be the father of implies the fatherhood relation.
The verb to be the mother of implies the motherhood relation.

Childhood relates a person (called the child) to a person (called the parent) when the parent is the father of the child or the parent is the mother of the child.

The verb to be the child of implies the childhood relation.

**Section 1.1 - Family tree**

Bestla is the mother of Odin.
Borr is the father of Odin.

Odin is the father of Balder.
Frig is the mother of Balder.
Odin is the father of Hod.
Odin is the father of Thor.

The example illustrates how Inform 7 allows clear and concise declarations in natural language, at the cost of some quite verbose definitions.

Of particular interest is the readability of the source code. Even to readers having no previous experience with the language, the code remains quite clear, simply because it mimics the English language. In comparison, the Prolog code may be a lot more concise, but it's also completely opaque to people without the necessary mathematical background.

In this respect, Inform 7 can be seen as a perfection of Donald Knuth's "literate programming" philosophy. Unlike most realizations (TEX, javadoc, etc.), where programming code and documentation are merely interleaved in the source files, Inform 7 conflates the two: the documentation *is* the code.

[Nelson05] discuss this and other Inform 7 design issues.

# 3   Lexical analysis

During execution, a natural language processing program receives a string of characters as input. Since few non-logographic languages are best described as a string of separate characters, the input characters must be grouped into words (*tokens* or *terminals*), before being sent to the parser. This procedure is known as lexical analysis, scanning or tokenization.

It is the first step that a program must undertake when dealing with textual input, and for programming languages is often quite simple.[1]

## 3.1   Words of natural language

Of course, nothing in natural language is simple. [Trask04] gives four markedly different definitions of what a word is, two of which are useful for our purposes.

**Orthographic words** are strings of letters, numbers and hyphens, delimited by spaces or other non-alphanumeric characters. Orthographic words are thus readily identifiable in a text, and a suitable basis for tokenization.

This division does not always match the semantics, however. "Ice cream" is semantically a single word, but consists of two orthographic words.

**Lexemes** are semantic units of the text. A lexeme like JUMP correspond to the ortographic word "jump", as well as all its inflections (jumps, jumped, jumping). The set of known lexemes is referred to as the *lexicon*.

A lexeme always belong to exactly one syntactic class. As such, BILL$_N$ (the noun "a bill"), BILL$_V$ (the verb "to bill") and BILL$_{PN}$ (the proper noun "Bill") are three different lexemes.

A lexeme may also correspond to a sequence of orthographic words, as in "ice cream" (ICE-CREAM$_N$).

Determining exactly which orthographic words match which lexemes is a complicated procedure which requires understanding of the sentence structure, something which we only have after the parsing step. Hence the parser is fed othographic words and not lexemes.

## 3.2   Compound nouns

The above mentioned "ice cream" is an example of a compound noun. Compound nouns are nouns that consist of two or more words, and the challenge is to properly identify them as a single lexeme.

---

[1] A notable exception is C++, where the lexical analysis is intertwined not only with the parsing, but also with the following semantic analysis.

Although it may seem desireable for the semantic step to analyse the underlying composition (ice + cream, ash + tray), this is very difficult to do properly even for humans, as the compound does not reveal how the parts relate to eachother: Ice cream is cream made out of ice, but an ashtray is not a tray made out of ash.

One may distinguish between three ways of writing compound nouns:

The **closed form** is when the words are joined together, as in "ashtray", and is thus quite easy to handle from a parsing perspective, as we have a one-to-one mapping between orthographic word and lexeme.

The **hyphenated form** is when the words are connected by a hyphen, and is equally easy to parse.

The **open form** is when the words are separated by spaces, as in "ice cream", and presents a difficult parsing problem, easily resulting in ambiguities.

Compound nouns are further complicated by the fact that while many are limited to one of the above forms (e.g. "ice age", not "iceage" or "ice-age"), others can be written in two or all three of the above forms (both "ice cream" and "ice-cream" are correct). Additionally, some forms are in widespread use, even if not sanctioned by dictionaries ("ash tray").

One might conclude that program should accept compound nouns in all three forms, but this causes problems on its own: A "green house" is quite different from a "greenhouse".

For general-purpose natural language processing, it is not practical to build a comprehensive catalog of compound nouns either, since their number is too high, new compounds are invented every day, and old compounds change form quickly – "data base" became "database" in less than two decades.

## 3.3    Classifying words

Determining which word classes (nouns, verbs, etc.) to divide words into is in itself a difficult task.

A single class for *proper nouns* and another for *common nouns* may for instance seem appropriate, but English have two different types of common nouns: *count nouns* and *mass nouns*. Count nouns are ordinary nouns (e.g. "cat") which are countable and have a singular and plural form. Mass nouns are nouns like "water", which are uncountable, and are neither singular nor plural.

To further complicate matters, mass nouns may in some circumstances be used as count nouns, e.g. "Which of these *oils* is better for cooking?" One may even have *a beef* with someone.

It seems reasonable then to distinguish between the count noun BEEF$_{CN}$ and the mass noun BEEF$_{MN}$ at the lexeme level.

## 3.4 Inflection and stemming

A lexeme has an associated *citation form* or *lemma*; the base form of the word. For count nouns, it is the singular, for verbs it's the infinitive, for mass nouns and proper nouns, it's the only form of the word (not counting the possessive case).

The lemma is inflected to produce the different forms of the lexeme. Regular words are inflected according to a simple set of rules, but English also contains a large number of irregular words, where the different forms cannot simply be computed from the citation form, and instead must be stored explicitly. (The worst offender is the verb *to be*.)

The regular inflection rules for plural count nouns can be described thus: If the word ends on "s", "sh", "x", "z", add "es". Otherwise, if the word ends on "y", replace that "y" with "ies". Otherwise, add "s".

In shorthand notation, the inflection rules become:

| | |
|---|---|
| *Plural* | s/sh/x/z → -es |
| | y → ies |
| | $\epsilon$ → s |
| *Posessive* | s/x/z → -' |
| | $\epsilon$ → 's |
| *Comparative* | &d → -der |
| | &g → -ger |
| | &t → -ter |
| | e → -r |
| | y → -ier |
| | $\epsilon$ → -er |
| *Superlative* | &d → -dest |
| | &g → -gest |
| | &t → -test |
| | e → -st |
| | y → -iest |
| | $\epsilon$ → -est |

(Here, & indicates any vowel.)

Once we can inflect lexemes, we can also *stem* orthographic words and obtain the corresponding lexeme along with its grammatical features. This is done by adding all lemmas and irregular inflections to a *multi map*, mapping to the corresponding lexeme. A multi map is required since an orthographic word may map to multiple lexemes (possibly in different forms).

To find the lexemes corresponding to a word, and the above rules are applied in reverse, and the result is looked up in the map. An implementation of this can be found in the `jp.lexer.Feature` and `jp.lexer.Lexicon` (appendix B.3 and B.3), and a sample program can be found in `jp.StemTest` (appendix B.4).

# 4 Formal languages

A formal language is a mathematical representation of language, in which a language $L$ is defined as a set of symbol strings.[2]

The symbols can be anything, and are often characters; in describing natural languages, it's useful to define symbols as words and punctuation marks. These symbols are denoted *terminal symbols* or *tokens*, and the set $\Sigma$ of allowed symbols must be finite.

The language set contains all strings that are valid in that language. An example of a finite language could thus be the language of twin Arabic digits,

$$L = \{00, 11, 22, 33, 44, 55, 66, 77, 88, 99\}$$

over the symbols

$$\Sigma = \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$$

Most languages, however, are infinite, thus requiring us to describe them in other ways than simply enumerating every legal string. Formal devices for this purpose are called *formal grammars*.

The Chomsky hierachy (figure 2) orders languages (and the grammars describing them) after the complexity of recognizing valid strings (and rejecting invalid strings), such that every category is a proper subset of the categories above it in the hierachy.

| Language | Automaton |
|---|---|
| Recursively enumerable | Turing machine |
| Recursive | Decider |
| Context-sensitive | Linear-bounded |
| Indexed | Nested stack |
| Context-free | Nondeterministic pushdown |
| Deterministic context-free | Deterministic pushdown |
| Regular | Finite |

Figure 2: *The augmented Chomsky hierarchy of formal languages* [Elder05]

The most well-known of these are *regular languages* (typically described by regular expressions, and recognizable by a finite automaton) and *context-free languages* (described by context-free grammars, such as BNF grammars, and recognized by a non-deterministic pushdown automaton.)

---

[2]These strings are sometimes referred to as *words*, but to avoid confusion (as we're discussing natural language), we'll use *word* in the senses discussed in section 3.1.

## 4.1   Context-free grammars

Context-free grammars (CFGs) can themselves be divided into several complexity categories:

$$\text{LL(k)} \subset \text{LR(1)} \subseteq \text{deterministic CFG} \subset \text{unambiguous CFG} \subset \text{all CFG}$$

A context-free grammar is said to be deterministic, if it can be implemented using a deterministic pushdown automaton (as opposed to a non-deterministic PDA).

It's preferable to deal with deterministic context-free languages, since they may be parsed using efficient $\text{LL}(k)$ and $\text{LR}(1)$ parsers. For this reason, most programming languages strive to be deterministic context-free.

A grammar is said to be ambiguous, if there's more than one way to derive the same string from the grammar. An example of an ambiguous CFG would be

$$
\begin{aligned}
\mathsf{S} &\rightarrow \mathsf{A\ A} \\
\mathsf{A} &\rightarrow \mathsf{1} \ \mid \ \epsilon
\end{aligned}
$$

which generates the language $\{\epsilon, 1, 11\}$ but has two ways of deriving the string 1 in that the 1 may be the expansion of either the first or the second $\mathsf{A}$.

A context-free grammar can easily become ambiguous, as seen in this highly simplifed expression grammar:

$$\mathsf{expr} \ \rightarrow \ \mathsf{number} \ \mid \ \mathsf{expr} + \mathsf{expr} \ \mid \ \mathsf{expr} * \mathsf{expr}$$

Here, ambiguity causes problems because the individual productions are assigned semantics.

A smaller problem with the above grammar is that the expression $\mathsf{1\ +\ 2\ +\ 3}$ can be parsed as either $(1 + 2) + 3$ or $1 + (2 + 3)$. Either way, we get the correct mathematical result, and we can ignore the superfluous parse by e.g. requiring the parser to proceed from left to right.[3]



Figure 3: *The ambiguous grammar allows for multiple derivations.*

---

[3] This parser-level "hack" is usually employed to handle the *dangling else* problem of many programming languages.

Much worse is that the expression 1 + 2 * 3 can be parsed as either $(1 + 2) \cdot 3$ or $1 + (2 \cdot 3)$. In mathematics, we use operator precedence to remove this ambiguity, and we can encode this in the grammar:

$$
\begin{aligned}
\textsf{multiplicative} \quad &\rightarrow \quad \textsf{number} \quad | \quad \textsf{number} * \textsf{multiplicative} \\
\textsf{additive} \quad &\rightarrow \quad \textsf{multiplicative} \quad | \quad \textsf{multiplicative} + \textsf{additive}
\end{aligned}
$$

In general, however, it might not be that easy. Any sufficiently complex grammar of natural language will almost certainly be ambiguous and thus non-deterministic, requiring the use of more complex parsers than e.g. those used in parsing programming languages.

# 5  Parsing

While many simple and highly efficient algorithms exist for parsing LL and LR grammars, parsing algorithms for arbitrary context free grammars are somewhat more complex.

Several algorithms require that the grammar be rewritten to Chomsky normal form (for the Cocke-Younger-Kasami algorithm) or Greibach normal form (after which the input can be recognized by a simple non-deterministic push-down automaton).

These algorithms are general since the required normalization can be performed on any CFG, without changing the recognized language. However, since we attach semantic meaning to the productions of our grammar, rewriting the grammar eases parsing at the cost of complicating the semantic analysis.

Instead, the Earley parser may be applied, which accepts any context-free grammar in Backus-Naur form.

## 5.1  A modified Earley parser

[Earley70] skims over the issue of constructing a parse tree, and the described algorithm contains unnecessary complexities (such as look-ahead[4]), as noted by later authors. The following formulation of a modified Earley parser is thus based not only on Earley's paper, but also later resources ([Aycock02] and [Dollin02]).

**Input and preparations**

The parser takes as input:

- A set of terminals $a, b, \ldots$
- A set of nonterminals $A, B, \ldots$, of which one (denoted $R$) is the *root*
- A set of BNF productions $A \to \alpha$, where $\alpha$ is a string consisting of symbols from the two sets just mentioned
- An $n$ symbol string $X_0 \cdots X_{n-1}$ to be parsed

As a preparation, we choose a new (unused) terminal symbol $\dashv$ (the *terminator* and a new nonterminal $\phi$, which is promoted to root, and given the following production:

$$\phi \to R \dashv$$

$\dashv$ is appended to the input string as well, becoming $X_n$.

---

[4][Earley70] suggests a look-ahead of 1; [Aycock02] and [Dollin02] dispenses with look-ahead altogether.

**Parse trees**

The goal of the parser is to construct one or more valid parse trees, describing ways to parse the input according to the given productions.

A parse tree is an ordered tree consisting of nodes, which may be either simple *terminal nodes* (without child nodes), or *production nodes* (which *may* have child nodes). Each node is assigned a value, which for terminal nodes is a terminal symbol $a$, and for production nodes is a production $A \to \alpha$.

A production node can thus also be seen as a tuple

$$\langle A \to \alpha, \langle c_1, c_2, \ldots, c_n \rangle \rangle$$

consisting of the production $A \to \alpha$ and the child nodes $c_1, c_2, \ldots, c_n$.

In the following, we use the short-hand notation $A(c_1 c_2 \cdots c_n)$ to denote such a production node.[5] This enables the following notation for whole parse trees:

$$\phi(E(a + E(a)))$$

**Parse states**

Being a chart parser, the Earley parser constructs a number of parse states during execution.

These states are maintained in $n + 2$ *ordered* sets $S_0 \ldots S_{n+1}$: one set for every position between input symbols (as well as at the beginning and end of the input).

A parse state is a tuple $\langle P \to \alpha, p, i, T \rangle$, where $P \to \alpha$ is a production we're currently parsing, $p$ is our current position in the string $\alpha$ (the expansion of that production), $i$ denotes the input position at which we began parsing $P$, and $T$ is a partial parse tree, which may end up in the resulting parse tree(s).

As a shorthand, we write

$$P \to \alpha \ . \ \beta \quad i \quad T \tag{1}$$

to denote the state $\langle P \to \alpha\beta, p, i, T \rangle$ where the position $p$ is between $\alpha$ and $\beta$, as denoted by the dot.

In a state set $S_j$, the state (1) represents the following facts:

- We're currently testing whether input characters starting with $X_i$ can be derived from $P \to \alpha\beta$.
- $\alpha \overset{*}{\Rightarrow} X_i \cdots X_{j-1}$, that is, the input symbols $X_i$ through $X_{j-1}$ have been verified to match the part of the production before the dot, $\alpha$.

---

[5] For the sake of simplicity, this notation dispenses with the right-hand side of the production, $\alpha$. An actual implementation would usually need to track the whole production $A \to \alpha$.

- $T$ is a partial parse tree for a parse starting at $X_i$ with root $P$.

If the dot is at the end of the production, we further have that

- $P \overset{*}{\Rightarrow} X_i \cdots X_{j-1}$ (the input symbols $X_i$ through $X_{j-1}$ can be derived from $P$.)
- $T$ is a complete parse tree for a parse of $X_i$ through $X_{j-1}$ with root $P$.

At the end of the parse, the following statements are equivalent:

- The state is in $S_{n+1}$.
- The state is $\phi \to R \dashv .$    $0$    $\phi(T \dashv)$.
- $\phi(T \dashv)$ is a complete parse tree from a parse of the entire input (including $\dashv$) with root $\phi$, from which the relevant parse tree $T$ with root $R$ can trivially be extracted.

### Algorithm

The algorithm starts out by adding a single state to $S_0$,

$$\phi \to . R \dashv \quad 0 \quad \phi()$$

indicating that we're currently testing whether our whole input can be derived from our root $\phi$ and that no input symbols have been verified to match so far.

The algorithm then iterates over the state sets $S_0$ through $S_n$ in ascending order, and for each set $S_j$, the algorithm processes the states of $S_j$ in order.

Processing a state in $S_j$ may cause new states to be added to $S_j$; these must be processed as well, thus the requirement that the sets be ordered.

Depending on the state to be processed, one of three actions may be taken.

**The scanner** applies when the symbol following the dot is a terminal, $a$:

$$P \to \alpha . a \beta \quad i \quad T$$

We compare $a$ to the next input symbol $X_j$, and, if it's a match, adds the state

$$P \to \alpha a . \beta \quad i \quad T$$

to $S_{j+1}$, indicating that we successfully parsed the $a$.

**The predictor** applies when the symbol following the dot is a nonterminal, $A$:

$$P \to \alpha . A \beta \quad i \quad T$$

The predictor adds a state to $S_j$ for every production $A \rightarrow \gamma$ of that nonterminal,

$$A \rightarrow \; . \; \gamma \quad j \quad A()$$

Thus, we begin a check of whether $A$ matches the input at position $j$.

**The completer** applies when the dot is at the end of the production:

$$A \rightarrow \gamma \; . \quad i \quad T$$

This means that we successfully parsed an $A$ at position $i$, so we go back to $S_i$, and for every state which predicted $A$,

$$P \rightarrow \alpha \; . \; A \; \beta \quad i' \quad P(s)$$

we add a new state to $S_j$

$$P \rightarrow \alpha \; A \; . \; \beta \quad i' \quad P(s \; T)$$

indicating that we successfully parsed $A$, with the resulting parse tree $T$.

**Example**

As an example, take the input $a + a$ and the following simple grammar:

$$
\begin{aligned}
\phi &\rightarrow E \dashv \\
E &\rightarrow a \quad | \quad E + E
\end{aligned}
$$

After parsing completes, we will have produced the following state sets:

|       |   |        |               |                |   |                               |
|-------|---|--------|---------------|----------------|---|-------------------------------|
|       |   | $\phi$ | $\rightarrow$ | $. \; E \; \dashv$ | 0 | $\phi()$                  |
| $S_0$ |   | $E$    | $\rightarrow$ | $. \; a$       | 0 | $E()$                         |
|       |   | $E$    | $\rightarrow$ | $. \; E \; + \; E$ | 0 | $E()$                    |
|       |   | $E$    | $\rightarrow$ | $a \; .$       | 0 | $E(a)$                        |
| $S_1$ | $*$ | $\phi$ | $\rightarrow$ | $E \; . \; \dashv$ | 0 | $\phi(E(a))$             |
|       |   | $E$    | $\rightarrow$ | $E \; . \; + \; E$ | 0 | $E(E(a))$                |
|       |   | $E$    | $\rightarrow$ | $E \; + \; . \; E$ | 0 | $E(E(a)+)$               |
| $S_2$ |   | $E$    | $\rightarrow$ | $. \; a$       | 2 | $E()$                         |
|       | $*$ | $E$    | $\rightarrow$ | $. \; E \; + \; E$ | 2 | $E()$                    |
|       |   | $E$    | $\rightarrow$ | $a \; .$       | 2 | $E(a)$                        |
|       |   | $E$    | $\rightarrow$ | $E \; + \; E \; .$ | 0 | $E(E(a) + E(a))$         |
| $S_3$ | $*$ | $E$    | $\rightarrow$ | $E \; . \; + \; E$ | 2 | $E(E(a))$                |
|       |   | $\phi$ | $\rightarrow$ | $E \; . \; \dashv$ | 0 | $\phi(E(E(a) + E(a)))$   |
|       | $*$ | $E$    | $\rightarrow$ | $E \; . \; + \; E$ | 0 | $E(E(E(a) + E(a)))$      |
| $S_4$ |   | $\phi$ | $\rightarrow$ | $E \; \dashv \; .$ | 0 | $\phi(E(E(a) + E(a)) \dashv)$ |

Blind alleys in the parsing process are marked with an $*$.

## 5.2   The epsilon problem

Earley briefly touches upon a problematic feature of many context-free grammars. Productions of the form $P \to \epsilon$, so-called epsilon productions, and more generally, any *nullable* nonterminal $E \overset{*}{\Rightarrow} \epsilon$.

Such a nullable nonterminal $E$ may be both predicted *and* completed in the same state set $S_j$, as no input symbols are scanned:

$$
\begin{aligned}
A &\to\ .\,E \quad j \quad A() \\
A &\to E\ . \quad\ j \quad A()
\end{aligned}
$$

Since we're still in the middle of processing $S_j$, we can't "go back" and iterate over all states of $S_j$, as a later predictor may yet add more states to $S_j$.

As a contrived example, take the input + and the grammar

$$
\begin{aligned}
\phi &\ \to\ \ S \dashv \\
S &\ \to\ \ E \ \mid\ P \\
P &\ \to\ \ Q + \\
Q &\ \to\ \ E \\
E &\ \to\ \ \epsilon
\end{aligned}
$$

Clearly, $\phi(S(P(Q(E())+)))$ is a valid parse. However, the algorithm rejects the input if we're not careful. During execution, we end up with the following states in $S_0$,

$$
\begin{array}{rclccl}
 & \phi & \to & .\,S \dashv & 0 & \phi() \\
* \ (1) & S & \to & .\,E & 0 & S() \\
 & S & \to & .\,P & 0 & S() \\
(2) & E & \to & . & 0 & E() \\
(3) & P & \to & .\,Q + & 0 & P()
\end{array}
$$

just as we're about to process state (2). As the dot is at the end of the production, we run the completer, and finding that (1) is the only state with $E$ to the right of the dot, the following state is added to $S_0$:

$$
* \qquad S \ \to\ E\ . \quad 0 \quad S(E())
$$

However, then running the predictor on (3) results in a new state being added to $S_0$:

$$
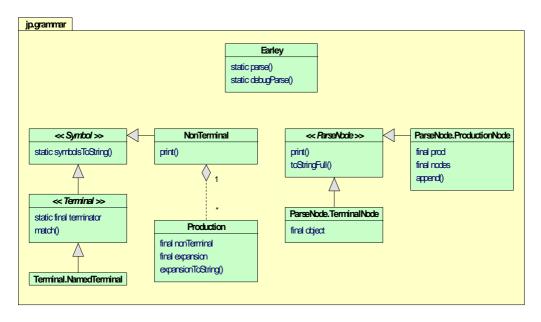Q \ \to\ \ .\,E \quad 0 \quad Q()
$$

This state too has $E$ to the right of the dot, but we missed it when we ran the completer earlier.

**Solutions**

[Aycock02] suggests that the parser keep track of the nullable productions, and preemptively complete them in the predictor step. This is certainly an excellent solution when implementing a recognizer, but there's no easy way to generate a valid parse tree, if nullable productions are simply skipped in this fashion.

A clean and simple solution, and the solution used by this author, is to alternate between running the predictor and completer, until neither has any more states to add to $S_j$. After this, the scanner may be run on all appropriate states in $S_j$ to construct $S_{j+1}$.

To speed this up, the predictor and completer are only run again if the predictor predicted a nullable production. (Which productions are nullable can trivially be determined before the parsing starts.)

# 6    Parser implementation

Figure 4: *The classes of the* jp.grammar *package.*

The modified Earley parser described above is implemented in the non-instantiable class jp.grammar.Earley, which provide two static methods:

```
static public Set<ParseNode.ProductionNode>
    parse(NonTerminal root, Object... input)

static public Set<ParseNode.ProductionNode>
    debugParse(NonTerminal root, Object... input)
```

18

`debugParse` is identical to `parse`, except that it also prints the generated parse states to standard output, for debugging.

The grammar used for parsing must be given in terms of `jp.grammar.Terminal`, `jp.grammar.NonTerminal` and `jp.grammar.Production` objects, which represent the corresponding mathematical objects.

A simple example of how to use the parser is given in `jp.ExprTest`, from which the following snippet is taken:

```
11          Terminal a = new Terminal.NamedTerminal("a");
12          Terminal plus = new Terminal.NamedTerminal("+");
13          NonTerminal E = new NonTerminal("E");
14
15          new Production(E, a);
16          new Production(E, E, plus, E);
17
18          Set<ParseNode.ProductionNode> result = Earley.debugParse(E, a, plus, a);
19
20          System.out.println("\nParse trees:");
21          for (ParseNode.ProductionNode l: result) l.print(0);
```

This implements the grammar described in the example of section 5.1 (page 16).

In the code above, `debugParse` is used to display the internal parse states, and the returned parse trees are simply printed. (A typical application would instead have to interpret the parse trees.)

## 6.1   Grammar compiler

A complex context-free grammar converted to Java code isn't a pretty sight, and rather tiresome to write. Instead, the `jp.GC` grammar compiler may be used to convert a grammar source file written in a BNF-like syntax.

The grammar compiler accepts the file names of one or more grammar specifications, and outputs the corresponding Java code to standard output.

The syntax is quite simple. Each line describes one or more alternatives for a given nonterminal, using the following syntax:

$$
\begin{aligned}
\text{line} \quad &\rightarrow \quad \text{nonterminal ``='' symbolString alternatives} \\
\text{alternatives} \quad &\rightarrow \quad \epsilon \\
&\mid \quad \text{``|'' symbolString alternatives} \\
\text{symbolString} \quad &\rightarrow \quad \epsilon \\
&\mid \quad \text{symbol symbolString}
\end{aligned}
$$

Here, nonterminal and symbol are arbitrary Java identifiers. The generated code will define local variables for the nonterminals. Presumably, symbols refer to either nonterminals defined by the grammar, or terminals defined by the host application; the grammar compiler does not verify their validity.

On all lines containing the comment symbol #, # and everything following it is ignored. Blank lines are ignored as well.

A line may be continued by indenting the following lines by one or more whitespace characters.

The grammar syntax can describe itself as follows:[6]

```
line = symbol EQUAL symbolString alternatives

alternatives = # epsilon
             | PIPE symbolString alternatives

symbolString = # epsilon
             | symbol symbolString
```

(Here, EQUAL and PIPE must be defined by the host application to represent the equal sign and pipe symbol respectively.) The result is shown in figure 6.1.

A larger example can be found in appendix B.5, and the generated code in appendix B.4.

```
NonTerminal symbolString = new NonTerminal("symbolString");
NonTerminal line = new NonTerminal("line");
NonTerminal alternatives = new NonTerminal("alternatives");
new Production(symbolString);
new Production(symbolString, symbol, symbolString);
new Production(line, symbol, EQUAL, symbolString, alternatives);
new Production(alternatives);
new Production(alternatives, PIPE, symbolString, alternatives);
```

Figure 5: *The Java code generated by the grammar compiler*

Since Java does not support include files, and since the generated code references symbols defined by the host application, the generated code must be pasted into the appropriate place in the host application source code, either manually or using some external tool to merge the compiler output into the host application code. The gram target of this project's makefile uses sed for this purpose.

---

[6]The parser of the grammar compiler is not actually implemented using the jp.grammar package, but instead employs a handwritten non-recursive parser.

# 7  Semantics

Semantics is the study of meaning and understanding, a field in which philosophy and science intersects. In natural language processing, *semantic analysis* translates (ambiguous) natural language into simpler, more well-defined (usually non-ambiguous) elements.

A central notion in semantics is the basic unit of knowledge, which I'll refer to as *idea*, a word used by John Locke to denote "whatever is the object of understanding when a man thinks" in his *Essay Concerning Human Understanding.* [Locke90] Ideas thus include all things, whether concrete physical objects (e.g. the Eiffel tower), physical properties (the color purple), abstract classes (mammals), as well as abstract concepts (good and evil).

[Locke90] also puts emphasis on the duality of knowledge and language, suggesting the conclusion that humans cannot think what they cannot put into words (nor, obviously, put into words what they cannot think).[7]

This serves as a caution against straying too far from the original natural language when doing the semantic processing.

## 7.1  Meanings of "to be"

In implementing semantics, it's important to recognize the different meanings of "to be" ("is the same as", "is a subset of" and "is characterized by"), as demonstrated in the following sentences:

*Cats are felines.* "cats" and "felines" are the same; both refer to animals of the *Felidae* family. (Although "cat" can also refer specifically to the domestic cat, a person, and a number of other things.) We say that cat = mammal.

*Cats are mammals.* Cats are a proper subset of mammals. We say that cat is-a mammal.

*Cats are cute.* Cats are characterized by being cute. This paper will use the notation cat is cute.

## 7.2  Semantic networks using is-a

A semantic network is a directed graph, in which each node represents an idea, and the edges are binary relations, typically is-a-relations.

Each idea is further assigned a number of properties. Although some semantic network implementations make a distinction between properties and nodes, in the

---

[7] As Wittgenstein put it in his 1922 *Tractatus Logico-Philosophicus*: "The limits of my language mean the limits of my world." Also, compare Orwell's *Nineteen Eighty-Four*, in which the goal of the *Newspeak* language is that only loyal thoughts be expressible, and thus, thinkable.
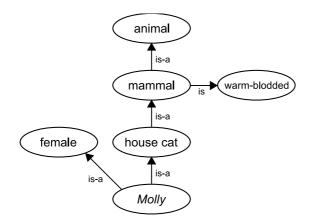
Figure 6: *A semantic network*

general case, properties (like "warm-blooded" in figure 6) are just ideas, linked
with a different type of relation (here denoted is). This way, the semantic network
can also make assertions about properties, and properties may be assigned other
properties.

The interesting feature of is-a is that it enables properties to be inherited. Once
we know that mammal is warm-blooded and Molly is-a house-cat is-a mammal, we
can assume that Molly is warm-blooded.

## 7.3 Cancellation

Property inheritance, however, is only an assumption. The classic example is
bird is capable-of-flight, which is usually true, but not always, e.g. in the case
of penguins.

To encode this in a semantic network, one explicitly asserts the opposite (penguin
is-not capable-of-flight), thus *cancelling* inheritance of the capable-of-flight property.

## 7.4 The duck test

The is-a inheritance can also be used in reverse, to determine the class of an idea.
For instance, assuming that all we know about Tom is that Tom is warm-blooded,
we might then assume that Tom is-a mammal.

Each is-a relation is thus accompanied by an implicit reverse relation, which this
paper will refer to as "implies".

[Brachman83] correctly argues that just because an entity has certain properties, all
of which are typical for a certain class, that entity is not necessarily a member of

that class.

However, inductive reasoning prescribes the opposite: If an entity has all or some of the known properties of a class, and these properties are unique to that class, we must assume that entity to belong to that class, at least until new information comes to light and suggests otherwise.

This boils down to an old philosophic discussion about the *identity of indiscernibles*, aptly formulated by James W. Riley in an oft-quoted poem:

> When I see a bird that walks like a duck and swims like a duck and quacks like a duck, I call that bird a duck.

On a mathematical level, the identity of indiscernibles is still debated, but in a physical reality, it's an indispensable rule, and humans certainly apply this so-called *duck test* all the time to determine identities.

## 7.5   Uncertainty

[Brachman83] has a point though: our assumption about Tom are based on a fragile premise: that no other animal (or generally, no other idea) could be warm-blooded. Indeed, Tom might turn out to be another warm-blooded animal, such as a white shark. To deal with this, relations in the semantic network can be assigned a numerical level of confidence.

To gauge the confidence of "warm-blooded implies mammal", however, we must assess the probability of new, contradictory information (such that the fact that non-mammals may be warm-blooded) appearing, which is of course difficult. 50% may be as good as any other value.

The remaining 50% must then be divided between mammal and white shark, recognizing that

- We know of one mammal, but of no white sharks, which arguably makes mammal more likely than shark.
- We don't know if the class of mammals and the class of white sharks overlap. Tom could be both.

Again, the best we can do is a wild guess.

Another problem is that "warm-blooded" isn't a scientifically well-defined term, and some may thus disagree with the claim that white sharks are warm-blooded. We may indicate this by assigning a confidence of (say) 50% to the white-shark is warm-blooded relation. Mammals, on the other hand, are most certainly warm-blooded.

The result is a semantic network looking like figure 7. At this point, all we need to do is to apply a feed-back mechanism that adjusts the confidence levels in response to user input, and we'll essentially have a neural network.
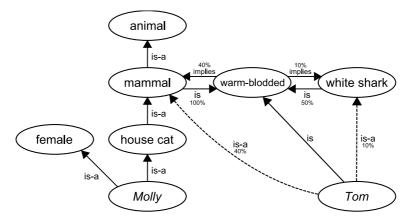


Figure 7: *The updated semantic network*

## 7.6 Individuals versus generics

[Brachman83] argues that one should treat *generics* (sets, e.g. humans) and *individuals* (set members, e.g. John) differently. This leads to a distinction between two kinds of is-a relationships, namely $\subseteq$ and $\in$.

But for many purposes, treating individuals as single-member sets is a usuable abstraction, such that only $\subseteq$ remains:

$$\mathsf{John} = \{john\} \subseteq \mathsf{humans}$$

In this representation, the value *john* is never directly referred to. As such we cannot enumerate the known elements of the set humans, but we *can* enumerate the known *subsets* (which include John).

## 7.7 Meta-semantics

Take the following sentence, "A person is either male or female." The sentence carries an implicit "but not both", which can be expressed using predicate logic as such ($\forall p$):

$$\mathsf{person}(p) \Rightarrow (\mathsf{male}(p) \vee \mathsf{female}(p)) \wedge \overline{\mathsf{male}(p) \wedge \mathsf{female}(p)}$$

– or equivalently using set logic:

$$\mathsf{persons} \subseteq (\mathsf{males} \cup \mathsf{females}) \cap \overline{\mathsf{males} \cap \mathsf{females}}$$

24

The similar sentence, if interpreted literally, "A person can be either male or female." is more vague. The use of "either" still carries an implicit "but not both", but the exchange of "is" for "can be" means that a person *also* could be something else entirely. The sentence thus boils down to "A person cannot *both* be male and female".

$$\mathsf{person}(p) \Rightarrow \overline{\mathsf{male}(p) \wedge \mathsf{female}(p)}$$

Finally, "A person can be male or female." dismisses with both "either" and "is", reducing the sentence to the tautology

$$\mathsf{person}(p) \Rightarrow \mathsf{male}(p) \vee \mathsf{female}(p) \vee \mathsf{true}$$

The sentence is thus devoid of semantics, but not *meta-semantics*: It tells us something very important, namely that some persons *may* possess the properties male and female.

Such semantics may be described using a "can-be" relation in the semantic network. For a general solution, a second-order predicate $\mathrm{can}(P, x, y)$ may instead be introduced, where $P$ is any first-order relation ("is", "is-a", "implies", etc.), $x$ the class of ideas that the relation *could* apply to, and $y$ the associated idea. For instance:

$$\mathrm{can}(\mathsf{is}, \mathsf{person}, \mathsf{female})$$

Note that in order to represent second-order predicates in a semantic network, the different types of relations must be represented as ideas as well. To represent *tertiary* relations like the above, it may be necessary to have *relation instances* (like person is female) represented as ideas as well. At this point, it becomes clear that the semantic network is reasoning about its own reasoning.

# 8 Challenges in the English grammar

The grammar of natural language is immensely complex, and hard to pin down as well. As [Winograd71] says, "the task of codifying an entire language in any formalism is so large that it would be folly to try in the course of a single research project."

Different grammars may accept the same sentences, but result in different parse trees, allowing different levels of insight into the semantics of the sentence.

Take the simple sentence "John ate the apple." Classic English grammar renders this as as sentence consisting of a noun phrase ("John") and verb phrase or *predicate* ("ate the apple").

A *systemic grammar*, as for instance employed by SHRDLU, would render the sentence as a clause consisting of a noun group ("John"), a verb group ("ate") and another noun group ("the apple"). This is illustrated in figure 8.
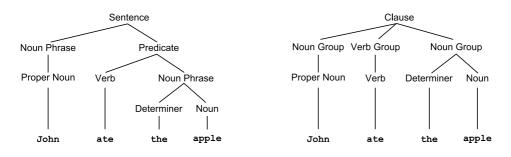


Figure 8: *Same sentence, different grammars*

## 8.1 Participial ambiguity

One difficult challenge is sentences that can only be interpreted in context of other sentences. [Allen95] gives the following example of this, in the form of a participial ambiguity:

```
Visiting relatives can be trying.
```

"Visiting" is the present participle of VISIT$_V$; however, "Visiting relatives" is a noun phrase that can be interpreted in two ways:

| Participal use | Interpretation | Head | Number |
|---|---|---|---|
| Gerund | the act of visiting relatives | "visiting" | singular |
| Verbal adjective | relatives that are visiting | "relatives" | plural |

As can be seen, the two interpretations have different heads and grammatical number, which would have allowed disambiguation had the sentence been `Visiting relatives are trying.` or `Visiting relatives is trying.`.

The first sentence remains ambiguous because "can be" does not indicate number.

## 8.2   Phrasal verbs

Phrasal verbs are compounds consisting of a verb together with one or more particles. Phrasal verbs are grammatical exceptions, in that they do not simply denote that the verb is modified by the particles, but have an entirely new meaning.

An example is the intransitive phrasal verb "give up", which has nothing to do with either giving nor the upwards direction. Additionally, the particle is fixed for this particular meaning of the verb; the syntactically equivalent "give down" is meaningless, for instance.

As such, "give", "give up", "give in", "give over", etc. are all different verbs, and all but the first are phrasal verbs.

One way to deal with phrasal verbs is to incorporate them into the grammar, but since a goal is often to design a simple, unchanging grammar (with only the lexicon growing over time), this is less than ideal.

# 9   Conclusion

The Earley parser is a simple and efficient parser for dealing with arbitrary context-free grammars, although Earley's original formulation has since been improved upon, as detailed in this paper. However, parsing is only the first step a program must take to gain understanding of natural language, and models of semantics and knowledge representation (such as semantic networks) have yet to reach a level of maturity, where they may be used for general purpose natural language processing.

Computers will not understand natural language any day soon. Not only is the rules of natural language incredibly complex, possibly more complex than today's computers can handle, but they're indeed too complex for even humans to fully comprehend them, cf. the continuing research into which kind of grammar is more suitable.

Still, even if they are few and far between, a number of projects have demonstrated that natural language may be viable for highly domain-specific programming and data entry, cf. SHRDLU and Inform 7, just like artificial intelligence can compete with humans in specific domains, such as chess.

# A  References

[Allen95]       James Allen. *Natural language understanding*. 2nd edition, 1995.

[Aycock02]      John Aycock and R. Nigel Horspool. *Practical Earley parsing*. The
                Computer Journal, 45, no. 6; 620–630, 2002.

[Brachman83]    Ronald J. Brachman. *What IS-A is and isn't: An analysis of tax-*
                *onomic links in semantic networks*. IEEE Computer, 16, no. 10;
                30–36, 1983.

[Codd70]        Edgar F. Codd. *A relational model of data for large shared data*
                *banks*. Communications of the ACM, 13, no. 6; 377–387, Jun 1970.

[Dollin02]      Chris Dollin. *Earley parser*. comp.programming, Apr 2002.

[Earley70]      Jay C. Earley. *An efficient context-free parsing algorithm*. Commu-
                nications of the ACM, 13, no. 2; 94–102, Feb 1970.

[Elder05]       Murray Elder. *G-automata, counter languages and the Chomsky*
                *hierarchy*. London Mathematical Society Lecture Note Series no.
                339, 2005.

[Locke90]       John Locke. *An Essay Concerning Human Understanding*, 1690.

[Nelson05]      Graham Nelson. *Natural language, semantic analysis and interactive*
                *fiction*, 2005.

[Trask04]       Larry Trask. *What is a word?*, Jan 2004.

[Winograd71]    Terry Winograd. *Procedures as a representation for data in a com-*
                *puter program for understanding natural language*, 1971.

# B  Source

## B.1  jp.util package

The jp.util package contains two utility classes, one (MultiMap) serving as the
bass for jp.lexer.Lexicon, the other (TablePrint) provides pretty debug output
for jp.grammar.Earley.

### jp.util.MultiMap

```
1   package jp.util;
2
3   import java.util.*;
4
5   public class MultiMap<K, V>
6   {
7       private Map<K, Set<V>> map = new HashMap<K, Set<V>>();
8
9       private class MultiMapEntry<K, V> implements Map.Entry<K, V>
10      {
11          private final K key;
12          private final V value;
13
14          public MultiMapEntry(K key, V value)
15          {
16              this.key = key;
17              this.value = value;
18          }
19
20          /** Compares the specified object with this entry for equality. */
21          public boolean equals(Object o)
22          {
23              if (!(o instanceof Map.Entry)) return false;
24
25              Map.Entry entry = (Map.Entry) o;
26              return (entry.getKey() == key && entry.getValue() == value);
27          }
28
29          /** Returns the key corresponding to this entry. */
30          public K getKey() { return key; }
31
32          /** Returns the value corresponding to this entry. */
33          public V getValue() { return value; }
34
35          /** Returns the hash code value for this map entry. */
36          public int hashCode() { return key.hashCode() ^ value.hashCode(); }
37
38          /**
39           *  Replaces the value corresponding to this entry with the specified
40           *  value (optional operation).
41           */
42          public V setValue(V value)
43          {
44              throw new UnsupportedOperationException();
45          }
46      }
47
48      /** Removes all mappings from this map. */
49      public void clear() { map.clear(); }
50
51      /** Returns true if this map contains a mapping for the specified key. */
52      public boolean containsKey(Object key) { return map.containsKey(key); }
53
```

```
54        /** Returns true if this map maps one or more keys to the specified value.*/
55        public boolean containsValue(Object value)
56        {
57            for (Set<V> values: map.values())
58                if (values.contains(value)) return true;
59
60            return false;
61        }
62
63        /** Returns a set view of the mappings contained in this map. */
64        public Set<Map.Entry<K,V>> entrySet()
65        {
66            Set<Map.Entry<K,V>> result = new HashSet<Map.Entry<K,V>>();
67            for (Map.Entry<K,Set<V>> entry: map.entrySet())
68            {
69                for (V v: entry.getValue())
70                    result.add(new MultiMapEntry<K, V>(entry.getKey(), v));
71            }
72            return result;
73        }
74
75        /** Compares the specified object with this map for equality. */
76        public boolean equals(Object o)
77        {
78            if (!(o instanceof MultiMap)) return false;
79            return map.equals(((MultiMap) o).map);
80        }
81
82        /** Returns the values to which this map maps the specified key. */
83        @SuppressWarnings("unchecked")
84        public Set<V> get(Object key)
85        {
86            Set<V> c = map.get(key);
87            if (c == null) return Collections.EMPTY_SET; // unchecked
88            return Collections.unmodifiableSet(c);
89        }
90
91        /** Returns the hash code value for this map. */
92        public int hashCode() { return map.hashCode(); }
93
94        /** Returns true if this map contains no key-value mappings. */
95        public boolean isEmpty() { return map.isEmpty(); }
96
97        /** Returns a set view of the keys contained in this map. */
98        public Set<K> keySet() { return map.keySet(); }
99
100       /** Associates the specified value with the specified key in this map. */
101       public void put(K key, V value)
102       {
103           Set<V> set = map.get(key);
104           if (set == null) map.put(key, set = new HashSet<V>());
105           set.add(value);
106       }
107
108       /** Associates the specified values with the specified key in this map. */
109       public void put(K key, Collection<V> values)
110       {
111           Set<V> set = map.get(key);
112           if (set == null) map.put(key, set = new HashSet<V>());
113           set.addAll(values);
114       }
115
116       /** Copies all of the mappings from the specified map to this map. */
117       public void putAll(Map<? extends K,? extends V> t)
118       {
119           for (Map.Entry<? extends K,? extends V> entry: t.entrySet())
120               put(entry.getKey(), entry.getValue());
121       }
122
123       /** Removes the mapping for this key from this map if it is present. */
```

```
124        public void remove ( Object key ) { map . remove ( key ) ; }
125
126        /** Returns the number of key mappings in this map. */
127        public int size () { return map . size () ; }
128
129        /**  Returns a set view of the values contained in this map. */
130        public Set <V> values ()
131        {
132            Set <V> result = new HashSet <V >() ;
133
134            for ( Map . Entry <K , Set <V >> entry : map . entrySet () )
135                result . addAll ( entry . getValue () ) ;
136
137            return result ;
138        }
139    }
```

## jp.util.TablePrint

```
 1   package jp . util ;
 2
 3   import java . util .*;
 4
 5   public class TablePrint
 6   {
 7       private List < String []> rows = new ArrayList < String []>() ;
 8       private int [] columnWidths ;
 9       private int cSpace ;
10
11       public TablePrint ( int columnCount , int cSpace )
12       {
13           columnWidths = new int [ columnCount ];
14           this . cSpace = cSpace ;
15       }
16
17       public void addRow ( String ... row )
18       {
19           assert ( row . length <= columnWidths . length ) ;
20           rows . add ( row ) ;
21           for ( int i = 0; i < row . length ; i ++)
22           {
23               if ( columnWidths [i] < row [i]. length () )
24                   columnWidths [i] = row [i]. length () ;
25           }
26       }
27
28       public void addRow ( Object ... row )
29       {
30           assert ( row . length <= columnWidths . length ) ;
31           String [] strs = new String [ row . length ];
32
33           for ( int i = 0; i < row . length ; i ++) strs [i] = row [i]. toString () ;
34           addRow ( strs ) ;
35       }
36
37       public void print ()
38       {
39           for ( String [] row : rows )
40           {
41               for ( int i = 0; i < row . length ; i ++)
42               {
43                   System . out . print ( row [i]) ;
44                   for ( int j = row [i]. length () ; j < columnWidths [i] + cSpace ; j ++)
45                       System . out . print ( " " ) ;
46               }
47               System . out . println () ;
48           }
49       }
50   }
```

## B.2    jp.grammar package

The jp.grammar package contains classes for managing a context-free grammar and an implementation of a modified Earley parser, as described in section 6. See also the class diagram on page 18.

### jp.grammar.Earley

```
1    package jp.grammar;
2
3    import java.util.*;
4    import jp.util.TablePrint;
5
6    public class Earley
7    {
8        private Earley() { }
9
10       private static class ParseState
11       {
12           Production prod; // production
13           int prodPos;      // position in the expansion of the production
14           int inputPos;     // position in input of the beginning of the production
15           ParseNode.ProductionNode parse;
16
17           public boolean equals(Object o)
18           {
19               ParseState s = (ParseState) o;
20               return s.prod == prod && s.prodPos == prodPos
21                   && s.inputPos == inputPos && s.parse.equals(parse);
22           }
23
24           ParseState(Production prod, int prodPos, int inputPos,
25               ParseNode.ProductionNode parse)
26           {
27               this.prod = prod;
28               this.prodPos = prodPos;
29               this.inputPos = inputPos;
30               this.parse = parse;
31           }
32
33           public String toString()
34           {
35               return prod.nonTerminal + " -> " +
36                   prod.expansionToString(prodPos) +
37                   "  " + inputPos;
38           }
39
40           boolean hasNextSymbol()
41           {
42               return (prodPos < prod.expansion.length);
43           }
44
45           Symbol nextSymbol()
46           {
47               return hasNextSymbol() ? prod.expansion[prodPos] : null;
48           }
49
50           // "Move the dot", that is, increase the prodPos.
51           ParseState moveOver(ParseNode.ProductionNode parse)
52           {
53               return new ParseState(prod, prodPos + 1, inputPos, parse);
54           }
55       }
56
57       // An ordered set, and thus not an ordinary Java Set.
```

```
58      private static class StateSet extends AbstractCollection<ParseState>
59      {
60          private List<ParseState> set = new ArrayList<ParseState>();
61
62          // For debugging
63          public void print()
64          {
65              TablePrint tbl = new TablePrint(4, 2);
66
67              for (ParseState ps: set)
68              {
69                  tbl.addRow(ps.prod.nonTerminal,
70                      " -> " + ps.prod.expansionToString(ps.prodPos),
71                      ps.inputPos,
72                      ps.parse.toStringFull());
73              }
74              tbl.print();
75          }
76
77          public Iterator<ParseState> iterator() { return set.iterator(); }
78
79          public boolean add(ParseState elm)
80          {
81              if (set.contains(elm)) return false;
82              set.add(elm);
83              return true;
84          }
85
86          public ParseState get(int i) { return set.get(i); }
87
88          public int size() { return set.size(); }
89      }
90
91      static private Set<ParseNode.ProductionNode> parseInternal(
92          NonTerminal root, Object[] input_,
93          Production phiProduction, StateSet[] stateSets)
94      {
95          int inputLength = input_.length; // 'n' in Earley70
96
97          // Append a terminator to the input.
98          Object[] input = new Object[input_.length + 1];
99          System.arraycopy(input_, 0, input, 0, input_.length);
100         input[input_.length] = Terminal.terminator;
101
102         stateSets[0].add(new ParseState(phiProduction, 0, 0,
103                 new ParseNode.ProductionNode(phiProduction)));
104
105         for (int inputPos = 0; inputPos <= inputLength; inputPos++)
106         {
107             StateSet set = stateSets[inputPos];
108
109             int oldSetSize = set.size();
110             boolean runAgain = false;
111
112             // Complete the set S[inputPos].
113             for (int i = 0; i < set.size(); i++)
114             {
115                 ParseState state = set.get(i);
116
117                 // If no next symbol in prod: Do the COMPLETER operation.
118                 if (!state.hasNextSymbol())
119                 {
120                     // Look at the input pos where we entered this parsestate.
121                     // (Might be the current set.)
122                     StateSet beginSet = stateSets[state.inputPos];
123                     for (int j = 0; j < beginSet.size(); j++)
124                     {
125                         ParseState ps = beginSet.get(j);
126                         if (state.prod.nonTerminal != ps.nextSymbol()) continue;
127
```

```
128                          set.add(ps.moveOver(ps.parse.append(state.parse)));
129                      }
130                  }
131
132                  // If next symbol is a NonTerminal: Do the PREDICTOR operation.
133                  else if (state.nextSymbol() instanceof NonTerminal)
134                  {
135                      NonTerminal nt = (NonTerminal) state.nextSymbol();
136                      if (nt.isNullable)
137                      {
138                          // This non-terminal may expand to the empty string.
139                          // We thus may need to run the completer on previous
140                          // entries in the set.
141                          runAgain = true;
142                      }
143
144                      for (Production prod: nt.alternatives)
145                      {
146                          ParseState s = new ParseState(prod, 0, inputPos,
147                              new ParseNode.ProductionNode(prod));
148                          set.add(s);
149                      }
150                  }
151                  else assert(state.nextSymbol() instanceof Terminal);
152              }
153
154              if (runAgain && set.size() > oldSetSize)
155              {
156                  // Run completer/predictor again on all entries in the set.
157                  inputPos--;
158                  continue;
159              }
160
161              // Construct the set S[inputPos + 1].
162              // We can use foreach here, because the current set doesn't change.
163              for (ParseState state: set)
164              {
165                  // If the next symbol is a Terminal: Do the SCANNER operation.
166                  if (state.nextSymbol() instanceof Terminal)
167                  {
168                      Terminal t = (Terminal) state.nextSymbol();
169                      ParseNode.TerminalNode node = t.match(input[inputPos]);
170                      if (node != null)
171                      {
172                          ParseState ps = state.moveOver(
173                              state.parse.append(node));
174
175                          stateSets[inputPos + 1].add(ps);
176                      }
177                  }
178              }
179          }
180
181          Set<ParseNode.ProductionNode> result
182              = new HashSet<ParseNode.ProductionNode>();
183
184          StateSet finalSet = stateSets[inputLength + 1];
185          for (ParseState ps: finalSet)
186          {
187              assert(ps.prod == phiProduction && ps.prodPos == 2 &&
188                  ps.inputPos == 0);
189              result.add((ParseNode.ProductionNode) ps.parse.nodes[0]);
190          }
191
192          return result;
193      }
194
195      static public Set<ParseNode.ProductionNode>
196          parse(NonTerminal root, Object... input)
197      {
```

```
198        StateSet [] stateSets = new StateSet [input . length + 2];
199        for (int i = 0; i < stateSets . length ; i++)
200            stateSets [i] = new StateSet ();
201
202        NonTerminal phi = new NonTerminal ("#");
203        Production phiProduction = new Production (
204            phi , root , Terminal . terminator );
205
206        phi . initialize (new HashSet < Symbol >());
207
208        return parseInternal (
209            root , input , phiProduction , stateSets );
210    }
211
212    static public Set < ParseNode . ProductionNode >
213        debugParse ( NonTerminal root , Object ... input )
214    {
215        StateSet [] stateSets = new StateSet [input . length + 2];
216        for (int i = 0; i < stateSets . length ; i++)
217            stateSets [i] = new StateSet ();
218
219        NonTerminal phi = new NonTerminal ("#");
220        Production phiProduction = new Production (
221            phi , root , Terminal . terminator );
222
223        phi . initialize (new HashSet < Symbol >());
224
225        Set < ParseNode . ProductionNode > result = parseInternal (
226            root , input , phiProduction , stateSets );
227
228        for (int i = 0; i < stateSets . length ; i++)
229        {
230            System . out . println ("============= S [" + i + "] =============");
231            stateSets [i]. print ();
232        }
233
234        return result ;
235    }
236 }
```

### jp.grammar.NonTerminal

```
 1  package jp . grammar ;
 2
 3  import java . util .*;
 4
 5  /**
 6   *  A NonTerminal is a symbol in the grammar that may match a varying number of
 7   *  symbols in an input string, depending on its associated alternatives (a list
 8   *  of productions .)
 9   */
10  public class NonTerminal extends Symbol
11  {
12      private final String name ;
13
14      /* package private */ List < Production > alternatives
15          = new ArrayList < Production >();
16
17      public NonTerminal (String name) { this . name = name ; }
18
19      public String toString () { return name ; }
20
21      public void print ()
22      {
23          for (Production p: alternatives )
24          {
25              System . out . println (p);
26          }
27      }
```

```
28
29        /* package private */ void initialize(Set<Symbol> done)
30        {
31            if (done.contains(this)) return;
32            done.add(this);
33
34            assert(alternatives.size() > 0) : name + " has no alternatives";
35            isNullable = false;
36
37            for (Production p: alternatives)
38            {
39                boolean thisAlternativeIsNullable = true;
40
41                for (Symbol s: p.expansion)
42                {
43                    s.initialize(done);
44                    if (!s.isNullable) thisAlternativeIsNullable = false;
45                }
46                if (thisAlternativeIsNullable) isNullable = true;
47            }
48        }
49    }
```

## jp.grammar.ParseNode

```
1    package jp.grammar;
2
3    import java.util.*;
4
5    abstract public class ParseNode
6    {
7        public void print(int indent)
8        {
9            while (indent-- > 0) System.out.print(" ");
10            System.out.println(this);
11        }
12
13        public String toStringFull() { return toString(); }
14
15        static public class TerminalNode extends ParseNode
16        {
17            public final Object symbol;
18
19            public TerminalNode(Object t) { symbol = t; }
20
21            public String toString() { return symbol.toString(); }
22        }
23
24        static public class ProductionNode extends ParseNode
25        {
26            public final ParseNode[] nodes;
27
28            public final Production prod;
29
30            public ProductionNode(Production prod)
31            {
32                this.prod = prod;
33                this.nodes = new ParseNode[0];
34            }
35
36            private ProductionNode(Production p, ParseNode[] nodes, ParseNode extra)
37            {
38                this.prod = p;
39                this.nodes = new ParseNode[nodes.length + 1];
40                System.arraycopy(nodes, 0, this.nodes, 0, nodes.length);
41                this.nodes[nodes.length] = extra;
42            }
43
44            public boolean equals(Object o)
```

```
45              {
46                  ParseNode . ProductionNode pnl = ( ParseNode . ProductionNode ) o ;
47                  return pnl . prod == prod && Arrays . equals ( pnl . nodes , nodes );
48              }
49
50              public ParseNode . ProductionNode append ( ParseNode node )
51              {
52                  return new ProductionNode ( prod , nodes , node );
53              }
54
55              public void print ( int indent )
56              {
57                  super . print ( indent );
58
59                  indent += 2;
60                  for ( ParseNode n : nodes ) n . print ( indent );
61              }
62
63              public String toString () { return prod . toString (); }
64
65              public String toStringFull ()
66              {
67                  StringBuilder sb = new StringBuilder ();
68                  sb . append ( prod . nonTerminal . toString ());
69                  sb . append ( " ( " );
70                  for ( ParseNode node : nodes )
71                  {
72                      sb . append ( " " );
73                      sb . append ( node . toStringFull ());
74                      sb . append ( " " );
75                  }
76                  sb . append ( " ) " );
77                  return sb . toString ();
78              }
79          }
80  }
```

## jp.grammar.Production

```
1   package jp . grammar ;
2
3   import java . util . Set ;
4
5   /**
6    *  A Production is a single expansion of an associated nonterminal.
7    */
8   public class Production
9   {
10      public final NonTerminal nonTerminal ;
11      public final Symbol [] expansion ;
12
13      public Production ( NonTerminal nonTerminal , Symbol ... expansion )
14      {
15          assert ( nonTerminal != null );
16          this . nonTerminal = nonTerminal ;
17          this . expansion = expansion ;
18
19          nonTerminal . alternatives . add ( this );
20      }
21
22      public String toString ()
23      {
24          return nonTerminal + " -> " + Symbol . symbolsToString ( expansion );
25      }
26
27      public String expansionToString ( int dotPos )
28      {
29          String result = "";
30
```

```
31          if (dotPos == 0) result += ".";
32          else result += " ";
33
34          int i = 0;
35          for (Symbol sym: expansion)
36          {
37              result += sym;
38              if (++i == dotPos) result += ".";
39              else result += " ";
40          }
41
42          return result;
43      }
44  }
```

## jp.grammar.Symbol

```
1   package jp.grammar;
2
3   import java.util.Set;
4
5   abstract public class Symbol
6   {
7       public boolean isNullable = false;
8
9       /** Recursively initializes symbols (determine isNullable) before parse. */
10      /* package private */ void initialize(Set<Symbol> done) { }
11
12      /** Converts an array of symbols to a string, for debug purposes. */
13      public static String symbolsToString(Symbol[] symbols)
14      {
15          String result = "";
16          for (Symbol sym: symbols)
17          {
18              if (result.length() > 0) result += " ";
19              result += sym;
20          }
21          return result;
22      }
23  }
```

## jp.grammar.Terminal

```
1   package jp.grammar;
2
3   import java.util.Set;
4
5   /**
6    *  A Terminal is a symbol in the grammar that may match exactly one symbol in
7    *  an input string. Typically, a Terminal will match a specific input symbol,
8    *  but it may also match a group of input symbols (e.g. any letter.)
9    */
10  abstract public class Terminal extends Symbol
11  {
12      /** A pre-defined NamedTerminal with name "|-" */
13      static final public Terminal terminator = new NamedTerminal("|-");
14
15      /** A NamedTerminal has a name and matches itself */
16      static public class NamedTerminal extends Terminal
17      {
18          private final String name;
19
20          public NamedTerminal(String name) { this.name = name; }
21
22          public String toString() { return name; }
23      }
24
```

```
25        /**
26         *  The match method of a Terminal determines if the terminal matches
27         *  the argument. If so, it returns a ParseNode.TerminalNode, otherwise
28         *  it returns null.
29         *  By default, we match ourselves.
30         */
31        public ParseNode.TerminalNode match(Object t)
32        {
33            if (t == this) return new ParseNode.TerminalNode(t);
34            return null;
35        }
36  }
```

## B.3 jp.lexer package

The jp.lexer package contains classes related to lexical analysis, including a tok-
enizer and classes for dealing with lexicons, lexemes, word classes and grammatical
features.

### jp.lexer.ClassResolver

```
1   package jp.lexer;
2
3   import jp.grammar.*;
4
5   import java.util.*;
6
7   /**
8    *  Interface used by ClassTerminal to test if a word belongs to a certain word
9    *  class.
10   *  The StaticCR implementation naively looks the word up in the lexicon.
11   *  LenientCR simply matches if the word is all-lowercase, and if so, adds the
12   *  word to the lexicon.
13   */
14  public interface ClassResolver
15  {
16      public Object resolve(String word);
17
18      public static class StaticCR implements ClassResolver
19      {
20          public final WordClass wordClass;
21          public final Lexicon lexicon;
22
23          public StaticCR(Lexicon lexicon, WordClass wordClass)
24          {
25              this.lexicon = lexicon;
26              this.wordClass = wordClass;
27          }
28
29          public Object resolve(String word)
30          {
31              Set<Lexeme> lexemes = lexicon.getLexemes(word, wordClass);
32              if (lexemes.size() > 0) return lexemes;
33              return null;
34          }
35
36          public String toString()
37          {
38              return "<" + wordClass.toString() + ">";
39          }
40      }
41
42      public static class LenientCR extends StaticCR
43      {
44          public LenientCR(Lexicon lexicon, WordClass wordClass)
45          {
46              super(lexicon, wordClass);
47          }
48
49          public Object resolve(String word)
50          {
51              for (int i = 0; i < word.length(); i++)
52                  if (Character.isUpperCase(word.charAt(i))) return null;
53
54              return lexicon.addLexeme(word, wordClass);
55          }
56      }
57  }
```

## jp.lexer.ClassTerminal

```
1   package jp.lexer;
2
3   import jp.grammar.*;
4
5   import java.util.*;
6
7   /**
8    *  A CFG terminal matching a whole class of words (e.g. all verbs).
9    *  The actual test is performed by a ClassResolver.
10   */
11  public class ClassTerminal extends Terminal
12  {
13      public ClassResolver resolver;
14
15      public ClassTerminal() { }
16
17      public ParseNode.TerminalNode match(Object t)
18      {
19          if (t instanceof String)
20          {
21              Object o = resolver.resolve((String) t);
22              if (o != null) return new ParseNode.TerminalNode(o);
23          }
24          return null;
25      }
26
27      public String toString() { return resolver.toString(); }
28  }
```

## jp.lexer.ExtendedLexicon

```
1   package jp.lexer;
2
3   import jp.grammar.*;
4
5   import java.util.*;
6
7   /**
8    *  A Lexicon that contains all words of a baseLexicon, plus possible additions.
9    *  This enables a lexicon segmented into a core vocabulary, and an extended
10   *  vocabulary.
11   */
12  public class ExtendedLexicon extends Lexicon
13  {
14      public final Lexicon baseLexicon;
15
16      public ExtendedLexicon(Lexicon baseLexicon)
17      {
18          this.baseLexicon = baseLexicon;
19      }
20
21      public Set<Lexeme> getLexemes(String word)
22      {
23          Set<Lexeme> set = new HashSet<Lexeme>(baseLexicon.getLexemes(word));
24          set.addAll(super.getLexemes(word));
25          return set;
26      }
27  }
```

## jp.lexer.Feature

```
1   package jp.lexer;
2
3   import java.util.*;
4
```

```
 5  /**
 6   *  Grammatical features, along with inflection and stemming rules.
 7   *  Each class is assigned a unique bit-index, to allow packing a set of
 8   *  features into a bit field.
 9   */
10
11  public enum Feature
12  {
13      // Person:
14      firstPerson,
15      secondPerson,
16      thirdPerson,
17
18      // Number
19      singular,
20      plural (new Inflx("s/sh/x/z", "-es"),
21              new Inflx("y", "ies"),
22              new Inflx("", "s")),
23
24      // Tense:
25      nonpast,
26      past,
27
28      // big, bigger, biggest
29      positive,
30      comparative(new Inflx("&d", "-der"), new Inflx("&g", "-ger"),
31          new Inflx("&t", "-ter"), new Inflx("e", "-r"),
32          new Inflx("y", "-ier"), new Inflx("", "-er")),
33      superlative(new Inflx("&d", "-dest"), new Inflx("&g", "-gest"),
34          new Inflx("&t", "-test"), new Inflx("e", "-st"),
35          new Inflx("y", "-iest"), new Inflx("", "est")),
36
37      // Etc:
38      possessive(new Inflx("s/x/z", "-'"), new Inflx("", "'s"))
39      ;
40
41      private static class Inflx
42      {
43          public final String[] oldEndings;
44          public final String newEnding;
45          public final boolean append;
46
47          /**
48           *  oldEndings: a slash ('/') separated list of old endings to match.
49           *  An & matches any vowel (a, e, i, o, u).
50           */
51          public Inflx(String oldEndings, String newEnding)
52          {
53              this.oldEndings = oldEndings.split("/");
54
55              append = newEnding.startsWith("-");
56
57              this.newEnding = append ? newEnding.substring(1) : newEnding;
58          }
59
60          public static boolean endingMatches(String word, String ending)
61          {
62              int wordIdx = word.length() - ending.length();
63              if (wordIdx < 0) return false;
64              for (int i = 0; i < ending.length(); i++, wordIdx++)
65              {
66                  char w = word.charAt(wordIdx);
67                  char e = ending.charAt(i);
68                  if (e == '&')
69                  {
70                      if (!(w == 'a' || w == 'e' || w == 'i' || w == 'o'
71                          || w == 'u')) return false;
72                  }
73                  else if (e != w) return false;
74              }
```

```
75              return true;
76          }
77
78          public String apply(String word)
79          {
80              for (String ending: oldEndings)
81              {
82                  if (!endingMatches(word, ending)) continue;
83                  if (append) return word + newEnding;
84                  return word.substring(0, word.length() - ending.length())
85                      + newEnding;
86              }
87              return null;
88          }
89
90          public boolean reverse(String word, List<String> results)
91          {
92              if (!word.endsWith(newEnding)) return false;
93
94              if (append)
95              {
96                  String oldWord =
97                      word.substring(0, word.length() - newEnding.length());
98
99                  for (String oldEnding: oldEndings)
100                 {
101                     if (!endingMatches(oldWord, oldEnding)) continue;
102
103                     results.add(oldWord);
104                     return true;
105                 }
106                 return false;
107             }
108
109             boolean anyWordsFound = false;
110             for (String oldEnding: oldEndings)
111             {
112                 String oldWord =
113                     word.substring(0, word.length() - newEnding.length())
114                     + oldEnding;
115
116                 if (!endingMatches(oldWord, oldEnding)) continue;
117
118                 results.add(oldWord);
119                 anyWordsFound = true;
120             }
121             return anyWordsFound;
122         }
123     }
124
125     public final int bit;
126     private final Inflx[] regularInflections;
127
128     private Feature()
129     {
130         bit = 1 << ordinal();
131         regularInflections = new Inflx[0];
132     }
133
134     private Feature(Inflx... regularInflections)
135     {
136         bit = 1 << ordinal();
137         this.regularInflections = regularInflections;
138     }
139
140     public static String toString(int features)
141     {
142         StringBuilder sb = new StringBuilder();
143
144         for (Feature f: Feature.values())
```

```
145            {
146                if ((f.bit & features) == 0) continue;
147
148                if (sb.length() != 0) sb.append("/");
149                sb.append(f);
150            }
151
152            return sb.toString();
153        }
154
155        public String inflect(String word)
156        {
157            for (Inflx inflx: regularInflections)
158            {
159                String result = inflx.apply(word);
160                if (result != null) return result;
161            }
162            return word;
163        }
164
165        public List<String> stem(String word)
166        {
167            List<String> results = new ArrayList<String>();
168
169            for (Inflx inflx: regularInflections)
170            {
171                inflx.reverse(word, results);
172            }
173
174            // Remove bogus results by checking with inflect().
175            Iterator<String> it = results.iterator();
176            while (it.hasNext())
177            {
178                if (!inflect(it.next()).equals(word)) it.remove();
179            }
180            return results;
181        }
182    }
```

## jp.lexer.Lexeme

```
 1  package jp.lexer;
 2
 3  import java.util.*;
 4
 5  public class Lexeme
 6  {
 7      /** A helper class coupling a lexeme with a set of grammatical features */
 8      public static class Inflected
 9      {
10          public Lexeme lexeme;
11          public int features;
12
13          public Inflected(Lexeme lexeme, int features)
14          {
15              this.lexeme = lexeme;
16              this.features = features;
17          }
18
19          public String toString()
20          {
21              return lexeme.toString() + "(" + Feature.toString(features) + ")";
22          }
23      }
24
25      public final WordClass wordClass;
26
27      public final String lemma; // citation form
28
```

```
29        public final Map<Integer, String> inflections
30            = new HashMap<Integer, String>();
31
32        public Lexeme(String lemma, WordClass wordClass)
33        {
34            this.wordClass = wordClass;
35            this.lemma     = lemma;
36        }
37
38        public String toString()
39        {
40            return lemma + "[" + wordClass.abbreviation + "]";
41        }
42
43        public String inflect(int features)
44        {
45            // Check for hardcoded inflection for this specific feature set.
46            String str = inflections.get(features);
47            if (str != null) return str;
48
49            str = lemma;
50
51            // Check for closest hardcoded inflection.
52            // Closeness is measured as number of bits in common.
53            // The match must not include features not specified in our argument.
54            int closestCount = 0;
55            int closestFeatures = 0;
56            for (Map.Entry<Integer, String> entry: inflections.entrySet())
57            {
58                int refFeatures = entry.getKey();
59                if ((refFeatures & ~features) != 0) continue;
60                int count = Integer.bitCount(refFeatures & features);
61
62                if (count <= closestCount) continue;
63
64                // A better match!
65                closestCount = count;
66                str = entry.getValue();
67                closestFeatures = refFeatures;
68            }
69            // Remove features already encoded
70            features &= ~closestFeatures;
71
72            for (Feature f: Feature.values())
73            {
74                if ((f.bit & features) == 0) continue;
75
76                str = f.inflect(str);
77            }
78            return str;
79        }
80
81        public Lexeme conjugate(int features, String value)
82        {
83            if (!inflections.containsKey(features) &&
84                inflect(features).equals(value)) return this;
85
86            inflections.put(features, value);
87            return this;
88        }
89
90        public Lexeme conjugate(Feature person, Feature number, Feature tense,
91            String value)
92        {
93            return conjugate(person.bit | number.bit | tense.bit, value);
94        }
95
96        public Lexeme conjugate(Feature person, Feature tense,
97            String singular, String plural)
98        {
```

```
 99              conjugate ( person , Feature . singular , tense , singular );
100              return conjugate ( person , Feature . plural , tense , plural );
101          }
102      }
```

## jp.lexer.LexemeMatch

```
 1    package jp . lexer ;
 2
 3    import jp . grammar .*;
 4
 5    import java . util .*;
 6
 7    /**
 8     *  A tuple of orthographic word , lexeme , and features .
 9     */
10    public class LexemeMatch
11    {
12        public final String word ;
13        public final Lexeme lexeme ;
14        public final int features ;
15
16        public LexemeMatch ( String word , Lexeme lexeme , int features )
17        {
18            this . word = word ;
19            this . lexeme = lexeme ;
20            this . features = features ;
21        }
22
23        public String toString ()
24        {
25            return lexeme . toString () + "[" + Feature . toString ( features ) + "]";
26        }
27    }
```

## jp.lexer.LexerException

```
1    package jp . lexer ;
2
3    public class LexerException extends Exception
4    {
5        public LexerException ( String msg )
6        {
7            super ( msg );
8        }
9    }
```

## jp.lexer.Lexer

```
 1    package jp . lexer ;
 2
 3    import jp . grammar .*;
 4
 5    import java . util .*;
 6
 7    /**
 8     *  A basic tokenizer , which splits a string into orthographic words ,
 9     *  punctuation and quoted strings .
10     */
11
12    public class Lexer
13    {
14        final private static int mNone = 0;
15        final private static int mWord = 1;
16        final private static int mString = 2;
```

```
17
18      static public class WordToken extends Terminal
19      {
20          public final String word;
21
22          public WordToken(String word) { this.word = word; }
23
24          public String toString() { return word; }
25      }
26
27      static public class StringToken extends Terminal
28      {
29          public final String string;
30
31          public StringToken(String string) { this.string = string; }
32
33          public String toString() { return '"' + string + '"'; }
34      }
35
36      static public class PunctToken extends Terminal
37      {
38          public final char punct;
39
40          public PunctToken(char punct) { this.punct = punct; }
41
42          public String toString() { return Character.toString(punct); }
43      }
44
45      protected Object wordToken(String token) { return new WordToken(token); }
46      protected Object punctToken(char token)  { return new PunctToken(token); }
47
48      public List<Object> tokenize(String string) throws LexerException
49      {
50          List<Object> result = new ArrayList<Object>();
51
52          int tokenIndex = 0;
53
54          int mode = mNone;
55
56          for (int i = 0; i < string.length(); i++)
57          {
58              char c = string.charAt(i);
59
60              switch (mode)
61              {
62              case mWord:
63                  if (Character.isLetter(c)) continue;
64
65                  result.add(wordToken(string.substring(tokenIndex, i)));
66                  mode = mNone;
67
68                  // fall through //
69
70              case mNone:
71                  switch (c)
72                  {
73                  case ' ': continue;
74                  case '"':
75                      mode = mString;
76                      tokenIndex = i;
77                      break;
78                  case '.':
79                  case ',':
80                  case ':':
81                  case '(':
82                  case ')':
83                  case ';':
84                      result.add(punctToken(c));
85                      break;
86
```

47

```
87                    default:
88                        if (Character.isLetter(c))
89                        {
90                            tokenIndex = i;
91                            mode = mWord;
92                        }
93                        else throw new LexerException("Invalid character at #" + i);
94                    }
95                    break;
96
97               case mString:
98                    if (c == '"')
99                    {
100                       mode = mNone;
101                       result.add(new StringToken(
102                           string.substring(tokenIndex + 1, i)));
103                   }
104                   break;
105            }
106        }
107
108        switch (mode)
109        {
110        case mWord:
111            result.add(wordToken(string.substring(tokenIndex)));
112            break;
113        case mNone: /* Do nothing */ break;
114        case mString: throw new LexerException("Unterminated string");
115        }
116        return result;
117    }
118 }
```

## jp.lexer.Lexicon

```
1   package jp.lexer;
2
3   import java.util.*;
4
5   import jp.util.*;
6
7   public class Lexicon implements Iterable<Lexeme>
8   {
9       private MultiMap<String, Lexeme> lexemes = new MultiMap<String, Lexeme>();
10      private MultiMap<String, Lexeme> irregular = new MultiMap<String, Lexeme>();
11
12      public Set<Lexeme> getLexemes(String word)
13      {
14          Set<Lexeme> set = new HashSet<Lexeme>(lexemes.get(word));
15          return set;
16      }
17
18      public Set<Lexeme> getLexemes(String word, WordClass wc)
19      {
20          Set<Lexeme> set = getLexemes(word);
21
22          Iterator<Lexeme> it = set.iterator();
23          while (it.hasNext())
24          {
25              if (it.next().wordClass != wc) it.remove();
26          }
27          return set;
28      }
29
30      public Set<Lexeme> addLexeme(String lemma, WordClass wordClass)
31      {
32          Set<Lexeme> lexes = getLexemes(lemma, wordClass);
33          if (lexes.size() == 0)
34          {
```

```
35                Lexeme lexeme = new Lexeme(lemma, wordClass);
36                lexes.add(lexeme);
37                lexemes.put(lemma, lexeme);
38            }
39            return lexes;
40        }
41
42        public Lexeme addLexeme(Lexeme lexeme)
43        {
44            lexemes.put(lexeme.lemma, lexeme);
45            for (String s: lexeme.inflections.values())
46                irregular.put(s, lexeme);
47
48            return lexeme;
49        }
50
51        public WordTerminal wordTerm(String word)
52        {
53            addLexeme(word, WordClass.particle);
54            return new WordTerminal(word);
55        }
56
57        public WordTerminal wordTerm(String word, WordClass... classes)
58        {
59            for (WordClass wc: classes) addLexeme(word, wc);
60            return new WordTerminal(word);
61        }
62
63        public Iterator<Lexeme> iterator() { return lexemes.values().iterator(); }
64
65        private void internalStem(String goal, String word, int features,
66            List<Lexeme.Inflected> results)
67        {
68            for (Feature f: Feature.values())
69            {
70                for (String stem: f.stem(word))
71                {
72                    assert(stem.length() < word.length());
73
74                    // For every result, recurse with the new feature set.
75                    internalStem(goal, stem, features | f.bit, results);
76                }
77            }
78
79            for (Lexeme lexeme: lexemes.get(word))
80            {
81                if (lexeme.inflect(features).equals(goal))
82                    results.add(new Lexeme.Inflected(lexeme, features));
83            }
84        }
85
86        public List<Lexeme.Inflected> stemAndLookUp(String word)
87        {
88            List<Lexeme.Inflected> results = new ArrayList<Lexeme.Inflected>();
89
90            for (Lexeme lexeme: irregular.get(word))
91            {
92                for (Map.Entry<Integer, String> e: lexeme.inflections.entrySet())
93                    if (e.getValue().equals(word))
94                        results.add(new Lexeme.Inflected(lexeme, e.getKey()));
95            }
96
97            internalStem(word, word, 0, results);
98            return results;
99        }
100    }
```

## jp.lexer.WordClass

```
 1  package jp.lexer;
 2
 3  /**
 4   *  The syntactic classes of words.
 5   *  Each class is assigned a unique bit-index, to allow packing a set of
 6   *  word-classes into a bit field.
 7   */
 8
 9  public enum WordClass
10  {
11      particle ("."),
12      countNoun ("CN"),
13      massNoun ("MN"),
14      properNoun ("PN"),
15      adjective ("Adj"),
16      verb ("V");
17
18      public final int bit;
19      public final String abbreviation;
20
21      WordClass(String abbreviation)
22      {
23          bit = 1 << ordinal();
24          this.abbreviation = abbreviation;
25      }
26
27      public static String toString(int classes)
28      {
29          StringBuilder sb = new StringBuilder();
30
31          for (WordClass wc: WordClass.values())
32          {
33              if ((wc.bit & classes) == 0) continue;
34
35              if (sb.length() != 0) sb.append("/");
36              sb.append(wc);
37          }
38
39          return sb.toString();
40      }
41  }
```

## jp.lexer.WordTerminal

```
 1  package jp.lexer;
 2
 3  import jp.grammar.*;
 4
 5  /** A CFG terminal matching a specific word (String). */
 6  public class WordTerminal extends Terminal
 7  {
 8      public final String word;
 9
10      public WordTerminal(String word) { this.word = word; }
11
12      public String toString() { return word; }
13
14      public ParseNode.TerminalNode match(Object t)
15      {
16          if (t instanceof String && t.equals(word))
17              return new ParseNode.TerminalNode(word);
18          return null;
19      }
20  }
```

## B.4   jp package

The `jp` package contains the grammar compiler (`GC`) described in section 6.1, and various test-classes, some of which implement examples found in this paper.

### jp.GC

```
 1   /*
 2    *    GC: Grammar Compiler
 3    *
 4    *    Usage: java GC [files]
 5    *
 6    *    Compiles input files containing BNF productions, then prints Java code
 7    *    suitable for use with the jp.grammar package to standard output.
 8    */
 9
10   package jp;
11
12   import java.util.*;
13   import java.io.*;
14
15   public class GC
16   {
17       public static class CompilerException extends Exception
18       {
19           public CompilerException(File filename, int lineNo, String message)
20           {
21               super(filename + ":" + lineNo + ": " + message);
22           }
23       }
24
25       public static class Grammar
26       {
27           private static class Production
28           {
29               private List<String> tokens = new ArrayList<String>();
30
31               public void addToken(String token) { tokens.add(token); }
32
33               public void print()
34               {
35                   for (int i = 0; i < tokens.size(); i++)
36                   {
37                       System.out.print(", ");
38                       System.out.print(tokens.get(i));
39                   }
40               }
41           }
42
43           // Maps nonTerminals to productions
44           private Map<String, List<Production>> productions
45               = new HashMap<String, List<Production>>();
46
47           private String currentNonTerminal;
48           private Production currentExpansion;
49
50           public void addExpansionToken(String token)
51           {
52               currentExpansion.addToken(token);
53           }
54
55           public void beginProduction(String nonTerminal)
56           {
57               finishProduction();
58               currentNonTerminal = nonTerminal;
59               currentExpansion = new Production();
```

```
60
61              if (!productions.containsKey(nonTerminal))
62                  productions.put(nonTerminal, new ArrayList<Production>());
63          }
64
65          public void finishExpansion()
66          {
67              if (currentExpansion != null)
68                  productions.get(currentNonTerminal).add(currentExpansion);
69              currentExpansion = new Production();
70          }
71
72          public void finishProduction()
73          {
74              if (currentNonTerminal == null) return;
75
76              finishExpansion();
77              currentNonTerminal = null;
78          }
79
80          public void print()
81          {
82              for (Map.Entry<String, List<Production>> entry:
83                  productions.entrySet())
84              {
85              System.out.print("    NonTerminal " + entry.getKey()
86                  + " = new NonTerminal(\"" + entry.getKey() + "\");\n");
87              }
88
89              for (Map.Entry<String, List<Production>> entry:
90                  productions.entrySet())
91              {
92                  for (Production p: entry.getValue())
93                  {
94                      System.out.print(
95                          "    new Production(" + entry.getKey());
96                      p.print();
97                      System.out.print(");\n");
98                  }
99              }
100         }
101     }
102
103     public static boolean isValidIdentifier(String s)
104     {
105         if (s.length() == 0 || !Character.isJavaIdentifierStart(s.charAt(0)))
106             return false;
107
108         for (int i = 1; i < s.length(); i++)
109         {
110             if (!Character.isJavaIdentifierPart(s.charAt(i))) return false;
111         }
112         return true;
113     }
114
115     public static void compile(File filename)
116         throws IOException, CompilerException
117     {
118         BufferedReader reader = new BufferedReader(new FileReader(filename));
119
120         Grammar grammar = new Grammar();
121
122         int lineNo = 0;
123         String line;
124         while ((line = reader.readLine()) != null)
125         {
126             lineNo++;
127
128             {
129                 int i = line.indexOf('#');
```

```
130            if (i != -1) line = line.substring(0, i);
131        }
132
133        String[] tokens = line.split("\\s+");
134        int tokenIndex = 0;
135
136        // Second check due to Sun JRE bug: split should return an empty
137        // list for an empty input, but doesn't.
138        if (tokens.length == 0 ||
139            (tokens.length == 1 && tokens[0].length() == 0)) continue;
140
141        if (tokens[0].length() > 0)
142        {
143            String nonTerminal = tokens[0];
144            if (!isValidIdentifier(nonTerminal))
145                throw new CompilerException(filename, lineNo,
146                    "Expected start of production, but \""
147                    + nonTerminal + "\" is not a valid identifier");
148
149            grammar.beginProduction(nonTerminal);
150
151            if (tokens.length < 2 || !tokens[1].equals("="))
152            {
153                throw new CompilerException(filename, lineNo,
154                    '"' + nonTerminal + "\" not followed by =");
155            }
156            tokenIndex = 2;
157        }
158        else if (grammar.currentNonTerminal == null)
159        {
160            throw new CompilerException(filename, lineNo,
161                "First line must not start with a space");
162        }
163
164        for (int i = tokenIndex; i < tokens.length; i++)
165        {
166            if (tokens[i].length() == 0) continue;
167
168            if (tokens[i].equals("|")) grammar.finishExpansion();
169            else grammar.addExpansionToken(tokens[i]);
170        }
171    }
172
173    grammar.finishProduction();
174
175    grammar.print();
176    }
177
178    public static void main(String[] args)
179        throws IOException
180    {
181        try
182        {
183            for (int i = 0; i < args.length; i++)
184            {
185                compile(new File(args[i]));
186            }
187        }
188        catch (CompilerException e)
189        {
190            System.err.println(e.getMessage());
191            System.exit(1);
192        }
193    }
194 }
```

## jp.EpsTest

```
1  package jp;
```

```
2
3   import jp.grammar.*;
4
5   import java.util.*;
6
7   /** An implementation of the epsilon problem */
8   public class EpsTest
9   {
10      public static void main(String[] args) throws Exception
11      {
12          Terminal plus = new Terminal.NamedTerminal("+");
13          NonTerminal S = new NonTerminal("S");
14          NonTerminal P = new NonTerminal("P");
15          NonTerminal Q = new NonTerminal("Q");
16          NonTerminal E = new NonTerminal("E");
17
18          new Production(S, E);
19          new Production(S, P);
20          new Production(P, Q, plus);
21          new Production(Q, E);
22          new Production(E);
23
24          Set<ParseNode.ProductionNode> result = Earley.debugParse(S, plus);
25
26          System.out.println("\nParse trees:");
27          for (ParseNode.ProductionNode l: result) l.print(0);
28      }
29  }
```

## jp.ExprTest

```
1   package jp;
2
3   import jp.grammar.*;
4
5   import java.util.*;
6
7   /** An implementation of the Parsing example */
8   public class ExprTest
9   {
10      public static void main(String[] args) throws Exception
11      {
12          Terminal a = new Terminal.NamedTerminal("a");
13          Terminal plus = new Terminal.NamedTerminal("+");
14          NonTerminal E = new NonTerminal("E");
15
16          new Production(E, a);
17          new Production(E, E, plus, E);
18
19          Set<ParseNode.ProductionNode> result = Earley.debugParse(E, a, plus, a);
20
21          System.out.println("\nParse trees:");
22          for (ParseNode.ProductionNode l: result) l.print(0);
23      }
24  }
```

## jp.StemTest

```
1   package jp;
2
3   import jp.grammar.*;
4   import jp.lexer.*;
5
6   import static jp.lexer.Lexicon.*;
7   import static jp.lexer.WordClass.*;
8   import static jp.lexer.Feature.*;
9
```

```
10    import java.util.*;
11
12    /** Runs a number of inflection and stemmer tests */
13    public class StemTest
14    {
15        public static Lexicon lexicon = new Lexicon();
16
17        public static void inflectAndStem(Lexeme lex, int features)
18        {
19            String word = lex.inflect(features);
20            System.out.println(word + " = " + lexicon.stemAndLookUp(word));
21        }
22
23        public static void main(String[] args) throws Exception
24        {
25            Lexeme weatherCN = new Lexeme("weather", countNoun);
26            Lexeme cloudCN   = new Lexeme("cloud", countNoun);
27            Lexeme rainV     = new Lexeme("rain", verb);
28            Lexeme bigAdj    = new Lexeme("big", adjective);
29
30            Lexeme beV = new Lexeme("be", verb)
31                .conjugate(firstPerson,  nonpast, "am",  "are")
32                .conjugate(secondPerson, nonpast, "are", "are")
33                .conjugate(thirdPerson,  nonpast, "is",  "are")
34                .conjugate(firstPerson,  past,    "was",  "were")
35                .conjugate(secondPerson, past,    "were", "were")
36                .conjugate(thirdPerson,  past,    "was",  "were");
37
38            lexicon.addLexeme(beV);
39            lexicon.addLexeme(weatherCN);
40            lexicon.addLexeme(rainV);
41            lexicon.addLexeme(bigAdj);
42
43            inflectAndStem(weatherCN, singular.bit);
44            inflectAndStem(weatherCN, plural.bit);
45            inflectAndStem(beV, firstPerson.bit | nonpast.bit | singular.bit);
46            inflectAndStem(rainV, nonpast.bit | firstPerson.bit | singular.bit);
47            inflectAndStem(bigAdj, superlative.bit);
48
49            for (String s: new String[] { "big", "large", "fair", "hot", "long" })
50            {
51                System.out.println(
52                    positive.inflect(s) + ", " +
53                    comparative.inflect(s) + ", " +
54                    superlative.inflect(s));
55            }
56
57            System.out.println("class   -> " + plural.stem("class"));
58            System.out.println("ties    -> " + plural.stem("ties"));
59            System.out.println("hottest -> " + superlative.stem("hottest"));
60        }
61    }
```

## jp.Test

```
 1    package jp;
 2
 3    import jp.grammar.*;
 4    import jp.lexer.*;
 5
 6    import static jp.lexer.Lexicon.*;
 7    import static jp.lexer.WordClass.*;
 8    import static jp.lexer.Feature.*;
 9
10    import java.util.*;
11
12    public class Test
13    {
14        /** A Lexer that return the same punctuation token for a given character. */
```

```
15        private static class MyLexer extends Lexer
16        {
17            private Map<Character, PunctToken> punctTokens
18                = new HashMap<Character, PunctToken>();
19
20            protected String wordToken(String token)
21            {
22                return token;
23            }
24
25            protected PunctToken punctToken(char token)
26            {
27                PunctToken t = punctTokens.get(token);
28                if (t != null) return t;
29
30                t = new PunctToken(token);
31                punctTokens.put(token, t);
32                return t;
33            }
34        }
35
36        /** A ClassResolver for proper nouns. Accepts new PNs if capitalized. */
37        private static class ProperNounCR implements ClassResolver
38        {
39            public final Lexicon lexicon;
40
41            public ProperNounCR(Lexicon lexicon)
42            {
43                this.lexicon = lexicon;
44            }
45
46            public Object resolve(String word)
47            {
48                Set<Lexeme> lexemes = lexicon.getLexemes(word, properNoun);
49                if (lexemes.size() > 0) return lexemes;
50                return Character.isUpperCase(word.charAt(0)) ?
51                        lexicon.addLexeme(word, properNoun) : null;
52            }
53
54            public String toString()
55            {
56                return "<properNoun*>";
57            }
58        }
59
60        public static void main(String[] args) throws Exception
61        {
62            MyLexer lexer = new MyLexer();
63
64            Terminal period = lexer.punctToken('.');
65
66            Lexicon baseLex = new Lexicon();
67
68            baseLex.addLexeme(new Lexeme("be", verb)
69                .conjugate(firstPerson,  nonpast, "am",  "are")
70                .conjugate(secondPerson, nonpast, "are", "are")
71                .conjugate(thirdPerson,  nonpast, "is",  "are")
72                .conjugate(firstPerson,  past,    "was", "were")
73                .conjugate(secondPerson, past,    "were", "were")
74                .conjugate(thirdPerson,  past,    "was", "were"));
75
76            baseLex.addLexeme(new Lexeme("star", countNoun));
77
78            WordTerminal shines = baseLex.wordTerm("shines", verb);
79            WordTerminal is     = new WordTerminal("is");
80            WordTerminal are    = new WordTerminal("are");
81
82            WordTerminal the   = baseLex.wordTerm("the");
83            WordTerminal it    = baseLex.wordTerm("it");
84            WordTerminal in    = baseLex.wordTerm("in");
```

56

```
85          WordTerminal at    = baseLex.wordTerm("at");
86          WordTerminal on    = baseLex.wordTerm("on");
87          WordTerminal under = baseLex.wordTerm("under");
88          WordTerminal above = baseLex.wordTerm("above");
89          WordTerminal beside = baseLex.wordTerm("beside");
90          WordTerminal behind = baseLex.wordTerm("behind");
91
92          ClassTerminal _countNoun_ = new ClassTerminal();
93          ClassTerminal _massNoun_ = new ClassTerminal();
94          ClassTerminal _adjective_ = new ClassTerminal();
95          ClassTerminal _properNoun_ = new ClassTerminal();
96
97          ClassResolver[] baseResolvers = new ClassResolver[] {
98              new ClassResolver.StaticCR(baseLex, countNoun),
99              new ClassResolver.StaticCR(baseLex, massNoun),
100             new ClassResolver.StaticCR(baseLex, adjective),
101             new ClassResolver.StaticCR(baseLex, properNoun)
102         };
103
104         ClassTerminal[] classTerminals = new ClassTerminal[]
105             { _countNoun_, _massNoun_, _adjective_, _properNoun_ };
106
107  // GRAMMAR BEGIN //
108      NonTerminal transitiveVerb = new NonTerminal("transitiveVerb");
109      NonTerminal sentence = new NonTerminal("sentence");
110      NonTerminal staticSpatialAdposition = new NonTerminal("
               staticSpatialAdposition");
111      NonTerminal input = new NonTerminal("input");
112      NonTerminal simpleSentence = new NonTerminal("simpleSentence");
113      NonTerminal adjectives = new NonTerminal("adjectives");
114      NonTerminal nounPhrase = new NonTerminal("nounPhrase");
115      NonTerminal adpositionalPhrase = new NonTerminal("adpositionalPhrase");
116      NonTerminal intransitiveVerb = new NonTerminal("intransitiveVerb");
117      NonTerminal predicate = new NonTerminal("predicate");
118      NonTerminal determiner = new NonTerminal("determiner");
119      NonTerminal sentences = new NonTerminal("sentences");
120      new Production(transitiveVerb, is);
121      new Production(transitiveVerb, are);
122      new Production(sentence, simpleSentence);
123      new Production(staticSpatialAdposition, in);
124      new Production(staticSpatialAdposition, on);
125      new Production(staticSpatialAdposition, at);
126      new Production(staticSpatialAdposition, under);
127      new Production(staticSpatialAdposition, above);
128      new Production(staticSpatialAdposition, behind);
129      new Production(staticSpatialAdposition, beside);
130      new Production(input, sentence);
131      new Production(input, sentences);
132      new Production(simpleSentence, nounPhrase, predicate);
133      new Production(adjectives, _adjective_, adjectives);
134      new Production(adjectives);
135      new Production(nounPhrase, determiner, _countNoun_);
136      new Production(nounPhrase, _properNoun_);
137      new Production(adpositionalPhrase, staticSpatialAdposition, nounPhrase);
138      new Production(intransitiveVerb, shines);
139      new Production(intransitiveVerb, is);
140      new Production(intransitiveVerb, are);
141      new Production(predicate, intransitiveVerb);
142      new Production(predicate, intransitiveVerb, adpositionalPhrase);
143      new Production(predicate, transitiveVerb, nounPhrase);
144      new Production(predicate, transitiveVerb, _adjective_);
145      new Production(determiner, the, adjectives);
146      new Production(sentences, sentence, period, sentences);
147      new Production(sentences, sentence, period);
148  // GRAMMAR END //
149
150         // Parse every argument as a separate sentence.
151         for (String s: args)
152         {
153             List<Object> tokens = lexer.tokenize(s);
```

57

```
154                 System.out.println(tokens);
155
156                 for (int i = 0; i < classTerminals.length; i++)
157                     classTerminals[i].resolver = baseResolvers[i];
158
159                 System.out.println("Old lexemes:");
160                 for (Lexeme lex: baseLex)
161                 {
162                     System.out.print(lex + "; ");
163                 }
164                 System.out.println();
165
166                 System.out.println("Parse:");
167                 Set<ParseNode.ProductionNode> result = Earley.parse(input,
168                     tokens.toArray(new Object[0]));
169
170                 if (result.size() == 0)
171                 {
172                     System.out.println("No results; trying new words.");
173
174                     Lexicon newLex = new ExtendedLexicon(baseLex);
175
176                     ClassResolver[] lenientResolvers = new ClassResolver[] {
177                         new ClassResolver.LenientCR(newLex, countNoun),
178                         new ClassResolver.LenientCR(newLex, massNoun),
179                         new ClassResolver.LenientCR(newLex, adjective),
180                         new ProperNounCR(newLex)
181                     };
182
183                     for (int i = 0; i < classTerminals.length; i++)
184                         classTerminals[i].resolver = lenientResolvers[i];
185
186                     result = Earley.parse(input, tokens.toArray(new Object[0]));
187
188                     if (result.size() != 0)
189                     {
190                         System.out.println("New lexemes:");
191                         for (Lexeme lex: newLex)
192                         {
193                             System.out.print(lex + "; ");
194                         }
195                         System.out.println();
196                     }
197                 }
198
199                 for (ParseNode.ProductionNode l: result)
200                 {
201                     l.print(0);
202                 }
203             }
204         }
205     }
```

58

## B.5   Other files

### Makefile

```
1   Test:    Test.java
2            javac -cp .. Test.java util/*.java grammar/*.java lexer/*.java semantics
                 /*.java
3   #        java -cp .. jp.Test "the sky is blue"
4   #        java -ea -cp .. jp.Test "the stars are shiny"
5            java -ea -cp .. jp.Test "the cat is wet"
6
7   StemTest:     StemTest.java
8            javac -cp .. StemTest.java util/*.java grammar/*.java lexer/*.java
9            java -ea  -cp .. jp.StemTest
10
11  ExprTest:     ExprTest.java
12           javac -cp .. ExprTest.java util/*.java grammar/*.java
13           java -ea  -cp .. jp.ExprTest
14
15  EpsTest:      EpsTest.java
16           javac -cp .. EpsTest.java util/*.java grammar/*.java
17           java -ea  -cp .. jp.EpsTest
18
19  SemTest:      SemTest.java
20           javac -cp .. SemTest.java util/*.java lexer/*.java semantics/*.java
21           java -ea  -cp .. jp.SemTest
22
23  rungc:  GC.class
24           java -ea  -cp .. jp.GC TestGrammar.gr
25
26  gram:   TestGrammar.gr GC.class
27           sed '/GRAMMAR BEGIN/,$$d' Test.java > Test.java.new
28           echo '// GRAMMAR BEGIN //' >> Test.java.new
29           java -ea  -cp .. jp.GC TestGrammar.gr >> Test.java.new
30           sed -n '/GRAMMAR END/,$$p' Test.java >> Test.java.new
31           mv Test.java Test.java.bak
32           mv Test.java.new Test.java
33
34  GC.class: GC.java
35           javac -cp .. GC.java
36
37  clean:
38           rm -f *.class **/*.class *.bak
```

### TestGrammar.gr

The grammar for jp.Test, for use with the grammar compiler.

```
1   #################
2   # Words         #
3   #################
4
5   adjectives      = _adjective_ adjectives |
6
7   determiner      = the adjectives
8
9   transitiveVerb = is | are
10
11  intransitiveVerb = shines | is | are
12
13  staticSpatialAdposition = in | on | at | under | above | behind | beside
14
15  #################
16  # Phrases       #
17  #################
18
```

59

```
19   nounPhrase         = determiner _countNoun_
20                      | _massNoun_
21                      | _properNoun_
22
23   adpositionalPhrase = staticSpatialAdposition nounPhrase
24
25   predicate          = intransitiveVerb
26                      | intransitiveVerb adpositionalPhrase
27                      | transitiveVerb nounPhrase
28                      | transitiveVerb _adjective_
29
30   ################
31   # Sentence      #
32   ################
33
34   simpleSentence  = nounPhrase predicate
35
36   sentence        = simpleSentence
37
38   sentences       = sentence period sentences
39                   | sentence period
40
41   input           = sentence
42                   | sentences
```