



DEPARTMENT OF INFORMATICS

TECHNISCHE UNIVERSITÄT MÜNCHEN

Master's Thesis in Informatics

Formal Verification of an Earley Parser

Martin Rau



DEPARTMENT OF INFORMATICS

TECHNISCHE UNIVERSITÄT MÜNCHEN

Master's Thesis in Informatics

Formal Verification of an Earley Parser

Formale Verifikation eines Earley Parsers

Author:	Martin Rau
Supervisor:	Tobias Nipkow
Advisor:	Tobias Nipkow
Submission Date:	15.06.2023

I confirm that this master's thesis in informatics is my own work and I have documented all sources and material used.

Munich, 15.06.2023

Martin Rau

Acknowledgments

I owe an enormous debt of gratitude to my family which always supported me throughout my studies. Thank you. I also would like to thank Prof. Tobias Nipkow for introducing me to the world of formal verification through Isabelle and for supervising both my Bachelor's and my Master's thesis. It was a pleasure to learn from and to work with you.

Abstract

TODO: Abstract

Contents

Acknowledgments	iii
Abstract	iv
1 Questions:	1
2 Introduction	2
2.1 Motivation	2
2.2 Related Work	2
2.3 Structure	3
2.4 Contributions	3
3 Earley's Algorithm	5
4 Earley's Algorithm Formalization	8
4.1 Context-free grammars and Isabelle/HOL	8
4.2 Earley's Algorithm	11
4.3 Well-formedness	15
4.4 Soundness	16
4.5 Completeness	19
4.6 Finiteness	23
5 Earley Recognizer Implementation	24
5.1 Draft	24
5.2 Definitions	25
5.3 Wellformedness	28
5.4 Soundness	30
5.5 Completeness	30
5.6 Main Theorem	32
6 Earley Parser Implementation	34
6.1 Draft	34
6.2 Pointer lemmas	34
6.3 Trees and Forests	35

Contents

6.4	A Single Parse Tree	36
6.5	All Parse Trees	38
6.6	A Word on Completeness	41
7	Usage	42
8	Conclusion	44
8.1	Summary	44
8.2	Future Work	44
	List of Figures	45
	List of Tables	46

1 Questions:

- How long? Min/Max?
- How to define Earley inductively?

2 Introduction

2.1 Motivation

some introduction about parsing, formal development of correct algorithms: an example based on earley's recogniser, the benefits of formal methods, LocalLexing and the Bachelor thesis.

2.2 Related Work

Tomita [Tomita:1987] presents an generalized LR parsing algorithm for augmented context-free grammars that can handle arbitrary context-free grammars.

Izmaylova *et al* [Izmaylova:2016] develop a general parser combinator library based on memoized Continuation-Passing Style (CPS) recognizers that supports all context-free grammars and constructs a Shared Packed Parse Forest (SPPF) in worst case cubic time and space.

Obua *et al* [Obua:2017] introduce local lexing, a novel parsing concept which interleaves lexing and parsing whilst allowing lexing to be dependent on the parsing process. They base their development on Earley's algorithm and have verified the correctness with respect to its local lexing semantics in the theorem prover Isabelle/HOL. The background theory of this Master's thesis is based upon the local lexing entry [LocalLexing-AFP] in the Archive of Formal Proofs.

Lasser *et al* [Lasser:2019] verify an LL(1) parser generator using the Coq proof assistant.

Barthwal *et al* [Barthwal:2009] formalize background theory about context-free languages and grammars, and subsequently verify an SLR automaton and parser produced by a parser generator.

Blaudeau *et al* [Blaudeau:2020] formalize the metatheory on Parsing expression grammars (PEGs) and build a verified parser interpreter based on higher-order parsing combinators for expression grammars using the PVS specification language and verification system. Koprowski *et al* [Koprowski:2011] present TRX: a parser interpreter formally developed in Coq which also parses expression grammars.

Jourdan *et al* [Jourdan:2012] present a validator which checks if a context-free grammar and an LR(1) parser agree, producing correctness guarantees required by verified

compilers.

Lasser *et al* [Lasser:2021] present the verified parser CoStar based on the ALL(*) algorithm. They proof soundness and completeness for all non-left-recursive grammars using the Coq proof assistant.

2.3 Structure

2.4 Contributions

SNIPPET:

Context-free grammars have been used extensively for describing the syntax of programming languages and natural languages. Parsing algorithms for context-free grammars consequently play a large role in the implementation of compilers and interpreters for programming languages and of programs which understand or translate natural languages. Numerous parsing algorithms have been developed. Some are general, in the sense that they can handle all context-free grammars, while others can handle only subclasses of grammars. The latter, restricted algorithms tend to be much more efficient. The algorithm described here seems to be the most efficient of the general algorithms, and also it can handle a larger class of grammars in linear time than most of the restricted algorithms.

SNIPPET:

The Computer Science community has been able to automatically generate parsers for a very wide class of context free languages. However, many parsers are still written manually, either using tool support or even completely by hand. This is partly because in some application areas such as natural language processing and bioinformatics we don not have the luxury of designing the language so that it is amendable to know parsing techniques, but also it is clear that left to themselves computer language designers do not naturally write LR(1) grammars. A grammar not only defines the syntax of a language, it is also the starting point for the definition of the semantics, and the grammar which facilitates semantics definition is not usually the one which is LR(1). Given this difficulty in constructing natural LR(1) grammars that support desired semantics, the general parsing techniques, such as the CYK Younger [Younger:1967], Earley [Earley:1970] and GLR Tomita [Tomita:1985] algorithms, developed for natural language processing are also of interest to the wider computer science community. When using grammars as the starting point for semantics definition, we distinguish between recognizers which simply determine whether or not a given string is in the language defined by a given grammar, and parser which also return some form of derivation of the string, if one exists. In their basic form the CYK and Earley

algorithms are recognizers while GLR-style algorithms are designed with derivation tree construction, and hence parsing, in mind.

There is no known linear time parsing or recognition algorithm that can be used with all context free grammars. In their recognizer forms the CYK algorithm is worst case cubic on grammars in Chomsky normal form and Earley's algorithm is worst case cubic on general context free grammars and worst case n^2 on non-ambiguous grammars. General recognizers must, by definition, be applicable to ambiguous grammars. Tomita's GLR algorithm is of unbounded polynomial order in the worst case. Expanding general recognizers to parser raises several problems, not least because there can be exponentially many or even infinitely many derivations for a given input string. A cubic recognizer which was modified to simply return all derivations could become an unbounded parser. Of course, it can be argued that ambiguous grammars reflect ambiguous semantics and thus should not be used in practice. This would be far too extreme a position to take. For example, it is well known that the if-else statement in the ANSI-standard grammar for C is ambiguous, but a longest match resolution results in a linear time parser that attach the else to the most recent if, as specified by the ANSI-C semantics. The ambiguous ANSI-C grammar is certainly practical for parser implementation. However, in general ambiguity is not so easily handled, and it is well known that grammar ambiguity is in fact undecidable Hopcroft *et al* [Hopcroft:2006], thus we cannot expect a parser generator simply to check for ambiguity in the grammar and report the problem back to the user. Another possibility is to avoid the issue by just returning one derivation. However, if only one derivation is returned then this creates problems for a user who wants all derivations and, even in the case where only one derivation is required, there is the issue of ensuring that it is the required derivation that is returned. A truly general parser will return all possible derivations in some form. Perhaps the most well known representation is the shared packed parse forest SPPF described and used by Tomita [Tomita:1985]. Tomita's description of the representation does not allow for the infinitely many derivations which arise from grammars which contain cycles, the source adapt the SPPF representation to allow these. Johnson [Johnson:1991] has shown that Tomita-style SPPFs are worst case unbounded polynomial size. Thus using such structures will also turn any cubic recognition technique into a worst case unbounded polynomial parsing technique. Leaving aside the potential increase in complexity when turning a recognizer into a parser, it is clear that this process is often difficult to carry out correctly. Earley gave an algorithm for constructing derivations of a string accepted by his recognizer, but this was subsequently shown by Tomita [Tomita:1985] to return spurious derivations in certain cases. Tomita's original version of his algorithm failed to terminate on grammars with hidden left recursion and, as remarked above, had no mechanism for constructing complete SPPFs for grammars with cycles.

3 Earley's Algorithm

We present a slightly simplified version of Earley's original recognizer algorithm [Earley:1970], omitting Earley's proposed look-ahead since it's primary purpose is to increase the efficiency of the resulting recognizer. Throughout this thesis we are working with a running example. The considered grammar is a tiny excerpt of a toy arithmetic expression grammar: $\mathcal{G} ::= S \rightarrow x \mid S \rightarrow S + S$ and the input is $\omega = x + x + x$.

Intuitively, Earley's recognizer works in principle like a top-down parser carrying along all possible parses simultaneously in an efficient manner. In detail, the algorithm works as follows: it scans the input $\omega = a_0, \dots, a_n$, constructing $n + 1$ Earley bins B_i that are sets of Earley items. An initial bin B_0 and one bin B_{i+1} for each symbol a_i of the input. In general, an Earley item $A \rightarrow \alpha \bullet \beta, i, j$ consists of four parts: a production rule of the grammar that we are currently scanning, a bullet signalling how much of the production's right-hand side we have recognized so far, an origin i describing the position of ω where we started scanning, and an end j indicating the position of ω we are currently considering next for the remaining right-hand side of the production rule. Note that there will be only one set of earley items or only one bin B and we say an item is conceptually part of bin B_j if it's end is the index j . Table 3.1 lists the items for our example grammar. Bin B_4 contains for example the item $S \rightarrow S + \bullet S, 2, 4$. Or, we are scanning the rule $S \rightarrow S + S$, have recognized the substring from 2 to 4 (the first index being inclusive the second one exclusive) of ω by $\alpha = S +$, and are trying to scan $\beta = S$ from position 4 in ω .

The algorithm initializes B by applying the *Init* operation. It then proceeds to execute the *Scan*, *Predict* and *Complete* operations listed in figure 3.1 until there are no more new items being generated and added to B . Next we describe these four operations in detail:

1. The *Init* operation adds items $S \rightarrow \bullet \alpha, 0, 0$ for each production rule containing the start symbol S on its left-hand side.

For our example *Init* adds the items $S \rightarrow \bullet x, 0, 0$ and $S \rightarrow \bullet S + S, 0, 0$.

2. The *Scan* operation applies if there is a terminal to the right side of the bullet, or items of the form $A \rightarrow \alpha \bullet a\beta, i, j$, and the j -th symbol of ω matches the terminal symbol following the bullet. We add one new item $A \rightarrow \alpha a \bullet \beta, i, j + 1$ to B

moving the bullet over the scanned terminal symbol.

Considering our example, bin B_3 contains the item $S \rightarrow S \bullet + S, 2, 3$, the third element of ω is the terminal $+$, so we add the item $S \rightarrow S + \bullet S, 2, 4$ to the conceptual bin B_4 .

3. The *Predict* operation is applicable to an item when there is a non-terminal to the right of the bullet or items of the form $A \rightarrow \alpha \bullet B\beta, i, j$. It adds one new item $B \rightarrow \bullet \gamma, j, j$ to the bin for each alternate $B \rightarrow \gamma$ of that non-terminal.

E.g. for the item $S \rightarrow S + \bullet S, 0, 2$ in B_2 we add the two items $S \rightarrow \bullet x, 2, 2$ and $S \rightarrow \bullet S + S, 2, 2$ corresponding to the two alternates of S . The bullet is set to the beginning of the right-hand side of the production rule, the origin and end are set to $j = 2$ to indicate that we are starting to scan in the current bin and have not scanned anything so far.

4. The *Complete* operation applies if we process an item with the bullet at the end of the right side of its production rule. For an item $B \rightarrow \gamma \bullet, j, k$ we have successfully scanned the substring $\omega[j..k)$ and are now going back to the origin bin B_j where we predicted this non-terminal. There we look for any item of the form $A \rightarrow \alpha \bullet B\beta, i, j$ containing a bullet in front of the non-terminal we completed, or the reason we predicted it on the first place. Since we scanned the predicted non-terminal successfully, we are allowed to move over the bullet, resulting in one new item $A \rightarrow \alpha B \bullet \beta, i, k$. Note in particular the origin and end indices.

Looking back at our example, we can add the item $S \rightarrow S + S \bullet, 0, 5$ for two different reasons corresponding conceptually to the two different ways we can derive ω . When processing $S \rightarrow x \bullet, 4, 5$ we find $S \rightarrow S + \bullet S, 0, 4$ in the origin bin B_4 which corresponds to recognizing $(x + x) + x$. We 'add' the same item again while applying the *Complete* operation to $S \rightarrow S + S \bullet, 2, 5$ and $S \rightarrow S + \bullet S, 0, 2$ which corresponds to recognizing the input as $x + (x + x)$.

If the algorithm encounters an item of the form $S \rightarrow \alpha \bullet, 0, |\omega| + 1$, it returns *true*, otherwise it returns *false*. For the tiny arithmetic expression grammar we generate the item $S \rightarrow S + S \bullet, 0, 5$ and return the correct answer *true*, since there exist derivations for $\omega = x + x + x$, e.g. $S \Rightarrow S + S \Rightarrow x + S \Rightarrow x + S + S \xRightarrow{*} x + x + x$ or $S \Rightarrow S + S \Rightarrow S + x \Rightarrow S + S + x \xRightarrow{*} x + x + x$.

To proof the correctness of Earley's recognizer algorithm we need to show the following theorem:

$$S \rightarrow \alpha \bullet, 0, |\omega| + 1 \in B \text{ iff } S \xRightarrow{*} \omega$$

It follows from the following three lemmas:

1. Termination: there only exist a finite number of Earley items
2. Soundness: for every generated item there exists an according derivation:
 $A \rightarrow \alpha \bullet \beta, i, j \in B$ implies $A \xRightarrow{*} \omega[i..j)\beta$
3. Completeness: for every derivation we generate an according item:
 $A \xRightarrow{*} \omega[i..j)\beta$ implies $A \rightarrow \alpha \bullet \beta, i, j \in B$

INIT	SCAN	PREDICT
$\frac{}{S \rightarrow \bullet \alpha, 0, 0}$	$\frac{A \rightarrow \alpha \bullet a \beta, i, j \quad \omega[j] = a}{A \rightarrow \alpha a \bullet \beta, i, j+1}$	$\frac{A \rightarrow \alpha \bullet B \beta, i, j \quad B \rightarrow \gamma \in \mathcal{G}}{B \rightarrow \bullet \gamma, j, j}$
	COMPLETE	
	$\frac{A \rightarrow \alpha \bullet B \beta, i, j \quad B \rightarrow \gamma \bullet, j, k}{A \rightarrow \alpha B \bullet \beta, i, k}$	

Figure 3.1: Earley inference rules

Table 3.1: Earley items for the grammar \mathcal{G} : $S \rightarrow x, S \rightarrow S + S$

B_0	B_1	B_2
$S \rightarrow \bullet x, 0, 0$ $S \rightarrow \bullet S + S, 0, 0$	$S \rightarrow x \bullet, 0, 1$ $S \rightarrow S \bullet + S, 0, 1$	$S \rightarrow S + \bullet S, 0, 2$ $S \rightarrow \bullet x, 2, 2$ $S \rightarrow \bullet S + S, 2, 2$
B_3	B_4	B_5
$S \rightarrow x \bullet, 2, 3$ $S \rightarrow S + S \bullet, 0, 3$ $S \rightarrow S \bullet + S, 2, 3$ $S \rightarrow S \bullet + S, 0, 3$	$S \rightarrow S + \bullet S, 2, 4$ $S \rightarrow S + \bullet S, 0, 4$ $S \rightarrow \bullet x, 4, 4$ $S \rightarrow \bullet S + S, 4, 4$	$S \rightarrow x \bullet, 4, 5$ $S \rightarrow S + S \bullet, 2, 5$ $S \rightarrow S + S \bullet, 0, 5$ $S \rightarrow S \bullet + S, 4, 5$ $S \rightarrow S \bullet + S, 2, 5$ $S \rightarrow S \bullet + S, 0, 5$

4 Earley’s Algorithm Formalization

In this chapter we shortly introduce the interactive theorem prover Isabelle/HOL [Nipkow:2002] used as the tool for verification in this thesis and recap some of the formalism of context-free grammars and its representation in Isabelle. Finally we formalize the simplified Earley recognizer algorithm presented in Chapter 3; discussing the implementation and the proofs for termination, soundness and completeness. Note that most of the definitions of Sections 4.1 and 4.2 are not our own work but only slightly adapted from [Obua:2017] [LocalLexing-AFP]. All of the proofs in this chapter are our own work.

4.1 Context-free grammars and Isabelle/HOL

Isabelle/HOL [Nipkow:2002] is an interactive theorem prover based on a fragment of higher-order logic. It supports the core concepts commonly known from functional programming languages. The notation $t :: \tau$ means that term t has type τ . Basic types include *bool*, *nat*; type variables are written $'a$, $'b$, etc. Pairs are written (a, b) ; triples and so forth are written (a, b, c) but are internally represented as nested pairs; the nesting is on the first component of a pair. Functions *fst* and *snd* return the first and second component of a pair; the operator (\times) represents pairs at the type level. Most type constructors are written postfix, e.g. $'a\ set$ and $'a\ list$; the function space arrow is \Rightarrow ; function *set* converts a list into a set. Type synonyms are introduced via the *type_synonym* command. Algebraic data types are defined with the keyword *datatype*. Non-recursive definitions are introduced with the *definition* keyword.

It is standard to define a language as a set of strings over a finite set of symbols. We deviate slightly by introducing a type variable $'a$ for the type of symbols. Thus a string corresponds to a list of symbols and a language is formalized as a set of lists of symbols. We represent a context-free grammar as the datatype *CFG*. An instance *cfg* consists of (1) a list of nonterminals ($\mathfrak{N}\ cfg$), (2) a list of terminals ($\mathfrak{T}\ cfg$), (3) a list of production rules ($\mathfrak{R}\ cfg$), and a start symbol ($\mathfrak{S}\ cfg$) where \mathfrak{N} , \mathfrak{T} , \mathfrak{R} and \mathfrak{S} are projections accessing the specific part of the instance *cfg* of the datatype *CFG*. Each rule consists of a left-hand side or *rule-head*, a single symbol, and a right-hand side or *rule-body*, a list of symbols. The productions with a particular nonterminal N on their left-hand sides are called the alternatives of N . We make the usual assumptions about

the well-formedness of a context-free grammar: the intersection of the set of terminals and the set of nonterminals is empty; the start symbol is a nonterminal; the rule head of a production is a nonterminal and its rule body consists of only symbols. Additionally, since we are working with a list of productions, we make the assumption that this list is distinct.

type-synonym 'a rule = 'a × 'a list

type-synonym 'a rules = 'a rule list

datatype 'a cfg =

CFG

(\mathfrak{N} : 'a list)

(\mathfrak{T} : 'a list)

(\mathfrak{R} : 'a rules)

(\mathfrak{S} : 'a)

definition rule-head :: 'a rule \Rightarrow 'a **where**

rule-head = fst

definition rule-body :: 'a rule \Rightarrow 'a list **where**

rule-body = snd

definition disjunct-symbols :: 'a cfg \Rightarrow bool **where**

disjunct-symbols cfg \equiv set (\mathfrak{N} cfg) \cap set (\mathfrak{T} cfg) = {}

definition wf-startsymbol :: 'a cfg \Rightarrow bool **where**

wf-startsymbol cfg \equiv \mathfrak{S} cfg \in set (\mathfrak{N} cfg)

definition wf-rules :: 'a cfg \Rightarrow bool **where**

wf-rules cfg \equiv $\forall (N, \alpha) \in$ set (\mathfrak{R} cfg). $N \in$ set (\mathfrak{N} cfg) \wedge ($\forall s \in$ set α . $s \in$ set (\mathfrak{N} cfg) \cup set (\mathfrak{T} cfg))

definition distinct-rules :: 'a cfg \Rightarrow bool **where**

distinct-rules cfg \equiv distinct (\mathfrak{R} cfg)

definition wf-cfg :: 'a cfg \Rightarrow bool **where**

wf-cfg cfg \equiv disjunct-symbols cfg \wedge wf-startsymbol cfg \wedge wf-rules cfg \wedge distinct-rules cfg

Furthermore, in Isabelle, lists are constructed from the empty list [] via the infix cons-operator (#); the operator (@) appends two lists. Sets follow the standard mathematical notation including the commonly found set builder notation or set comprehensions $\{x \mid P x\}$, and can also be defined inductively using the keyword *inductive_set*.

Next we formalize the concept of a derivation. We use lowercase letters a, b, c indicating terminal symbols; capital letters A, B, C denote nonterminals; lists of symbols are represented by greek letters: α, β, γ , occasionally also by lowercase letters u, v, w .

The empty list in the context of a language is ϵ . A sentential is a list consisting of only symbols. A sentence is a sentential if it only contains terminal symbols. We first define a predicate $\text{derives1 } \text{cfg } u \ v$ which expresses that we can derive v from u in a single step or the predicate holds if there exist α, β, N and γ such that $u = \alpha @ [N] @ \beta$, $v = \alpha @ \gamma @ \beta$ and (N, γ) is a production rule. We then can define the set of single-step derivations using derives1 , and subsequently the set of all derivations given a particular grammar is the reflexive-transitive closure of the set of single-step derivations. Finally, we say v can be derived from u given a grammar cfg , or $\text{derives } \text{cfg } u \ v$ if $(u, v) \in \text{derivations } \text{cfg}$.

type-synonym $'a \text{ sentential} = 'a \text{ list}$

definition $\text{is-terminal} :: 'a \text{ cfg} \Rightarrow 'a \Rightarrow \text{bool}$ **where**
 $\text{is-terminal } \text{cfg } s \equiv s \in \text{set } (\mathcal{T} \text{ cfg})$

definition $\text{is-nonterminal} :: 'a \text{ cfg} \Rightarrow 'a \Rightarrow \text{bool}$ **where**
 $\text{is-nonterminal } \text{cfg } s \equiv s \in \text{set } (\mathcal{N} \text{ cfg})$

definition $\text{is-symbol} :: 'a \text{ cfg} \Rightarrow 'a \Rightarrow \text{bool}$ **where**
 $\text{is-symbol } \text{cfg } s \equiv \text{is-terminal } \text{cfg } s \vee \text{is-nonterminal } \text{cfg } s$

definition $\text{wf-sentential} :: 'a \text{ cfg} \Rightarrow 'a \text{ sentential} \Rightarrow \text{bool}$ **where**
 $\text{wf-sentential } \text{cfg } s \equiv \forall x \in \text{set } s. \text{is-symbol } \text{cfg } x$

definition $\text{is-sentence} :: 'a \text{ cfg} \Rightarrow 'a \text{ sentential} \Rightarrow \text{bool}$ **where**
 $\text{is-sentence } \text{cfg } s \equiv \forall x \in \text{set } s. \text{is-terminal } \text{cfg } x$

definition $\text{derives1} :: 'a \text{ cfg} \Rightarrow 'a \text{ sentential} \Rightarrow 'a \text{ sentential} \Rightarrow \text{bool}$ **where**
 $\text{derives1 } \text{cfg } u \ v \equiv$
 $\quad \exists \alpha \ \beta \ N \ \gamma.$
 $\quad \quad u = \alpha @ [N] @ \beta$
 $\quad \quad \wedge v = \alpha @ \gamma @ \beta$
 $\quad \quad \wedge (N, \gamma) \in \text{set } (\mathcal{R} \text{ cfg})$

definition $\text{derivations1} :: 'a \text{ cfg} \Rightarrow ('a \text{ sentential} \times 'a \text{ sentential}) \text{ set}$ **where**
 $\text{derivations1 } \text{cfg} = \{ (u, v) \mid u \ v. \text{derives1 } \text{cfg } u \ v \}$

definition $\text{derivations} :: 'a \text{ cfg} \Rightarrow ('a \text{ sentential} \times 'a \text{ sentential}) \text{ set}$ **where**
 $\text{derivations } \text{cfg} = (\text{derivations1 } \text{cfg})^*$

definition $\text{derives} :: 'a \text{ cfg} \Rightarrow 'a \text{ sentential} \Rightarrow 'a \text{ sentential} \Rightarrow \text{bool}$ **where**
 $\text{derives } \text{cfg } u \ v \equiv (u, v) \in \text{derivations } \text{cfg}$

Potentially recursive but provably total functions that may make use of pattern matching are defined with the *fun* and *function* keywords; partial functions are defined

via *partial_function*. Take for example the function *slice* defined below. Term *slice i j xs* computes the slice of a list *xs* between indices *i* (inclusive) and *j* (exclusive), e.g. *slice 2 4 [a, b, c, d, e]* evaluates to *[c, d]*.

```
fun slice :: nat  $\Rightarrow$  nat  $\Rightarrow$  'a list  $\Rightarrow$  'a list where
  slice - - [] = []
| slice - 0 (x#xs) = []
| slice 0 (Suc b) (x#xs) = x # slice 0 b xs
| slice (Suc a) (Suc b) (x#xs) = slice a b xs
```

Lemmas, theorems and corollaries are presented using the keywords *lemma*, *theorem*, *corollary* respectively, followed by their names. They consist of zero or more assumptions marked by *assumes* keywords and one conclusion indicated by *shows*. E.g. we can proof a simple lemma about the interaction between the *slice* function and the append operator (@), stating the conditions under which we can split one slice into two.

```
lemma slice-append:
  assumes  $i \leq j \leq k$ 
  shows slice i j xs @ slice j k xs = slice i k xs
```

4.2 Earley's Algorithm

Next we formalize the algorithm presented in Chapter 3. First we define the datatype *item* representing Earley items. For example, the item $S \rightarrow S + \bullet S, 2, 4$ consists of four parts: a production rule (*item-rule*), a natural number (*item-bullet*) indicating the position of the bullet in the production rule, and two natural numbers (*item-origin* inclusive, *item-end* exclusive) representing the portion of the input string ω that has been scanned by the item. Additionally we introduce a few useful abbreviations: the functions *item-rule-head* and *item-rule-body* access the *rule-head* respectively *rule-body* of an item. Functions *item- α* and *item- β* split the production rule body at the bullet, e.g. $S \rightarrow \alpha \bullet \beta$. We call an item *complete* if the bullet is at the end of the production rule body. The next symbol (*next-symbol*) of an item is either *None* if it is complete, or *Some s* where *s* is the symbol in the production rule body following the bullet. An item is finished if the item rule head is the start symbol, the item is complete, and the whole input has been scanned or *item-origin item* = 0 and *item-end item* = $|\omega|$. Finally, we call a set of items *recognizing* if it contains at least one finished item, indicating that this set of items recognizes the input ω .

```
datatype 'a item =
  Item
    (item-rule: 'a rule)
    (item-bullet : nat)
    (item-origin : nat)
```

(*item-end* : nat)

type-synonym 'a items = 'a item set

definition *item-rule-head* :: 'a item \Rightarrow 'a **where**
item-rule-head x = rule-head (item-rule x)

definition *item-rule-body* :: 'a item \Rightarrow 'a sentential **where**
item-rule-body x = rule-body (item-rule x)

definition *item- α* :: 'a item \Rightarrow 'a sentential **where**
item- α x = take (item-bullet x) (item-rule-body x)

definition *item- β* :: 'a item \Rightarrow 'a sentential **where**
item- β x = drop (item-bullet x) (item-rule-body x)

definition *is-complete* :: 'a item \Rightarrow bool **where**
is-complete x \equiv item-bullet x \geq length (item-rule-body x)

definition *next-symbol* :: 'a item \Rightarrow 'a option **where**
next-symbol x \equiv if is-complete x then None else Some ((item-rule-body x) ! (item-bullet x))

definition *is-finished* :: 'a cfg \Rightarrow 'a sentential \Rightarrow 'a item \Rightarrow bool **where**
is-finished cfg ω x \equiv
item-rule-head x = \mathcal{S} cfg \wedge
item-origin x = 0 \wedge
item-end x = length ω \wedge
is-complete x

definition *recognizing* :: 'a items \Rightarrow 'a cfg \Rightarrow 'a sentential \Rightarrow bool **where**
recognizing I cfg $\omega \equiv \exists x \in I. \text{is-finished } \text{cfg } \omega \ x$

Normally we don't construct items directly via the *Item* constructor but use two auxiliary constructors: the function *init-item* is used by the *Init* and *Predict* operations. It sets the *item-bullet* to 0 or the beginning of the production rule body, initializes the *item-rule*, and indicates that this is an initial item by assigning *item-origin* and *item-end* to the current position in the input. The function *inc-item* returns a new item moving the bullet over the next symbol (assuming there is one) and setting the *item-end* to the current position in the input, leaving the item rule and origin untouched. It is utilized by the *Scan* and *Complete* operations.

definition *init-item* :: 'a rule \Rightarrow nat \Rightarrow 'a item **where**
init-item r k = Item r 0 k

definition $\text{inc-item} :: 'a \text{ item} \Rightarrow \text{nat} \Rightarrow 'a \text{ item}$ **where**
 $\text{inc-item } x \ k = \text{Item } (\text{item-rule } x) \ (\text{item-bullet } x + 1) \ (\text{item-origin } x) \ k$

There are different approaches of defining the set of Earley items in accordance with the rules of Figure 3.1. We can take an abstract approach and define the set inductively using Isabelle's inductive sets, or a more operational point of view. We take the latter approach and discuss the reasoning for this decision end the end of this section.

Note that, as mentioned previously, even though we are only constructing one set of Earley items, conceptually all items with the same item end form one Earley bin. Our operational approach is then the following: we generate Earley items bin by bin in ascending order, starting from the 0-th bin which contains all initial items computed by the *Init* operation. The three operations *Scan*, *Predict*, and *Complete* all take as arguments the index of the current bin and the current set of Earley items. For the k -th bin the *Scan* operation initializes the $k + 1$ -th bin, whereas the *Predict* and *Complete* operations only generate items belonging to the k -th bin. We then define a function *E-step* (short for Earley step) that returns the union of applying the three operations to a set of Earley items and this set itself. We complete the k -th bin and initialize the $k + 1$ -th bin by iterating *E-step* until the set of items stabilizes, captured by the *E-bin* definition. The function \mathcal{E} then generates the bins up to the n -th bin by applying the *E-bin* function first to the initial set of items *Init* and continuing in ascending order bin by bin. Finally, we compute the set of Earley items by applying \mathcal{E} to the length of the input.

definition $\text{bin} :: 'a \text{ items} \Rightarrow \text{nat} \Rightarrow 'a \text{ items}$ **where**
 $\text{bin } I \ k = \{ x . x \in I \wedge \text{item-end } x = k \}$

definition $\text{Init} :: 'a \text{ cfg} \Rightarrow 'a \text{ items}$ **where**
 $\text{Init } \text{cfg} = \{ \text{init-item } r \ 0 \mid r . r \in \text{set } (\mathfrak{R} \text{ cfg}) \wedge \text{fst } r = (\mathfrak{S} \text{ cfg}) \}$

definition $\text{Scan} :: \text{nat} \Rightarrow 'a \text{ sentential} \Rightarrow 'a \text{ items} \Rightarrow 'a \text{ items}$ **where**
 $\text{Scan } k \ \omega \ I =$
 $\{ \text{inc-item } x \ (k+1) \mid x . a .$
 $x \in \text{bin } I \ k \wedge$
 $\omega!k = a \wedge$
 $k < \text{length } \omega \wedge$
 $\text{next-symbol } x = \text{Some } a \}$

definition $\text{Predict} :: \text{nat} \Rightarrow 'a \text{ cfg} \Rightarrow 'a \text{ items} \Rightarrow 'a \text{ items}$ **where**
 $\text{Predict } k \ \text{cfg} \ I =$
 $\{ \text{init-item } r \ k \mid r . x .$
 $r \in \text{set } (\mathfrak{R} \text{ cfg}) \wedge$
 $x \in \text{bin } I \ k \wedge$
 $\text{next-symbol } x = \text{Some } (\text{rule-head } r) \}$

definition *Complete* :: $\text{nat} \Rightarrow 'a \text{ items} \Rightarrow 'a \text{ items}$ **where**

Complete $k \ I =$
 $\{ \text{inc-item } x \ k \mid x \ y.$
 $x \in \text{bin } I \ (\text{item-origin } y) \wedge$
 $y \in \text{bin } I \ k \wedge$
 $\text{is-complete } y \wedge$
 $\text{next-symbol } x = \text{Some } (\text{item-rule-head } y) \}$

definition *E-step* :: $\text{nat} \Rightarrow 'a \text{ cfg} \Rightarrow 'a \text{ sentential} \Rightarrow 'a \text{ items} \Rightarrow 'a \text{ items}$ **where**

E-step $k \ \text{cfg} \ \omega \ I = I \cup \text{Scan } k \ \omega \ I \cup \text{Complete } k \ I \cup \text{Predict } k \ \text{cfg} \ I$

fun *funpower* :: $('a \Rightarrow 'a) \Rightarrow \text{nat} \Rightarrow ('a \Rightarrow 'a)$ **where**

funpower $f \ 0 \ x = x$
 $\mid \text{funpower } f \ (\text{Suc } n) \ x = f \ (\text{funpower } f \ n \ x)$

definition *natUnion* :: $(\text{nat} \Rightarrow 'a \text{ set}) \Rightarrow 'a \text{ set}$ **where**

natUnion $f = \bigcup \{ f \ n \mid n. \text{True} \}$

definition *limit* :: $('a \text{ set} \Rightarrow 'a \text{ set}) \Rightarrow 'a \text{ set} \Rightarrow 'a \text{ set}$ **where**

limit $f \ x = \text{natUnion } (\lambda n. \text{funpower } f \ n \ x)$

definition *E-bin* :: $\text{nat} \Rightarrow 'a \text{ cfg} \Rightarrow 'a \text{ sentential} \Rightarrow 'a \text{ items} \Rightarrow 'a \text{ items}$ **where**

E-bin $k \ \text{cfg} \ \omega \ I = \text{limit } (\text{E-step } k \ \text{cfg} \ \omega) \ I$

fun \mathcal{E} :: $\text{nat} \Rightarrow 'a \text{ cfg} \Rightarrow 'a \text{ sentential} \Rightarrow 'a \text{ items}$ **where**

$\mathcal{E} \ 0 \ \text{cfg} \ \omega = \text{E-bin } 0 \ \text{cfg} \ \omega \ (\text{Init } \text{cfg})$
 $\mid \mathcal{E} \ (\text{Suc } n) \ \text{cfg} \ \omega = \text{E-bin } (\text{Suc } n) \ \text{cfg} \ \omega \ (\mathcal{E} \ n \ \text{cfg} \ \omega)$

definition *earley* :: $'a \text{ cfg} \Rightarrow 'a \text{ sentential} \Rightarrow 'a \text{ items}$ **where**

earley $\text{cfg} \ \omega = \mathcal{E} \ (\text{length } \omega) \ \text{cfg} \ \omega$

We followed the operational approach of defining the set of Earley items primarily for two reasons: first of all, we adapted most of the definitions of this chapter from the work on Local Lexing [Obua:2017] [LocalLexing-AFP], which takes the more operational approach and already defines useful lemmas, for example on function iteration. Secondly, the operational approach maps more easily to the list-based implementation of the next chapter which necessarily takes an ordered approach to generating Earley items. Nonetheless, in hindsight, defining the set of Earley items inductively seems to be not only the more elegant approach but also might simplify some of the proofs of this chapter.

4.3 Well-formedness

Due to the operational view of generating the set of Earley items, the proofs of, not only, well-formedness, but also soundness and completeness follow the same structure: we first proof a property about the basic building blocks, the *Init*, *Scan*, *Predict*, and *Complete* operations. Then, we proof that this property is maintained iterating the function *E-step*, and thus holds for the *E-bin* operation. Finally, we show that the function \mathcal{E} maintains this property for all conceptual bins and thus for the *earley* definition, or the set of Earley items.

Before we start to proof soundness and completeness of the generated set of Earley items, especially the completeness proof is more involved, we highlight the general proof structure once in detail, for a simpler property: well-formedness of the items, allowing us to concentrate only on the core aspects for the soundness and completeness proofs.

An Earley item is well-formed (*wf-item*) if the item rule is a rule of the grammar; the item bullet is bounded by the length of the item rule body; the item origin does not exceed the item end, and finally the item end is at most the length of the input.

definition *wf-item* :: 'a cfg \Rightarrow 'a sentential \Rightarrow 'a item \Rightarrow bool **where**

wf-item cfg ω x \equiv
 item-rule x \in set (\mathfrak{R} cfg) \wedge
 item-bullet x \leq length (item-rule-body x) \wedge
 item-origin x \leq item-end x \wedge
 item-end x \leq length ω

definition *wf-items* :: 'a cfg \Rightarrow 'a sentential \Rightarrow 'a items \Rightarrow bool **where**

wf-items cfg ω I $\equiv \forall x \in I. \text{wf-item cfg } \omega x$

lemma *wf-Init*:

shows *wf-items* cfg ω (Init cfg)

lemma *wf-Scan-Predict-Complete*:

assumes *wf-items* cfg ω I

shows *wf-items* cfg ω (Scan k ω I \cup Predict k cfg I \cup Complete k I)

lemma *wf-E-step*:

assumes *wf-items* cfg ω I

shows *wf-items* cfg ω (E-step k cfg ω I)

Lemmas *wf-Init*, *wf-Scan-Predict-Complete*, and *wf-E-step* follow trivially by definition of the respective operations.

lemma *wf-funpower*:

assumes *wf-items* cfg ω I

shows $wf\text{-}items\ cfg\ \omega\ (funpower\ (E\text{-}step\ k\ cfg\ \omega)\ n\ I)$

lemma $wf\text{-}E\text{-}bin$:

assumes $wf\text{-}items\ cfg\ \omega\ I$

shows $wf\text{-}items\ cfg\ \omega\ (E\text{-}bin\ k\ cfg\ \omega\ I)$

lemma $wf\text{-}E\text{-}bin0$:

shows $wf\text{-}items\ cfg\ \omega\ (E\text{-}bin\ 0\ cfg\ \omega\ (Init\ cfg))$

We proof the lemma $wf\text{-}funpower$ by induction on n using $wf\text{-}E\text{-}step$, and lemmas $wf\text{-}E\text{-}bin$ and $wf\text{-}E\text{-}bin0$ follow immediately using additionally the fact that $x \in limit\ f\ X \equiv \exists n. x \in funpower\ f\ n\ X$ and lemma $wf\text{-}Init$.

lemma $wf\text{-}\mathcal{E}$:

shows $wf\text{-}items\ cfg\ \omega\ (\mathcal{E}\ n\ cfg\ \omega)$

lemma $wf\text{-}earley$:

shows $wf\text{-}items\ cfg\ \omega\ (earley\ cfg\ \omega)$

Finally, lemma $wf\text{-}\mathcal{E}$ is proved by induction on n using lemmas $wf\text{-}E\text{-}bin$ and $wf\text{-}E\text{-}bin0$; lemma $wf\text{-}earley$ follows by definition.

4.4 Soundness

Next, we proof the soundness of the generated items, or: $A \rightarrow \alpha \bullet \beta, i, j \in B$ implies $A \xRightarrow{*} \omega[i..j)\beta$ which is stated in terms of our formalization by the *sound-item* definition below. As mentioned previously, the general proof structure follows the proof for well-formedness. Thus, we only highlight one slightly more involved lemma stating the soundness of the *Complete* operation while stating the remaining lemmas without explicit proof. Additionally, proving lemma *sound-Complete* provides some insight into working with and proving properties about derivations.

definition $sound\text{-}item :: 'a\ cfg \Rightarrow 'a\ sentential \Rightarrow 'a\ item \Rightarrow bool$ **where**

$sound\text{-}item\ cfg\ \omega\ x = derives\ cfg\ [item\text{-}rule\text{-}head\ x]\ (slice\ (item\text{-}origin\ x)\ (item\text{-}end\ x)\ \omega\ @\ item\text{-}\beta\ x)$

definition $sound\text{-}items :: 'a\ cfg \Rightarrow 'a\ sentential \Rightarrow 'a\ items \Rightarrow bool$ **where**

$sound\text{-}items\ cfg\ \omega\ I \equiv \forall x \in I. sound\text{-}item\ cfg\ \omega\ x$

Obua [Obua:2017] [LocalLexing-AFP] defines derivations at two different abstraction levels. The first representation is as the reflexive-transitive closure of the set of one-step derivations as introduced earlier in this chapter. The second representation is again more operational. He defines a predicate $Derives1\ cfg\ u\ i\ r\ v$ that is conceptually analogous to the predicate $derives1\ cfg\ u\ v$ but also captures the rule r used for a single rewriting step and the position i in u where the rewriting occurs.

definition *Derives1* :: 'a cfg \Rightarrow 'a sentential \Rightarrow nat \Rightarrow 'a rule \Rightarrow 'a sentential \Rightarrow bool **where**

Derives1 cfg u i r v \equiv
 $\exists \alpha \beta N \gamma.$
 $u = \alpha @ [N] @ \beta$
 $\wedge v = \alpha @ \gamma @ \beta$
 $\wedge (N, \gamma) \in \text{set } (\mathfrak{R} \text{ cfg})$
 $\wedge r = (N, \gamma) \wedge i = \text{length } \alpha$

We then define the type of a *derivation* as a list of pairs representing precisely the positions and rules used to apply each rewrite step. The predicate *Derivation* is defined recursively as follows: $\alpha \Rightarrow \beta$ only holds for the empty derivation or list. If the derivation consists of at least one rewrite pair (i, r) , or *Derivation* cfg $\alpha ((i, r) \# D) \beta$, then there must exist a γ such that *Derives1* cfg $\alpha i r \gamma$ and *Derivation* cfg $\gamma D \beta$. Obua then proves that both notions of a derivation are equivalent (lemma *derives-equiv-Derivation*)

type-synonym 'a derivation = (nat \times 'a rule) list

fun *Derivation* :: 'a cfg \Rightarrow 'a sentential \Rightarrow 'a derivation \Rightarrow 'a sentential \Rightarrow bool **where**

Derivation - $\alpha [] \beta = (\alpha = \beta)$
 $| \text{Derivation} \text{ cfg } \alpha (d \# D) \beta = (\exists \gamma. \text{Derives1} \text{ cfg } \alpha (\text{fst } d) (\text{snd } d) \gamma \wedge \text{Derivation} \text{ cfg } \gamma D \beta)$

lemma *derives-equiv-Derivation*:

shows *derives* cfg $\alpha \beta \equiv \exists D. \text{Derivation} \text{ cfg } \alpha D \beta$

Next, we state a small but useful lemma about rewriting derivations using the more operational definition of derivations defined above without explicit proof.

lemma *Derivation-append-rewrite*:

assumes *Derivation* cfg $\alpha D (\beta @ \gamma @ \delta)$
assumes *Derivation* cfg $\gamma E \gamma'$
shows $\exists F. \text{Derivation} \text{ cfg } \alpha F (\beta @ \gamma' @ \delta)$

And finally, we proof soundness of the *Complete* operation:

lemma *sound-Complete*:

assumes wf: wf-items cfg ωI
assumes sound: sound-items cfg ωI
shows sound-items cfg $\omega (\text{Complete } k I)$

Proof. Let z denote an arbitrary but fixed item of *Complete* $k I$. By the definition of the *Complete* operation there exist items x and y such that: $x \in \text{bin } I (\text{item-origin } y)$, $y \in \text{bin } I k$, *is-complete* y , *next-symbol* $x = \text{Some } (\text{item-rule-head } y)$, and $z = \text{inc-item } x k$.

Since y is in bin k , it is complete and the set I is sound (assumption *sound*), there exists a derivation E such that

$$\text{Derivation} \text{ cfg } [\text{item-rule-head } y] E (\text{slice } (\text{item-origin } y) (\text{item-end } y) \omega)$$

by lemma *derives-equiv-Derivation*. Similarly, since x is in bin *item-origin* y and due to assumption *sound*, there exists a derivation D such that

$$\text{Derivation cfg [item-rule-head } x] D (\text{slice (item-origin } x) (\text{item-origin } y) \omega @ \text{item-}\beta x)$$

Note that $\text{item-}\beta x = \text{item-rule-head } y \# \text{tl (item-}\beta x)$ since the next symbol of x is equal to the item rule head of y . Thus, by lemma *Derivation-append-rewrite*, and the definition of D and E , there exists a derivation F such that

$$\begin{aligned} \text{Derivation cfg [item-rule-head } x] F (\text{slice (item-origin } x) (\text{item-origin } y) \omega) @ \\ \text{slice (item-origin } y) (\text{item-end } y) \omega @ \text{tl (item-}\beta x) \end{aligned}$$

Additionally, we know that x and y are well-formed items due to the facts that x is in bin *item-origin* y , y is in bin k , and the assumption *wf-items cfg* ωI . Thus, we can discharge the assumptions of lemma *slice-append* (*item-origin* $x \leq \text{item-origin } y$ and *item-origin* $y \leq \text{item-end } y$) and have

$$\text{Derivation cfg [item-rule-head } x] F (\text{slice (item-origin } x) (\text{item-end } y) \omega @ \text{tl (item-}\beta x))$$

Moreover, since $z = \text{inc-item } x k$ and the next symbol of x is the item rule head of y , it follows that $\text{tl (item-}\beta x) = \text{item-}\beta z$, and ultimately *sound-item cfg* ωz , again by the definition of z and lemma *derives-equiv-Derivation*. □

lemma *sound-Init*:

shows *sound-items cfg* ω (*Init cfg*)

lemma *sound-Scan*:

assumes *wf-items cfg* ωI

assumes *sound-items cfg* ωI

shows *sound-items cfg* ω (*Scan* $k \omega I$)

lemma *sound-Predict*:

assumes *sound-items cfg* ωI

shows *sound-items cfg* ω (*Predict* $k \text{ cfg } I$)

lemma *sound-E-step*:

assumes *wf-items cfg* ωI

assumes *sound-items cfg* ωI

shows *sound-items cfg* ω (*E-step* $k \text{ cfg } \omega I$)

lemma *sound-funpower*:

assumes *wf-items* *cfg* ω *I*
assumes *sound-items* *cfg* ω *I*
shows *sound-items* *cfg* ω (*funpower* (*E-step* *k* *cfg* ω) *n* *I*)

lemma *sound-E-bin*:
assumes *wf-items* *cfg* ω *I*
assumes *sound-items* *cfg* ω *I*
shows *sound-items* *cfg* ω (*E-bin* *k* *cfg* ω *I*)

lemma *sound-E-bin0*:
shows *sound-items* *cfg* ω (*E-bin* 0 *cfg* ω (*Init* *cfg*))

lemma *sound-E*:
shows *sound-items* *cfg* ω (\mathcal{E} *k* *cfg* ω)

lemma *sound-earley*:
shows *sound-items* *cfg* ω (*earley* *cfg* ω)

Finally, using *sound-earley* and the definitions of *sound-item*, *recognizing*, *is-finished* and *is-complete* the final theorem follows: if the generated set of Earley items is *recognizing*, or contains a *finished* item then there exists a derivation from the start symbol of the grammar to the input ω .

theorem *soundness*:
assumes *recognizing* (*earley* *cfg* ω) *cfg* ω
shows *derives* *cfg* [\mathcal{S} *cfg*] ω

4.5 Completeness

Next, we prove completeness and consequently obtain a concluded correctness proof using theorem *soundness*. The completeness proof is by far the most involved proof of this chapter. Thus, we present it in greater detail, and also slightly deviate from the proof structure of the well-formedness and soundness proofs presented previously. We directly start to prove three properties of the \mathcal{E} function that correspond conceptually to the three different operations that can occur while generating the bins.

We need three simple lemmas concerning the *E-bin* function, stated without explicit proof: (1) *E-bin* *k* *cfg* ω *I* only (potentially) changes bins *k* and *k* + 1 (lemma *E-bin-bin-idem*); (2) the *E-step* operation is subsumed by the *E-bin* operation, since it computes the limit of *E-step* (lemma *E-step-sub-E-bin*); and (3) the function *E-bin* is idempotent (lemma *E-bin-idem*).

lemma *E-bin-bin-idem*:
assumes $i \neq k$

assumes $i \neq k+1$
shows $\text{bin } (E\text{-bin } k \text{ cfg } \omega \text{ } I) \text{ } i = \text{bin } I \text{ } i$

lemma *E-step-sub-E-bin*:
shows $E\text{-step } k \text{ cfg } \omega \text{ } I \subseteq E\text{-bin } k \text{ cfg } \text{inp } I$

lemma *E-bin-idem*:
shows $E\text{-bin } k \text{ cfg } \omega \text{ } (E\text{-bin } k \text{ cfg } \omega \text{ } I) = E\text{-bin } k \text{ cfg } \omega \text{ } I$

Next, we proof lemma *Scan- \mathcal{E}* in detail: it describes under which assumptions the function \mathcal{E} generates a 'scanned' item:

lemma *Scan- \mathcal{E}* :
assumes $i+1 \leq k$
assumes $x \in \text{bin } (\mathcal{E} \text{ } k \text{ cfg } \omega) \text{ } i$
assumes $\text{next-symbol } x = \text{Some } a$
assumes $k \leq \text{length } \omega$
assumes $\omega[i] = a$
shows $\text{inc-item } x \text{ } (i+1) \in \mathcal{E} \text{ } k \text{ cfg } \omega$

Proof. The proof is by induction in k for arbitrary i , x , and a :

The base case $k = 0$ is trivial, since we have the assumption $i + 1 \leq 0$.

For the induction step we can assume $i + 1 \leq k + 1$, $k + 1 \leq |\omega|$, $x \in \text{bin } (\mathcal{E} \text{ } (k + 1) \text{ cfg } \omega) \text{ } i$, $\text{next-symbol } x = \text{Some } a$, and $\omega[i] = a$. Assumptions $x \in \text{bin } (\mathcal{E} \text{ } (k + 1) \text{ cfg } \omega) \text{ } i$ and $i + 1 \leq k + 1$ imply that $x \in \text{bin } (\mathcal{E} \text{ } k \text{ cfg } \text{inp}) \text{ } i$ by lemma *E-bin-bin-idem*.

We then consider two cases:

- $i + 1 \leq k$: We can apply the induction hypothesis using assumptions $k + 1 \leq |\omega|$, $\text{next-symbol } x = \text{Some } a$, $\omega[i] = a$ and additionally $x \in \text{bin } (\mathcal{E} \text{ } k \text{ cfg } \text{inp}) \text{ } i$ and have $\text{inc-item } x \text{ } (i + 1) \in \mathcal{E} \text{ } k \text{ cfg } \omega$. The statement to proof follows by lemma *E-step-sub-E-bin* and the definition of *E-step*.
- $k < i + 1$: We have $i = k$, since $i + 1 \leq k + 1$. Thus, we have $\text{inc-item } x \text{ } (i + 1) \in \text{Scan } k \text{ } \omega \text{ } (\mathcal{E} \text{ } k \text{ cfg } \omega)$ using assumptions $k + 1 \leq |\omega|$, $\text{next-symbol } x = \text{Some } a$, $\omega[i] = a$, and additionally $x \in \text{bin } (\mathcal{E} \text{ } k \text{ cfg } \text{inp}) \text{ } i$ by the definition of the *Scan* operation. This in turn implies $\text{inc-item } x \text{ } (i + 1) \in E\text{-step } k \text{ cfg } \omega \text{ } (\mathcal{E} \text{ } k \text{ cfg } \omega)$ by lemma *E-step-sub-E-bin* and the definition of *E-step*. Since the function *E-bin* is idempotent (lemma *E-bin-idem*), we have $\text{inc-item } x \text{ } (i + 1) \in \mathcal{E} \text{ } k \text{ cfg } \omega$ by definition of \mathcal{E} . And again, the final statement follows by lemma *E-step-sub-E-bin* and the definition of *E-step*.

□

lemma *Predict- \mathcal{E}* :

assumes $i \leq k$
assumes $x \in \text{bin } (\mathcal{E} \text{ k cfg } \omega) \ i$
assumes $\text{next-symbol } x = \text{Some } N$
assumes $(N, \alpha) \in \text{set } (\mathfrak{R} \text{ cfg})$
shows $\text{init-item } (N, \alpha) \ i \in \mathcal{E} \text{ k cfg } \omega$

lemma *Complete- \mathcal{E}* :

assumes $i \leq j$
assumes $j \leq k$
assumes $x \in \text{bin } (\mathcal{E} \text{ k cfg } \omega) \ i$
assumes $\text{next-symbol } x = \text{Some } N$
assumes $(N, \alpha) \in \text{set } (\mathfrak{R} \text{ cfg})$
assumes $y \in \text{bin } (\mathcal{E} \text{ k cfg } \omega) \ j$
assumes $\text{item-rule } y = (N, \alpha)$
assumes $i = \text{item-origin } y$
assumes $\text{is-complete } y$
shows $\text{inc-item } x \ j \in \mathcal{E} \text{ k cfg } \omega$

The proof of lemmas *Predict- \mathcal{E}* and *Complete- \mathcal{E}* are similar in structure to the proof of lemma *Scan- \mathcal{E}* with the exception of the base case that is in both cases non-trivial but can be proven with the help of lemmas *E-step-sub-E-bin* and *E-bin-idem*, the definition of *E-bin* and the definitions of *Predict* and *Complete*, respectively.

The core idea for the completeness proof is the following: if there exists an item of the form $N \rightarrow \alpha \bullet \beta, i, j$ in a set of items I and there exists a derivation $\beta \xRightarrow{D} \omega[j..k)$, then I also contains the complete item $N \rightarrow \alpha \beta \bullet, i, k$. Note that this statement (lemma *partially-completed-up-to* below) holds only if $j \leq k$, $k \leq |\omega|$, all items of I are well-formed and most importantly I must be *partially-completed* up to the length of the derivation D .

Intuitively, a set I is *partially-completed* if for each non-complete item x in I , the existence of a derivation D from the next symbol of x implies that the item that can be obtained by moving the bullet over the next symbol of x is also present in I . The full definition of *partially-completed* is slightly more involved since we need to keep track of the validity of the indices and the fact that we sometimes want to limit derivations to a certain depth as done above.

definition *partially-completed* :: $\text{nat} \Rightarrow 'a \text{ cfg} \Rightarrow 'a \text{ sentential} \Rightarrow 'a \text{ items} \Rightarrow ('a \text{ derivation} \Rightarrow \text{bool}) \Rightarrow \text{bool}$ **where**

$\text{partially-completed } k \text{ cfg } \omega \ I \ P \equiv$
 $\forall i \ j \ x \ a \ D.$
 $i \leq j \wedge j \leq k \wedge k \leq \text{length } \omega \wedge$
 $x \in \text{bin } I \ i \wedge \text{next-symbol } x = \text{Some } a \wedge$
 $\text{Derivation cfg } [a] \ D \ (\text{slice } i \ j \ \omega) \wedge P \ D \longrightarrow$
 $\text{inc-item } x \ j \in I$

To proof lemma *partially-completed-upto*, we need an auxiliary lemma about derivations: a derivation $\alpha\beta \xRightarrow{D} \gamma$, can be split into two derivations E and F whose length is bounded by the length of D , and there exist α' and β' such that $\alpha \xRightarrow{E} \alpha'$, $\beta \xRightarrow{F} \beta'$ and $\gamma = \alpha' @ \beta'$.

lemma *Derivation-append-split*:

assumes *Derivation cfg* $(\alpha @ \beta) D \gamma$

shows $\exists E F \alpha' \beta'. \text{Derivation cfg } \alpha E \alpha' \wedge \text{Derivation cfg } \beta F \beta' \wedge \gamma = \alpha' @ \beta' \wedge \text{length } E \leq \text{length } D \wedge \text{length } F \leq \text{length } D$

Proof.

□

lemma *partially-completed-upto*:

assumes *wf-items cfg* ωI

assumes $j \leq k$

assumes $k \leq \text{length inp}$

assumes $x = \text{Item } (N, \alpha) b i j$

assumes $x \in I$

assumes *Derivation cfg* $(\text{item-}\beta x) D (\text{slice } j k \omega)$

assumes *partially-completed k cfg* $\omega I (\lambda D'. \text{length } D' \leq \text{length } D)$

shows *Item* $(N, \alpha) (\text{length } \alpha) i k \in I$

lemma *partially-completed- \mathcal{E}* :

assumes *wf-cfg cfg*

shows *partially-completed k cfg inp* $(\mathcal{E} k \text{cfg } \omega) (\lambda-. \text{True})$

lemma *slice-singleton*:

assumes $b \leq \text{length } xs$

assumes $[x] = \text{slice } a b xs$

shows $b = a + 1$

lemma *slice-nth*:

assumes $a < \text{length } xs$

shows $\text{slice } a (a+1) xs = [xs!a]$

lemma *partially-completed-earley*:

assumes *wf-cfg cfg*

shows *partially-completed* $(\text{length } \omega) \text{cfg } \omega (\text{earley cfg } \omega) (\lambda-. \text{True})$

lemma *Derivation- $\mathfrak{S}1$* :

assumes *wf-cfg cfg*

assumes *is-sentence cfg* ω

assumes *Derivation cfg* $[\mathfrak{S} \text{cfg}] D \omega$

shows $\exists \alpha E. \text{Derivation cfg } \alpha E \omega \wedge (\mathfrak{S} \text{cfg}, \alpha) \in \text{set } (\mathfrak{R} \text{cfg})$

theorem *completeness*:

assumes *wf-cfg* *cfg*
assumes *is-sentence* *cfg* ω
assumes *derives* *cfg* $[\mathfrak{S} \text{ } \text{cfg}] \omega$
shows *recognizing* (*earley* *cfg* ω) *cfg* ω

corollary *correctness*:

assumes *wf-cfg* *cfg*
assumes *is-sentence* *cfg* ω
shows *recognizing* (*earley* *cfg* *inp*) *cfg* $\omega \longleftrightarrow \text{derives } \text{cfg } [\mathfrak{S} \text{ } \text{cfg}] \omega$

4.6 Finiteness

At last, we prove that the generated set of Earley items is finite. In Chapter 5 we are using this result to prove the termination of an executable version of the algorithm.

Since *earley* *cfg* ω only generates well-formed items (lemma *wf-earley*) it suffices to prove that there only exists a finite number of well-formed items. Define

$$T = \text{set } (\mathfrak{R} \text{ } \text{cfg}) \times \{0..m\} \times \{0..|\omega|\} \times \{0..|\omega|\}$$

where $m = \text{Max } \{|rule-body \text{ } r| \mid r \in \text{set } (\mathfrak{R} \text{ } \text{cfg})\}$. The set T is finite and $\{x \mid \text{wf-item } \text{cfg } \omega \text{ } x\}$ is a subset of mapping the *Item* constructor over T (strictly speaking we need to first unpack the quadruple).

lemma *finiteness-UNIV-wf-item*:

shows *finite* $\{x \mid x. \text{wf-item } \text{cfg } \omega \text{ } x\}$

theorem *finiteness*:

shows *finite* (*earley* *cfg* ω)

5 Earley Recognizer Implementation

5.1 Draft

- introduce auxiliary definitions: `filter_with_index`, `pointer`, `entry` in more detail most everything else in text
- overview over earley implementation with linked list and pointers and the mapping into a functional setting
- introduce `Init_it`, `Scan_it`, `Predict_it` and `Complete_it`, compare them with the set notation and discuss performance improvements (Grammar in more specific form) Why do they all return a list?!
- discuss `bin(s)_upd(s)` functions. Why `bin_upds` like this -> easier than fold for proofs!
- discuss `pi_it` and why it is a partial function -> only terminates for valid input and foreshadow how this is done in isabelle
- introduce remaining definitions (analog to sets)
- discuss wf proofs quickly and go into detail about isabelle specifics about termination and the custom induction scheme using finiteness
- outline the approach to proof correctness aka subsumption in both directions
- discuss list to set proofs
- discuss soundness proofs (maybe omit since obvious)

- discuss completeness proof focusing on the complete case shortly explaining scan and predict which don't change via iteration and order does not matter
- highlight main theorems

Table 5.1: Earley items with pointers for the grammar $\mathcal{G}: S \rightarrow x, S \rightarrow S + S$

	B_0	B_1	B_2
0	$S \rightarrow \bullet x, 0, 0;$	$S \rightarrow x \bullet, 0, 1; 0$	$S \rightarrow S + \bullet S, 0, 2; 1$
1	$S \rightarrow \bullet S + S, 0, 0;$	$S \rightarrow S \bullet + S, 0, 1; (0, 1, 0)$	$S \rightarrow \bullet x, 2, 2;$
2			$S \rightarrow \bullet S + S, 2, 2;$
	B_3	B_4	B_5
0	$S \rightarrow x \bullet, 2, 3; 1$	$S \rightarrow S + \bullet S, 2, 4; 2$	$S \rightarrow x \bullet, 4, 5; 2$
1	$S \rightarrow S + S \bullet, 0, 3; (2, 0, 0)$	$S \rightarrow S + \bullet S, 0, 4; 3$	$S \rightarrow S + S \bullet, 2, 5; (4, 0, 0)$
2	$S \rightarrow S \bullet + S, 2, 3; (2, 2, 0)$	$S \rightarrow \bullet x, 4, 4;$	$S \rightarrow S + S \bullet, 0, 5; (4, 1, 0), (2, 0, 1)$
3	$S \rightarrow S \bullet + S, 0, 3; (0, 1, 1)$	$S \rightarrow \bullet S + S, 4, 4;$	$S \rightarrow S \bullet + S, 4, 5; (4, 3, 0)$
4			$S \rightarrow S \bullet + S, 2, 5; (2, 2, 1)$
5			$S \rightarrow S \bullet + S, 0, 5; (0, 1, 2)$

5.2 Definitions

fun *filter-with-index'* :: $\text{nat} \Rightarrow ('a \Rightarrow \text{bool}) \Rightarrow 'a \text{ list} \Rightarrow ('a \times \text{nat}) \text{ list}$ **where**
filter-with-index' - - [] = []
filter-with-index' i P (x#xs) = (
 if P x then (x,i) # *filter-with-index'* (i+1) P xs
 else *filter-with-index'* (i+1) P xs)

definition *filter-with-index* :: $('a \Rightarrow \text{bool}) \Rightarrow 'a \text{ list} \Rightarrow ('a \times \text{nat}) \text{ list}$ **where**
filter-with-index P xs = *filter-with-index'* 0 P xs

datatype *pointer* =
 Null
 | Pre nat
 | PreRed nat \times nat \times nat (nat \times nat \times nat) list

datatype $'a$ *entry* =
 Entry

(item : 'a item)
(pointer : pointer)

type-synonym 'a bin = 'a entry list
type-synonym 'a bins = 'a bin list

definition items :: 'a bin \Rightarrow 'a item list **where**
items b = map item b

definition pointers :: 'a bin \Rightarrow pointer list **where**
pointers b = map pointer b

definition bins-eq-items :: 'a bins \Rightarrow 'a bins \Rightarrow bool **where**
bins-eq-items bs0 bs1 \longleftrightarrow map items bs0 = map items bs1

definition bins-items :: 'a bins \Rightarrow 'a items **where**
bins-items bs = $\bigcup \{ \text{set } (\text{items } (bs ! k)) \mid k. k < \text{length } bs \}$

definition bin-items-upto :: 'a bin \Rightarrow nat \Rightarrow 'a items **where**
bin-items-upto b i = $\{ \text{items } b ! j \mid j. j < i \wedge j < \text{length } (\text{items } b) \}$

definition bins-items-upto :: 'a bins \Rightarrow nat \Rightarrow nat \Rightarrow 'a items **where**
bins-items-upto bs k i = $\bigcup \{ \text{set } (\text{items } (bs ! l)) \mid l. l < k \} \cup \text{bin-items-upto } (bs ! k) i$

definition wf-bin-items :: 'a cfg \Rightarrow 'a sentential \Rightarrow nat \Rightarrow 'a item list \Rightarrow bool **where**
wf-bin-items cfg inp k xs $\equiv \forall x \in \text{set } xs. \text{wf-item } \text{cfg } \text{inp } x \wedge \text{item-end } x = k$

definition wf-bin :: 'a cfg \Rightarrow 'a sentential \Rightarrow nat \Rightarrow 'a bin \Rightarrow bool **where**
wf-bin cfg inp k b $\equiv \text{distinct } (\text{items } b) \wedge \text{wf-bin-items } \text{cfg } \text{inp } k (\text{items } b)$

definition wf-bins :: 'a cfg \Rightarrow 'a list \Rightarrow 'a bins \Rightarrow bool **where**
wf-bins cfg inp bs $\equiv \forall k < \text{length } bs. \text{wf-bin } \text{cfg } \text{inp } k (bs ! k)$

definition nonempty-derives :: 'a cfg \Rightarrow bool **where**
nonempty-derives cfg $\equiv \forall N. N \in \text{set } (\mathfrak{N} \text{ cfg}) \longrightarrow \neg \text{derives } \text{cfg } [N] []$

definition Init-list :: 'a cfg \Rightarrow 'a sentential \Rightarrow 'a bins **where**
Init-list cfg inp \equiv
let rs = filter ($\lambda r. \text{rule-head } r = \mathfrak{S} \text{ cfg}$) ($\mathfrak{R} \text{ cfg}$) in
let b0 = map ($\lambda r. (\text{Entry } (\text{init-item } r 0) \text{Null})$) rs in
let bs = replicate (length inp + 1) ([]) in
bs[0 := b0]

definition Scan-list :: nat \Rightarrow 'a sentential \Rightarrow 'a \Rightarrow 'a item \Rightarrow nat \Rightarrow 'a entry list **where**

Scan-list $k \text{ inp } a \ x \ pre \equiv$
 if $\text{inp}!k = a$ then
 let $x' = \text{inc-item } x \ (k+1)$ in
 [Entry $x' \ (\text{Pre } pre)$]
 else []

definition Predict-list $:: \text{nat} \Rightarrow 'a \text{ cfg} \Rightarrow 'a \Rightarrow 'a \text{ entry list where}$

Predict-list $k \text{ cfg } X \equiv$
 let $rs = \text{filter } (\lambda r. \text{rule-head } r = X) \ (\mathfrak{R} \text{ cfg})$ in
 map $(\lambda r. (\text{Entry } (\text{init-item } r \ k) \ \text{Null})) \ rs$

definition Complete-list $:: \text{nat} \Rightarrow 'a \text{ item} \Rightarrow 'a \text{ bins} \Rightarrow \text{nat} \Rightarrow 'a \text{ entry list where}$

Complete-list $k \ y \ bs \ red \equiv$
 let $\text{orig} = bs \ ! \ (\text{item-origin } y)$ in
 let $is = \text{filter-with-index } (\lambda x. \text{next-symbol } x = \text{Some } (\text{item-rule-head } y)) \ (\text{items orig})$ in
 map $(\lambda (x, pre). (\text{Entry } (\text{inc-item } x \ k) \ (\text{PreRed } (\text{item-origin } y, pre, red) \ []))) \ is$

fun bin-upd $:: 'a \text{ entry} \Rightarrow 'a \text{ bin} \Rightarrow 'a \text{ bin where}$

bin-upd $e' \ [] = [e']$
 | bin-upd $e' \ (e\#es) =$
 case (e', e) of
 $(\text{Entry } x \ (\text{PreRed } px \ xs), \text{Entry } y \ (\text{PreRed } py \ ys)) \Rightarrow$
 if $x = y$ then $\text{Entry } x \ (\text{PreRed } py \ (px\#xs@ys)) \# es$
 else $e \# \text{bin-upd } e' \ es$
 | - \Rightarrow
 if $\text{item } e' = \text{item } e$ then $e \# es$
 else $e \# \text{bin-upd } e' \ es$

fun bin-upds $:: 'a \text{ entry list} \Rightarrow 'a \text{ bin} \Rightarrow 'a \text{ bin where}$

bin-upds $[] \ b = b$
 | bin-upds $(e\#es) \ b = \text{bin-upds } es \ (\text{bin-upd } e \ b)$

definition bins-upd $:: 'a \text{ bins} \Rightarrow \text{nat} \Rightarrow 'a \text{ entry list} \Rightarrow 'a \text{ bins where}$

bins-upd $bs \ k \ es = bs[k := \text{bin-upds } es \ (bs!k)]$

partial-function (tailrec) E-list' $:: \text{nat} \Rightarrow 'a \text{ cfg} \Rightarrow 'a \text{ sentential} \Rightarrow 'a \text{ bins} \Rightarrow \text{nat} \Rightarrow 'a \text{ bins where}$

E-list' $k \text{ cfg inp } bs \ i =$
 if $i \geq \text{length } (\text{items } (bs \ ! \ k))$ then bs
 else
 let $x = \text{items } (bs!k) \ ! \ i$ in
 let $bs' =$
 case $\text{next-symbol } x$ of
 Some $a \Rightarrow$
 if $\text{is-terminal } \text{cfg } a$ then

```

    if k < length inp then bins-upd bs (k+1) (Scan-list k inp a x i)
    else bs
  else bins-upd bs k (Predict-list k cfg a)
  | None ⇒ bins-upd bs k (Complete-list k x bs i)
  in E-list' k cfg inp bs' (i+1))

```

definition $E\text{-list} :: \text{nat} \Rightarrow 'a \text{ cfg} \Rightarrow 'a \text{ sentential} \Rightarrow 'a \text{ bins} \Rightarrow 'a \text{ bins}$ **where**
 $E\text{-list } k \text{ cfg inp bs} = E\text{-list}' k \text{ cfg inp bs } 0$

fun $\mathcal{E}\text{-list} :: \text{nat} \Rightarrow 'a \text{ cfg} \Rightarrow 'a \text{ sentential} \Rightarrow 'a \text{ bins}$ **where**
 $\mathcal{E}\text{-list } 0 \text{ cfg inp} = E\text{-list } 0 \text{ cfg inp (Init-list cfg inp)}$
 $|\ \mathcal{E}\text{-list (Suc n) cfg inp} = E\text{-list (Suc n) cfg inp (\mathcal{E}\text{-list } n \text{ cfg inp})$

definition $\text{earley-list} :: 'a \text{ cfg} \Rightarrow 'a \text{ sentential} \Rightarrow 'a \text{ bins}$ **where**
 $\text{earley-list cfg inp} = \mathcal{E}\text{-list (length inp) cfg inp}$

5.3 Wellformedness

lemma wf-bin-bin-upd :

assumes $\text{wf-bin cfg inp } k \text{ b wf-item cfg inp (item e) item-end (item e) = k}$
shows $\text{wf-bin cfg inp } k \text{ (bin-upd e b)}$

lemma wf-bin-bin-upds :

assumes $\text{wf-bin cfg inp } k \text{ b distinct (items es)}$
assumes $\forall x \in \text{set (items es)}. \text{wf-item cfg inp } x \wedge \text{item-end } x = k$
shows $\text{wf-bin cfg inp } k \text{ (bin-upds es b)}$

lemma wf-bins-bins-upd :

assumes $\text{wf-bins cfg inp bs distinct (items es)}$
assumes $\forall x \in \text{set (items es)}. \text{wf-item cfg inp } x \wedge \text{item-end } x = k$
shows $\text{wf-bins cfg inp (bins-upd bs k es)}$

lemma wf-bins-Init-list :

assumes wf-cfg cfg
shows $\text{wf-bins cfg inp (Init-list cfg inp)}$

lemma wf-bins-Scan-list :

assumes $\text{wf-bins cfg inp bs } k < \text{length bs } x \in \text{set (items (bs!k)) } k < \text{length inp next-symbol } x \neq \text{None}$
shows $\forall y \in \text{set (items (Scan-list k inp a x pre))}. \text{wf-item cfg inp } y \wedge \text{item-end } y = k+1$

lemma $\text{wf-bins-Predict-list}$:

assumes $\text{wf-bins cfg inp bs } k < \text{length bs } k \leq \text{length inp wf-cfg cfg}$
shows $\forall y \in \text{set (items (Predict-list k cfg X))}. \text{wf-item cfg inp } y \wedge \text{item-end } y = k$

lemma $\text{wf-bins-Complete-list}$:

assumes $\text{wf-bins cfg inp bs } k < \text{length bs } y \in \text{set (items (bs!k))}$
shows $\forall x \in \text{set (items (Complete-list k y bs red))}. \text{wf-item cfg inp } x \wedge \text{item-end } x = k$

fun $\text{earley-measure} :: \text{nat} \times 'a \text{ cfg} \times 'a \text{ sentential} \times 'a \text{ bins} \Rightarrow \text{nat} \Rightarrow \text{nat}$ **where**

$\text{earley-measure } (k, \text{cfg}, \text{inp}, \text{bs}) \text{ i} = \text{card } \{ x \mid x. \text{wf-item cfg inp } x \wedge \text{item-end } x = k \} - i$

definition *wf-earley-input* :: (nat × 'a cfg × 'a sentential × 'a bins) set **where**

$$\begin{aligned} \text{wf-earley-input} = \{ \\ & (k, \text{cfg}, \text{inp}, \text{bs}) \mid k \text{ cfg inp bs.} \\ & k \leq \text{length inp} \wedge \\ & \text{length bs} = \text{length inp} + 1 \wedge \\ & \text{wf-cfg cfg} \wedge \\ & \text{wf-bins cfg inp bs} \\ & \} \end{aligned}$$

lemma *wf-earley-input-Init-list*:

assumes $k \leq \text{length inp}$ *wf-cfg* *cfg*
shows $(k, \text{cfg}, \text{inp}, \text{Init-list } \text{cfg } \text{inp}) \in \text{wf-earley-input}$

lemma *wf-earley-input-Complete-list*:

assumes $(k, \text{cfg}, \text{inp}, \text{bs}) \in \text{wf-earley-input} \neg \text{length (items (bs!k))} \leq i$
assumes $x = \text{items (bs!k)!i next-symbol } x = \text{None}$
shows $(k, \text{cfg}, \text{inp}, \text{bins-upd bs k (Complete-list k x bs red)}) \in \text{wf-earley-input}$

lemma *wf-earley-input-Scan-list*:

assumes $(k, \text{cfg}, \text{inp}, \text{bs}) \in \text{wf-earley-input} \neg \text{length (items (bs!k))} \leq i$
assumes $x = \text{items (bs!k)!i next-symbol } x = \text{Some } a$
assumes *is-terminal* *cfg* $a < \text{length inp}$
shows $(k, \text{cfg}, \text{inp}, \text{bins-upd bs (k+1) (Scan-list k inp a x pre)}) \in \text{wf-earley-input}$

lemma *wf-earley-input-Predict-list*:

assumes $(k, \text{cfg}, \text{inp}, \text{bs}) \in \text{wf-earley-input} \neg \text{length (items (bs!k))} \leq i$
assumes $x = \text{items (bs!k)!i next-symbol } x = \text{Some } a \neg \text{is-terminal } \text{cfg } a$
shows $(k, \text{cfg}, \text{inp}, \text{bins-upd bs k (Predict-list k cfg a)}) \in \text{wf-earley-input}$

lemma *wf-earley-input-E-list'*:

assumes $(k, \text{cfg}, \text{inp}, \text{bs}) \in \text{wf-earley-input}$
shows $(k, \text{cfg}, \text{inp}, \text{E-list' k cfg inp bs i}) \in \text{wf-earley-input}$

lemma *wf-earley-input-E-list*:

assumes $(k, \text{cfg}, \text{inp}, \text{bs}) \in \text{wf-earley-input}$
shows $(k, \text{cfg}, \text{inp}, \text{E-list k cfg inp bs}) \in \text{wf-earley-input}$

lemma *wf-earley-input- \mathcal{E} -list*:

assumes $k \leq \text{length inp}$ *wf-cfg* *cfg*
shows $(k, \text{cfg}, \text{inp}, \mathcal{E}\text{-list k cfg inp}) \in \text{wf-earley-input}$

lemma *wf-earley-input-earley-list*:

assumes $k \leq \text{length inp}$ *wf-cfg* *cfg*
shows $(k, \text{cfg}, \text{inp}, \text{earley-list } \text{cfg } \text{inp}) \in \text{wf-earley-input}$

lemma *wf-bins-E-list'*:

assumes $(k, \text{cfg}, \text{inp}, \text{bs}) \in \text{wf-earley-input}$
shows *wf-bins* *cfg* *inp* (*E-list'* *k* *cfg* *inp* *bs* *i*)

lemma *wf-bins-E-list*:

assumes $(k, \text{cfg}, \text{inp}, \text{bs}) \in \text{wf-earley-input}$
shows *wf-bins* *cfg* *inp* (*E-list* *k* *cfg* *inp* *bs*)

lemma *wf-bins- \mathcal{E} -list*:
assumes $k \leq \text{length } \text{inp } \text{wf-cfg } \text{cfg}$
shows $\text{wf-bins } \text{cfg } \text{inp } (\mathcal{E}\text{-list } k \text{ cfg } \text{inp})$
lemma *wf-bins-earley-list*:
assumes $\text{wf-cfg } \text{cfg}$
shows $\text{wf-bins } \text{cfg } \text{inp } (\text{earley-list } \text{cfg } \text{inp})$

5.4 Soundness

lemma *Init-list-eq-Init*:
shows $\text{bins-items } (\text{Init-list } \text{cfg } \text{inp}) = \text{Init } \text{cfg}$
lemma *Scan-list-sub-Scan*:
assumes $\text{wf-bins } \text{cfg } \text{inp } \text{bs } \text{bins-items } \text{bs } \subseteq I \ x \in \text{set } (\text{items } (\text{bs } ! \ k))$
assumes $k < \text{length } \text{bs } k < \text{length } \text{inp } \text{next-symbol } x = \text{Some } a$
shows $\text{set } (\text{items } (\text{Scan-list } k \text{ inp } a \ x \ \text{pre})) \subseteq \text{Scan } k \text{ inp } I$
lemma *Predict-list-sub-Predict*:
assumes $\text{wf-bins } \text{cfg } \text{inp } \text{bs } \text{bins-items } \text{bs } \subseteq I \ x \in \text{set } (\text{items } (\text{bs } ! \ k)) \ k < \text{length } \text{bs}$
assumes $\text{next-symbol } x = \text{Some } X$
shows $\text{set } (\text{items } (\text{Predict-list } k \text{ cfg } X)) \subseteq \text{Predict } k \text{ cfg } I$
lemma *Complete-list-sub-Complete*:
assumes $\text{wf-bins } \text{cfg } \text{inp } \text{bs } \text{bins-items } \text{bs } \subseteq I \ y \in \text{set } (\text{items } (\text{bs } ! \ k)) \ k < \text{length } \text{bs}$
assumes $\text{next-symbol } y = \text{None}$
shows $\text{set } (\text{items } (\text{Complete-list } k \ y \ \text{bs } \text{red})) \subseteq \text{Complete } k \ I$
lemma *E-list'-sub-E*:
assumes $(k, \text{cfg}, \text{inp}, \text{bs}) \in \text{wf-earley-input}$
assumes $\text{bins-items } \text{bs } \subseteq I$
shows $\text{bins-items } (\text{E-list}' k \text{ cfg } \text{inp } \text{bs } i) \subseteq E \ k \text{ cfg } \text{inp } I$
lemma *E-list-sub-E*:
assumes $(k, \text{cfg}, \text{inp}, \text{bs}) \in \text{wf-earley-input}$
assumes $\text{bins-items } \text{bs } \subseteq I$
shows $\text{bins-items } (\text{E-list } k \text{ cfg } \text{inp } \text{bs}) \subseteq E \ k \text{ cfg } \text{inp } I$
lemma *\mathcal{E} -list-sub- \mathcal{E}* :
assumes $k \leq \text{length } \text{inp } \text{wf-cfg } \text{cfg}$
shows $\text{bins-items } (\mathcal{E}\text{-list } k \text{ cfg } \text{inp}) \subseteq \mathcal{E} \ k \text{ cfg } \text{inp}$
lemma *earley-list-sub-earley*:
assumes $\text{wf-cfg } \text{cfg}$
shows $\text{bins-items } (\text{earley-list } \text{cfg } \text{inp}) \subseteq \text{earley } \text{cfg } \text{inp}$

5.5 Completeness

lemma *impossible-complete-item*:
assumes $\text{wf-cfg } \text{cfg } \text{wf-item } \text{cfg } \text{inp } x \text{ sound-item } \text{cfg } \text{inp } x$

assumes *is-complete* x *item-origin* $x = k$ *item-end* $x = k$ *nonempty-derives* *cfg*

shows *False*

lemma *Complete-Un-eq-terminal:*

assumes *next-symbol* $z = \text{Some } a$ *is-terminal* *cfg* a *wf-items* *cfg* *inp* I *wf-item* *cfg* *inp* z *wf-cfg* *cfg*

shows *Complete* k $(I \cup \{z\}) = \text{Complete } k I$

lemma *Complete-Un-eq-nonterminal:*

assumes *next-symbol* $z = \text{Some } a$ *is-nonterminal* *cfg* a *sound-items* *cfg* *inp* I *item-end* $z = k$

assumes *wf-items* *cfg* *inp* I *wf-item* *cfg* *inp* z *wf-cfg* *cfg* *nonempty-derives* *cfg*

shows *Complete* k $(I \cup \{z\}) = \text{Complete } k I$

lemma *Complete-sub-bins-Un-Complete-list:*

assumes *Complete* $k I \subseteq \text{bins-items } bs$ $I \subseteq \text{bins-items } bs$ *is-complete* z *wf-bins* *cfg* *inp* bs *wf-item* *cfg* *inp* z

shows *Complete* k $(I \cup \{z\}) \subseteq \text{bins-items } bs \cup \text{set } (\text{items } (\text{Complete-list } k z bs \text{ red}))$

lemma *E-list'-mono:*

assumes $(k, \text{cfg}, \text{inp}, bs) \in \text{wf-earley-input}$

shows *bins-items* $bs \subseteq \text{bins-items } (E\text{-list}' k \text{ cfg } \text{inp } bs i)$

lemma *E-step-sub-E-list':*

assumes $(k, \text{cfg}, \text{inp}, bs) \in \text{wf-earley-input}$

assumes *E-step* $k \text{ cfg } \text{inp}$ $(\text{bins-items-upto } bs k i) \subseteq \text{bins-items } bs$

assumes *sound-items* *cfg* *inp* $(\text{bins-items } bs)$ *is-sentence* *cfg* *inp* *nonempty-derives* *cfg*

shows *E-step* $k \text{ cfg } \text{inp}$ $(\text{bins-items } bs) \subseteq \text{bins-items } (E\text{-list}' k \text{ cfg } \text{inp } bs i)$

lemma *E-step-sub-E-list:*

assumes $(k, \text{cfg}, \text{inp}, bs) \in \text{wf-earley-input}$

assumes *E-step* $k \text{ cfg } \text{inp}$ $(\text{bins-items-upto } bs k 0) \subseteq \text{bins-items } bs$

assumes *sound-items* *cfg* *inp* $(\text{bins-items } bs)$ *is-sentence* *cfg* *inp* *nonempty-derives* *cfg*

shows *E-step* $k \text{ cfg } \text{inp}$ $(\text{bins-items } bs) \subseteq \text{bins-items } (E\text{-list } k \text{ cfg } \text{inp } bs)$

lemma *E-list'-bins-items-eq:*

assumes $(k, \text{cfg}, \text{inp}, as) \in \text{wf-earley-input}$

assumes *bins-eq-items* as *bs* *wf-bins* *cfg* *inp* as

shows *bins-eq-items* $(E\text{-list}' k \text{ cfg } \text{inp } as i) (E\text{-list}' k \text{ cfg } \text{inp } bs i)$

lemma *E-list'-idem:*

assumes $(k, \text{cfg}, \text{inp}, bs) \in \text{wf-earley-input}$

assumes $i \leq j$ *sound-items* *cfg* *inp* $(\text{bins-items } bs)$ *nonempty-derives* *cfg*

shows *bins-items* $(E\text{-list}' k \text{ cfg } \text{inp } (E\text{-list}' k \text{ cfg } \text{inp } bs i) j) = \text{bins-items } (E\text{-list}' k \text{ cfg } \text{inp } bs i)$

lemma *E-list-idem:*

assumes $(k, \text{cfg}, \text{inp}, bs) \in \text{wf-earley-input}$

assumes *sound-items* *cfg* *inp* $(\text{bins-items } bs)$ *nonempty-derives* *cfg*

shows *bins-items* $(E\text{-list } k \text{ cfg } \text{inp } (E\text{-list } k \text{ cfg } \text{inp } bs)) = \text{bins-items } (E\text{-list } k \text{ cfg } \text{inp } bs)$

lemma *funpower-E-step-sub-E-list:*

assumes $(k, \text{cfg}, \text{inp}, bs) \in \text{wf-earley-input}$

assumes *E-step* $k \text{ cfg } \text{inp}$ $(\text{bins-items-upto } bs k 0) \subseteq \text{bins-items } bs$ *sound-items* *cfg* *inp* $(\text{bins-items } bs)$

assumes *is-sentence* *cfg* *inp* *nonempty-derives* *cfg*

shows *funpower* $(E\text{-step } k \text{ cfg } \text{inp}) n (\text{bins-items } bs) \subseteq \text{bins-items } (E\text{-list } k \text{ cfg } \text{inp } bs)$

lemma *E-sub-E-list:*

assumes $(k, \text{cfg}, \text{inp}, \text{bs}) \in \text{wf-earley-input}$
assumes $E\text{-step } k \text{ cfg inp } (\text{bins-items-upto } \text{bs } k \ 0) \subseteq \text{bins-items } \text{bs sound-items } \text{cfg inp } (\text{bins-items } \text{bs})$
assumes $\text{is-sentence } \text{cfg inp nonempty-derives } \text{cfg}$
shows $E \ k \ \text{cfg inp } (\text{bins-items } \text{bs}) \subseteq \text{bins-items } (E\text{-list } k \ \text{cfg inp } \text{bs})$
lemma $\mathcal{E}\text{-sub-}\mathcal{E}\text{-list}$:
assumes $k \leq \text{length } \text{inp wf-cfg } \text{cfg}$
assumes $\text{is-sentence } \text{cfg inp nonempty-derives } \text{cfg}$
shows $\mathcal{E} \ k \ \text{cfg inp} \subseteq \text{bins-items } (\mathcal{E}\text{-list } k \ \text{cfg inp})$
lemma $\text{earley-sub-earley-list}$:
assumes $\text{wf-cfg } \text{cfg is-sentence } \text{cfg inp nonempty-derives } \text{cfg}$
shows $\text{earley } \text{cfg inp} \subseteq \text{bins-items } (\text{earley-list } \text{cfg inp})$

5.6 Main Theorem

definition $\text{recognizing-list} :: 'a \text{ bins} \Rightarrow 'a \text{ cfg} \Rightarrow 'a \text{ sentential} \Rightarrow \text{bool}$ **where**
 $\text{recognizing-list } I \ \text{cfg inp} \equiv \exists x \in \text{set } (\text{items } (I \ ! \ \text{length } \text{inp})). \text{is-finished } \text{cfg inp } x$

theorem $\text{recognizing-list-iff-earley-recognized}$:

assumes $\text{wf-cfg } \text{cfg is-sentence } \text{cfg inp nonempty-derives } \text{cfg}$
shows $\text{recognizing-list } (\text{earley-list } \text{cfg inp}) \ \text{cfg inp} \longleftrightarrow \text{recognizing } (\text{earley } \text{cfg inp}) \ \text{cfg inp}$

corollary correctness-list :

assumes $\text{wf-cfg } \text{cfg is-sentence } \text{cfg inp nonempty-derives } \text{cfg}$
shows $\text{recognizing-list } (\text{earley-list } \text{cfg inp}) \ \text{cfg inp} \longleftrightarrow \text{derives } \text{cfg } [\mathcal{E} \ \text{cfg}] \ \text{inp}$

SNIPPET:

It is this latter possibility, adding items to S_i while representing sets as lists, which causes grief with epsilon-rules. When Completer processes an item $A \rightarrow \text{dot}, j$ which corresponds to the epsilon-rule $A \rightarrow \epsilon$, it must look through S_j for items with the dot before an A . Unfortunately, for epsilon-rule items, j is always equal to i . Completer is thus looking through the partially constructed set S_i . Since implementations process items in S_i in order, if an item $B \rightarrow \alpha \text{dot } A \beta$, k is added to S_i after Completer has processed $A \rightarrow \text{dot}, j$, Completer will never add $B \rightarrow \alpha A \text{dot } \beta, k$ to S_i . In turn, items resulting directly and indirectly from $B \rightarrow \alpha A \text{dot } \beta, k$ will be omitted too. This effectively prunes potential derivation paths which might cause correct input to be rejected. (EXAMPLE) Aho *et al* [Aho:1972] propose the stay clam and keep running the Predictor and Completer in turn until neither has anything more to add. Earley himself suggest to have the Completer note that the dot needed to be moved over A , then looking for this whenever future items were added to S_i . For efficiency's sake the collection of on-terminals to watch for should be stored in a data structure which allows fast access. Neither approach is very satisfactory. A third solution [Aycoack:2002] is a simple modification of the Predictor based on the idea of nullability. A non-terminal A

is said to be nullable if A derives star epsilon. Terminal symbols of course can never be nullable. The nullability of non-terminals in a grammar may be precomputed using well-known techniques [Appel:2003] [Fischer:2009] Using this notion the Predictor can be stated as follows: if $A \rightarrow \alpha \dot{B} \beta$, j is in S_i , add $B \rightarrow \dot{\gamma}$, i to S_i for all rules $B \rightarrow \gamma$. If B is nullable, also add $A \rightarrow \alpha B \dot{\beta}$, j to S_i . Explanation why I decided against it. Involves every grammar can be rewritten to not contain epsilon productions. In other words we eagerly move the dot over a nonterminal if that non-terminal can derive epsilon and effectively disappear. The source implements this precomputation by constructing a variant of a LR(0) deterministic finite automata (DFA). But for an earley parser we must keep track of which parent pointers and LR(0) items belong together which leads to complex and inelegant implementations [McLean:1996]. The source resolves this problem by constructing split epsilon DFAs, but still need to adjust the classical earley algorithm by adding not only predecessor links but also causal links, and to construct the split epsilon DFAs not the original grammar but a slightly adjusted equivalent grammar is used that encodes explicitly information that is crucial to reconstructing derivations, called a grammar in nihilist normal form (NNF) which might increase the size of the grammar whereas the authors note empirical results that the increase is quite modest (a factor of 2 at most).

Example: $S \rightarrow AAAA$, $A \rightarrow a$, $A \rightarrow E$, $E \rightarrow \text{epsilon}$, input a S_0 $S \rightarrow \dot{A} AAAA, 0$, $A \rightarrow \dot{a}, 0$, $A \rightarrow \dot{E}, 0$, $E \rightarrow \dot{\epsilon}, 0$, $A \rightarrow E \dot{\epsilon}, 0$, $S \rightarrow A \dot{A} AA, 0$ S_1 $A \rightarrow a \dot{\epsilon}, 0$, $S \rightarrow A \dot{A} AA, 0$, $S \rightarrow AA \dot{A} A, 0$, $A \rightarrow \dot{a}, 1$, $A \rightarrow \dot{E}, 1$, $E \rightarrow \dot{\epsilon}, 1$, $A \rightarrow E \dot{\epsilon}, 1$, $S \rightarrow AAA \dot{A}, 0$

6 Earley Parser Implementation

6.1 Draft

6.2 Pointer lemmas

definition *predicts* :: 'a item \Rightarrow bool **where**
predicts $x \equiv \text{item-origin } x = \text{item-end } x \wedge \text{item-bullet } x = 0$

definition *scans* :: 'a sentential \Rightarrow nat \Rightarrow 'a item \Rightarrow 'a item \Rightarrow bool **where**
scans $\text{inp } k \ x \ y \equiv y = \text{inc-item } x \ k \wedge (\exists a. \text{next-symbol } x = \text{Some } a \wedge \text{inp}!(k-1) = a)$

definition *completes* :: nat \Rightarrow 'a item \Rightarrow 'a item \Rightarrow 'a item \Rightarrow bool **where**
completes $k \ x \ y \ z \equiv y = \text{inc-item } x \ k \wedge \text{is-complete } z \wedge \text{item-origin } z = \text{item-end } x \wedge$
 $(\exists N. \text{next-symbol } x = \text{Some } N \wedge N = \text{item-rule-head } z)$

definition *sound-null-ptr* :: 'a entry \Rightarrow bool **where**
sound-null-ptr $e \equiv \text{pointer } e = \text{Null} \longrightarrow \text{predicts } (\text{item } e)$

definition *sound-pre-ptr* :: 'a sentential \Rightarrow 'a bins \Rightarrow nat \Rightarrow 'a entry \Rightarrow bool **where**
sound-pre-ptr $\text{inp } bs \ k \ e \equiv \forall \text{pre. pointer } e = \text{Pre } \text{pre} \longrightarrow$
 $k > 0 \wedge \text{pre} < \text{length } (bs!(k-1)) \wedge \text{scans } \text{inp } k \ (\text{item } (bs!(k-1)!\text{pre})) \ (\text{item } e)$

definition *sound-prered-ptr* :: 'a bins \Rightarrow nat \Rightarrow 'a entry \Rightarrow bool **where**
sound-prered-ptr $bs \ k \ e \equiv \forall p \ ps \ k' \ \text{pre } \text{red. pointer } e = \text{PreRed } p \ ps \wedge (k', \text{pre}, \text{red}) \in \text{set } (p\#ps) \longrightarrow$
 $k' < k \wedge \text{pre} < \text{length } (bs!k') \wedge \text{red} < \text{length } (bs!k) \wedge \text{completes } k \ (\text{item } (bs!k'!\text{pre})) \ (\text{item } e) \ (\text{item } (bs!k!\text{red}))$

definition *sound-ptrs* :: 'a sentential \Rightarrow 'a bins \Rightarrow bool **where**
sound-ptrs $\text{inp } bs \equiv \forall k < \text{length } bs. \forall e \in \text{set } (bs!k). \text{sound-null-ptr } e \wedge$
 $\text{sound-pre-ptr } \text{inp } bs \ k \ e \wedge$
 $\text{sound-prered-ptr } bs \ k \ e$

definition *mono-red-ptr* :: 'a bins \Rightarrow bool **where**
mono-red-ptr $bs \equiv \forall k < \text{length } bs. \forall i < \text{length } (bs!k). \forall k' \ \text{pre } \text{red } ps. \text{pointer } (bs!k!i) = \text{PreRed } (k', \text{pre}, \text{red}) \ ps \longrightarrow \text{red} < i$

lemma *sound-ptrs-bin-upd*:
assumes *sound-ptrs inp bs k < length bs es = bs!k distinct (items es)*
assumes *sound-null-ptr e sound-pre-ptr inp bs k e sound-prered-ptr bs k e*
shows *sound-ptrs inp (bs[k := bin-upd e es])*

lemma *mono-red-ptr-bin-upd*:
assumes *mono-red-ptr bs k < length bs es = bs!k distinct (items es)*
assumes $\forall k' \text{ pre red ps. pointer } e = \text{PreRed } (k', \text{pre}, \text{red}) \text{ ps} \longrightarrow \text{red} < \text{length es}$
shows *mono-red-ptr (bs[k := bin-upd e es])*

lemma *sound-mono-ptrs-bin-upds*:
assumes *sound-ptrs inp bs mono-red-ptr bs k < length bs b = bs!k distinct (items b) distinct (items es)*
assumes $\forall e \in \text{set es. sound-null-ptr } e \wedge \text{sound-pre-ptr inp bs k e} \wedge \text{sound-prered-ptr bs k e}$
assumes $\forall e \in \text{set es. } \forall k' \text{ pre red ps. pointer } e = \text{PreRed } (k', \text{pre}, \text{red}) \text{ ps} \longrightarrow \text{red} < \text{length b}$
shows *sound-ptrs inp (bs[k := bin-upds es b]) \wedge mono-red-ptr (bs[k := bin-upds es b])*

lemma *sound-mono-ptrs-E-list'*:
assumes $(k, \text{cfg}, \text{inp}, \text{bs}) \in \text{wellformed-bins}$
assumes *sound-ptrs inp bs sound-items cfg inp (bins-items bs)*
assumes *mono-red-ptr bs*
assumes *nonempty-derives cfg wf-cfg cfg*
shows *sound-ptrs inp (E-list' k cfg inp bs i) \wedge mono-red-ptr (E-list' k cfg inp bs i)*

lemma *sound-mono-ptrs-E-list*:
assumes $(k, \text{cfg}, \text{inp}, \text{bs}) \in \text{wellformed-bins}$
assumes *sound-ptrs inp bs sound-items cfg inp (bins-items bs)*
assumes *mono-red-ptr bs*
assumes *nonempty-derives cfg wf-cfg cfg*
shows *sound-ptrs inp (E-list k cfg inp bs) \wedge mono-red-ptr (E-list k cfg inp bs)*

lemma *sound-ptrs-Init-list*:
shows *sound-ptrs inp (Init-list cfg inp)*

lemma *mono-red-ptr-Init-list*:
shows *mono-red-ptr (Init-list cfg inp)*

lemma *sound-mono-ptrs-E-list*:
assumes $k \leq \text{length inp wf-cfg cfg nonempty-derives cfg wf-cfg cfg}$
shows *sound-ptrs inp (E-list k cfg inp) \wedge mono-red-ptr (E-list k cfg inp)*

lemma *sound-mono-ptrs-earley-list*:
assumes *wf-cfg cfg nonempty-derives cfg*
shows *sound-ptrs inp (earley-list cfg inp) \wedge mono-red-ptr (earley-list cfg inp)*

6.3 Trees and Forests

datatype *'a tree* =
Leaf 'a
| *Branch 'a 'a tree list*

```

fun yield-tree :: 'a tree  $\Rightarrow$  'a sentential where
  yield-tree (Leaf a) = [a]
| yield-tree (Branch ts) = concat (map yield-tree ts)

fun root-tree :: 'a tree  $\Rightarrow$  'a where
  root-tree (Leaf a) = a
| root-tree (Branch N _) = N

fun wf-rule-tree :: 'a cfg  $\Rightarrow$  'a tree  $\Rightarrow$  bool where
  wf-rule-tree - (Leaf a)  $\longleftrightarrow$  True
| wf-rule-tree cfg (Branch N ts)  $\longleftrightarrow$  (
  ( $\exists r \in \text{set } (\mathcal{R} \text{ cfg}). N = \text{rule-head } r \wedge \text{map root-tree ts} = \text{rule-body } r$ )  $\wedge$ 
  ( $\forall t \in \text{set ts}. \text{wf-rule-tree cfg } t$ ))

fun wf-item-tree :: 'a cfg  $\Rightarrow$  'a item  $\Rightarrow$  'a tree  $\Rightarrow$  bool where
  wf-item-tree cfg - (Leaf a)  $\longleftrightarrow$  True
| wf-item-tree cfg x (Branch N ts)  $\longleftrightarrow$  (
   $N = \text{item-rule-head } x \wedge \text{map root-tree ts} = \text{take } (\text{item-bullet } x) (\text{item-rule-body } x) \wedge$ 
  ( $\forall t \in \text{set ts}. \text{wf-rule-tree cfg } t$ ))

definition wf-yield-tree :: 'a sentential  $\Rightarrow$  'a item  $\Rightarrow$  'a tree  $\Rightarrow$  bool where
  wf-yield-tree inp x t  $\equiv \text{yield-tree } t = \text{slice } (\text{item-origin } x) (\text{item-end } x) \text{ inp}$ 

datatype 'a forest =
  FLeaf 'a
| FBranch 'a 'a forest list list

fun combinations :: 'a list list  $\Rightarrow$  'a list list where
  combinations [] = [[]]
| combinations (xs#xss) = [ x#cs . x <- xs, cs <- combinations xss ]

fun trees :: 'a forest  $\Rightarrow$  'a tree list where
  trees (FLeaf a) = [Leaf a]
| trees (FBranch N fss) = (
  let tss = (map ( $\lambda f$ . concat (map ( $\lambda f$ . trees f) fs)) fss) in
  map ( $\lambda ts. \text{Branch } N \text{ ts}$ ) (combinations tss)
  )

```

6.4 A Single Parse Tree

```

partial-function (option) build-tree' :: 'a bins  $\Rightarrow$  'a sentential  $\Rightarrow$  nat  $\Rightarrow$  nat  $\Rightarrow$  'a tree option where
  build-tree' bs inp k i = (
    let e = bs!k!i in (

```

```

case pointer e of
  Null  $\Rightarrow$  Some (Branch (item-rule-head (item e)) [])
| Pre pre  $\Rightarrow$  (
  do {
    t  $\leftarrow$  build-tree' bs inp (k-1) pre;
    case t of
      Branch N ts  $\Rightarrow$  Some (Branch N (ts @ [Leaf (inp!(k-1))]))
    | -  $\Rightarrow$  None
  })
| PreRed (k', pre, red) -  $\Rightarrow$  (
  do {
    t  $\leftarrow$  build-tree' bs inp k' pre;
    case t of
      Branch N ts  $\Rightarrow$ 
        do {
          t  $\leftarrow$  build-tree' bs inp k red;
          Some (Branch N (ts @ [t]))
        }
    | -  $\Rightarrow$  None
  })
))

```

definition build-tree :: 'a cfg \Rightarrow 'a sentential \Rightarrow 'a bins \Rightarrow 'a tree option **where**
 build-tree cfg inp bs \equiv
 let k = length bs - 1 in (
 case filter-with-index ($\lambda x. \text{is-finished } \text{cfg } \text{inp } x$) (items (bs!k)) of
 [] \Rightarrow None
 | (-, i)#- \Rightarrow build-tree' bs inp k i)

fun build-tree'-measure :: ('a bins \times 'a sentential \times nat \times nat) \Rightarrow nat **where**
 build-tree'-measure (bs, inp, k, i) = foldl (+) 0 (map length (take k bs)) + i

definition wf-tree-input :: ('a bins \times 'a sentential \times nat \times nat) set **where**
 wf-tree-input = {
 (bs, inp, k, i) | bs inp k i.
 sound-ptrs inp bs \wedge
 mono-red-ptr bs \wedge
 k < length bs \wedge
 i < length (bs!k)
 }

lemma wf-tree-input-pre:
assumes (bs, inp, k, i) \in wf-tree-input
assumes e = bs!k!i pointer e = Pre pre

shows $(bs, inp, (k-1), pre) \in wf-tree-input$

lemma *wf-tree-input-prered-pre*:

assumes $(bs, inp, k, i) \in wf-tree-input$

assumes $e = bs!k!i$ pointer $e = PreRed (k', pre, red) ps$

shows $(bs, inp, k', pre) \in wf-tree-input$

lemma *wf-tree-input-prered-red*:

assumes $(bs, inp, k, i) \in wf-tree-input$

assumes $e = bs!k!i$ pointer $e = PreRed (k', pre, red) ps$

shows $(bs, inp, k, red) \in wf-tree-input$

lemma *build-tree'-termination*:

assumes $(bs, inp, k, i) \in wf-tree-input$

shows $\exists N ts. build-tree' bs inp k i = Some (Branch N ts)$

lemma *wf-item-tree-build-tree'*:

assumes $(bs, inp, k, i) \in wf-tree-input$

assumes $wf-bins\ cfg\ inp\ bs$

assumes $k < length\ bs\ i < length\ (bs!k)$

assumes $build-tree' bs inp k i = Some\ t$

shows $wf-item-tree\ cfg\ (item\ (bs!k!i))\ t$

lemma *wf-ylid-tree-build-tree'*:

assumes $(bs, inp, k, i) \in wf-tree-input$

assumes $wf-bins\ cfg\ inp\ bs$

assumes $k < length\ bs\ i < length\ (bs!k)\ k \leq length\ inp$

assumes $build-tree' bs inp k i = Some\ t$

shows $wf-ylid-tree\ inp\ (item\ (bs!k!i))\ t$

theorem *wf-rule-root-ylid-tree-build-tree*:

assumes $wf-bins\ cfg\ inp\ bs\ sound-ptrs\ inp\ bs\ mono-red-ptr\ bs\ length\ bs = length\ inp + 1$

assumes $build-tree\ cfg\ inp\ bs = Some\ t$

shows $wf-rule-tree\ cfg\ t \wedge root-tree\ t = \mathfrak{S}\ cfg \wedge ylid-tree\ t = inp$

corollary *wf-rule-root-ylid-tree-build-tree-earley-list*:

assumes $wf-cfg\ cfg\ nonempty-derives\ cfg$

assumes $build-tree\ cfg\ inp\ (earley-list\ cfg\ inp) = Some\ t$

shows $wf-rule-tree\ cfg\ t \wedge root-tree\ t = \mathfrak{S}\ cfg \wedge ylid-tree\ t = inp$

theorem *correctness-build-tree-earley-list*:

assumes $wf-cfg\ cfg\ is-sentence\ cfg\ inp\ nonempty-derives\ cfg$

shows $(\exists t. build-tree\ cfg\ inp\ (earley-list\ cfg\ inp) = Some\ t) \longleftrightarrow derives\ cfg\ [\mathfrak{S}\ cfg]\ inp$

6.5 All Parse Trees

fun *insert-group* :: $('a \Rightarrow 'k) \Rightarrow ('a \Rightarrow 'v) \Rightarrow 'a \Rightarrow ('k \times 'v\ list)\ list \Rightarrow ('k \times 'v\ list)\ list$ **where**

insert-group $K\ V\ a\ [] = [(K\ a, [V\ a])]$

insert-group $K\ V\ a\ ((k, vs)\#xs) =$

 if $K\ a = k$ then $(k, V\ a\ \# vs)\ \# xs$

 else $(k, vs)\ \# insert-group\ K\ V\ a\ xs$

)

fun group-by :: ('a ⇒ 'k) ⇒ ('a ⇒ 'v) ⇒ 'a list ⇒ ('k × 'v list) list **where**
 group-by K V [] = []
 | group-by K V (x#xs) = insert-group K V x (group-by K V xs)

partial-function (option) build-trees' :: 'a bins ⇒ 'a sentential ⇒ nat ⇒ nat ⇒ nat set ⇒ 'a forest list option **where**

build-trees' bs inp k i I = (
 let e = bs!k!i in (
 case pointer e of
 Null ⇒ Some ([FBranch (item-rule-head (item e)) []])
 | Pre pre ⇒ (
 do {
 pres ← build-trees' bs inp (k-1) pre {pre};
 those (map (λf.
 case f of
 FBranch N fss ⇒ Some (FBranch N (fss @ [[FLeaf (inp!(k-1))]]))
 | - ⇒ None
) pres)
 })
 | PreRed p ps ⇒ (
 let ps' = filter (λ(k', pre, red). red ∉ I) (p#ps) in
 let gs = group-by (λ(k', pre, red). (k', pre)) (λ(k', pre, red). red) ps' in
 map-option concat (those (map (λ((k', pre), reds).
 do {
 pres ← build-trees' bs inp k' pre {pre};
 rss ← those (map (λred. build-trees' bs inp k red (I ∪ {red})) reds);
 those (map (λf.
 case f of
 FBranch N fss ⇒ Some (FBranch N (fss @ [concat rss]))
 | - ⇒ None
) pres)
 })
) gs))
)
))

definition build-trees :: 'a cfg ⇒ 'a sentential ⇒ 'a bins ⇒ 'a forest list option **where**

build-trees cfg inp bs ≡
 let k = length bs - 1 in
 let finished = filter-with-index (λx. is-finished cfg inp x) (items (bs!k)) in
 map-option concat (those (map (λ(-, i). build-trees' bs inp k i {i}) finished))

fun build-forest'-measure :: ('a bins × 'a sentential × nat × nat × nat set) ⇒ nat **where**
 build-forest'-measure (bs, inp, k, i, I) = foldl (+) 0 (map length (take (k+1) bs)) - card I

definition wf-trees-input :: ('a bins × 'a sentential × nat × nat × nat set) set **where**
 wf-trees-input = {
 (bs, inp, k, i, I) | bs inp k i I.
 sound-ptrs inp bs ∧
 k < length bs ∧
 i < length (bs!k) ∧
 I ⊆ {0..<length (bs!k)} ∧
 i ∈ I
 }

lemma wf-trees-input-pre:

assumes (bs, inp, k, i, I) ∈ wf-trees-input
assumes e = bs!k!i pointer e = Pre pre
shows (bs, inp, (k-1), pre, {pre}) ∈ wf-trees-input

lemma wf-trees-input-prered-pre:

assumes (bs, inp, k, i, I) ∈ wf-trees-input
assumes e = bs!k!i pointer e = PreRed p ps
assumes ps' = filter (λ(k', pre, red). red ∉ I) (p#ps)
assumes gs = group-by (λ(k', pre, red). (k', pre)) (λ(k', pre, red). red) ps'
assumes ((k', pre), reds) ∈ set gs
shows (bs, inp, k', pre, {pre}) ∈ wf-trees-input

lemma wf-trees-input-prered-red:

assumes (bs, inp, k, i, I) ∈ wf-trees-input
assumes e = bs!k!i pointer e = PreRed p ps
assumes ps' = filter (λ(k', pre, red). red ∉ I) (p#ps)
assumes gs = group-by (λ(k', pre, red). (k', pre)) (λ(k', pre, red). red) ps'
assumes ((k', pre), reds) ∈ set gs red ∈ set reds
shows (bs, inp, k, red, I ∪ {red}) ∈ wf-trees-input

lemma build-trees'-termination:

assumes (bs, inp, k, i, I) ∈ wf-trees-input
shows ∃fs. build-trees' bs inp k i I = Some fs ∧ (∀f ∈ set fs. ∃N fss. f = FBranch N fss)

lemma wf-item-tree-build-trees':

assumes (bs, inp, k, i, I) ∈ wf-trees-input
assumes wf-bins cfg inp bs
assumes k < length bs i < length (bs!k)
assumes build-trees' bs inp k i I = Some fs
assumes f ∈ set fs
assumes t ∈ set (trees f)
shows wf-item-tree cfg (item (bs!k!i)) t

lemma wf-forest-build-trees':

assumes (bs, inp, k, i, I) ∈ wf-trees-input

assumes $wf\text{-}bins\ cfg\ inp\ bs$
assumes $k < length\ bs\ i < length\ (bs!k)\ k \leq length\ inp$
assumes $build\text{-}trees'\ bs\ inp\ k\ i\ I = Some\ fs$
assumes $f \in set\ fs$
assumes $t \in set\ (trees\ f)$
shows $wf\text{-}yield\text{-}tree\ inp\ (item\ (bs!k!i))\ t$
theorem $wf\text{-}rule\text{-}root\text{-}yield\text{-}tree\text{-}build\text{-}trees$:
assumes $wf\text{-}bins\ cfg\ inp\ bs\ sound\text{-}ptrs\ inp\ bs\ length\ bs = length\ inp + 1$
assumes $build\text{-}trees\ cfg\ inp\ bs = Some\ fs\ f \in set\ fs\ t \in set\ (trees\ f)$
shows $wf\text{-}rule\text{-}tree\ cfg\ t \wedge root\text{-}tree\ t = \S\ cfg \wedge yield\text{-}tree\ t = inp$
corollary $wf\text{-}rule\text{-}root\text{-}yield\text{-}tree\text{-}build\text{-}trees\text{-}earley\text{-}list$:
assumes $wf\text{-}cfg\ cfg\ nonempty\text{-}derives\ cfg$
assumes $build\text{-}trees\ cfg\ inp\ (earley\text{-}list\ cfg\ inp) = Some\ fs\ f \in set\ fs\ t \in set\ (trees\ f)$
shows $wf\text{-}rule\text{-}tree\ cfg\ t \wedge root\text{-}tree\ t = \S\ cfg \wedge yield\text{-}tree\ t = inp$
theorem $soundness\text{-}build\text{-}trees\text{-}earley\text{-}list$:
assumes $wf\text{-}cfg\ cfg\ is\text{-}sentence\ cfg\ inp\ nonempty\text{-}derives\ cfg$
assumes $build\text{-}trees\ cfg\ inp\ (earley\text{-}list\ cfg\ inp) = Some\ fs\ f \in set\ fs\ t \in set\ (trees\ f)$
shows $derives\ cfg\ [\S\ cfg]\ inp$
theorem $termination\text{-}build\text{-}tree\text{-}earley\text{-}list$:
assumes $wf\text{-}cfg\ cfg\ nonempty\text{-}derives\ cfg\ derives\ cfg\ [\S\ cfg]\ inp$
shows $\exists fs.\ build\text{-}trees\ cfg\ inp\ (earley\text{-}list\ cfg\ inp) = Some\ fs$

6.6 A Word on Completeness

SNIPPET:

A shared packed parse forest SPPF is a representation designed to reduce the space required to represent multiple derivation trees for an ambiguous sentence. In an SPPF, nodes which have the same tree below them are shared and nodes which correspond to different derivations of the same substring from the same non-terminal are combined by creating a packed node for each family of children. Nodes can be packed only if their yields correspond to the same portion of the input string. Thus, to make it easier to determine whether two alternates can be packed under a given node, SPPF nodes are labelled with a triple (x,i,j) where $a_{j+1} \dots a_i$ is a substring matched by x . To obtain a cubic algorithm we use binarised SPPFs which contain intermediate additional nodes but which are of worst case cubic size. (EXAMPLE SPPF running example???)

We can turn earley's algorithm into a correct parser by adding pointers between items rather than instances of non-terminals, and labelling the pointers in a way which allows a binarised SPPF to be constructed by walking the resulting structure. However, in order to construct a binarised SPPF we also have to introduce additional nodes for grammar rules of length greater than two, complicating the final algorithm.

7 Usage

definition $\varepsilon\text{-free} :: 'a \text{ cfg} \Rightarrow \text{bool}$ **where**
 $\varepsilon\text{-free } \text{cfg} \longleftrightarrow (\forall r \in \text{set } (\mathfrak{N} \text{ cfg}). \text{rule-body } r \neq [])$

lemma $\varepsilon\text{-free-impl-non-empty-deriv}$:
 $\varepsilon\text{-free } \text{cfg} \Longrightarrow N \in \text{set } (\mathfrak{N} \text{ cfg}) \Longrightarrow \neg \text{derives } \text{cfg } [N] \quad []$

datatype $t = x \mid \text{plus}$

datatype $n = S$

datatype $s = \text{Terminal } t \mid \text{Nonterminal } n$

definition $\text{nonterminals} :: s \text{ list}$ **where**
 $\text{nonterminals} = [\text{Nonterminal } S]$

definition $\text{terminals} :: s \text{ list}$ **where**
 $\text{terminals} = [\text{Terminal } x, \text{Terminal plus}]$

definition $\text{rules} :: s \text{ rule list}$ **where**
 $\text{rules} = [$
 $(\text{Nonterminal } S, [\text{Terminal } x]),$
 $(\text{Nonterminal } S, [\text{Nonterminal } S, \text{Terminal plus}, \text{Nonterminal } S])$
 $]$

definition $\text{start-symbol} :: s$ **where**
 $\text{start-symbol} = \text{Nonterminal } S$

definition $\text{cfg} :: s \text{ cfg}$ **where**
 $\text{cfg} = \text{CFG nonterminals terminals rules start-symbol}$

definition $\text{inp} :: s \text{ list}$ **where**
 $\text{inp} = [\text{Terminal } x, \text{Terminal plus}, \text{Terminal } x, \text{Terminal plus}, \text{Terminal } x]$

lemma wf-cfg :

shows $\text{wf-cfg } \text{cfg}$

lemma is-sentence-inp :

shows $\text{is-sentence } \text{cfg } \text{inp}$

lemma nonempty-derives :

shows $\text{nonempty-derives } \text{cfg}$

lemma correctness :

shows *recognizing-list (earley-list cfg inp) cfg inp \longleftrightarrow derives cfg [\mathfrak{S} cfg] inp*

8 Conclusion

8.1 Summary

8.2 Future Work

Different approaches:

(1) SPPF style parse trees as in Scott et al -> need Imperative/HOL for this

Performance improvements:

(1) Look-ahead of k or at least 1 like in the original Earley paper. (2) Optimize the representation of the grammar instead of single list, group by production, ... (3) Keep a set of already inserted items to not double check item insertion. (4) Use a queue instead of a list for bins. (5) Refine the algorithm to an imperative version using a single linked list and actual pointers instead of natural numbers.

Parse tree disambiguation:

Parser generators like YACC resolve ambiguities in context-free grammars by allowing the user to specify precedence and associativity declarations restricting the set of allowed parses. But they do not handle all grammatical restrictions, like 'dangling else' or interactions between binary operators and functional 'if'-expressions.

Grammar rewriting:

Adams *et al* [Adams:2017] describe a grammar rewriting approach reinterpreting CFGs as the tree automata, intersectiong them with tree automata encoding desired restrictions and reinterpreting the results back into CFGs.

Afroozeh *et al* [Afroozeh:2013] present an approach to specifying operator precedence based on declarative disambiguation rules basing their implementation on grammar rewriting.

Thorup [Thorup:1996] develops two concrete algorithms for disambiguation of grammars based on the idea of excluding a certain set of forbidden sub-parse trees.

Parse tree filtering:

Klint *et al* [Klint:1997] propose a framework of filters to describe and compare a wide range of disambiguation problems in a parser-independent way. A filter is a function that selects from a set of parse trees the intended trees.

List of Figures

3.1 Earley inference rules	7
--------------------------------------	---

List of Tables

3.1	Earley items running example	7
5.1	Earley items with pointers running example	25