

A Drop-in Replacement for LR(1) Table-Driven Parsing

Michael Oudshoorn

Founding Dean, High Point University, School of Engineering

One University Parkway, High Point, NC 27268, USA

Email: {moudshoo@highpoint.edu}

Received on, 27 April 2021 - Accepted on, 10 July 2021 - Published on, 15 November 2021

ABSTRACT

This paper presents a construction method for a deterministic one-symbol look-ahead LR parser which allows non-terminals in the parser look-ahead. This effectively relaxes the requirement of parsing the reverse of the right-most derivation of a string/sentence. This is achieved by replacing the deterministic push down automata of LR parsing by a two-stack automata. The class of grammars accepted by the two-stack parser properly contains the LR(k) grammars. Since the modification to the table-driven LR parsing process is relatively minor and mostly impacts the creation of the *goto* and *action* tables, a parser modified to adopt the two-stack process should be comparable in size and performance to LR parsers.

Keywords: LR Parsing, Parsing Algorithm.

1. Introduction

The LR parser [1], characterized as a one-symbol look-ahead parser [2], is based on a deterministic push-down automata (DPDA). A strength of this technique is that the same parser engine can be used with parse tables generated by any of the increasingly general LR(0), SLR(1), LALR(1) and LR(1) algorithms. Parsing is limited to LR(1) by the parser mechanism. Although LR parsers accept the deterministic context-free languages (DCFLs) [3, 4] some context-free grammars (CFGs) for DCFLs require significant transformation before they can be used in LR parser construction. Seite [5] observes that such transformations are tedious and may reduce readability of the grammar, or generate a grammar that is not exactly equivalent to the original. However, a parser generator accepting a wider class of grammars evades many of these disadvantages by simply avoiding the need to transform the grammar.

This paper discusses a one-symbol look-ahead parser employing two stacks as a practical replacement for the LR parsing engine. The engine is based on a deterministic two-stack machine [6] and other two-stack (bottom-up) parsers [3, 7, 8]. A parser construction method for the two-stack machine which produces efficient, deterministic parsers for a broad class of grammars properly including LR(k) grammars is presented. For non- LR(k) grammars, the parser can defer parsing decisions while an arbitrary sequence of input symbols is processed.

Unbounded forward parsing activity is converted into non-terminals and is not repeated. This allows the parsing of a class of non- $LR(k)$ grammars using a single [terminal or non-terminal] symbol look-ahead.

Without formal proof, it is asserted that the generated parsers preserve the desirable properties of LR parsing, namely acceptance of a broad class of practical grammars, error detection at the first invalid symbol (hence guaranteed termination) [8] and $O(n)$ parsing efficiency (both time and space).

Section 2 introduces the terminology and notation used throughout the paper. Section 3 introduces an example that highlights some of the deficiencies of LR parser construction. Section 4 examines the two-stack parsing process and Section 5 presents the basic parser construction. Section 6 extends the parser generator for increased generality Results are presented in Section 7 and conclusions are offered in Section 8.

2. Terminology

Formally, a context-free grammar is a tuple $G = \langle N, T, P, S \rangle$ where N and T are disjoint sets of, respectively, non-terminal and terminal symbols (hence $N \cap T = \emptyset$), P is a finite set of productions, and $S \in N$ is the goal symbol. $V = N \cup T$ is defined as the vocabulary of the grammar. A production is a pair, written $A \rightarrow \alpha$, where $A \in N$ and $\alpha \in V^*$. The following conventions are used in this paper:

$$\begin{aligned} S, A, B, C \dots &\in N \\ \dots X, Y, Z &\in V \\ \alpha, \beta, \gamma \dots &\in V^* \\ a, b, c, \dots &\in T \\ \dots, x, y, z &\in T^* \end{aligned}$$

The notation, ε , represents an empty sequence and the length of a sequence α is $|\alpha|$, hence $|\varepsilon| = 0$.

The relation \Rightarrow , pronounced "directly produces", is defined on V^* such that $\alpha A \gamma \Rightarrow \alpha \beta \gamma$ where $A \rightarrow \beta \in P$. The relation closures \Rightarrow^* and \Rightarrow^+ , both pronounced "produces", are used to define the sentential forms, $\{\alpha \mid S \Rightarrow^* \alpha\}$, of the grammar G including the (probably infinite) language generated, $L(G) = \{x \mid S \Rightarrow^+ x\}$.

3. An exercise in parsing

Backus-Naur form (BNF) is a notation [9] describing the productions of a grammar. BNF is a small language that provides parsing problems for conventional parsing technologies.

The grammar G_1 [Figure 1] is a natural syntax definition of BNF notation without options in the right-hand side of productions. The language $L(G_1) = \text{BNF}$ is used

as the primary language example in this paper. In grammar G_1 , left recursion is used to describe a non-empty sequence S of productions. Each production P is a sequence $n ::= R$ where R is a possibly empty sequence of terminal t and non-terminal n symbols. A typical parsing system will distinguish terminal and non-terminal symbols in a separate lexical analysis phase.

An attempt to build a LALR(1) or LR(1) parser from this grammar fails, reporting a shift/reduce conflict. For some grammars, resolving such a conflict in favor of the shift action effectively disambiguates the grammar [10]. This strategy fails to produce a BNF parser based on G_1 because G_1 is not (suitably) ambiguous. BNF, though not strictly a programming language, is a simple and relevant language which shares characteristics with many computer languages. The failure of LR(1) parser generation for G_1 is disappointing given the expectations built up in the literature:

- “Most unambiguous grammars for programming languages are SLR(1) – hence LALR(1) and LR(1) – in the form in which they are originally written.” [11]
- “This class of grammars [LR(1)] has great importance for compiler design, since they are broad enough to include the syntax of almost all programming languages, yet restrictive enough to have efficient parsers that are essentially DPDA's.” [6]
- “Theoretically, LR(k) parsers are of interest in that they are the most powerful class of deterministic bottom-up parsers using at most k look-ahead symbols. ... if a grammar, G , can be parsed by any deterministic bottom-up parser using k look-ahead symbols, a LR(k) parser can be built for G .” [12]

$$\begin{aligned}
 (1) \quad S &\rightarrow S P \\
 (2) \quad S &\rightarrow P \\
 (3) \quad P &\rightarrow n ::= R \\
 (4) \quad R &\rightarrow \varepsilon \\
 (5) \quad R &\rightarrow R n \\
 (6) \quad R &\rightarrow R t
 \end{aligned}$$

Figure 1. Numbered productions of grammar G_1

More precisely, it is well known that all, and only, the deterministic CFL's have LR(1) grammars [6] While parsers can be developed for a wide class of languages, but they are not always based on a grammar that relates naturally to the semantics. Knuth [4] made it clear from the start that:

“... LR(k) is a property of the grammar, not of the language alone. The distinction between grammar and language is extremely important when semantics is being considered as well as syntax.”

BNF is an example of such a language. BNF is a regular language (RL), so $L(G_1) \in RL \subset DCFL$ and can be described by the equivalent regular expressions

in Figure 2. From the latter regular expression, $G_2 \in LR(1)$ where $L(G_1) \equiv L(G_2)$ (see Figure 3) can be constructed.

$$\begin{aligned} (n ::= (n / t)^*) + &\equiv n ::= (n / t)^* (n ::= (n / t)^*)^* \\ &\equiv n ::= ((n / t)^* n ::=)^* (n / t)^* \end{aligned}$$

Figure 2. Regular expressions for BNF.

$$\begin{aligned} S &\rightarrow P R \\ P &\rightarrow n ::= / P R n ::= \\ R &\rightarrow \varepsilon / R n / R t \end{aligned}$$

Figure 3. Grammar $G_2 \in LR(1)$ and $L(G_2) \equiv$ BNF.

The link between the syntax and the semantics of G_2 is less obvious than with G_1 : namely, each R is associated with the last $n ::=$ in the preceding P . Grammar G_2 is arguably a less natural expression, albeit an equivalent, of the grammar for BNF.

A language designer may also alter the language as well as the grammar to suit a class of parser generator, e.g. by adding terminators and/or separators [13, 14] as in grammars G_{3a} , G_{3b} and G_{3c} (Figure 4). This syntactic notation is reminiscent of statement syntax in programming languages such as Algol-9 [9], C [15], Pascal [16], Modula-2 [17], Modula-3 [18], Ada [19, 20], etc. It is also argued that such notation aids error recovery and is simply good language design practice. Other languages, e.g. BCPL [21] and Haskell [22, 23], adopt this solution but try to use layout to mask the consequences.

$\begin{aligned} S &\rightarrow S ; P / P \\ P &\rightarrow n ::= R \\ R &\rightarrow \varepsilon / R n / R t \end{aligned}$	$\begin{aligned} S &\rightarrow S P / P \\ P &\rightarrow n ::= R ; \\ R &\rightarrow \varepsilon / R n / R t \end{aligned}$	$\begin{aligned} S &\rightarrow S' semi_{opt} \\ S' &\rightarrow S' ; P / P \\ P &\rightarrow n ::= R \\ R &\rightarrow \varepsilon / R n / R t \\ semi_{opt} &\rightarrow \varepsilon / ; \end{aligned}$
G_{3a} : separator	G_{3b} : terminator	G_{3c} : separator/terminator

Figure 4. Grammars for BNF with terminators and/or separators.

Since G_1 cannot be parsed using $LR(1)$ parsers, the more general $LR(k)$ techniques can be considered. This immediately meets with success since $G_1 \in LR(2)$. An experienced language design could manipulate the grammar to produce grammar $G_{4a} \in LR(1)$ (Figure 5) by treating $n ::=$ as the single lexical item, *lhs*. Alternatively, using the forward context capabilities of a lexical analysis tool such as [24, 25], the symbol n followed by $::=$ can be distinguished from n followed by some other symbol, before it reaches the parser.

A general parser generator could allow the separation of a language into lexical and syntactic elements as a matter of style or taste rather than necessity. A comparison of G_{4a} with G_2 suggests that the lexical/syntactic division can be important in the preservation of semantic structure during grammar development. Grammar G_{4b} (Figure 5) for Algol-60 BNF [9] is based on another lexical/syntactic separation and is as "natural" as G_1 . Unfortunately, G_{4b} is not LR(2) nor even LR(k); from such a starting point, a lexical solution to the parsing problem may be more difficult to discover.

$S \rightarrow SP / P$ $P \rightarrow lhs R$ $R \rightarrow \varepsilon / R n / R t$	$S \rightarrow SP / P$ $P \rightarrow N ::=$ $P \rightarrow P N / P t$ $N \rightarrow < C >$ $C \rightarrow c / C c$
--	--

Grammar $G_{4a} \in LR(1)$. Grammar $G_{4b} \in LR(k)$.

Figure 5. Alternate syntax/lexical separations.

The essence of this difficulty is that LR grammars have poor closure behavior under homomorphism [26,6] and substitution with regular expressions. Syntactic difficulties detected in parser generation can be treated as semantic issues or transferred to the lexical analyzer. The limited capabilities of most lexical analyzers and the increased costs (space and/or time) diminish the applicability of the latter strategy. The two-stack parser generator presented in this paper is less sensitive to this division of responsibility between the lexical analysis and parsing modules.

Grammars G_{5a} and G_{5b} (Figure 6) represents grammars that can be parsed with a two-stack mechanism where $L(G_{5a}) \equiv L(G_{5b}) \equiv BNF$. LR(1) parser generators, working exclusively with terminal-symbol look-ahead fail for G_{5a} and G_{5b} . With G_{5b} , the sequence $R n$ can be derived from either R or $R L$ resulting in a non-deterministic parser. Further, the LR parser cannot reverse a derivation step, $P L \Rightarrow L ::= R L$, as the trailing L would be overwritten in the parsing engine. This illustrates the differences between LR parsing and the two-stack technique. By using non-terminal look-ahead, the two-stack parser can reverse both G_{5b} derivation sequences:

$S \rightarrow P S / P$ $P \rightarrow n ::= R$ $R \rightarrow \varepsilon / R n / R t$	$S \rightarrow S P / P$ $P \rightarrow L ::= R$ $R \rightarrow \varepsilon / R n / R t$ $L \rightarrow n$
Grammar G_{5a}	Grammar G_{5b}

Figure 6. BNF grammars for two-stack parsing.

1. $P L ::= \dots \Rightarrow L ::= R L ::= \dots \Rightarrow L ::= R n ::= \dots$
2. $P \dots \Rightarrow ::= R \dots \Rightarrow L ::= R n \dots \Rightarrow L ::= R n n \dots$

The technique described in Section 5 develops a parser using two-stack automata via a construction generalized from LR parser generation. The parser generator presented fails for G_1 , but succeeds for G_{5a} , G_{5b} and G_6 (Figure 7). Each of these grammars might be described as natural relative to G_2 . Grammar G_6 , like G_{4b} , is a BNF grammar without empty productions, and is the basis for the example two-stack BNF parser developed in Sections 4 and 5.

- (1) $S \rightarrow S P$
- (2) $S \rightarrow P$
- (3) $P \rightarrow n ::=$
- (4) $P \rightarrow P n$
- (5) $P \rightarrow P t$

Figure 7. Grammar G_6 with numbered productions.

The parser generator described below is more general than the LR(1) algorithm and more practical than for the deterministic regular parsable (DRP) grammars [8]. The two-stack parser is simpler than parsers for other developments from LR(1), e.g. LR(k) and LR(k, t) [4], LRR [26, 5], LAR(m) [28] and DRP.. Figure 8 shows the capabilities of the parsers for the BNF-describing grammars given above.

	G_1	G_2	G_{4a}	G_{4b}	G_{5a}	G_{5b}	G_6
LR(1)		√	√				
LR(2)/LR(k)	√	√	√		√	√	√
Two-stack		√	√	√	√	√	√

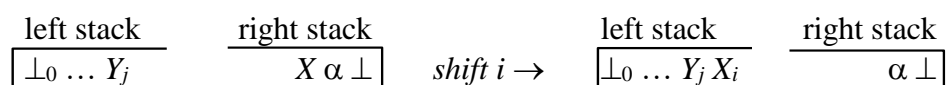
Figure 8. Grammar class membership.

4. A Two-Stack Parser

The two-stack parser has a left-stack containing integer-subscripted symbols and a right-stack containing symbols. The set of symbols, V , comes from a grammar and the integers from the set, K , of node identifiers from the parser construction. The stacks are delimited by a special \perp symbol and are illustrated below with stack-tops toward the center.

The parser has four actions:

- *Shift i* where $i \in K$, pops a symbol, X , from the right-stack and pushes the subscripted symbol, X_i , into the left-stack. Pictorially:



- *Reduce $\alpha \rightarrow \beta$* pops β' , an annotated version of β , from the left-stack and pushes into the right-stack. The notation *Reduce $n + m$* is used for *Reduce $A\mu \rightarrow \alpha\mu$* where $P_n = A \rightarrow \alpha$ and $m = |\mu|$, or simply *Reduce n* when $\mu = \epsilon$ ($m = 0$).

$$\begin{array}{c} \text{left stack} \\ \boxed{\perp_0 \dots X_I \alpha' \mu'} \end{array} \quad \begin{array}{c} \text{right stack} \\ \boxed{\gamma \perp} \end{array} \quad \text{reduce } n + m \rightarrow \begin{array}{c} \text{left stack} \\ \boxed{\perp_0 \dots X} \end{array} \quad \begin{array}{c} \text{right stack} \\ \boxed{A \mu \gamma \perp} \end{array}$$

- *Accept* causes successful termination.
- *Error* causes failure termination.

The parser uses the subscript and symbol at the top of the stacks to access the parse table, a function $PT: K \times V \rightarrow action$, to determine the action to be taken at each step. This table is represented with *Error* entries left blank, e.g. Figure 9 shows a parser table derived from grammar G_6 .

Node	\perp	n	t	$::=$	S	P
0		S3			S2	S1
1	R2	S4	S5			R2
2	Acc	S3				S6
3				S8		
4	R4	R4	R4	S7		R4
5	R5	R5	R5			R5
6	R1	S4	S5			R1
7	R3	R3	R3			R1
8	R3	R3	R3			

Figure 9. Parse table for grammar G_6 .

Initially, the left-stack holds the \perp_0 symbol and the right-stack holds the input string followed by the \perp marker. Valid input is accepted in a state where the left-stack contains the goal symbol and the right-stack is empty.

$$\begin{array}{c} \text{left stack} \\ \boxed{\perp_0} \end{array} \quad \begin{array}{c} \text{right stack} \\ \boxed{\text{input } \perp} \end{array} \quad \text{action} \rightarrow \dots$$

$$\rightarrow \text{action} \rightarrow \begin{array}{c} \text{left stack} \\ \boxed{\perp_0 S_i} \end{array} \quad \begin{array}{c} \text{right stack} \\ \boxed{\perp} \end{array} \quad \text{Accept}$$

With invalid input, the parser reaches an *Error* action. The complete sequence of parser actions for a short string from $L(G_6)$ is shown in Figure 10.

Provided the left-stack is correctly conditioned for each reduce (as in LR parsing), partial correctness of the parser can be shown by reverse induction. At termination, the stacks contain a \perp delimited sentential form. Shift preserves a sentential form in the stacks, hence it can be ignored. Reduce is effectively the inverse of \Rightarrow ; a *reduce* action resulting in a sentential form necessarily starts from a sentential form. A successful parse, a sequence of *shift* and *reduce* actions terminating with *accept*, must commence with a sentential form in the stacks. Responsibility for ensuring that the parser moves forward from an accessible parser state rests with the parser construction method.

<i>left stack</i>	<i>right stack</i>	<i>action</i>
\perp_0	$n ::= n t n ::= n \perp$	Shift 3
$\perp_0 n_3$	$::= n t n ::= n \perp$	Shift 8
$\perp_0 n_3 ::= 7$	$n t n ::= n \perp$	Reduce 3
\perp_0	$P n t n ::= n \perp$	Shift 1
$\perp_0 P_1$	$n t n ::= n \perp$	Shift 4
$\perp_0 P_1 n_4$	$t n ::= n \perp$	Reduce 4
\perp_0	$P t n ::= n \perp$	Shift 1
$\perp_0 P_1$	$t n ::= n \perp$	Shift 5
$\perp_0 P_1 t_5$	$n ::= n \perp$	Reduce 5
\perp_0	$P n ::= n \perp$	Shift 1
$\perp_0 P_1$	$n ::= n \perp$	Shift 4
$\perp_0 P_1 n_4$	$::= n \perp$	Shift 7
$\perp_0 P_1 n_4 ::= 7$	$n \perp$	Reduce 3
$\perp_0 P_1$	$P n \perp$	Reduce 2
\perp_0	$S P n \perp$	Shift 2
$\perp_0 S_2$	$P n \perp$	Shift 6
$\perp_0 S_2 P_6$	$n \perp$	Shift 4
$\perp_0 S_2 P_6 n_4$	\perp	Reduce 3
$\perp_0 S_2$	$P \perp$	Shift 6
$\perp_0 S_2 P_6$	\perp	Reduce 1
\perp_0	$S \perp$	Shift 2
$\perp_0 S_2$	\perp	Accept

Figure 10. Parsing $n ::= n t n ::= n$ with grammar G_6 .

5. Parse table construction

A parser is constructed from a grammar augmented with the parser symbols $\perp \in T$ and $S' \in N$ used in a new production $P_0 = S' \rightarrow \perp S \in P$. The following relations based on the grammar are defined:

$$\begin{aligned}
 FRONT(A\alpha) &= FRONT(A) = \{X \mid A \Rightarrow X\beta\} \\
 FRONT_0(X\alpha) &= \{X\} \\
 FRONT(\epsilon) &= \emptyset \\
 FWD^+(\alpha) &= \cup\{X \mid \alpha \Rightarrow^+ X\beta\} \\
 FWD^*(\alpha) &= \cup\{X \mid \alpha \Rightarrow^* X\beta\} \\
 FIRST_k(\alpha) &= \begin{cases} \epsilon, k = 0 \\ \alpha, \alpha = \perp \vee \alpha = \epsilon \\ x, \alpha \Rightarrow^* \alpha\beta \wedge x \in \alpha FIRST_{k-1}(\beta) \end{cases}
 \end{aligned}$$

An item is a tuple $\langle n, \mu, \beta, \gamma \rangle$, alternatively represented as $A\mu \rightarrow \alpha \bullet \beta; \gamma$, where $P_n = A \rightarrow \delta \in P$, $\delta\mu = \alpha\beta$, and $\gamma \in V^*$ is an element of a [finite] set of contexts. Informally, an item reads: having found α , look for β and [an initial segment of] context γ before reducing $\alpha\beta$ to $A\mu$. An initial item has the form $A\mu \rightarrow \bullet \delta\mu; \gamma$ or $\langle n, \mu, \delta, \gamma \rangle$ where $P_n = A \rightarrow \delta$. A final item has the form $A\mu \rightarrow \delta\mu \bullet; \gamma$ or $\langle n, \mu, \epsilon, \gamma \rangle$. The construction described in this section maintains $\mu = \epsilon$; the construction in Section 6 produces non-empty values of μ . Items with differing contexts may be displayed as a group, i.e:

$$A\mu \rightarrow \alpha \bullet \beta; \gamma, \dots, A\mu \rightarrow \alpha \bullet \beta; \delta \equiv A\mu \rightarrow \alpha \bullet \beta; \gamma, \dots, \delta$$

Initially, a directed graph (digraph) is constructed where nodes are sets of items and edges are labelled with symbols. The nodes in the digraph are uniquely numbered. The initial node is the singleton item set, $I_0 = \{S' \rightarrow \perp \bullet S; \perp\} \equiv \{<0, \varepsilon, S, \perp>\}$. The digraph for G_6 is shown in Figure 11.

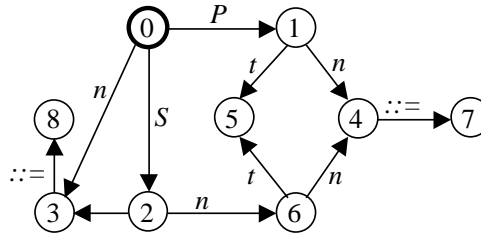


Figure 11. Digraph for grammar G_6 .

The parser table develops from the digraph. Each edge in the graph represents a shift action in the table. The reduce and accept entries in the table are based on the final items at a node. The node items for grammar G_6 are shown in Figure 12.

Node	Node items
0	$S' \rightarrow \perp \bullet S; \perp$
1	$S \rightarrow P \bullet; \perp, P \perp, P n$ $P \rightarrow P \bullet n; \perp, P \perp, P n, n, t$ $P \rightarrow P \bullet t; \perp, P \perp, P n, n, t$
2	$S' \rightarrow \perp S \bullet; \perp$ $S \rightarrow S \bullet P; \perp, P \perp, P n$
3	$P \rightarrow n \bullet ::=; \perp, P \perp, P n, n, t$
4	$P \rightarrow n \bullet ::=; \perp, n, t$ $P \rightarrow P n \bullet; \perp, P \perp, P n, n, t$
5	$P \rightarrow P t \bullet; \perp, P \perp, P n, n, t$
6	$S \rightarrow S P \bullet; \perp, P \perp, P n$ $P \rightarrow P \bullet n; \perp, P \perp, P n, n, t$ $P \rightarrow P \bullet t; \perp, P \perp, P n, n, t$
7	$P \rightarrow n ::= \bullet; \perp, P \perp, P n, n, t$
8	$P \rightarrow n ::= \bullet; \perp, n, t$

Figure 12. Item sets of nodes from grammar G_6 .

5.1 Digraph Construction

Each node is uniquely numbered and represents a set of items similar to basis items in LR parsing literature. The digraph is constructed by starting from node I_0 and constructing edges and nodes until the reachable nodes have been identified.

A node I_n has a (possibly empty) set of shift symbols, $SS(I_n)$. Two groups of symbols contribute to the shift symbols. Basic shift symbols, $BS(I_n)$, are derived from non-final items and are the symbols that must be shifted by the parser to progress towards a reduce action. Reduce conflict symbols, $RC(I_n)$, are symbols which lead to multiple reduce entries, i.e. a reduce/reduce conflict, in an LR parser generation. Two-stack parser construction treats these symbols as shift symbols.

In two-stack parsing, the symbols in $RC(I_n)$ are shifted with the expectation that the conflict will be resolved once a non-terminal is recognized. $RC(I_n)$ is the set of symbols which appear in more than one reduce symbol set for an item group, $RS(I_n, n, \mu, \alpha)$. Consider a node containing two final items, $\langle i, \mu, \varepsilon, \alpha \rangle$, $\langle i, \mu, \varepsilon, \alpha \rangle$, $\langle j, \eta, \varepsilon, \beta \rangle \in I_k$, where $\alpha \Rightarrow^* X \gamma \Rightarrow Z \dots$ and $\beta \Rightarrow^* Y \delta \Rightarrow Z \dots$. The parser generator avoids potential parse table conflicts:

$$\begin{aligned} \{reduce\ i, reduce\ j\} &\subseteq PT(I_k, W) \\ \text{where } W &\in FWD^*(Z) \subset RS(I_k, i, \mu, \varepsilon) \cap RS(I_k, j, \eta, \varepsilon) \\ \text{hence } FWD^*(Z) &\subseteq RC(I_k) \subseteq SS(I_k) \end{aligned}$$

As a result, items $X \rightarrow \bullet Z \dots; \gamma$ and $Y \rightarrow \bullet Z \dots; \delta$ and their derivatives are included in the node closure. The result is a parser which tries to recognize either X or Y starting with Z . Having reduced X or Y onto the top of the right-stack, the parser is able to decide between reduce i and reduce j .

The same strategy deals with conflicts between shift and reduce actions. The presence of empty productions, productions of the form $A \rightarrow \varepsilon$, in the grammar requires that non-final items be included in the definition of RC .

$$\begin{aligned} SS(I_n) &= BS(I_n) \cup RC(I_n) - \{\perp\} \\ BS(I_n) &= \cup \{FRONT^*(\alpha) \mid \langle i, \mu, \alpha, \beta \rangle \in I_i\} \\ RC(I_n) &= \cup \{RS(I_n, i, \mu, \alpha) \cap RS(I_n, j, \eta, \beta) \mid \langle i, \mu, \alpha \rangle \neq \langle j, \eta, \beta \rangle\} \\ RS(I_i, n, \mu, \alpha) &= \cup \{FWD^*(\alpha \beta) - FRONT^*(\alpha) \mid \langle n, \mu, \alpha, \beta \rangle \in I_i\} \end{aligned}$$

Every digraph node has a labelled edge leaving it for each shift symbol associated with that node, hence:

$$\begin{aligned} Edges(I_n) &= \{\langle I_n, X, Next(I_n, X) \rangle \mid X \in SS(I_n)\} \\ Next(I_i, X) &= \{\langle n, \mu, \alpha, \beta \rangle \mid \langle n, \mu, X \alpha, \beta \rangle \in Close(I_i)\} \end{aligned}$$

Each node has a (possibly empty) set of closure items. These are computed by recursively extending the set of items that may be applicable in the (local) parsing position represented by the node. A complete set of closure items is identified, prior to the relevant items being selected, as this simplifies development of item contexts.

$$\begin{aligned}
Close(I_n) &= I_n \cup Relevant(SS(I_n), Derive^+(I_n)) \\
Derive(I_i) &= \cup \left\{ \begin{array}{l} CI(A, \delta), \langle n, \mu, \alpha A \beta, \gamma \rangle \in I_i \wedge \alpha \Rightarrow^* \varepsilon \wedge \delta \in Context(\beta, \gamma) \\ CI(A, \gamma, \delta), \langle n, \mu, \alpha A \beta, \beta A \gamma \rangle \in I_i \wedge \alpha \beta \Rightarrow^* \varepsilon \end{array} \right. \\
CI(A, \beta) &= \{ \langle n, \varepsilon, \alpha, \rangle \mid P_n = A \rightarrow \alpha \} \\
Context(\alpha, \beta) &= \left\{ \begin{array}{l} \beta, \alpha = \varepsilon \\ \gamma \alpha, \gamma \alpha \delta = \alpha b \wedge \gamma \in N^* \wedge b \in FIRST_1(\beta) \end{array} \right. \\
Relevant(R, Q) &= \left\{ \begin{array}{l} i, i = A \rightarrow \alpha \bullet X \beta; \gamma \in Q \wedge X \in R \\ i, i = X \rightarrow \bullet; \alpha \in Q \wedge X \in R \end{array} \right.
\end{aligned}$$

The Derive function adds initial items which may be applicable in the parsing position represented by an item set I_i . Closure items derive from all non-terminals immediately proceeded by the mark (\bullet) in the remaining item right-hand sides and contexts. Once these items, with their contexts, have been derived, Relevant items are selected. Final items lead to reduce entries in the parse table (see below), non-final items contribute to the items of other nodes (via Next in the edge construction). Figure 13 shows the digraph construction for grammar G_6 . Figure 14 shows a computed closure for node I_1 of grammar G_6 and indicates which items are relevant to parser construction.

For a given grammar, the set of possible items is finite since all items $A \varepsilon \rightarrow \alpha \bullet \beta; \gamma$ are constructed from a finite production $A \rightarrow \alpha \beta \in P$ where P is a finite set, and a m -bounded context $\gamma \in V^* \times T$ where $|\gamma| < m = \max\{|\alpha| \mid A \rightarrow \alpha \in P\}$. Since the items are finite, the nodes are finite and a finite algorithm for digraph construction exists.

5.2 Filling the Parse Table

Parser generation is completed by filling in the parse table, PT . A table entry is determined for each node I_i and symbol X .

$$PT(I_i, X) = \begin{cases} Shift j, & \langle I_i, X, I_j \rangle \in Edges(I_i) \\ Reduce n + |\mu|, \langle n, \mu, \varepsilon, X \alpha \rangle \in Close(I_i) \\ Reduce n + |\mu|, & X \in RS'(I_i, n, \mu) - SS(I_i) \end{cases}$$

$$RS'(I_i, n, \mu) = \cup \{ FWD^+(\gamma) \mid \langle n, \mu, \varepsilon, \gamma \rangle \in Close(I_i) \}$$

The single instance of *Reduce 0* (at $PT(Next(0, S), \perp)$ in the table) is replaced with *Accept*. This is consistent with the previous definition of *Accept* for the two-stack parser. Undefined entries are interpreted as *Error*. A table with at most one entry

in each table position is a deterministic parser. This is the case for a significant class of grammars including LR(1) as is shown informally in Section 6.

Node	Node items	Closure items	Edges
0	$S' \rightarrow \perp \bullet S; \perp$	$S \rightarrow \bullet S P; \perp, P \perp, P n$ $S \rightarrow \bullet P; \perp, P \perp, P n$ $P \rightarrow \bullet n ::=; \perp, P \perp, P n, n, t$ $P \rightarrow \bullet P n; \perp, P \perp, P n, n, t$ $P \rightarrow \bullet P t; \perp, P \perp, P n, n, t$	$\langle 0, P, 1 \rangle$ $\langle 0, S, 2 \rangle$ $\langle 0, n, 3 \rangle$
1	$S \rightarrow P \bullet; \perp, P \perp, P n$ $P \rightarrow P \bullet n; \perp, P \perp, P n, n, t$ $P \rightarrow P \bullet t; \perp, P \perp, P n, n, t$	$P \rightarrow \bullet n ::=; \perp, n, t$	$\langle 1, n, 4 \rangle$ $\langle 1, t, 5 \rangle$
2	$S' \rightarrow \perp S \bullet; \perp$ $S \rightarrow S \bullet P; \perp, P \perp, P n$	$P \rightarrow \bullet n ::=; \perp, P \perp, P n, n, t$ $P \rightarrow \bullet P n; \perp, P \perp, P n, n, t$ $P \rightarrow \bullet P t; \perp, P \perp, P n, n, t$	$\langle 2, n, 3 \rangle$ $\langle 2, P, 6 \rangle$
3	$P \rightarrow n \bullet ::=; \perp, P \perp, P n, n, t$		$\langle 3, ::=, 8 \rangle$
4	$P \rightarrow n \bullet ::=; \perp, n, t$ $P \rightarrow P n \bullet; \perp, P \perp, P n, n, t$		$\langle 4, ::=, 7 \rangle$
5	$P \rightarrow P t \bullet; \perp, P \perp, P n, n, t$		
6	$S \rightarrow S P \bullet; \perp, P \perp, P n$ $P \rightarrow P \bullet n; \perp, P \perp, P n, n, t$ $P \rightarrow P \bullet t; \perp, P \perp, P n, n, t$	$P \rightarrow \bullet n ::=; \perp, n, t$	$\langle 6, n, 4 \rangle$ $\langle 4, t, 5 \rangle$
7	$P \rightarrow n ::= \bullet; \perp, P \perp, P n, n, t$		
8	$P \rightarrow n ::= \bullet; \perp, n, t$		

Figure 13. Digraph and closure from grammar G_6 .

	Items	Relevant
node I_1	$S \rightarrow P \bullet; \perp, P \perp, P n$ $P \rightarrow P \bullet n; \perp, P \perp, P n, n, t$ $P \rightarrow P \bullet t; \perp, P \perp, P n, n, t$	√
$Derive(I_1)$	$P \rightarrow \bullet n ::=; \perp, n$ $P \rightarrow \bullet P n; \perp, n$ $P \rightarrow \bullet P t; \perp, n$ $P \rightarrow \bullet n ::=; t$	

Figure 14. Closure construction for grammar G_6 , node I_1 .

6. Extended two-stack parsing: using linear backup

The parser generator above will fail to produce a deterministic parser when a node contains items such that:

$$\begin{aligned}
 A \mu &\rightarrow \alpha \bullet; \beta \\
 C \eta &\rightarrow \gamma \bullet; X \delta
 \end{aligned}$$

where $\beta \Rightarrow^* X \dots$. The algorithm fails with a shift/reduce conflict when $\beta \Rightarrow^+ X \dots$ and a reduce/reduce conflict when $\beta = X \dots$.

This occurs in grammar G_{7a} (Figure 15), a LR(0, 2) grammar [4]. The basic two-stack parser generator, described in Section 5, produces a parser (Figure 16) from the modified G_{7b} grammar. The grammar transformation from G_{7a} to G_{7b} introduces a form of backup. The amount of backup introduced is proportional to the length of the input so parsing remains efficient.

$ \begin{aligned} S &\rightarrow A C c \mid B C d \\ A &\rightarrow a \\ B &\rightarrow a \\ C &\rightarrow C b \mid b \\ C_1 &\rightarrow C \\ C_2 &\rightarrow C \end{aligned} $ <p>Grammar G_{7a}</p>	$ \begin{aligned} S &\rightarrow A C_1 c \mid B C_2 d \\ A &\rightarrow a \\ B &\rightarrow a \\ C &\rightarrow C b \mid b \end{aligned} $ <p>Grammar G_{7b}</p>
--	--

Figure 15. LR(0, 2) grammars.

Node	\perp	a	b	c	d	S	A	B	C	C_1	C_2
0		S4				S3	S1	S2			
1			S7						S5	S6	
2			S10						S8		S9
3	Acc										
4			S12						S11	R3	R4
5			S13	R7							
6				S14							
7			R6	R6							
8			S15		R8						
9					S16						
10			R6		R6						
11			S17	R7	R8						
12			R6	R6	R6						
13			R5	R5							
14	R1										
15			R5		R5						
16	R2										
17			R5	R5	R5						

Figure 16. Parser table for grammar G_{7b} .

When applicable, a comparable transformation is performed automatically by the parser generator. The following modification to the parser generator adheres to the

policy of extending node closure allowing the parser to shift conflicting symbols until additional input indicates which reduce action is appropriate.

For grammar G_{7a} , immediate context is insufficient information on which to base parsing decisions. With extended context, a decision can be reached. The primary change to the parser generator is the additional *Extend* step in the new closure construction. By replacing *Close* with *Close'* in *Next*, a parse table (Figure 17) can be constructed.

$$Close'(I_n) = \{Extend(i, SS(I_n)) \mid i \in Close(I_n)\}$$

$$Extend(i, SS) = \begin{cases} A\mu X \rightarrow \alpha \bullet X; \beta\alpha, & i = A\mu \rightarrow \alpha \bullet; X\beta\alpha \wedge X \in SS \\ i, & \text{otherwise} \end{cases}$$

Node	\perp	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>S</i>	<i>A</i>	<i>B</i>	<i>C</i>
0		S4				S3	S1	S2	
1			S6						S5
2			S8						S7
3	Acc								
4			S10						S9
5			S11	S12					
6			R6	R6					
7			S13		S14				
8			R6		R6				
9			S15	R3+1	R4+1				
10			R6	R6	R6				
11			R5	R5					
12	R1								
13			R5		R5				
14	R2								
15			R5	R5	R5				

Figure 17. Parse table from grammar G_{7a} .

Figure 18 shows the nodes and the closure for the digraph nodes where *Extend* has an impact. Though this mechanism may introduce items which appear to be context sensitive, all the reduce actions are based on the context free rules of the grammar. The generator uses items with extended context to develop determinism during parsing.

This construction may produce redundant or unreachable nodes, e.g. nodes 6, 8, 11 and 13 in Figure 17 cannot be used since a sequence of *b*'s will have already been recognized as (or reduced to) a *C* when the parser reaches nodes 1, 2, 5 and 7. The set of items constructed by *Extend* is finite since each item $\langle n, \epsilon, \alpha, \beta\alpha \rangle$ introduced by *Derive* can be *Extended* to at most $|\beta|$ different items. Combined with node merging, this construction produces parsers comparable to other published techniques [29].

Node	Items	Close'
4	$A \rightarrow a \bullet; C c$ $B \rightarrow a \bullet; C d$	$A C \rightarrow a \bullet C; c$ $A C \rightarrow a \bullet C; d$ $C \rightarrow \bullet C b; b, c, d$ $C \rightarrow \bullet b; b, c, d$
9	$A C \rightarrow a C \bullet; c$ $B C \rightarrow a C \bullet; d$ $C \rightarrow C \bullet b; b, c, d$	$A C \rightarrow a C \bullet; c$ $B C \rightarrow a C \bullet; d$ $C \rightarrow C \bullet b; b, c, d$

Figure 18. Close' for selected nodes from the G_{7a} digraph.

7. Discussion

The parser construction technique described above is more complex than LR parser construction. Extra effort goes into constructing items, some of which eventually lead to non-terminal look-ahead in the parser, while others are not Relevant and are subsequently discarded. The closure items in nodes 1 and 6 of Figure 12 are examples of items that lead to non-terminal look-ahead.

The result of this extra work is the ability to resolve some of the conflicts that arise in LR parser generation. While the amount of effort has not yet been quantified, the overhead is anticipated to be relatively small in general and the method is expected to be practical.

7.1 Other Example Languages

Grammar G_{8a} (Figure 19) is a relatively natural grammar which defines a series of data records with option fields. A parse table and its derivation from G_{8a} are shown in Figures 20 and 21, respectively. Like BNF, $L(G_{8a})$ is a regular language, and hence a DCFL. In practice, the language is usually changed for LR parsing: either using a field terminator (rather than an initiator), or adding semantic constraints to an unordered syntax. Grammars G_{8b} and G_{8c} (Figure 19) are LR grammars for this language.

$S \rightarrow \epsilon \mid S R$	$S \rightarrow \epsilon \mid S R$	$S \rightarrow \epsilon \mid S R$
$R \rightarrow \text{hdr } f_1 f_2 f_3$	$R \rightarrow \text{hdr } F_1$	$R \rightarrow \text{hdr} \mid F_1 \mid F_2 \mid F_3$
$f_1 \rightarrow \epsilon \mid \text{sep data}_1$	$F_1 \rightarrow F_2 \mid \text{sep data}_1 F_2$	$F_1 \rightarrow \text{hdr sep data}_1$
$f_2 \rightarrow \epsilon \mid \text{sep data}_2$	$F_2 \rightarrow F_2 \mid \text{sep data}_2 F_3$	$F_2 \rightarrow \text{hdr sep data}_2 \mid F_1 \text{ sep data}_2$
$f_3 \rightarrow \epsilon \mid \text{sep data}_3$	$F_3 \rightarrow \epsilon \mid \text{sep data}_3$	$F_3 \rightarrow \text{hdr sep data}_3 \mid F_1 \text{ sep data}_3 \mid F_2 \text{ sep data}_3$
Grammar G_{8a}	Grammar G_{8b}	Grammar G_{8c}

Figure 19. Syntax records with optional fields.

Node	\perp	sep	hdr	data ₁	data ₂	data ₃	S	R	f ₁	f ₂	f ₃
0	R1		R1				S1	R1			
1	Acc		S3					S2			
2	R2		R2					R2			
3	R4	S5	R4						S4	R4	R4
4	R6	S7	R6							S6	R6
5				S8	S9	S10					
6	R8	S12	R8					R8			S11
7					S9	S10					
8	R5	R5	R5							R5	R5
9	R7	R7	R7								R7
10	R9		R9					R9			
11	R3		R3					R3			
12						S10					

Figure 20. Parse table for grammar G_{8a} .

Node	Node items	Closure items
0	$S' \rightarrow \perp \bullet S; \perp$	$S \rightarrow \bullet; \perp, R \perp, R \text{ hdr}$ $S \rightarrow \bullet S R; \perp, R \perp, R \text{ hdr}$
1	$S' \rightarrow \perp S \bullet; \perp$ $S \rightarrow S \bullet R; \perp, R \perp, R \text{ hdr}$	$R \rightarrow \bullet \text{ hdr } f_1 f_2 f_3; \perp, R \perp, R \text{ hdr}$
2	$S \rightarrow S R \bullet; \perp, R \perp, R \text{ hdr}$	
3	$R \rightarrow \text{hdr} \bullet f_1 f_2 f_3; \perp, R \perp, R \text{ hdr}$	$f_1 \rightarrow \bullet; f_2 f_3 \perp, f_2 f_3 \text{ hdr}$ $f_1 \rightarrow \bullet \text{ sep } data_1; f_2 f_3 \perp, f_2 f_3 \text{ hdr}$ $f_2 \rightarrow \bullet \text{ sep } data_2; f_3 \perp, f_3 \text{ hdr}$ $f_3 \rightarrow \bullet \text{ sep } data_3; \perp, R \perp, R \text{ hdr}$
4	$R \rightarrow \text{hdr } f_1 \bullet f_2 f_3; \perp, R \perp, R \text{ hdr}$	$f_2 \rightarrow \bullet; f_3 \perp, f_3 \text{ hdr}$ $f_2 \rightarrow \bullet \text{ sep } data_2; f_3 \perp, f_3 \text{ hdr}$ $f_3 \rightarrow \bullet \text{ sep } data_3; \perp, R \perp, R \text{ hdr}$
5	$f_1 \rightarrow \text{sep} \bullet data_1; f_2 f_3 \perp, f_2 f_3 \text{ hdr}$ $f_2 \rightarrow \text{sep} \bullet data_2; f_3 \perp, f_3 \text{ hdr}$ $f_3 \rightarrow \text{sep} \bullet data_3; \perp, R \perp, R \text{ hdr}$	
6	$R \rightarrow \text{hdr } f_1 f_2 \bullet f_3; \perp, R \perp, R \text{ hdr}$	$f_3 \rightarrow \bullet; \perp, R \perp, R \text{ hdr}$ $f_3 \rightarrow \bullet \text{ sep } data_3; \perp, R \perp, R \text{ hdr}$
7	$f_2 \rightarrow \text{sep} \bullet data_2; f_3 \perp, f_3 \text{ hdr}$ $f_3 \rightarrow \text{sep} \bullet data_3; \perp, R \perp, R \text{ hdr}$	
8	$f_1 \rightarrow \text{sep } data_1 \bullet; f_2 f_3 \perp, f_2 f_3 \text{ hdr}$	
9	$f_2 \rightarrow \text{sep } data_2 \bullet; f_3 \perp, f_3 \text{ hdr}$	
10	$f_3 \rightarrow \text{sep } data_3 \bullet; \perp, R \perp, R \text{ hdr}$	
11	$R \rightarrow \text{hdr } f_1 f_2 f_3 \bullet; \perp, R \perp, R \text{ hdr}$	
12	$f_3 \rightarrow \text{sep} \bullet data_3; \perp, R \perp, R \text{ hdr}$	

Figure 21. Nodes and items for grammar G_{8a} .

Grammar G_9 (Figure 22) describes BNF for context sensitive grammars (CSGs), where each rule is expected to start on a new line. This use of layout in the syntax of a language often leads to parsing problems.

$$\begin{aligned}
 S &\rightarrow \varepsilon \mid S R \\
 R &\rightarrow LHS ::= RHS \text{ eol} \\
 LHS &\rightarrow SS \\
 RHS &\rightarrow SS_{opt} \mid RHS \text{ eol} SS_{opt} \\
 SS &\rightarrow id \mid SS id \\
 SS_{opt} &\rightarrow \varepsilon \mid SS
 \end{aligned}$$

Figure 22. Grammar G_9 . BNF with layout for CSG productions.

A final example is drawn from compiler construction. In (revised) Pascal [16] statements are separated or terminated by semi-colons and *if* statements complicate the syntax by having optional else parts. Grammar G_{10a} (Figure 23) is a stylized unambiguous grammar describing this syntax. Inexperienced programmers may have problems learning where to put semi-colons in Pascal so a parser based on G_{10b} is more useful in a practical compiler.

$ \begin{aligned} BS &\rightarrow \textit{begin} SL s_{opt} \textit{end} \\ SL &\rightarrow ST \mid SL \textit{semi} ST \\ ST &\rightarrow MS \mid US \\ MS &\rightarrow \textit{stmt} \mid BS \mid \textit{if} MS \textit{else} MS \\ US &\rightarrow \textit{if} ST \mid \textit{if} MS \textit{else} US \\ s_{opt} &\rightarrow \textit{semi} \mid \varepsilon \end{aligned} $ <p>Grammar G_{10a}</p>	$ \begin{aligned} BS &\rightarrow \textit{begin} SL s_{opt} \textit{end} \\ SL &\rightarrow ST \mid SL s_{req} ST \\ ST &\rightarrow MS \mid US \\ MS &\rightarrow \textit{stmt} \mid BS \mid \textit{if} MS s_{err} \textit{else} MS \\ US &\rightarrow \textit{if} ST \mid \textit{if} MS s_{err} \textit{else} US \\ s_{opt} &\rightarrow \textit{semi} \mid \varepsilon \\ s_{req} &\rightarrow \textit{semi} \mid \varepsilon \\ s_{err} &\rightarrow \varepsilon \mid \textit{semi} \end{aligned} $ <p>Grammar G_{10b}</p>
--	--

Figure 23. Stylized Pascal syntax.

The error productions present problems for the LR parser generators but are acceptable to the two-stack parser generator. Figure 24 shows the number of states/nodes in parsers based on these grammars.

7.2 Choosing a Context Construction

The design of function Context determines the grammar class producing deterministic parsers based on the above parser construction. Our goal is to develop a terminating parser generator for a broad and decidable class of grammars. Since membership of $LR(k)$, LRR and DRP are all undecidable,

Parser	$G7a$	$G7b$	$G8a$	$G8b$	$G8c$	$G9$	$G10b$	$G10a$
LALR(1)	-	-	-	15	16	-	19	-
Two-stack	-	18	13	15	16	18	33	37
Extended two-stack	16	18	13	15	16	18	33	37

Figure 24. Parser state/node size.

considerable care is required. Informally, the context construction above has the required properties:

- item contexts at a node include contexts for LR items and states reached via the same symbols sequence;
- no context produces an empty string;
- the set of contexts is (practically) finite.

Other context constructions were considered and discarded. These are outlined in Figure 25.

Context construction	Problems
$Context(\alpha, \beta) = \alpha \beta$	infinite contexts and/or nodes
$Context(\alpha, \beta) = \begin{matrix} \beta, \alpha = \epsilon \\ \alpha, \text{ otherwise} \end{matrix}$	when $\alpha \Rightarrow^+ \epsilon$ may not include LR(1)
$Context(\alpha, X \beta) = \alpha X$	when $\alpha X \Rightarrow^+ \epsilon$ may not include LR(1)
<i>FOLLOW</i>	error detection may be delayed (works for G_I)
$Context(\alpha, \beta) = \begin{matrix} \beta, \alpha = \epsilon \\ \alpha x, x \in FIRST_1(\beta) \end{matrix}$	too many nodes and little advantage over the Context function presented in Section 5.1
$Context_k(\alpha, \beta) = \begin{matrix} \beta, \alpha = \epsilon \\ \alpha x, x \in FIRST_k(\beta) \end{matrix}$	too many nodes (but includes LR(k))

Figure 25. Alternative item context constructions.

Others have used or considered non-terminals in the look-ahead during parser construction. The LR(k) construction of Ancona et. al [30] is of considerable interest and deserves further investigation.

Parser generation can be extended using other context constructions. Parsers construction can use restricted regular expressions in item contexts. Closure for a left recursive symbols is treated specially, i.e. when $A \Rightarrow^+ A \alpha$ via rules $A \rightarrow A \alpha$ and $A \rightarrow \lambda$, then $Derive(I_i)$ where $B \mu \rightarrow \beta \bullet A \gamma; \delta \in I_i$, contains the items:

$$\begin{aligned} A \epsilon &\rightarrow \bullet A \alpha; \alpha^* \kappa \\ A \epsilon &\rightarrow \bullet \lambda; \alpha^* \kappa \end{aligned}$$

where $\kappa \in Context(\gamma, \delta)$. When $\alpha \Rightarrow^* \epsilon$, the grammar is ambiguous so parser construction will fail. The function Next in the Edge construction is extended to treat a compound context as an equivalent item group:

$$\begin{aligned} A\mu &\rightarrow \alpha \bullet; (\beta | \dots | \gamma)^* \delta \\ &\equiv A\mu \rightarrow \alpha \bullet; \delta, \beta (\beta | \dots | \gamma)^* \delta, \dots, \gamma (\beta | \dots | \gamma)^* \delta \end{aligned}$$

In the digraph construction for grammar G_I (Figure 26), node 6 includes the final item $P \rightarrow n ::= R \bullet; P^* \perp$ which is interpreted as the equivalent item groups:

$$\begin{aligned}
 &P \rightarrow n ::= R \bullet; P^* \perp \\
 \equiv &P \rightarrow n ::= R \bullet; \perp, P^+ \perp \\
 \equiv &P \rightarrow n ::= R \bullet; \perp, P P^* \perp
 \end{aligned}$$

Node	Node items	Closure items	Edges
0	$S' \rightarrow \perp \bullet S; \perp$	$S \rightarrow \bullet S P; P^* \perp$ $S \rightarrow \bullet P; P^* \perp$ $P \rightarrow \bullet n ::= R; P^* \perp$	$\langle 0, S, 1 \rangle$ $\langle 0, P, 2 \rangle$ $\langle 0, n, 3 \rangle$
1	$S' \rightarrow \perp S \bullet; \perp$ $S \rightarrow S \bullet P; P^* \perp$	$P \rightarrow \bullet N ::= R; P^* \perp$	$\langle 1, P, 4 \rangle$ $\langle 1, n, 3 \rangle$
2	$S \rightarrow P \bullet; P^* \perp$		
3	$P \rightarrow n \bullet ::= R; P^* \perp$		$\langle 3, ::=, 5 \rangle$
4	$S \rightarrow S P \bullet; P^* \perp$		
5	$P \rightarrow n ::= \bullet R; P^* \perp$	$R \rightarrow \bullet; P^* \perp, n, t$ $R \rightarrow \bullet R n; P^* \perp, n, t$ $R \rightarrow \bullet R t; P^* \perp, n, t$	$\langle 5, R, 6 \rangle$
6	$P \rightarrow n ::= R \bullet; P^* \perp$ $R \rightarrow R \bullet n; P^* \perp, n, t$ $R \rightarrow R \bullet t; P^* \perp, n, t$	$P \rightarrow \cdot N ::= R; P^* \perp$	$\langle 6, n, 7 \rangle$ $\langle 6, t, 8 \rangle$
7	$R \rightarrow R n \bullet; P^* \perp, n, t$ $P \rightarrow n \bullet ::= R; P^* \perp$		$\langle 7, ::=, 5 \rangle$
8	$R \rightarrow R t \bullet; P^* \perp, n, t$		

Figure 26. Digraph construction for grammar G_1 .

The Closure items for node 6 include (the *Relevant* part of) $CI(P, P^* \perp) = \{P \rightarrow \bullet n ::= R; P^* \perp\}$ since $n \in SS(I_6)$. The resulting parse table is shown in Figure 27.

Node	\perp	n	t	$::=$	S	P	R
0		S3			S1	S2	
1	Acc	S3				S3	
2	R2	R2				R2	
3				S5			
4	R1	R1				R1	
5	R4	R4	R4				S6
6	R3	S7	S8			R3	
7	R5	R5	R5	S5		R5	
8	R6	R6	R6			R6	

Figure 27. Parser table for grammar G_1 .

A parser created in this manner may use a large stack during parsing. Replacing the grammar rule $A \rightarrow A\alpha$ with a rule $A \rightarrow \alpha A$ (e.g. G_1 becomes G_{5b}) may achieve the same result and has the advantage of requiring an implementation choice with the acceptance of its consequences.

The use of regular expression in item contexts relates to LRR parsing [26]. The basic technique described here is not regarded as sufficiently general to justify implementation. Several extensions are under investigation.

7.3 Parser Generator Performance

Initial tests with small near- LR(1) grammars produced parsers comparable in size, i.e. number of states/nodes and number of table entries, to LR parsers. The parser produced for G_{8a} is smaller than a LR parser for the same language. In the observed cases, the technique appears to scale better than LR(k).

Practical parsers can be produced by adapting optimization techniques developed for LR parsing [31]. Practical LALR(1) construction can merge compatible LR(1) nodes during [32] or after [1] parser construction. The parse table for G_6 (Figure 9) can be compacted by creating a default reduce action for relevant nodes (i.e. all except 0, 2, 3) then merging nodes 7 and 8.

The controlled use of ambiguous grammars [10] allows small and efficient parsers to be constructed from compact language specifications. This technique, like the elimination of (semantically null) unit productions, appears to be applicable in two-stack parsing.

There is no problem with a parser generator producing a full LR(1) parser before merging states. Holub [32] indicates that the primary disadvantage of an LR(1) state machine is that it is typically twice the size of a corresponding LR(0) machine. Spector [33] points out that minimal-state full LR(1) tables are not significantly larger than LALR(1) tables but does not address the issue of whether the class of LR(1) grammars is practically larger than the LALR(1) grammars. Language designers tend to target LALR(1) acceptability rather than LR(1) parsers. Recent work by McPeak [34] and Chen and Pager [35] explores the effective implementation of LR(1) parsers and addresses the concerns noted earlier.

The two-stack parser generator needs "correct" context information in the items of a node so that it can avoid shifting error (terminal) symbols. The avoidance of node merging during parser generation is a simple (perhaps conservative) way to ensure the preservation of the valid prefix property in the parser.

7.4 Processing LR(1) and LR(k) Grammars

LR(1) grammars have special properties in the two-stack parser construction. The following function maps each digraph node I_n to a LR(1) state:

$$State(I_n) = \{A \rightarrow \alpha \bullet \beta; \alpha \mid A \rightarrow \alpha \bullet \beta; \gamma \in I_n \wedge \alpha \in FIRST_I(\gamma)\}$$

For each digraph node I_n , it is necessary that:

$$BS(I_n) \cap RS(I_n) = RC(I_n) = \emptyset$$

or the corresponding LR(1) state is inadequate. Two-stack closure and edge construction, working on LR(1) item sets, produce the same result as the corresponding LR(1) constructions.

For a LR(1) grammar, the terminal symbols associated with a particular reduce action at a node are the same as the LR(1) parser symbols at the corresponding

"state". The two-stack parser construction produces extra reduce table entries where the LR(1) parser says "don't care" and this may produce extra nodes. Merging compatible nodes produces a parser with the same number of states/nodes. Typically, these extra entries allow a two-stack parser based on a LR(1) grammar to recognize any sentential form thus are applicable in incremental parsing [27].

The two-stack construction algorithm is converted to LR(k) parser construction by replacing the Context function with the simpler $Context(\alpha, \beta) = FIRST_k(\alpha\beta)$. The $Context_k$ function mentioned in Section 7.2 lets the parser generator cover LR(k) in the same way as LR(1) is covered by the algorithm in Sections 5 and 6.

8. Conclusions

The two-stack parser can be regarded as a "drop-in replacement" for the LR parsing engine. Using LR(0), SLR(1), LALR(1) or LR(1) tables, the size and performance of the two parsing algorithms are the same. The two-stack parser has the advantage that it can be used with parse tables developed from a wider class of grammars.

The outline above shows that the LR(1) grammars, and their subclasses, are included in the class of grammars accepted by this two-stack parser generator. Given a LR(1) grammar, there is little advantage in using the two-stack parser generator. Given a non-LR(1) grammar, this parser generator can help in several ways. Most obviously, this generator could produce the tables for a practical parser. Alternatively, by preserving non-terminals in the item contexts, the generator may provide better diagnostic support for grammar manipulation. This parser construction extends the class of grammars for which efficient one-symbol look-ahead parsers can be constructed and provides an alternative to LR(k) when LR(1) proves inadequate.

The error handling capabilities of the parser and the parser generator deserve investigation. The error recovery mechanism used in YACC [13, 25] can be built into the two-stack parser, and the ability of the parser generator to resolve LR parser conflicts arising with grammars augmented for error handling can only be beneficial. More predictable error handling may be achieved through specific support in the parser generator.

An unoptimized parser table accepts many sentential forms as input. This property is expected to prove practical in incremental parsing/compiling, language sensitive editing or anywhere a parser could find non-terminals in its input [27].

The formal properties of grammar classes associated with the parser generator need to be identified, e.g. context and closure mechanisms described above construct parsers for a class of grammars properly including LR(k) for a given k.

Different context structures and closure mechanisms characterize different grammar classes.

9. REFERENCES

- [1] Aho, A.V. and Johnson, S.C., "LR Parsing," *Computing Surveys*, 6(2):99-124, June 1974.
- [2] DeRemer, F. and Pennello, T., "Efficient Computation of LALR(1) Look-Ahead Sets", *ACM Transactions on Programming Languages and Systems*, 4(4):615-649, October 1982.
- [3] Aho, A.V. and Ullman, J.D., *The Theory of Parsing, Translation and Compiling*, Volume. 1, Prentice-Hall, Englewood Cliffs, N.J., 1972.
- [4] Knuth, D.E., "On the Translation of Languages From Left to Right", *Information and Control*, 8(6):607-639, 1965.
- [5] Seite, B., "A YACC Extension for LRR Grammar Parsing", *Theoretical Computer Science*, 52:91-143, 1987.
- [6] Hopcroft, J.E. and Ullman, J.D., *Introduction to Automata Theory, Languages and Computation*, Addison-Wesley, 1979.
- [7] Harris, L.A., "SLR(1) and LALR(1) Parsing for Unrestricted Grammars", *Acta Informatica*, 24:191-209, 1987.
- [8] Turnbull, C.J.M. and Lee, E.S., "Generalized Deterministic Left to Right Parsing", *Acta Informatica*, 12:187-207, 1979.
- [9] Naur, P., [ed], "Report on the Algorithmic Language ALGOL 60", *Communications of the ACM*, 3(5):299-314, 1960. Revised in *Communications of the ACM*, 6:1-20, 1963.
- [10] Aho, A.V., Lam, M., Sethi, R., and Ullman, J.D., *Compilers: Principles, Techniques and Tools*, 2nd edition, Addison-Wesley, 2006.
- [11] Horning, J.J., "LR Grammars and Analyzers" in *Compiler Construction, an Advanced Course*, 2nd edition, F. L. Bauer and J. Eikel (eds), pp. 85-107, Springer, 1976.
- [12] Horning, J.J., "LR Grammars and Analyzers" in *Compiler Construction, an Advanced Course*, 2nd edition, F. L. Bauer and J. Eikel (eds), pp. 85-107, Springer, 1976.
- [13] Johnson, S.C., "YACC: Yet Another Compiler-Compiler", *Computer Science Technical Report 32*, Bell Laboratories, Murray Hill, NJ, 1975.
- [14] Wirth, N., "What Can we do About the Unnecessary Diversity of Notation for Syntactic Definitions", *Communications of the ACM*, 20(11):822-823, November 1977.
- [15] Kernighan, B.W. and Ritchie, D.M., *The C Programming Language*, 2nd edition, Prentice-Hall, Englewood Cliffs, N.J., 1988.
- [16] Jensen, K. and Wirth, N., "Pascal - User Manual and Report", 2nd edition, *Lecture Notes in Computer Science No. 18*, Springer, New York, 1975.
- [17] Wirth, N., *Programming in Modula-2*, 2nd edition, Springer-Verlag, New York, 3rd edition, 1985.

- [18] Cardelli, L., Donahue, J., Glassman, L., Jordan, M., Kalsow, B. and Nelson, G., "Modula-3 Language Definition", ACM SIGPLAN Notices, 27(8):42-82, 1992.
- [19] Barnes, J.G.P., Programming in Ada 2012, Cambridge, Cambridge Press, 2014.
- [20] Taft, S. Tucker. et al. Ada 2012 Reference Manual. Language and Standard Libraries International Standard ISO/IEC 8652/2012 (E) . 1st ed. 2013. Berlin, Heidelberg: Springer Berlin Heidelberg, 2013. Web.
- [21] Richards, M. and Whitby-Stevens, C., BCPL - the Language and its Compiler, Cambridge University Press, 1981.
- [22] Hudak, P., Jones, S.P., and Wadler, P.L., [eds], "Report on the Functional Programming Language Haskell, version 1.2", ACM SIGPLAN Notices, 27(5), May 1992.
- [23] Marlow, S., et al. Haskell 2010 Language Report, Available online <http://www.haskell.org/>(May 2011), 2010.
- [24] Lesk, M.E. and Schmidt, E., "Lex-a Lexical Analyzer Generator" in UNIX Programmer's Manual 2, AT&T Bell Labs., Murray Hill, N.J., 1975.
- [25] Mason, T. and Brown, D. Lex & Yacc, O'Reilly & Associates, 1990.
- [26] Culik, K. and Cohen, R., "LR-Regular Grammars - an Extension of LR(k) Grammars", Computer and System Science, 7:66-96, 1973.
- [27] DeRemer, F.L., "Simple LR(k) Grammars", Communications of the ACM, 14(7):453-460, July 1971.
- [28] Bermudez, M.E. and Schimpf, K.M., "A Practical Arbitrary Look-Ahead LR Parsing Technique", ACM SIGPLAN Notices, 21(7):136-144, July 1986.
- [29] Harford, A.G., Heuring, V.P., and Main, M. G., "A New Parsing Method for Non-LR(1) Grammars", Software - Practice and Experience, 22(5):419-437, Wiley, May 1992.
- [30] Ancona, M., Doderio, G., Gianuzzi, V., and Morgavi, M., "Efficient Construction of LR(k) States and Tables", ACM Transactions on Programming Languages and Systems, 13(1)150-178, January 1991.
- [31] Sippu, S. and Soisalon-Soininen, E., Parsing Theory, Volume II: LR(k) and LL(k) Parsing, Springer, Berlin, 1990.
- [32] Holub, A.I., Compiler Design in C, Prentice-Hall, Englewood Cliffs, N.J., 1990.
- [33] Spector, D., "Efficient Full LR(1) Parser Generation", SIGPLAN Notices, 23(12):143-150, December 1988.
- [34] McPeak, S., and Necula, G. C. "Elkhound: A fast, practical GLR parser generator." Proceedings of Compiler Construction, pp. 73-88, 2004. Retrieved from <https://escholarship.org/uc/item/8559j464>.
- [35] Chen, X. and Pager, D.. "Full LR(1) Parser Generator Hyacc and Study on the Performance of LR(1) algorithms". In Proceedings of The Fourth International C* Conference on Computer Science and Software Engineering (C3S2E '11). Association for Computing Machinery, New York, NY, USA, 83-92, 2011.