

## An efficient computation of right context for LR-based error repair

Min-Soo Jung \*, Kwang-Moo Choe, Taisook Han

*Programming Languages Laboratory, Department of Computer Science, Korea Advanced Institute of Science and Technology, 373-1,  
Kusong-Dong, Yusong-Gu, Taejon 305-701, South Korea*

(Communicated by K. Ikeda)

(Received 19 May 1993)

(Revised 26 August 1993 and 15 October 1993)

---

### Abstract

The *left context* in LR-based parsing is the sequence of states in the parsing stack. The *right context* is the vocabulary strings to appear for a given left context. We propose an efficient method of computing right context for LR-based syntax error repair. The efficiency of our method is achieved from removing some redundancies in the method of previous work.

*Key words:* Compilers; LR parsing; Syntax error repair; Kernel item; Right context

---

### 1. Introduction

The *left context* in LR-based parsing is the sequence of states in the parsing stack. The *right context* is the expected vocabulary strings to appear for a given left context. One of the important applications of right context is syntax error repair. When a syntax error occurs during a parsing, right context might be used to get an insertion string [4,5] which patches up the syntax error. LR-based parser, however, requires additional space and time in computing right context since an LR parsing stack contains vocabulary strings that have been seen till now, whereas an LL parsing stack contains expected vocabulary strings

from now. This overhead in LR-based parser may be one of the reasons that several error repair methods including FMQ [5] have selected LL as a target parser.

Choe and Chang [2] proposed a method of computing right context for LR-based error repair at parsing time, using a new LALR formalism [7]. Their method examined the kernel items only by means of factoring Closure relations [1] on states into Path relations [2,7] on nonterminals. One of the difficulties of computing right context in LR-based parsing is that there are multiple kernel items in a state. In the work of Choe and Chang, the right context of an LR(0) state was computed as follows: getting a kernel item from the state and then computing the right context of the kernel item, getting a next kernel item from the state and then computing the right

---

\* Corresponding author.

context of the kernel item, and so on. When the right contexts of all kernel items in the state have been computed one by one, then they are united to make the right context of the state. Their method, however, computes the right context redundantly since the right contexts of each kernel item in the same state may have common suffixes.

In this paper, we present an efficient method of computing right context for LR-based syntax error repair. Informally, our method is summarized as follows. We show that the suffixes of right context for each kernel item in the same state are the same in many cases, although their prefixes are normally different. In order not to compute the common suffixes redundantly, we should know what the common suffixes of a given state is. We present a method of getting common suffixes mostly at parser generation time.

The following section reviews the notation and the terminology used. Section 3 contains a method of analysing kernel items in a given LR(0) state to get common suffixes. Section 4, an extension of Section 3, contains a method of analysing addition states as well as the given state. Section 5 contains an experimental result to validate our method. Conclusion is given in Section 6.

## 2. Notation and terminology

The reader is assumed to be familiar with the notion and the notation related to context-free grammar and LR parsing theory in [1]. A context-free grammar  $G$  is a quadruple  $G = (N, \Sigma, P, S)$ , where  $N$  is a finite set of nonterminals,  $\Sigma$  is a finite set of terminals,  $P$  is a finite set of productions, and  $S$  is the start symbol in  $N$ . Given a CFG  $G = (N, \Sigma, P, S)$ , the corresponding augmented grammar is  $G' = (N', \Sigma, P', S')$ , where  $S'$  is the new start symbol not in  $N$ ,  $N' = N \cup \{S'\}$ , and  $P' = P \cup \{S' \rightarrow S\}$ . The vocabulary of  $G$  is denoted by  $V$  ( $V = N \cup \Sigma$ ). The concatenation over vocabulary strings or sets of vocabulary strings is defined as follows: Let  $X$  and  $Y$  be sets of vocabulary strings and  $\gamma$  be a vocabulary string, then  $XY = \{\alpha\beta \mid \alpha \in X, \beta \in Y\}$ , and  $\gamma X = \{\gamma\alpha \mid \alpha \in X\}$ . The length of string  $\alpha$ , denoted by  $|\alpha|$ , is the number of symbols in  $\alpha$ .

The prefix of  $\alpha$  of length  $k$  is denoted by  $k:\alpha$ , and the suffix of  $\alpha$  of length  $k$  is denoted by  $\alpha:k$ .

Let  $p$  be an LR(0) state. Then the set of all kernel items in the state  $p$  is denoted by  $\text{KERNEL}(p)$ , and the set of all mark symbols of kernel items in  $p$  is denoted by  $\text{MARK}(p)$ , where a mark symbol is the symbol immediately to the right of the dot in an LR(0) item. The  $\alpha$ -predecessor of  $p$  is defined by  $\text{PRED}(p, \alpha) = \{r \mid \text{GOTO}(r, \alpha) = p\}$ . Recall that a left context is a sequence of states. Then the  $\alpha$ -predecessor of a left context  $\sigma$  is defined by  $\text{PRED}(\sigma, \alpha) = \{\sigma - |\alpha| : \sigma$ , i.e. the  $\text{PRED}(\sigma, \alpha)$  is computed by removing the suffix of  $\sigma$  of length  $|\alpha|$  from  $\sigma$ .

**Definition 2.1** (Lookback [3]). Let  $\sigma$  be a left context and  $p$  be the last symbol of  $\sigma$ , i.e.,  $\sigma:1 = p$ , and let  $p$  have a kernel item  $[A \rightarrow \alpha \cdot \beta]$ . Then the lookback of  $\sigma$  and  $[A \rightarrow \alpha \cdot \beta]$  is defined as follows:

$$\text{LOOKBACK}(\sigma, [A \rightarrow \alpha \cdot \beta]) = \text{PRED}(\sigma, \alpha)q,$$

where  $q = \text{GOTO}(r, A)$ ,  $r = \text{PRED}(\sigma, \alpha):1$ .

**Example 2.2** (Predecessor and lookback). Assume that LR(0) states  $s, r, q$  and  $p$  are depicted in Fig. 1. Let left context  $\sigma = s \cdots r \cdots p$ . Then  $\text{PRED}(\sigma, \alpha) = s \cdots r$  and  $\text{LOOKBACK}(\sigma, [A \rightarrow \alpha \cdot \beta]) = s \cdots rq$ . This example shows the usefulness of the lookback operation in computing right context. For the given  $\sigma$ , we get  $\beta$  from the state  $p$ . And then, we get  $\delta$  from the state  $q$ . Note that  $q$  is the last symbol of  $\text{LOOKBACK}(\sigma, [A \rightarrow \alpha \cdot \beta])$ .

## 3. Computation of right context for LR-based error repair

In this section, we define the right context of a given left context. We show some redundancies of the method of Choe and Chang in computing

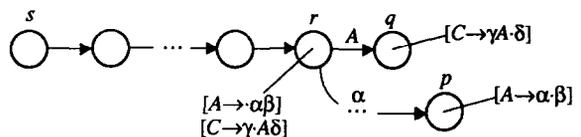


Fig. 1.

right context. Through an analysis of the kernel items in an LR(0) state, we show the rationale of removing the redundancies.

**Definition 3.1.** Let  $\sigma$  be a left context and  $p$  be the last symbol of  $\sigma$ . The *right context* of a given left context  $\sigma$ , denoted by  $\text{RC}(\sigma)$ , is defined as follows:

$$\text{RC}(\sigma) = \bigcup_{[A \rightarrow \alpha \cdot \beta] \in \text{KERNEL}(p)} \{\text{RCI}(\sigma, [A \rightarrow \alpha \cdot \beta])\},$$

where  $\text{RCI}(\sigma, [A \rightarrow \alpha \cdot \beta]) = \beta \text{RC}(\tau)$ ,  $\tau = \text{LOOKBACK}(\sigma, [A \rightarrow \alpha \cdot \beta])$ . As a terminating condition, we given  $\text{RCI}(\sigma', [S' \rightarrow S \cdot]) = \varepsilon$  (the empty string), where  $\sigma'$  is a left context.

Definition 3.1 shows that right context of an LR(0) state is the union of right contexts of all kernel items in the state. Now we are to give the key idea of this paper. Let  $\sigma$  be a left context and  $p$  be the last symbol of  $\sigma$ . Assume that  $\text{KERNEL}(p) = \{[A_1 \rightarrow \alpha_1 \cdot \beta_1], \dots, [A_n \rightarrow \alpha_n \cdot \beta_n]\}$  and  $\text{RCI}(\sigma, [A_1 \rightarrow \alpha_1 \cdot \beta_1]) = \beta_1 \gamma_1 \delta, \dots, \text{RCI}(\sigma, [A_n \rightarrow \alpha_n \cdot \beta_n]) = \beta_n \gamma_n \delta$ , where  $\gamma_i$  ( $1 \leq i \leq n$ ) and  $\delta$  are some vocabulary string or set of vocabulary strings. In the method of Choe and Chang,  $\delta$ , the common suffixes of the right contexts of every kernel items in  $p$ , is computed  $n$ -times redundantly. To get rid of this redundancy, we factor out the common suffixes  $\delta$ . The rest part of this section presents a method of finding common suffixes.

**Lemma 3.2.** Let  $p$  be an LR(0) state and  $\sigma$  be any left context whose last symbol is  $p$ . If state  $p$  has only one kernel item, say  $\text{KERNEL}(p) = \{[A \rightarrow \alpha \cdot \beta]\}$ , then  $\text{RC}(\sigma) = \beta \text{RC}(\tau)$ , where  $\tau = \text{LOOKBACK}(\sigma, [A \rightarrow \alpha \cdot \beta])$ .

According to Lemma 3.2, if the last symbol of  $\sigma$  has only one kernel item, then  $\text{RC}(\sigma)$  can be computed from vocabulary string  $\beta$ ,  $|\alpha|$ ,  $A$ , the left context in parsing stack and parsing table. The left context and parsing table are assumed to be given in the remaining part of this paper since they are necessary even for a “normal” parsing.

Since  $\beta$  is a prefix of  $\text{RC}(\sigma)$ , and  $\text{RC}(\tau)$  is the rest suffixes of  $\text{RC}(\sigma)$ , we store  $\beta$ ,  $|\alpha|$  and  $A$  into the state  $p$ , where  $|\alpha|$  and  $A$  are necessary to get  $\tau$ . Recall the definition of lookback,  $\tau$  can be computed from  $\sigma$ ,  $|\alpha|$  and  $A$ . And then computation of  $\text{RC}(\tau)$  will be started with  $\tau$ , recursively.

**Lemma 3.3.** Let  $p$  be an LR(0) state and  $\sigma$  be any left context whose last symbol is  $p$ . If  $\text{KERNEL}(p) = \{[A \rightarrow A \cdot \alpha], [B \rightarrow \gamma \cdot \beta]\}$ , then  $\text{RC}(\sigma) = \alpha^* \beta \text{RC}(\tau)$ , where  $\tau = \text{LOOKBACK}(\sigma, [B \rightarrow \gamma \cdot \beta])$ .

**Proof.** Since  $\text{LOOKBACK}(\sigma, [A \rightarrow A \cdot \alpha]) = \sigma$ ,

$$\begin{aligned} \text{RC}(\sigma) &= \beta \text{RC}(\tau) \cup \alpha \text{RC}(\sigma) \\ &= \beta \text{RC}(\tau) \cup \alpha (\beta \text{RC}(\tau) \cup \alpha \text{RC}(\sigma)) \\ &= \dots = \{\varepsilon, \alpha, \alpha\alpha, \dots\} \beta \text{RC}(\tau) \\ &= \alpha^* \beta \text{RC}(\tau). \quad \square \end{aligned}$$

Note the form of the kernel item  $[A \rightarrow A \cdot \alpha]$ . Left-hand side nonterminal in the kernel item,  $A$ , and the vocabulary string to the left part of dot in the kernel item,  $A$ , are the same. The kernel item  $[A \rightarrow A \cdot \alpha]$  does not contribute to computing  $\tau$ , where  $\text{RC}(\tau)$  will be common suffixes. But the kernel item  $[B \rightarrow \gamma \cdot \beta]$  contributes to getting  $|\gamma|$  and  $B$  to compute  $\tau$ . We say the kernel item  $[B \rightarrow \gamma \cdot \beta]$  is the *representative kernel item* of the state  $p$  in the sense that this kernel item has the informations for computing the common suffixes. To compute the common suffixes, we store  $\alpha^* \beta$ ,  $|\gamma|$  and  $B$  into the state  $p$ . Lemma 3.3 can be extended to an LR(0) state  $p$  whose kernel items are  $[A \rightarrow A \cdot \alpha_1], \dots, [A \rightarrow A \cdot \alpha_n]$ , and  $[B \rightarrow \gamma \cdot \beta]$ , where  $[B \rightarrow \gamma \cdot \beta]$  will be the representative kernel item of  $p$ .

**Lemma 3.4.** Let  $p$  be an LR(0) state and  $\sigma$  be any left context whose last symbol is  $p$ . If  $\text{KERNEL}(p) = \{[A \rightarrow \gamma \cdot \alpha], [A \rightarrow \gamma \cdot \beta]\}$ , then  $\text{RC}(\sigma) = \{\alpha, \beta\} \text{RC}(\tau)$ , where  $\tau = \text{LOOKBACK}(\sigma, [A \rightarrow \gamma \cdot \alpha]) = \text{LOOKBACK}(\sigma, [A \rightarrow \gamma \cdot \beta])$ .

Note the form of two kernel items  $[A \rightarrow \gamma \cdot \alpha]$  and  $[A \rightarrow \gamma \cdot \beta]$ . Both kernel items have the same left-hand side nonterminal,  $A$ , and the same vo-

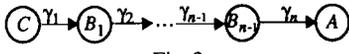


Fig. 2.

cabulary string to the left part of dot,  $\gamma$ . Not both kernel items  $[A \rightarrow \gamma \cdot \alpha]$  and  $[A \rightarrow \gamma \cdot \beta]$  are necessary for computing  $\tau$ , where  $RC(\tau)$  will be common suffixes; one of two is sufficient to compute  $\tau$ . We say the kernel item  $[A \rightarrow \gamma \cdot \alpha]$  (or  $[A \rightarrow \gamma \cdot \beta]$ ) is the representative kernel item of the state  $p$ . Lemma 3.4 can be extended to an LR(0) state  $p$  whose kernel items are  $[A \rightarrow \gamma \cdot \alpha_1], \dots, [A \rightarrow \gamma \cdot \alpha_n]$ , where any kernel item of them can be the representative kernel item of  $p$ .

**Definition 3.5** (Immediately represented state). If all kernel items in an LR(0) state  $p$  can be represented by one kernel item using Lemma 3.3 and/or Lemma 3.4, then we say that the state  $p$  is an *immediately represented state*.

**Example 3.6** (Immediately represented state). Let  $p$  be an LR(0) state and  $\sigma$  be any left context whose last symbol is  $p$ . If  $KERNEL(p) = \{(S \rightarrow iE \cdot tS), [S \rightarrow iE \cdot tSeS], [E \rightarrow E \cdot +E], [E \rightarrow E \cdot -E]\}$ , then

- $[S \rightarrow iE \cdot tS]$  } (one of two is sufficient to compute the common suffixes),
- $[S \rightarrow iE \cdot tSeS]$  }
- $[E \rightarrow E \cdot +E]$  (is not necessary to compute the common suffixes),
- $[E \rightarrow E \cdot -E]$  (is not necessary to compute the common suffixes).

Hence  $RC(\sigma) = \{+E, -E\}^* \{tS, tSeS\} RC(\tau)$ , where  $\tau = LOOKBACK(\sigma, [S \rightarrow iE \cdot tS])$  (or  $\tau = LOOKBACK(\sigma, [S \rightarrow iE \cdot tSeS])$ ), and  $[S \rightarrow iE \cdot tS]$  is the representative kernel item of  $p$ .

**Theorem 3.7.** Let  $p$  be an LR(0) state and  $\sigma$  be any left context whose last symbol is  $p$ . If  $p$  is an immediately represented state with the representative kernel item  $[A \rightarrow \alpha \cdot \beta]$ , then  $RC(\sigma) = VS RC(\tau)$ , where  $VS$  is a set of vocabulary strings which is computable from the kernel items in the state  $p$  and  $\tau = LOOKBACK(\sigma, [A \rightarrow \alpha \cdot \beta])$ .

**4. Analysis of non-immediately represented state**

In this section, we deal with non-immediately represented states which are not represented by one kernel item through the analysis presented in Section 3. However, by the analysis of additional predecessor states of a given non-immediately represented state  $p$ , we may get common suffixes for the state  $p$ . The following definitions are introduced in order to analyse a non-immediately represented state.

**Definition 4.1** (L-graph and L-relation [7]). Let  $G = (N, \Sigma, P, S)$  be a context-free grammar. Then the L-graph for  $G$  is a directed graph whose vertices are the same as the nonterminals in  $G$ , and the edges are defined by  $Edge(A, B) = \{\beta \mid A \rightarrow B\beta \in P\}$ . The nonterminal  $A$  has an L-relation to the nonterminal  $B$ , written  $A L B$ , if  $Edge(A, B) \neq \emptyset$ . The nonterminal  $A$  has a transitive L-relation to the nonterminal  $B$ , written  $A L^+ B$ , if there is a path from  $A$  to  $B$  in L-graphs.

**Definition 4.2** (Path [7]). Let  $G = (N, \Sigma, P, S)$  be a context-free grammar,  $A$  and  $C$  be nontermi-

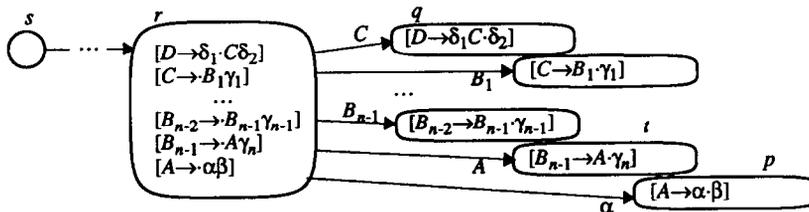


Fig. 3.

nals. Then the *paths* from  $C$  to  $A$ , denoted  $\text{PATH}(C, A)$ , is defined as follows:

$$\begin{aligned} \text{PATH}(C, A) = \{ & \gamma_n \dots \gamma_1 | B_0 \rightarrow B_1 \gamma_1, B_1 \rightarrow B_2 \gamma_2, \\ & \dots, B_{n-1} \rightarrow B_n \gamma_n \in P, \\ & B_0 = C, B_n = A, n \geq 1 \}, \end{aligned}$$

where the string  $\gamma_1 \dots \gamma_n$  describes a path from  $C$  to  $A$  in the L-graph of the form shown in Fig. 2, where  $C = B_0$  and  $A = B_n$ .

**Example 4.3** (The usefulness of L-graph, L-relation and PATH in computing right context). Fig. 3 shows an LR(0) automaton and let left context  $\sigma = s \dots r \dots p$ . One can get  $\beta$ , as a prefix of right context of  $\sigma$ , from the state  $p$ . Then go to state  $t$  using lookback operation, and get  $\gamma_n$ , which follows  $\beta$ . By the repeated computations using lookback operations, one can get  $\beta \gamma_n \gamma_{n-1} \dots \gamma_1 \delta_2 \dots$ , as a right context of  $\sigma$ . But Park, Choe and Chang computed  $\gamma_n \dots \gamma_1$  from  $\text{PATH}(C, A)$ , where paths between two nonterminals can be computed from traversing L-graph in reverse order. Also they point out that paths between two nonterminals is a property of grammar not a property of a specific LR(0) state.

**Definition 4.4.** Let  $r$  be an LR(0) state, then the *subgraph of L-graph related to  $r$* , denoted by  $\text{L-GRAPH}(r)$ , is defined as follows:

$$\begin{aligned} \text{L-GRAPH}(r) &= (V(r), E(r)), \\ V(r) &= \{A | C L^* A, C \in \text{MARK}(r)\}, \\ E(r) &= \{(A, B) | A L B, A, B \in V(r)\}. \end{aligned}$$

An example of L-graph and L-graph( $r$ ) for an actual grammar are given in Example 4.9. We will use the term nonterminal and the term vertex in L-graph interchangeably.

**Definition 4.5.** For a given  $\text{L-GRAPH}(r) = (V(r), E(r))$ , the *cut symbols* of an LR(0) state  $r$  and nonterminal  $A \in V(r)$ , denoted by  $\text{CUT}(r, A)$ , is defined as follows:

$$\text{CUT}(r, A) = \{A\} \cup \{B | B \in V(r), \text{deletion of vertex } B \text{ and its connected edges results that there is no path from any vertex in } \text{MARK}(r) \text{ to } A \text{ in } \text{L-GRAPH}(r)\}.$$

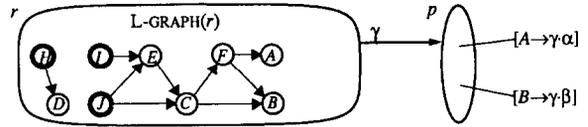


Fig. 4.

**Example 4.6** (Cut symbols). Let  $p$  be an LR(0) state and  $r \in \text{PRED}(p, \gamma)$ . Let  $\text{MARK}(r) = \{H, I, J\}$  and  $H \rightarrow D\gamma_1, I \rightarrow E\gamma_2, J \rightarrow E\gamma_3, J \rightarrow C\gamma_4, E \rightarrow C\gamma_5, C \rightarrow F\gamma_6, C \rightarrow B\gamma_7, F \rightarrow A\gamma_8$  and  $F \rightarrow B\gamma_9$  are productions over a CFG  $G$ . Assume that  $\text{L-GRAPH}(r)$  is depicted in Fig. 4.

From the  $\text{L-GRAPH}(r)$ , we have  $\text{CUT}(r, A) = \{A, F, C\}$ ,  $\text{CUT}(r, B) = \{B, C\}$  and  $\text{CUT}(r, A) \cap \text{CUT}(r, B) = \{C\}$ . To compute the common suffixes of  $[A \rightarrow \gamma \cdot \alpha]$  and  $[B \rightarrow \gamma \cdot \beta]$  in  $p$ , we examine the states  $r \in \text{PRED}(p, \gamma)$ . Recall that L-graph replaces the lookback operations in computing right context as is explained in Example 4.3. From the  $\text{L-GRAPH}(r)$ , the strings obtainable from  $A$  to  $C$ ,  $\text{PATH}(C, A)$ , and the strings obtainable from  $B$  to  $C$ ,  $\text{PATH}(C, B)$ , are different, but the strings obtainable from  $C$  to any vertex in  $\text{MARK}(r)$  are all the same with respect to  $A$  and  $B$ .

For the purpose of getting common suffixes for a given non-immediately represented state, we first classify non-immediately represented states into two classes. Let non-immediately represented state  $p$  have the kernel items  $[A_1 \rightarrow \alpha_1 \cdot \beta_1], [A_2 \rightarrow \alpha_2 \cdot \beta_2], \dots, [A_n \rightarrow \alpha_n \cdot \beta_n]$ . If  $\alpha_1 = \alpha_2 = \dots = \alpha_n$ , then we call the state  $p$  as *class 1*; otherwise we call the state  $p$  as *class 2*.

**Lemma 4.7** (Non-immediately represented state of class 1). *Let  $p$  be an LR(0) state and  $\sigma$  be any left context whose last symbol is  $p$ , and  $p$  be a non-immediately represented state of class 1. Let  $\text{KERNEL}(p) = \{[A \rightarrow \gamma \cdot \alpha], [B \rightarrow \gamma \cdot \beta]\}$ . For all  $r \in \text{PRED}(p, \gamma)$ , if there exists a common nonterminal  $C$  such that  $C \in \text{CUT}(r, A) \cap \text{CUT}(r, B)$ , then  $\text{RC}(\sigma) = \text{VS } \text{RC}(\tau)$ , where  $\text{VS}$  is a set of vocabulary strings,  $\text{RC}(\tau)$  is the common suffixes of right context of kernel items  $[A \rightarrow \gamma \cdot \alpha]$  and  $[B \rightarrow \gamma \cdot \beta]$  in  $p$ .*

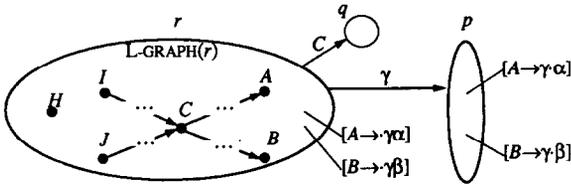


Fig. 5.

**Proof.** Assume that LR(0) states  $p$ ,  $q$ ,  $r$  and  $L\text{-GRAPHS}(r)$  are depicted in Fig. 5. For all  $r \in \text{PRED}(p, \gamma)$ , assume that  $C$  is a common nonterminal such that  $C \in \text{CUT}(r, A) \cap \text{CUT}(r, B)$ . Then  $\text{RCI}(\sigma, [A \rightarrow \gamma \cdot \alpha]) = \alpha \text{ PATH}(C, A) \text{ RC}(\tau)$ , and  $\text{RCI}(\sigma, [B \rightarrow \gamma \cdot \beta]) = \beta \text{ PATH}(C, B) \text{ RC}(\tau)$ , where  $\tau = \text{PRED}(\sigma, \gamma) q$  and  $q = \text{GOTO}(\text{PRED}(\sigma, \gamma):1, C)$ .

$$\begin{aligned} \text{RC}(\sigma) &= \text{RCI}(\sigma, [A \rightarrow \gamma \cdot \alpha]) \\ &\cup \text{RCI}(\sigma, [B \rightarrow \gamma \cdot \beta]) \\ &= \{\alpha \text{ PATH}(C, A), \beta \text{ PATH}(C, B)\} \\ &\times \text{RC}(\tau). \quad \square \end{aligned}$$

Note that we want to get a condition which holds for all left context whose last symbol is  $p$ , not for a specific left context. In other words, the condition presented in Lemma 4.7 is independent on history of parsing in the sense that a left context records a history of parsing. Hence we examine all  $\gamma$ -predecessor states of  $p$  and get a condition which holds commonly for all these states.

According to Lemma 4.7, for any left context  $\sigma$  whose last symbol  $p$  holds the condition in Lemma 4.7,  $\text{RC}(\sigma)$  is computable from  $\{\alpha \text{ PATH}(C, A), \beta \text{ PATH}(C, B)\}, |\gamma|, C$ . For the left context  $\sigma$  whose last symbol  $p$  does not hold the condition,  $\text{RC}(\sigma)$  can be computed by uniting  $\alpha \text{ RC}(\sigma')$  and  $\beta \text{ RC}(\tau')$ , where  $\sigma' = \text{LOOKBACK}(\sigma, [A \rightarrow \gamma \cdot \alpha])$  and  $\tau' = \text{LOOKBACK}(\sigma, [B \rightarrow \gamma \cdot \beta])$ . Let  $p$  be a non-immediately represented state of class 1 and  $p$  hold the condition in Lemma 4.7, then we call the state  $p$  as a *represented state of class 1*.

**Lemma 4.8** (Non-immediately represented state of class 2). *Let  $p$  be an LR(0) state and  $\sigma$  be any left context whose last symbol is  $p$ , and  $p$  be a non-immediately represented state of class 2. Let*

$\text{KERNEL}(p) = \{[A \rightarrow \delta \gamma \cdot \alpha], [B \rightarrow \gamma \cdot \beta]\}$ , where  $\delta \neq \epsilon$ . For all  $r \in \text{PRED}(p, \gamma)$ , if there exists a common set of kernel items whose mark symbol has a transitive L-relation to  $B$ , and all kernel items in this common set can be represented by the kernel item  $[A \rightarrow \delta \cdot \gamma \alpha]$ , then  $\text{RC}(\sigma) = \text{VS RC}(\tau)$ , where VS is a set of vocabulary strings,  $\text{RC}(\tau)$  is the common suffixes of right context of kernel items  $[A \rightarrow \delta \gamma \cdot \alpha]$  and  $[B \rightarrow \gamma \cdot \beta]$  in  $p$ .

**Proof.** Assume that LR(0) state  $p$ ,  $r$  and  $L\text{-GRAPHS}(r)$  are depicted in Fig. 6. For all  $r \in \text{PRED}(p, \gamma)$ , assume that the set of kernel items whose mark symbol has a transitive L-relation to  $B$  is  $\{[A \rightarrow \delta \cdot \gamma \alpha], [A \rightarrow \delta \cdot \phi]\}$ . So to speak,  $1:\gamma$  and  $1:\phi$  have a transitive L-relation to  $B$ . Since the kernel items  $[A \rightarrow \delta \cdot \gamma \alpha]$  and  $[A \rightarrow \delta \cdot \phi]$  in  $r$  are represented by the kernel item  $[A \rightarrow \delta \cdot \gamma \alpha]$ , we get

$$\begin{aligned} \text{RCI}(\sigma, [B \rightarrow \gamma \cdot \beta]) &= \{\beta \text{ PATH}(1:\gamma, B) \gamma : (|\gamma| - 1) \alpha, \\ &\beta \text{ PATH}(1:\phi, B) \phi : (|\phi| - 1) \text{RC}(\tau)\}, \end{aligned}$$

where  $\tau = \text{LOOKBACK}(\sigma, [A \rightarrow \delta \gamma \cdot \alpha])$ . And we get  $\text{RCI}(\sigma, [A \rightarrow \delta \gamma \cdot \alpha]) = \alpha \text{ RC}(\tau)$ .

$$\begin{aligned} \text{RC}(\sigma) &= \text{RCI}(\sigma, [A \rightarrow \delta \gamma \cdot \alpha]) \\ &\cup \text{RCI}(\sigma, [B \rightarrow \gamma \cdot \beta]) \\ &= \{\alpha, \beta \text{ PATH}(1:\gamma, B) \gamma \dots (|\gamma| - 1) \alpha, \\ &\beta \text{ PATH}(1:\phi, B) \phi : (|\phi| - 1)\} \text{RC}(\tau). \quad \square \end{aligned}$$

According to Lemma 4.8, for any left context  $\sigma$  whose last symbol  $p$  holds the condition in Lemma 4.8,  $\text{RC}(\sigma)$  is computable from  $\{\alpha, \beta \text{ PATH}(1:\gamma, B) \gamma \cdot (|\gamma| - 1) \alpha, \beta \text{ PATH}(1:\phi, B) \phi : (|\phi| - 1), |\delta \gamma|, A$ . Let  $p$  be a non-immediately represented state of class 2 and  $p$  holds the condition in Lemma 4.8, then we call the state  $p$  as a *represented state of class 2*.

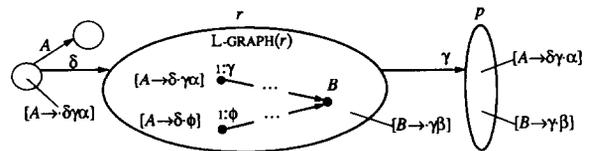


Fig. 6.

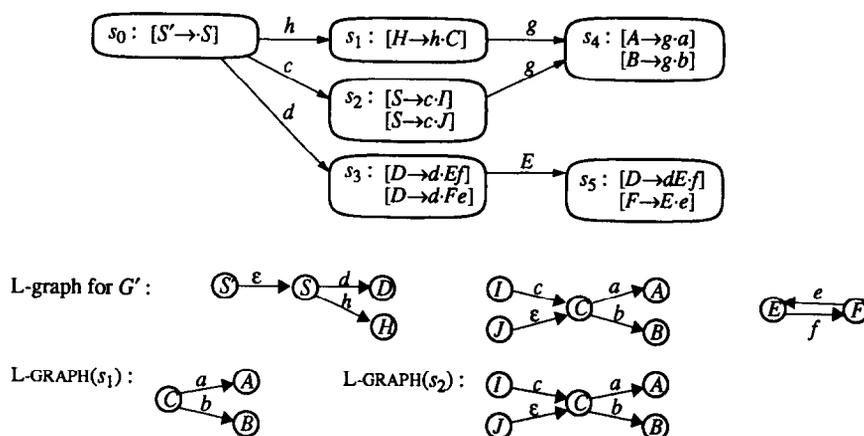


Fig. 7.

**Example 4.9** (Represented state of class 1 and class 2). As a short example, we give a CFG,  $G = (\{S, A, B, C, D, E, F, H, I, J\}, \{a, b, c, d, e, f, g, h\}, \{S \rightarrow cI, S \rightarrow cJ, s \rightarrow Dd, S \rightarrow Hh, I \rightarrow Cc, J \rightarrow C, H \rightarrow hC, C \rightarrow Aa, C \rightarrow Bb, A \rightarrow ga, B \rightarrow gb, D \rightarrow dEf, D \rightarrow dFe, E \rightarrow Ff, F \rightarrow Ee\}, S)$ .

LR(0) automaton for  $G'$ , the augmented grammar of  $S$  is shown in Fig. 7 (we do not give some LR(0) states and transitions).

From Lemma 4.7, we can show that  $s_4$  is a represented state of class 1 as follows: The state  $s_1$  and  $s_2$  are the  $g$ -predecessor states of  $s_4$ . From the  $L\text{-GRAPH}(s_1)$  and  $L\text{-GRAPH}(s_2)$ ,  $C$  can be the common cut symbols which holds the condition of

Lemma 4.7. Then the state  $s_4$  stores  $\{a \text{ PATH}(C, A), b \text{ PATH}(C, B)\}$ ,  $|g|$  and  $C$ , where  $\text{PATH}(C, A) = a, \text{PATH}(C, B) = b$ .

Also from Lemma 4.8, we can show that  $s_5$  is a represented state of class 2 as follows: The state  $s_3$  is the  $E$ -predecessor state of  $s_5$ . Since  $E$  and  $F$  have a transitive L-relation to  $F$ , and the kernel items  $[D \rightarrow d \cdot Ef]$  and  $[D \rightarrow d \cdot Fe]$  in  $s_3$  are represented by the kernel item  $[D \rightarrow d \cdot Ef]$ , the set of kernel items  $\{[D \rightarrow d \cdot Ef], [D \rightarrow d \cdot Fe]\}$  holds the condition of Lemma 4.8. Then the state  $s_5$  stores  $\{f, e \text{ PATH}(E, F)f, e \text{ PATH}(F, F)e, |dE|$  and  $D$ , where  $\text{PATH}(E, F) = \{f, fef, fefef, \dots\} = f\{ef\}^*$ ,  $\text{PATH}(F, F) = \{\epsilon, fe, fefe, \dots\} = \{fe\}^*$ .

Table 1

Experimental result with the programming languages PASCAL, C, ADA and CHILL

	PASCAL	C	ADA	CHILL
number of nonterminals	74	64	238	176
number of terminals	61	85	97	146
number of productions	186	214	459	484
number of LR(0) states ( $\#q$ )	341	367	860	875
number of states with multiple kernel items ( $\#m$ )	72	112	174	282
number of states with one kernel item ( $\#s$ )	269	255	686	593
(ratio 1 = $\#s/\#q$ )	(0.79)	(0.69)	(0.80)	(0.68)
number of immediately represented states with multiple kernel items ( $\#i$ )	63	97	90	242
(ratio 2 = $(\#s + \#i)/\#q$ )	(0.97)	(0.96)	(0.90)	(0.95)
number of represented states of class 1 or class 2 ( $\#r$ )	9	11	64	29
(ratio 3 = $(\#s + \#i + \#r)/\#q$ )	(1.00)	(0.99)	(0.98)	(0.99)
number of non-represented states	0	3	20	11
(ratio 4 = $(\#i + \#r)/\#m$ )	(1.00)	(0.97)	(0.89)	(0.96)

## 5. Experimental result and comparison

We summarize the experimental result in Table 1, which shows how many LR(0) states might be immediately represented states and how many states might be represented states of class 1 or class 2. We also compare the space and time requirement with the work of Choe and Chang [2].

Our experiment was carried out using *YACC* under *UNIX*. The values of  $\#q$ ,  $\#m$  and  $\#s$  were counted from *y.output* file which has the kernel item of all LR(0) state for a given grammar. The values of  $\#i$  was computed from the analysis method presented in Section 3. The values of  $\#r$  was computed from the analysis method presented in Section 4. Note the ratio 3 in Table 1. The ratio 3 is the portion of the states which have only one kernel item( $\#s$ ) or immediately represented states with multiple kernel items( $\#i$ ) or represented states( $\#r$ ) among the all states( $\#q$ ). The ratio 3 means that 98 to 100% of LR(0) states have a *pre-determinable* (determinable in parser generation time) common suffixes, hence it is possible to avoid redundant computation of the common suffixes during parsing time.

We compare the time requirement to compute right context at parsing time with the work of Choe and Chang. Let  $\sigma$  be a left context and  $p$  be the last symbol of  $\sigma$ . We can get prefixes of  $RC(\sigma)$  from the set of vocabulary strings stored in the state  $p$ . To get the rest suffixes of  $RC(\sigma)$ , backtracking to some LR(0) state through the lookback operation is required. The more backtracking leads to the longer right context as is explained in Example 2.2. It is application oriented how long the right context would be. Let  $n$  be the number of backtracking just mentioned, and  $\Phi$  is the ratio 3 in Table 1, and  $|\bar{K}|$  be the average number of kernel items in an LR(0) state. The method of Choe and Chang examines right context with respect to each kernel item, hence examines  $|\bar{K}|$  times for one state. But our method examines only one kernel item (the representative kernel item) for the 98% to 100% of states (states with single kernel item, immediately represented or represented states), and examines  $|\bar{K}|$  times for the remaining 0% to 2% of states

(non-represented states). Hence, in the method of Choe and Chang, the required time for computing right context is  $|\bar{K}|^n$  ( $|\bar{K}|$  is evaluated 1.43 to 1.71 in our experiment). In our method, however, the required time for computing right context is  $(\Phi \times 1 + (1 - \Phi) \times |\bar{K}|)^n$  ( $\Phi$  is evaluated 0.98 to 1.00 as is in Table 1). As an example for a comparison, assume that  $|\bar{K}|$  is 1.57,  $\Phi$  is 0.99, and  $n$  is 3. Then the required time of the method of Choe and Chang is 3.87, whereas that of our method is 1.02. The testing of conditions in Lemma 4.7 or Lemma 4.8 is time-consuming task, but this testing is accomplished during the parser generation time.

We compare the space requirement is computing right context with the work of Choe and Chang. Let  $|Q|$  and  $|\bar{R}|$  be the number of LR(0) states, and the average number of symbols in a production of grammar, respectively. In the work of Choe and Chang, the required space is  $|Q| \times |\bar{K}| \times |\bar{R}|$  for every kernel item in every LR(0) state, and  $|P| \times |\bar{R}|$  for L-graph. L-graph is needed to get paths between two nonterminals. In our method, the required space is  $\Phi \times |Q| \times |VS| + (1 - \Phi) \times |Q| \times |\bar{K}| \times |\bar{R}|$ , and  $|P| \times |\bar{R}|$  for L-graph. Note that  $|VS|$  is the vocabulary string or set of vocabulary strings that have been shown in Lemma 3.2, Theorem 3.7, Lemma 4.7 and Lemma 4.8. Table 2 shows the comparison of space requirement for the programming languages PASCAL, C, ADA and CHILL. The space for L-graph which is common to both methods is excluded.

The values in Table 2 are computed manually through the analysis of the grammar and LR(0) automaton for each grammar [6]. Space reduction is not so remarkable, but our emphasis is on time requirement.

Table 2  
Comparison of the space requirement (unit: elementary space for a short integer or a character)

	PASCAL	C	ADA	CHILL
Method of Choe and Chang	1756	2000	4347	5267
Our method	1687	1975	4474	5170

## 6. Conclusion

We have proposed an efficient method of computing right context for LR-based error repair. To avoid the redundancies in the work of Choe and Chang, we have presented a method of analysing the kernel items in an LR(0) state. By means of pre-determining the common suffixes of right context during parser generation time, it is possible to avoid redundant computation of the common suffixes during parsing time.

The experimental result summarized in Section 5 shows how many states might have a pre-determinable common suffixes. According to the result, 98 to 100% of states might have a pre-determinable common suffixes. A comparison of the time and space requirement with the work of Choe and Chang was also given. Our method has the advantage of time (at parsing time) in computing right context, and does not require more space than that of Choe and Chang.

## References

- [1] A.V. Aho, R. Sethi and J.D. Ullman, *Compilers: Principles, Techniques, and Tools* (Addison-Wesley, New York, 1986).
- [2] K.-M. Choe and C.-H. Chang, Efficient computation of the locally least-cost insertion string for the LR error repair, *Inform. Process. Lett.* **23** (6) (1986) 311–316.
- [3] F. DeRemer and T. Pennello, Efficient computation of LALR(1) lookahead sets, *ACM Trans. Programming Language Systems* **4** (4) (1982) 615–649.
- [4] C.N. Fischer, B.A. Dion and J. Mauney, A locally least cost LR-error corrector, Report 363, Computer Science Dept., University of Wisconsin, Madison, 1979.
- [5] C.N. Fischer, D.R. Milton and S.B. Quiring, Efficient LL(1) error correction and recovery using only insertions, *Acta Inform.* **13** (3) (1980) 141–154.
- [6] M.-S. Jung and K.-M. Choe, On the efficient computation of right context for LR-based repair, Report 77, Computer Science Dept., Korea Advanced Institute of Science and Technology, 1993.
- [7] J.C.H. Park, K.-M. Choe and C.-H. Chang, A new analysis of LALR formalism, *ACM Trans. Programming Language Systems* **7** (1) (1985) 159–175.