

영-한 기계 번역에서 조건 단일화 기반 변환 문법 해석기

김 남수*, 전 현경*, 박 영진**, 이 용석*

* 전북대학교 컴퓨터학과 ** 정인대학 전산정보처리과

Transfer Grammar Compiler Based on Conditional Unification for English-Korean Machine Translation

Nam-Su Kim*, Hyun-Gyung Chon*, Young-Jun Park**, Yong-Seok Lee*

*Dept. of Computer Science, Chonbuk National University

**Dept. of Computer & Information Processing, ChongIn College

요 약

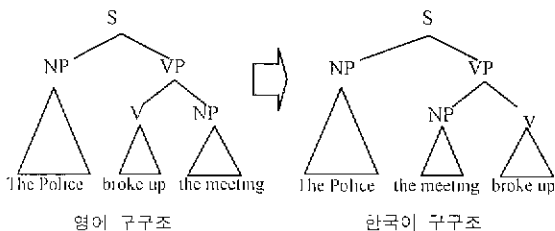
영-한 기계번역에 있어서 변환(Transfer)은 상이한 두 언어의 어순을 생성될 언어의 어순으로 결정하기 위한 변환 규칙에 의해서 영어의 구문 분석된 트리를 생성될 한국어의 구조에 맞게 재구성한다. 이러한 변환은 개발 단계 중 또는 그 후에 언어의 다양한 차이를 수용하기 위해 프로그램 수준에서 잦은 수정을 하게 된다. 이를 극복하기 위해서 본 논문에서는 변환 문법을 이용하여 좀더 체계적이고 확정이 쉬운 문법을 이용하여 변환을 수행하고자 한다. 이를 위해 영어의 구구조 자질들을 쉽게 검사 및 탐색하면서 조건에 맞는 자질들을 나누어 주는 기능을 하는 조건 단일화 연산자를 기반으로 일반적인 문맥자유문법(Context Free Grammar)을 이용한 변환 문법을 제안한다. 또한 변환 문법은 변환 문법 해석기에 의해 변환 프로그램으로 컴파일 됨을 보인다. 이러한 변환 방법은 영-한 변환에 있어서 문법 수준에서 변환 처리를 할 수 있어 변환 프로그램의 개발 및 유지보수에 많은 도움을 준다.

1. 서 론

영-한 기계 번역에 있어서 변환(transfer) 방식이란 분석(analysis), 변환(transfer), 생성(generation)의 세 단계를 통해 번역이 진행 되는데 변환 단계에서는 영어와 한국어의 어휘적 차이 및 구조적 차이를 해소하여 영어에 대한 중간 표현을 한국어의 생성에 적합한 중간표현으로 대응시킨다.

연환방식의 기계번역은 변환 과정에서 두 언어간의 용법이나 의미적 차이를 처리해 줄 수 있기 때문에 언어제계가 상이한 두 언어(예를 들어 영어와 한국어, 영어와 일본어)간의 기계번역에서 용이하다[1,2].

예를 들어 "The police broke up the meeting."의 분석된 영어 구조와 한국어 구조간의 어순을 트리 형태로 살펴보면 다음과 같다.



[그림 1] 예문에 대한 영어-한국어 구구조

변환 단계에서 이러한 상이한 언어군의 어휘와 구조적 차이점을 조절하여 자연스러운 한국어를 얻기 위하여 규칙 기반과

예문 기반, 속어 기반 등의 영한 기계번역이 시도되었다[2,3,4]. 규칙 기반 기계 번역은 대량의 복잡한 변환 규칙의 개발과 복잡성의 어려움을 가지며, 예문에 기반한 기계번역은 예문 데이터 베이스로 부터 원시 언어의 구구조와 가장 근사한 예문을 선택하고, 선택된 예문에 표현된 변환 정보를 이용하여 변환시키는 방법이다. 이러한 방법은 특정 영역(domain)의 번역, 많은 예문에 따른 속도 저하 및 자동 예문 추출 과정의 어려움을 가진다[5]. 또한 언어의 다양한 차이점들을 수용하여 구조 변환을 하기 위해 변환의 개발 단계 또는 개발 후 프로그램 수준에서 많은 수정을 하게 된다.

따라서 본 논문에서는 이러한 변환 처리를 좀더 체계적으로 수행하기 위해 정규화된 톨로서 지원할 수 있는 조건 단일화 기반의 변환 문법 및 변환 문법 해석기를 제안한다. 즉 변환 처리를 위해 변환 문법을 기술한 후 해석기에 의해 생성된 변환 프로그램을 이용하여 변환이 수행되도록 하는 방법이다.

이러한 방법은 변환 시스템에서 어순의 차이를 프로그램 수준에서 처리하지 않고 단지 어순을 고려한 변환 문법만 개발하여 변환 문법 해석기에 의해 생성된 프로그램을 실행함으로써 구조 변환이 가능하므로 기계번역에서 변환 모듈의 개발 및 유지보수에 많은 도움을 준다.

2. 조건 단일화 기반 변환 문법

본 논문에서 제안한 변환 문법은 일반적인 CFG(Context Free Grammar)문법 영식에 조건 단일화 연산자들의 기능들을 확장(augment)시켜 다음과 같은 정형화된 틀을 사용하고 있다.

(<A>->(<C>)

((조건 단일화식)
(조건 단일화식)
...))

<변환 문법 해석기를 위한 문법 표현 형식>

문법기술의 편의를 위해 규칙의 왼쪽 구성소부터 차례로 x0, x1, x2, 등으로 표시한다 위의 문법의 처음 세 요소(<A> → <C>)는 문맥자유형식의 문법 규칙이다 이러한 변환 문법은 입력된 트리의 자질 구조를 어떻게 문법적인 구조로 나누어야 하는지를 기술한 것이다

2.1 조건 단일화 연산자

단일화란 두 자질구조 안의 정보를 포함하는 최소한의 정보를 가지는 자질구조를 찾는 연산을 말한다 또한 조건 단일화는 자질구조 안의 특정 속성값을 동적으로 검사한 후 각 문장 성분이 가지는 문법적 관계를 검사하여 조건에 맞는 단일화만을 수행한다는 전략이다[6]

조건 단일화 표현 방식은 아래와 같으며, 조건 단일화를 수행하기 위한 몇 가지 연산자들의 기능은 다음과 같다

조건 단일화의 예	실제 예
(1) PATH=PATH	(x1 = (x0 np))
PATH=ATOMS	((x1 func) = subj)
(2) PATH==PATH	(x1 == (x0 vp np))
(3) PATH=cATOMS	((x0 subcat) = c T1)
(4) *DEFINED*	((x0 np) = *DEFINED*) ..
(5) *UNDEFINED*	((x0 np) = *UNDEFINED*)
(6) *OR*	*OR*
	((x0 subcat) = c T1)
	((x0 subcat) = c I0))

[그림 2] 조건 단일화의 표현방식

[그림 2]의 실제 예에 있는 x0는 [그림 1]과 같이 분석된 영어 구구조 NP와 VP의 자질 구조로 구성되어 있다.

- (1) = : x1에 값이 존재하면 새로운 값과 단일화를 수행한다 [그림 2]에 나열된 예제 중 '(x1 = (x0 np))'는 변환에서 x0의 자질 값 'NP'를 x1에 복사하는 역할을 하게 된다.
- (2) == : 값을 복사 후 삭제(이동)의 기능을 한다. '(x1 = (x0 vp np))'는 'VP'내의 'NP'를 x1에 복사 후 x0의 해당 'NP'를 삭제한다 결과적으로 x0의 'VP'에는 'V'만 존재하게 된다
- (3) =c : 해당 값이 존재하면 조건은 성공하여 다음에 기술된 연산자들을 수행하게 된다
- (4) *DEFINED* : 해당 값이 존재하면 성공하게 된다. 위의 예제에서 '(x0 np)'가 존재하면 참이 되어 다음에 기술된 단일화 연산자들을 실행하게 된다.
- (5) *UNDEFINED* : 해당 값이 존재하지 않으면 성공하게 된다. 위의 예제에서 '(x0 np)'가 존재하지 않으면 참이 되어 다음에 기술된 단일화 연산자들을 실행하게 된다
- (6) *OR* : 단일화 제약의 이점기술은 단일화를 선택적으로 실행하도록 하는 연산자이다 예를 들이, 위의 예제에서 x0의 subcat(동사의 하위범주)이 자동사(I0) 또는 타동사(T1)인지의 여부에 따라 서로 다른 조건 단일화 연산자들이 실행된다.

이러한 조건 단일화 연산자는 분석된 영어 구구조의 자질들을 쉽게 체크 및 탐색하면서 조건에 맞는 자질들을 나누어 주는 기능을 가지고 있다

2.2 변환 문법

다음은 앞에서 설명한 조건 단일화 연산자와 CFG 문법 형식으로 서론에서 기술된 예제 "The police broke up the meeting"을 처리하기 위한 변환 문법이다.

```

(<SENTENCE> → (<SUBJ> <OBJ> <PRED>))
  ((x0 subcat) =c T1)
    (x1 = (x0 np))
    (x2 == (x0 vp np))
    (x3 = (x0 vp v)) ))

(<SUBJ> → (<DET> <N>))
  ((x1 = (x0 det))
   (x2 = (x0 np)) ))

(<DET> → (@DET))
  ((x1 = x0))

(<N> → (@N))
  ((x1 = x0))

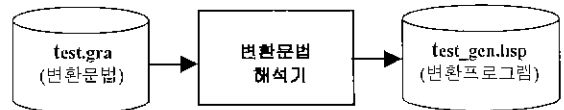
(<PRED> → (@PRED))
  ((x1 = x0))
    
```

[그림 3] 예문의 변환을 위한 문법

위의 문법은 목적으로 명사구를 갖는 타동사구를 변환하는 문법이다 변환의 입력인 구문분석된 결과는 x0인 <SENTENCE>에 입력되며 문장의 주어가 되는 자질인 (x0 np)가 x1으로 복사되고, 목적어가 되는 자질인 (x0 vp np)는 x2로 복사된 후 (x0 vp np)는 삭제된다. 또한 동사인 (x0 vp)에 남아 있는 v를 x3에 복사하게 된다. 이러한 조건 단일화 연산자 및 변환 문법에 의해 영어의 구구조를 자연스럽게 한국어의 구구조로 바꾸게 되는 것이다

3. 변환문법 해석기

[그림 4]와 같이 변환문법 해석기는 변환 문법들을 임파일하여 직접적으로 호출하여 실행할 수 있는 LISP 함수들로 생성한다



[그림 4] 변환문법 해석기

[그림 3]의 문맥자유형식의 문법규칙을 위한 번역 방법은 기존에 여러 연구들에서 발표되어 왔다[7, 8] 그리고 조건 단일화 연산자식의 해석은 다음과 같은 함수들을 통해 변환 프로그램으로 생성되어 진다.

```

(setq equation (caddr statements))
(defun compile-equation (equation)
  (dolist (process unify_equations)
    (let ()
      (case (cadr process)
        (=
          ((pathp (caddr process))
            (P=P ,(list (car process)) ,(list (caddr process))))
            ((or (atomp (caddr process)) (atomp (car process)))
              (P=A ,(list (car process)) ,(list (caddr process))))
            (=C
              (P=CA ,(list (car process)) ,(list (caddr process))))
          )))
      (defmacro P=P (path1 path2) 'setq x (path=path ',path1 ',path2)))
      (defmacro P=A (path1 path2) 'setq x (path=atoms ',path1 ',path2)))
      (defmacro P=CA (path1 path2) 'setq x (path=caloms ',path1 ',path2)))
    )
  )
    
```

[그림 5] 문법 해석을 위한 함수의 예

위의 [그림 5]에서 'statements'는 주어진 하나의 문법 표현이고, 'equation'은 조건 단일화식이 할당될 때, 그 안의 각 단일화 식 중 두 번째 요소가 어떠한 종류의 단일화 연산자 인지 식별하여 해석함수가 생성하게 된다. 여기서 LISP의 'defmacro'

