

표현 패턴에 의한 한국어-영어 기계 번역을 위한 개념 구성

이 호 석
뉴미디어학과 공과대학 호서대학교
hslee@office.hoseo.ac.kr

A Conceptual Framework for Korean-English Machine Translation using Expression Patterns

Ho Suk Lee
New Media Dept. College of Engineering Hoseo University

요 약

본 논문은 표현 패턴(expression patterns)을 이용한 한국어-영어 기계 번역 방법에 대하여 논의한다. 표현 패턴은 한국어-영어 기계 번역을 위하여 의미적이고 표현적인 관점에서 한국어 표현에 적합한 영어 표현을 대응하여 정의하였다. 그리고 새로운 한국어 파싱 방법을 개발할 것을 제안하였다. 한국어 파싱 방법은 한국어의 교착어로서의 특징, 표현 패턴 개념, 문장 분할 개념, 그리고 파싱 과정에 의미 구조를 포함하는 개념 등을 포함하여 개발할 것을 주장하였다. 논문의 마지막 부분에는 간단한 한국어 문법을 정의하여 새로운 한국어 구문 분석 방법의 가능성을 제시하였다.

Abstract

This paper discusses a Korean-English machine translation method using expression patterns. The expression patterns are defined for the purpose of aligning Korean expressions with appropriate English expressions in semantic and expressive senses. This paper also argues to develop a new Korean syntax analysis method using agglutinative characteristics of Korean language, expression pattern concept, sentence partition concept, and incorporation of semantic structures as well in the parsing process. We defined a simple Korean grammar to show the possibility of new Korean syntax analysis method.

1. 서론

한국어-영어 기계 번역을 위해서는 번역의 품질을 결정하는 주요 문장 요소를 파악하는 것이 중요하다. 논문에서 우선 한국어의 문법에 대하여 조사하였다. 참고 문헌 [1]은 한자어 통사 구조, 통사 규칙, 그리고 통사적 사실들의 지배 원리를 규명하고 기술하였다. 참고 문헌 [2]는 한국어의 응용사에 대하여 형태론, 의미/화용론, 통사론적인 관점에서 고찰하여 논의하였다. 참고 문헌 [3]은 한국어의 부사절에 대하여 상세하게 논의하였다. 참고 문헌 [4]는 한국어의 종결 어미와 조동사와 한국어의 복합문에 대하여 논의하였다. 참고 문헌 [5]는 외국인(일본인)의 관점에서 한국어 단어의 형태와 문법에 대하여 흥미를 들여서 간결하면서도 쉽게 설명하고 있다. 참고 문헌 [6]은 한국어 문법 체계에 대하여 폭넓게 설명하고 있다. 참고 문헌 [7]은 한국어와 만주 통구어 특히 그 중에서 어휘 구조 대조 연구가 상세히 설명되어 있다. 참고 문헌 [8]은 우리말의 형태에 대하여 형태소론, 단어의 구획(소용)과 파생법, 단어 합성법, 굴곡법(어미 활용), 추밀굴곡법(조사 활용) 등으로 나누어서 설명하고 있다. 한국어 단어와 어미의 형태를 이해하기 위해서는 반드시 참고하여야 될 책으로 생각된다. 참고 문헌 [9]는 10세기 말엽까지의 연대를 그 한국어의 시기로 정하고 이 시기 이전의 국어에 대하여 논의하고 있다. 1.1.3 삼국시대의 언어에서는 여러 단어와 지명들 예로써서 삼국의 언어들 사이에는 방언적인 차이 이외에 그다지 차이를 보이지 않았다고 설명하고 있다. 참고 문헌 [10]은 한국어의 문법에 대하여 음운론과 자소론, 형태론, 통사론과 의미론, 텍스트론과 타노는 개념의 분류 구조에 대하여 논의한다. 이 책의 내용들은 의미 표지(semantic feature)를 작성하는데 좋은 참고가 될 것으로 보인다. 참고 문헌 [12]는 한국어 문장 종류별로 문장의 구조, 화용론, 의미, 문장의 논리성, 문장의 유형에 따른 기능 등을 논의하고 있다. 참고 문헌 [13]은 한국어 문법에 대하여 체계적으로 설명하여 쉽게 살펴 볼 수 있도록 하였으며, 참고 문헌 [14]는 대조언어학의 관점에서 한국어와 스페인어 문장들을 대조하여 문법, 문장 형태, 그리고 의미론을 논의하고 있다. 참고 문헌 [15]는 한국어에 대하여 단어 형태의 구획, 굴곡법, 통사론, 통사론적 통사론(단어 자체로서 조사의 활용)이 어떻게 나타내는 방법으로 나누어서 논의하고 있다. 참고 문헌 [16]은 한국어 문법에 대하여 형태론과 통사론으로 구분하여 논의를 하고 있다. 참고 문헌 [17]은 통사론적 통사론과 통사론적 통사론(단어 자체로서 조사의 활용)이 어떻게 나타내는 방법으로 나누어서 논의를 하고 있다. 참고 문헌 [18]은 국어 명

사의 의미 체계, 의미 구조, 의미 구조의 구성, 명사 연결 구성과 의미 구조의 관계에 대하여 논의하고 있다. 우선 의미 영역에 대하여 분류하였는데, 의미 영역을 실제와 양식으로 분류하였고, 양식은 다시 사태와 관계로 분류하였으며, 사태는 다시 사건과 상태로 구분하였다. 의미 구조의 구성에서는 의미 네트워크(semantic network)에 대하여 논하였으며, 명사의 의미 구조인 속성 구조(qualia structure) 그리고 의미 자질(semantic feature)에 대하여도 논의하였다. 명사의 의미 구조 참고에 유용한 도서이다. 참고 문헌 [19]는 국어의 자동사를 비행위성 자동사와 행위성 자동사로 분류하여 논의하였다(80쪽 <표 1>). 비행위성 자동사는 대상자동사, 소재자동사, 비교자동사, 변성자동사, 대칭자동사, 피동자동사, 심리자동사로 구분하였다. 대상자동사의 예로는 “끓다”, “눕다” 등을 제시하였다. 소재 자동사의 예로는 “끓다”, “속하다” 등을 제시하였다. 행위성 자동사는 행위자동사, 위치자동사, 이동자동사, 대칭자동사로 구분하였다. 행위자동사의 예로는 “날다”, “뛰다” 등을 제시하였다. 위치자동사의 예로는 “이르다”, “달다” 등을 제시하였다. 그 밖에 자동과 타동, 행위성과 비행위성이 미분화된 동사들도 있는데 이러한 동사들은 중립동사로 분류하였다. 중립동사의 예로는 “가시다”, “들생거리다” 등을 들었다. 마지막으로 결론 부분에는 논의한 자동사들의 자세한 분류표가 제시되어 있다(367쪽). 참고 문헌 [20]은 비교적 전통적인 관점에서 국어의 문법에 대하여 논의하였다. 국어 문법학 형태론, 품사론, 명사론, 불규칙활용론, 의존용언론, 서법과 양태론 등으로 설명하였다. 참고 문헌 [21]은 한국어의 화용론(pragmatics)에 대하여 서술하였다. 화용론은 언어의 사용에 대한 것으로 음운론적, 형태론적, 통사론적, 의미론적, 담화/텍스트론적, 사회언어학적, 인지심리학적으로 화용론을 분류하여 논의하였다. 3.5절에는 대화의 유형이 분류되어 있으며, 7장에는 화용론에 대한 종합적 분석의 예가 제시되어 있는데, MBC 100분 토론(2006년 3월 6일) “과자 유해 논란” 토론 중에서 처음 30분 동안 대화 내용이 제시되어 있다. 이 대화 텍스트에 대하여 음운론적, 형태론적, 통사론적, 담화론적 분석이 제시되어 있다. 담화론적 분석에서는 95개의 담화 표지가 사용되었다고 제시되어 있다. 참고 문헌 [22]는 국어 조사 중에서 “-에”와 “-로”의 용법에 대하여 매우 상세하게 설명되어 있다. 한국어의 번역을 위해서 참고할 만한 도서이다. 참고 문헌 [23]은 한국어 문장 분석론적 관점의 분석을 위하여 거시구조의 틀과 미시구조의 틀을 제시하고 있다. 거시구조의 틀에는 어미 구조, 어미와 어미의 상관성, 한정조사와 어미의 상관관계, 부사기능어(문장 부사, 양태 부사)와 어미의 상관관계 등이 포함된다. 미시구조의 틀에는 언어구문성, 관용구 구성, 구동사(phrasal verb), 술어와 논항 구조(argument structure), 문법관계표지(격조사, 격어미) 등이 포함된다. 한국어의 기계 처리를 위해서는 참고할 만한 도서이

if (V = 선행동사(그밖의동사)) &
 부사(이미),
 영어 번역 : N₁ wonder
 whether N₂ already V N₃
 (sent (np(subj(N₁)) vp(vt(wonder)
 clause(whether(whether)
 np(subj(N₂)) vp(adv(already)
 v(V) obj(N₃)))

따라서 위의 예들을 보면 적절한 영어 번역을 위해서는 한국어 동사의 하위범주화 패턴이 매우 중요하며, 이러한 하위범주화 패턴 정보가 반드시 사전에 등재되어야 한다. 결국 하위범주화 패턴은 동사나 형용사 활용의 모든 패턴을 나타내는 것이다. 따라서 하위범주화 패턴은 절문(clause)을 포함할 수 있다. 이런 경우에는 절을 미리 영어로 번역한 다음에 천주적으로 만들어진 영어 문장에 삽입하여 전체 영어 문장을 구성하는 경우도 있다.

4. 보문

한국어-영어 기계 번역을 위해서는 보문에 대한 이해가 필요하다. 예를 들어, 다음의 문장들을 보자. 이 문장들은 보문에 대한 좋은 예이다.

- (문장 4) 그는 학자가 되었다.
- (문장 5) 코끼리는 코가 길다.
- (문장 6) 이 나무에 새가 많이 뜬다.

(문장 4)를 영어로 번역하면 "He became a scholar." 가 되고, (문장 5)를 영어로 번역하면 "The nose of elephant is long." 이 되고, (문장 6)을 영어로 번역하면 "A new leaf is coming up on this tree." 가 될 것이다. 따라서 한국어에서는 동일한 구문 구조를 가지고 있지만, 영어로 번역하면 완전히 다른 구문 구조로 바뀌어 나타난다. 따라서 한국어 보어 문장의 적절한 번역에는 동사와 그리고 동사 앞에 위치하는 두 개의 명사 사이의 관계가 중요하다고 하겠다. 따라서 관련된 언어적 정보는 동사와 명사 두 부분에 모두 등재되어야 한다. 이 경우에도 파서는 문장의 구문 구조를 파악하는 역할을 수행하고, 실질적인 문장의 의미 구조는 사전에 등재되어 있는 동사나 형용사의 하위범주화 패턴과 문장의 주격 명사와 보격 명사의 관계에 의하여 결정되어야 할 것이다. 즉, 위의 (문장 4)에서는 명사 "그"와 명사 "학자"의 관계를 알아야 하며, (문장 5)에서는 명사 "코끼리"와 명사 "코"의 관계를 알아야 한다. 이 두 예에서는 명사의 조사는 주격 조사라는 것을 주의하여야 한다. 따라서 명사는 영어로 번역된 문장에서도 주어가 된다. (문장 6)에서도 명사 "나무"와 명사 "잎"의 관계를 알아야 한다. 그리고 이 문장에서는 명사의 조사는 처소격임을 알아야 한다. 이 경우에는 영어로 번역하였을 경우에 전치사구가 될 수 있다. 이 예는 문장의 구문과 의미 정보이외에도 보어 문장에서 조사의 역할이 중요함을 나타낸다. (문장 4)의 예에서 알 수 있듯이, 동사 "되다"에는 사전에 다음 정보가 등재되어 있어야 한다.

- 되다 : 동사
- 하위범주화 패턴1 : (N₁M_{subj}, N₂M_{subj})
- 영어 번역 : N₁ become N₂

그리고 (문장 5)의 형용사 "길다"에는 다음 정보가 등재되어 있어야 한다.

- 길다 : 형용사
- 하위범주화 패턴1 : (N₁M_{subj})
- 영어 번역 : N₁ be long.
- 하위범주화 패턴2 : (N₁M_{subj}, N₂M_{subj})
- if(yes=Noun(N₁,N₂)), 영어 번역 : N₂ of N₁ be long.

(문장 6)의 동사 "뜬다"에는 다음 정보가 등재되어 있어야 한다.

- 뜬다 : 동사
- 하위범주화 패턴1 : (N₁M_{subj})
- 영어 번역 : N₁ come up.
- 하위범주화 패턴2 : (N₁M_{adv(에)), N₂M_{subj})}
- 영어 번역 : N₂ come up (Prep_{adv(에)) N₁}

보문에 대한 하위범주화 정보도 사전에 등재되어 있어야 문장에 대한 구문 분석이 끝나고 필요한 구문 정보를 사전에서 추출하여 영어로의 번역에 사용할 수가 있다.

5. 부사

다음에 한국어 부사절에 대하여 논의해 보자. 우선 다음 문장을 보자.

- (문장 7) 물이 너무 맑으면 고기가 적게 뜬다.
- (문장 8) 그는 글도 잘하며, 말도 잘한다.
- (문장 9) 그 아이가 형과는 달리 사교에 능하다.

(문장 7)을 영어로 번역하면, "If water is too clear, fishes little gather." 가 될 수 있다. 이런 번역을 가능하게 하는 문장 요소는 부사적 어미 "-면"이다. "-면"에 의하여 "if"가 도입된다. 동사 "뜬다"에는 다음 언어 정보가 등재되어 있어야 한다. 그리고 명사 "고기"에는 다음 언어 정보가 등재되어 있어야 한다.

- 뜬다 : 동사
- 하위범주화 패턴1 : (N₁M_{subj})
- 영어 번역 : N₁ gather.
- 고기 : 명사
- 관련명사 : 물
- 관련영어 번역 : fish
- 관련명사 : 고기
- 영어 번역 : meat

(문장 8)의 영어 번역은 "He writes well and talks well." 이 될 수 있다. 이 경우에는 "잘하다" 동사에 다음 언어 정보가 등재되어 있어야 한다.

- 잘하다 : 동사
- 하위범주화 패턴1 : (N₁M_{obj}, N₂M_{obj})
- if (N₂ = 글), 영어 번역 : N₁ write well
- if (N₂ = 말), 영어 번역 : N₁ talk well
- if (N₂ = 피아노), 영어 번역 : N₁ play piano well
- if (N₂ = 노래), 영어 번역 : N₁ sing well

(문장 9)는 "The boy is sociable unlike his old brother." 번역이 될 수 있다. 이 문장의 경우에는 동사 "능하다"와 "다르다"는 다음과 같이 되어야 한다.

- 능하다 : 동사
- 하위범주화 패턴1 : (N₁M_{subj}, N₂M_{place})
- if (N₂ = 사교), 영어 번역 : N₁ be sociable
- if (N₂ = 사격), 영어 번역 : N₁ shoot well
- if (N₂ = 일), 영어 번역 : N₁ be good at N₂
- 다르다 : 동사
- 하위범주화 패턴1 : (N₁M_{subj}, N₂M_{adv})
- if (문장 = 단문), 영어 번역 : N₁ be different from N₂
- if (문장 = 부사절), 영어 번역 : N₁ unlike N₂

이제까지 한국어-영어 기계 번역에 있어서 한국어의 조사, 동사, 부사, 그리고 (문장 7)에서는 어미가 중요한 역할을 한다. 이것을 제시하였다.

6. 한국어 기본 문장 문형

다음은 한국어의 기본 문장 문형들이다[6].

1. S -> NP VP
2. NP -> (Det) N₁M_{subj}
 (Det) N₂M_{obj}
 (Det) N₂M_{com}
 (Det) N₂M_{obj} N₃M_{com}
3. VP -> (Adb) ± V AUX
 (Det) N₁M_o, AUX
 (Det) N₂M_{obj} (Adb) ± V AUX
 (Det) N₂M_{com} (Adb) ± V AUX
 (Det) N₂M_{obj} (Det) N₃M_{com} Adb + V AUX
4. S -> (Det) N₁M_{subj} ± V AUX
 (Det) N₁M_{subj} (Det) N₂M_o AUX
 (Det) N₁M_{subj} (Det) N₂M_{obj} (Adb) ± V AUX
 (Det) N₁M_{subj} (Det) N₂M_{com} (Adb) ± V AUX
 (Det) N₁M_{subj} (Det) N₂M_{obj} (Det) N₃M_{com} (Adb) ± V AUX

여기서, S는 문장, NP는 명사구, VP는 동사구, Det는 관형어, Adb는 부사어, M은 지표, M_{subj}는 주격지표, M_{obj}은 목적격지표, M_{com}은 보격지표, M_o는 서술격지표, N₁은 첫 번째 명사(주어), N₂는 두 번째 명사(목적어나 보어), N₃는 세 번째 명사(보어), +V는 동사(+행위), -V는 동사(-행위), AUX는 어미(Auxiliary), ()는 임의적인 꾸밈말 등을 의미한다.

한국어의 기본 문형은 구절구조 문법에 입각하여 제시되었다. 문형을 보면 한국어의 문장은 명사구와 동사구와 구성되어 있다. 그리고 명사구는 (1) 관형어 + 주어 역할을 하는 명사 + 주격지표(조사), (2) 관형어 + 목적어나 보어 역할을 하는 명사 + 목적격이나 보격 지표(조사), 혹은 (3) 관형어 + 목적어나 보어 역할을 하는 명사 + 보어 역할을 하는 명사 + 보격지표(조사)로 구성되는 것을 알 수 있다. 동사구는 (1) 부사어 + 동사 + 어미, 혹은 (2) 관형어 + 목적어나 보어 역할을 하는 명사 + 부사어 + 동사 + 어미가 반복되어 구성되는 것을 알 수 있다. 한국어 파서를 구성할 때에는 이러한 관계를 파악하여 구성하여야 할 것이다.

7. 표현 패턴에 의한 한국어-영어 기계 번역

지금까지는 하위범주화 패턴 방식의 한국어-영어 기계 번역에 대하여 논의하였다. 즉, 한국어 표현 문장을 동일한 의미를 나타내는 영어 표현 문장으로 변환시키기 위하여 트랜스퍼 사전에 수록된 한국어 동사와 형용사의 하위범주화 패턴을 활용하였다.

그러나 한국어 표현 패턴을 의미가 동일한 영어 표현 패턴과 1:1로 대응시켜서 표현 패턴 사전을 구성하고 이를 기계 번역에 활용하는 방법을 생각할 수 있다. 다음의 한국어 문장의 예를 보자.

- (문장 10) 그날 아침을 먹고 출근하였다.
- (문장 11) 그날 아침을 거르고 출근하였다.
- (문장 12) 그날 아침 식사에 대하여 논의하였다.
- (문장 13) 그날 아침 개개에 대하여 조사하였다.
- (문장 14) 그날 아침 식사를 조사하였다.
- (문장 15) 그날 더위를 먹었다.
- (문장 16) 그날 더위를 심하게 먹었다.

(문장 10)을 영어로 번역하면 "He went to office after having breakfast."가 될 것이다. 그러나 한국어 문장에는 실제로 "사무실(office)"이라는 단어도 없고 "아침밥(breakfast)"이라는 단어도 없다. 한국어에서 "출근하였다"라는 단어는 "사무실에 일하러 간다"는 의미이며, "아침을 먹었다"라는 동사구 속에는 "아침 식사를 하였다"가 포함되어 있다. 따라서 동일한 의미를 나타내는 한국어 표현과 영어 표현을 표현 패턴이 다른 것이다. 따라서 이 경우엔 "go to office after having breakfast"라는 영어 표현 패턴 전체로 번역하면 될 것이다. (문장 11)의 경우에는 Today, he went to office without having breakfast."로 번역이 될 것이다. 한국어 타동사 "거르다"를 영어 전치사 "without"으로 번역하면 된다. (문장 12)의 경우에는 역사에 대하여 "about"으로 번역할 수 있다. 이 경우에는 "에 대하여"라는 구문에서 "about"이라는 표현 패턴이 있다. 이 표현 패턴은 영어로 "about"이라는 하나의 전치사로 번역이 될 수 있다. 그러나 (문장 13)의 경우에 "We checked the cause of accident."으로 번역할 수 있다. 이 경우에는 "에 대하여"로 번역하면 안 된다. (문장 14)의 경우도 영어로 번역할 수 있다. 따라서 한국어의 "에 대하여" 표현 패턴이 항상 영어의 "about"으로 번역되는 것은 아니라는 것을 알 수 있다. 영어로 번역되는 동사가 타동사인 경우엔 표현 패턴이 지배하는 문장성분의 문법적 성분과 기능에 따라 표현 패턴과 함께 표시함으로써 해결할 수 있다. 사전 내용의 비교를 위하여 하위범주화 패턴 방식과 표현 패턴 방식을 함께 시하였다. 영어에 한국어의 "조사하다"와 유사한 의미의 단어에는 check, investigate, examine, inquire into, probe 이 있다고 할 수 있다. 아래의 예에서는 이들 단어들 모두 두 문장을 생성하는가 중요한 문제가 될 수 있다. 이것을 문법적 성분과 관련이 있다고 할 수 있다. 예를 들어, 문체를 평범한 문체와 격조를 갖춘 문체 등으로 구분하여 생각할 수 있다. 영어 표현 "inquire into"는 다른 영어 표현에 비하여 비교적 격조 있는 표현을 갖춘 표현이라고 한다. 아래 사건의 내용에는 표현 패턴에 (N₂(noun obj))라고 표시하여 표현 패턴의 지배를 받는 N₂가 명사(noun)이고 목적어(obj)라는 것을 함께 나타내었다. 이것은 단수와 복수 일치와 시제 때문에 필요하다.

- 논의하다 : 동사
 - . 하위범주화 패턴1 : ((N₁M_{subj} N₂M_{adv(에)} V(대하여)))
 - . 영어 번역 : N₁ discuss about N₂
 - . 표현 패턴1 : ((N₁M_{subj} N₂Pattern(에 대하여)))
 - . 영어 번역 : N₁ discuss about N₂
- 조사하다 : 동사
 - . 하위범주화 패턴1 : ((N₁M_{subj} N₂M_{adv(에)} V(대하여)))
 - . 영어 번역 : N₁ check N₂
 - . 하위범주화 패턴2 : ((N₁M_{subj} N₂M_{obj(=)})))
 - . 영어 번역 : N₁ check N₂
 - . 표현 패턴1 : ((N₁M_{subj} N₂Pattern(에 대하여)))
 - . 영어 번역 : N₁ check N₂
 - . 표현 패턴2 : ((N₁M_{subj} N₂Pattern(에 대하여)))
 - . 영어 번역 : N₁ investigate N₂
 - . 표현 패턴3 : ((N₁M_{subj} N₂Pattern(에 대하여)))
 - . 영어 번역 : N₁ examine N₂
 - . 표현 패턴4 : ((N₁M_{subj} N₂Pattern(에 대하여)))
 - . 영어 번역 : N₁ inquire into N₂
 - . 표현 패턴5 : ((N₁M_{subj} N₂Pattern(에 대하여)))
 - . 영어 번역 : N₁ probe N₂

(문장 15)는 앞에서 트랜스퍼 사전 방식으로 다루었다. 그러나 "더위를 먹다" 전체를 표현 패턴으로 간주하여 "suffer from a hot weather" 표현으로 표현 패턴 방식으로 번역할 수 있다. (문장 16)의 경우는 "더위를 심하게 먹다" 전체를 표현 패턴으로 간주하여 "suffer from a hot weather severely"로 번역할 수 있다.

표현 패턴은 명사구, 동사구, 혹은 그 밖의 문장 표현에 존 재할 수 있으며 관형어도 넓은 의미에서는 표현 패턴으로 간주할 수 있다. 따라서 기존의 하위범주화 패턴에 의한 기계 번역 방식을 확장하여 표현 패턴에 의한 기계 번역 방법을 구성할 수 있다. 그러나 표현 패턴에 의한 기계 번역 방법의 효용성은 대하에서는 추후 연구와 조사가 필요하다. 표현 패턴 사전은 기존의 트랜스퍼 사전에 구성할 수가 있으며 또한 별도의 사전으로 구성할 수도 있다. 트랜스퍼 사전의 하위범주화 패턴은 문법적 관점에서 좀 더 치중하여 두 언어 간의 차이점을 다룬 것이라면, 표현 패턴은 표현과 의미적인 관점에서 좀 더 치중하여 두 언어 간의 차이점을 다룬 것이라고 할 수 있다. 표현 패턴의 장점은 번역에 있어서의 효과도 매우 크다는 것이다. 특히, 한국어에 있어서의 조사와 어미 활용이 상당히 복잡하고 표현 패턴은 이러한 한국어의 독특한 언어적 현상을 하나의 단위로 인식하여 다루므로, 영어를 대상으로 개발된 하위범주화 패턴보다도 기계 번역에 있어서 한 한국어에 더욱 적합한 개념이라고 할 수 있다. 예를 들어, 앞의 (문장 29)에서 사용한 "생각하다"도 이와 같은 한국어 표현 명사와 조사와 이루어진 간단한 표현인데도 영어로 번역하면 "although thinking hard"로 수어를 포함하면, "although I think hard"가 된다. 이 표현 앞에서는 아진 표현 패턴도 사실은 상 황에서 편의상 관형어만 표시한 것으로 다루었지만 사실은 한국어에서 관형어는 아니므로, 따라서 이러한 한국어 표현은 하나의 표현 패턴으로 간주하여 다루면 좋을 것이다. 그러나 문장에 의한 방법의 단점도 최종적으로 한국어 문장과 영어 문장을 구성하는 어휘의 차이 혹은 그 밖에 문체의 문제 등을 해결해 줄 문법 검사 및 변환에 대한 설명이다.

[표현 패턴의 정의] : 표현 패턴은 한국어-영어 기계 번역에 있어서 표현과 의미의 차이를 나타내는 한국어 단어 문자열의 최소 단위로 정의한다. 즉, 한국어에서 중요한 문법적 기능을 수행하는 조사와 어미 등을 다른 단어들과 함께 표시하여 동사의 하위범주화 패턴 형태로 나타낸다. 표현 패턴은 표현의 차이를 나타내는 단위로서 기계 번역에서 구문 변환과 의미 전달의 기능을 수행한다.

표현 사전을 사용하여 기계 번역을 수행한다는 관점에서는 기계 번역 수행 과정을 다음과 같이 생각해 볼 수 있다.

- (1) 한국어 문장을 구문 분석하여 문법을 검사하고 구문 트리를 구성한다. 문장의 표현 문자열을 표현 패턴 사전에서 찾아서 한국어 표현에 대응하는 영어 표현을 찾는다. 표현 패턴 사전은 동사, 명사, 형용사, 부사 순서로 찾는다. 표현 패턴 사전에서 찾을 수 없는 표현은 일반 한영 사전을 찾아서 한국어 단어에 대응하는 영어 단어를 찾는다. 최종적인 영어 문장 생성을 위하여 생성된 영어 표현 패턴의 문법적인 기능 정보는 유지하여야 한다.
- (2) 영어 생성기는 어순 문제, 시제 문제, 그리고 단수와 복수 일치 문제 등을 해결한다.
- (3) 영어 생성기는 최종적인 영어 문장을 생성한다.

또한 표현 패턴에 의한 기계 번역 방식을 기존의 트랜스퍼 방식을 확장하여 구성할 수 있다. 처음에는 트랜스퍼 사전에 입 각하여 한국어 문어의 분석을 시도하고 기계 번역을 수행한다. 만약 트랜스퍼 사전에서 적절한 영어 번역이 이루어지지 않으면, 다음에는 표현 패턴 사전에 의하여 다시 표면 문장의 문자열 비교를 통하여 기계 번역을 시도한다. 반면에 트랜스퍼 사전에 의한 기계 번역 방식을 지양하고 전적으로 표현 패턴 사전에 의한 기계 번역을 시도할 수 있다. 그러나 표현 패턴 방식에 의한 기계 번역이 성공을 거두기 위해서는 상당한 규모의 한국어 표현 패턴과 영어 표현 패턴 사이에 대한 대응 연구가 선행되어야 할 것이다.

8. 최근의 기계 번역 연구

근래에 중국어-영어 기계 번역 연구가 중국어권에서 비교적 활발히 진행되고 있다[33][34][35][36]. 국내에서는 2002년에도 전자통신연구소(ETRI)에서 한국어-중국어 기계 번역 시스템이 개발되기도 하였다[37]. 또한 영어-일본어 기계 번역 영역도 진척되고 있다[38]. 미국에서는 2005년도에 통계적인 방법을 사용한 한국어-영어 기계 번역 시스템이 개발되기도 하였다[39].

상품화된 기계 번역 시스템은 IBM[40], Microsoft[41], Systran[42] 등의 기관에서 제공하고 있다. 그리고 그 밖의 기관에서 기계 번역 관련 자료와 연구용 소프트웨어를 제공하고 있다[43][44][45]. 영어의 경우에는 [46][47]에서 통계적 방법을 사용하여 구현된 구문 분석기(statistical syntax analyzer)

에 대한 자료와 소프트웨어를 제공한다. 본 논문에서는 논의한 패턴에 의한 기계 번역은 새로운 시도로서 단어 간의 문맥 의존성(context-sensitivity)을 최대한으로 감소시키고 문맥 자유성(context-freeness)을 증가시키려는 시도가 있다. 문맥 자유성이 증가되면 자연언어가 문맥 자유(context-free) 언어가 되어 기계에 의한 처리가 용이해지게 된다. 그러나 텍스트 전반에 걸친 문맥 구조, 의미 구조, 담론 구조 그리고 텍스트의 배경이 되는 의식 구조 등을 기반한 기계 번역에 대한 연구도 논의될 필요가 있다.

9. 문장 분할 개념

문장 분할(sentence partition 혹은 segmentation) 개념은 한국어 문장 구성의 특징상 “주어+동사” 구문이 여러 번 반복되어 나타나는 경우가 많기 때문에 필요하다고 생각한다. 이 것의 비교적 평면적인 구문 구조에 기인한다고 할 수 있다. 예 들어 “나는 아침에 일어나서 세수를 하고 밥을 먹고 차를 타고 학교에 가서 공부를 하고 숙제를 하고 공부가 끝난 후에 집에 가서 공부를 하고 숙제를 하고 저녁을 먹고 잠을 잤다.”와 “공부가 끝난 후에 집에 와서 깨닫는 것이 씩었다. 그 다음에 숙제를 하고 저녁을 먹고 하루도 쉬지 않았다. 하지만 분할하는 과정에서 생략되는 요소들이 있고, 반면에 추가되는 요소들은 의미 구조에 포함하여 나타나고 다루어 할 것이다. 위 4개의 문장을 영어로 번역을 하면 “I got up early in the morning, washed my face, and had a breakfast.”, “Then, I went to school by car, and studied and took part in sports activity.”, “After finishing school, I went back home and took a shower.”, “I did my homework, had supper, and went to sleep.” 이 예에서도 (1) “차를 타고 학교에 가다-go to school by car”, (2) “운동에 참가하다-take part in sports activity”, (3) “공부가 끝난 후에-after finishing school”, (4) “깨끗이 씻다-take a shower”, (5) “잠을 자다-go to sleep” 과 같은 한국어와 영어 사이에 대응하는 표현 패턴을 발견할 수 있다. 문장 분할 개념과 관련된 것으로 부분 파싱(partial parsing) 개념이 있다. 이 개념도 문장의 길이가 길어져서 구조가 복잡한 경우에는 문장 전체를 파싱할 수가 없고 문장의 부분만을 파싱할 수 있는 것이다. 따른서 부분 파싱 개념은 부분 파싱만이 가능한 문장을 짧은 여러 개의 문장으로 분할하여 완전하게 파싱하는 기법과 관련이 있다. 문장 분할 개념은 파싱의 복잡도를 낮추는데도 도움이 될 것으로 생각한다.

10. 새로운 한국어 구문 분석 시스템

참고 문헌 [2][3][4][6][10][13][15][17][18][19]와 [21]~[28]에는 현대 한국어의 문법이 상세하게 논의되어 있다. 현대 한국어 문법에서는 한국어의 교차어적인 특징을 분명히 이해하고 이를 말동지를 통하여 밝히는 방향으로 연구가 진행되는 경향이 있다. 참고 문헌 [29]는 유용한 자료로서 연구의 목표로서 “한국어 문장 분석에 설정되어야 할 기본적인 통사 단위와 그것이 결합되어 형성되는 동사 구성을 범주화하고 체계화하는 것이다.”라고 하였다. 5장에서 조사의 역할과 의미의 활용 등에 대하여 거의 완벽한 자료 정리 결과를 제시하고 있다. 참고 문헌 [48]~[53]에는 90년대의 한국어 구문 분석 시스템에 대한 연구가 제시되어 있다. 참고 문헌 [52][53][54][60]는 의존 문법에 기반을 둔 한국어 구문 분석에 대한 연구가 논의되었다. 참고 문헌 [54]~[65]는 2000년대의 연구 결과를 제시하고 있다. 특히 최근 들어서는 한국어 구문 분석에 대한 연구가 다시 진행되고 있다는 것을 주목할 만하다[58]~[65]. 참고 문헌 [62]는 한국어의 병렬 명명 구조의 구성을 나타내는 확률적 모델을 제시하였다. 모델은 병렬 명명 구조의 대칭성(symmetry)과 상호교환성(interchangeability)에 근거하여 구성하였다. 참고 자료 [63]의 한국어 구문 분석 시스템은 주목할 만하다. 이 시스템은 의존 문법에 기반을 두고 개발되었다. 시스템은 문장 단위로 구문 분석을 실행하여 결과를 보여준다. 참고 문헌 [64][65]는 시스템 [63]에 관련된 논

문들로서 확률적 의존 문법을 사용하여 한국어 구문 분석기를 구현한 내용을 기술하고 있다. 최근에는 어휘 정보(lexical information)가 구문 분석의 정확성 향상에 일반적으로 미치는 것만큼은 기여하는 바가 없다는 연구 결과도 있다[64][66][67]. 심한지어 참고 문헌 [68]에서는 깊고 복잡한 문법 분석이 간단한 문법 분석에 비하여 생각보다는 효과적이라는 논의도 있다. 아마 이것은 어차피 자연언어는 의미에 의하여 좌우되는 부분이 많기 때문에 깊은 문법 분석 자체가 완벽할 수 없고 또한 큰 역할을 하지 못한다는 의미일 것이다. 참고 문헌 [29]의 22쪽에는 한국어 구문 분석에 있어서 HPSG나 LFG 관련 연구가 줄어든 원인으로 “HPSG나 LFG 한국어에서 조사나 어미가 가진 행의 정격을 부각시키기 어렵다.”고 설명하고 있다. 이것은 한국어는 교차어이기 때문에 서구 언어(영어)를 설명하기 위하여 개발된 문법 이론은 한국어의 고를 현상 설명하는 데는 적합하지 않다는 것을 의미한다. 따라서, 새로 개발하는 한국어 구문 분석 시스템의 요구 사항은 다음과 같이 정리할 수 있다. 즉, 새로운 시스템은 언어학적인 측면에서는 한국어의 교차어적인 특징(참고 문헌 [24]에서 제시한 특징 포함)을 충실히 고려하여 개발되어야 할 것이다. 교차어적인 특징은 보통 수서가 자유로운(free order) 그리고 핵이 뒤에 있는(head-final) 등의 용어로 표현된다. 이것은 당연한 것이다. 왜냐하면, 문법 기능을 나타내는 조사와 어미가 뒤에 있음으로서 수서가 자유로운 언어가 될 수 있는 것이기 때문이다. 참고 문헌 [29]의 연구 결과 자료는 충분한 자료로서 고려되어야 할 것이다. 여기에 참고 문헌 [21][27]의 확률적 개념도 고려되어야 할 것이다. 즉, 저자는 한국어의 문장 구성 체계는 의미 표현 단위(MEU, meaning expression unit)의 구성과 이들 간의 의존(연관) 관계로 되어 있다고 생각한다. 그리고 MEU들의 문법적인 기능은 조사나 어미의 활용이 나타낸다고 생각한다. MEU는 간단하게는 하나의 단어에서부터 구절 구조(phrase structure), 그리고 복잡한 문장까지 될 수 있다고 생각한다. 이러한 MEU를 참고 문헌 [54]에서는 단순히 단위(chunk)라고 표현하였다. 그러나 논문 [54]에서 설명한 단위(chunk)보다는 훨씬 복잡한 다양한 구조를 나타낼 수 있다. 따라서 MEU의 구성과 MEU들 간의 의존(연관) 관계를 밝히면 한국어에 대한 문법 체계 구성할 수 있을 것임은 가정해 본다. 단어의 품사(불완전)를 중요하게 고려하지 않는다는 것이다. 그리고 의존(불완전)을 마다간 단위로 대 한도 매우 중요하다. 예를 들어서 다음의 문법을 마다간 단위로 나타낼 수 있는지를 알 수 있다. DEP는 의존 관계만을 나타내며 표현 패턴과 문장 분할을 구현하기가 용이할 수 있다.

MEU -> Det	//	관사	조사
MEU -> Njosa	//	명사	조사
MEU -> Vomi	//	동사	어미
MEU -> ADJomi	//	형용사	어미
MEU -> ADV + Vomi	//	부사 + 동사	어미
MEU -> MEU + MEU	//	문법적 결합	
DEP -> MEU + MEU	//	전체 문장	

다음은 위의 문법을 적용하여 문장의 구조를 나타낸 것이다.

(문장 10) 그는 아침을 먹고 출근하였다.

MEU₁ -> MEU(아침을) + MEU(먹고)
 MEU₂ -> MEU(그는) + MEU₁
 DEP -> MEU₂ + MEU(출근하였다)

((그는) (아침을 먹고)) (출근하였다))

(문장 16) 그는 더위를 심하게 먹었다.

MEU₁ -> MEU(심하게) + MEU(먹었다)
 MEU₂ -> MEU(더위를) + MEU₁
 DEP -> MEU(그는) + MEU₂

((그는) ((더위를) (심하게 먹었다)))

(문장 5) 코끼리는 코가 길다.

MEU₁ -> MEU(코가) + MEU(길다)
 DEP -> MEU(코끼리는) + MEU₁

((코끼리는) (코가 길다))

(문장 17) 나는 그것을 알 수가 있다.

MEU₁ -> MEU(그것을) + MEU(알)
 MEU₂ -> MEU₁ + MEU(수가)
 MEU₃ -> MEU₂ + MEU(있다)
 DEP -> MEU(나는) + MEU₃

((((나는) (그것을 알)) (수가)) (있다))

