

병렬 코퍼스를 이용한 한중 기계번역 오류 탐지 방법

김윤^o 김영길

한국전자통신연구원 언어처리연구팀
wkim1019@etri.re.kr, kimyk@etri.re.kr

Method for Detecting Errors of Korean-Chinese MT Using Parallel Corpus

Yun Jin^o Young-Kil Kim

ETRI Natural Language Processing Research Team

요 약

본 논문에서는 패턴기반 자동번역시스템의 효율적인 번역 성능 향상을 위해 병렬 코퍼스(parallel corpus)를 이용한 오류 자동 탐지 방법을 제안하고자 한다. 번역시스템에 존재하는 대부분 오류는 크게 지식 오류와 엔진 오류로 나눌 수 있는데 통상 이런 오류는 이중 언어가 가능한 훈련된 언어학자가 대량의 자동번역 된 결과 문장을 읽음으로써 오류를 탐지하고 분석하여 번역 지식을 수정/확장하거나 또는 엔진을 개선하게 된다. 하지만, 이런 작업은 많은 시간과 노력을 필요로 하게 된다. 따라서 본 논문에서는 병렬 코퍼스 중의 목적 언어(Target Language) 문장 즉, 정답 문장과 자동번역 된 결과 문장을 다양한 방법으로 비교하면서 번역시스템에 존재하고 있는 지식 및 엔진 오류를 자동으로 탐지하는 방법을 제안한다. 제안한 방법은 한-중 자동번역시스템에 적용하여 그 정확률과 재현률을 측정하였으며, 자동적으로 오류를 탐지하여 추출 할 수 있음을 증명하였다.

1. 서 론

최근 국제화의 교류가 날로 증가함에 따라 이중 언어 간의 교류를 활발하게 하는 자동번역 활용이 날로 증가하고 있으며, 이와 같은 자동번역에 있어 번역 성능을 효율적으로 높이는 것이 중요한 과제로 떠오르고 있다.

효율적으로 번역 성능을 향상시키기 위해서는 번역시스템에 존재하고 있는 문제점을 파악하여야 하는데 이런 문제점 파악은 오류 탐지와 분석을 통하여 이루어진다. 특히 패턴 기반 번역시스템의 대부분 오류는 지식오류와 엔진 오류인데, 그 중 지식오류는 대량 지식(미등록어, 패턴 등) 구축 과정에서 지식 작업자의 지식 한계로 인해 잘못 구축된 지식, 새로운 지식 부재, 및 도메인 특화 지식 오류 등을 말하며, 엔진 오류는 이렇게 구축된 지식을 이용하여 번역시스템에서 논리적 처리 못하는 오류를 말한다.

통상 이와 같은 오류 탐지 및 분석 작업은 이중 언어가 가능한 훈련된 언어학자가 대량의 자동번역 된 결과 문장을 읽음으로써 오류를 탐지 및 분석 후 오류 유형을 분류하여 지식을 수정/확장하거나 엔진을 개선하게 된다. 하지만 이런 이중 언어학자는 구하기 힘들뿐만 아니라 대량의 번역문장을 읽어야 하는 시간적 노력을 필요로 하며, 또한 시스템적으로 분석할 수 있는 능력이 필요하다.

본 논문에서 제안하고자 하는 방법은 대량의 병렬 코퍼스 중 목적언어 문장을 정답 문장으로 간주하고 이 정답 문장과 원시언어(Source Language) 문장을 자동 번역한 결과 문장을 서로 비교하여 자동번역시스템의 오류

를 자동으로 탐지하고 분류하고자 한다. 이를 위해 본 논문에서는 병렬코퍼스 상의 각 문서 쌍에 대하여 먼저 문장 단위로 정렬(Alignment)을 수행한 다음, 정렬된 정답 문장과 원시언어 문장을 자동 번역한 결과문장을 단어 단위 정렬을 위해 형태소 단위로 분절을 수행하였다. 단어 단위의 정렬은 형태소 단위로 분절된 두 문장을 입력받아 단어 단위 매핑, 매핑 결과를 이용한 수평적 수직적 비교를 수행하였으며, 그 결과를 이용하여 디코딩 정렬 테이블(Decoding align table)을 구성하였다. 수평적 비교의 목적은 단어가 서로 매핑의 안 되면 신조어거나, 대역어 불일치, 또는 분절 등 오류를 탐지하기 위함이며, 수직적 비교의 목적은 구문 오류, 용언구 오류 등을 탐지하기 위해서다. 이렇게 생성된 디코딩 테이블을 이용하여 오류 탐지 및 분류를 수행하였다.

본 논문에서 제안한 방법은 다음과 같은 장점을 갖고 있다.

첫째, 이중언어 가능자만 오류를 찾을 수 있던 기존의 한계점을 극복할 수 있으며, 오류를 찾는데 투자하는 시간과 노력을 획기적으로 줄일 수 있다.

둘째, 발견된 대량의 오류를 통해 번역시스템 문제점 파악이 쉽고, 목적언어 문장을 참조할 수 있어 효과적인 문제점 개선책을 찾을 수 있다.

셋째, 새로운 도메인으로 번역시스템을 확장하는데 소요되는 많은 인적, 재적 비용과 시간적인 비용을 절약할 수 있다.

2. 관련 연구

번역시스템의 잠재적인 오류를 알기 위해 BLUE

(BiLingual Evaluation Understudy)[1]와 같은 자동 평가 도구를 광범위하게 이용하고 있지만, 오류여부를 평가한다고 해서 어떤 오류인지 알 수 없으며, 오류 여부에 대한 평가와 오류를 탐지하고 분석하여 찾아내는 것은 서로 다른 문제이며, 실질적으로 오류를 평가하는 것보다 오류를 찾아내는 것이 훨씬 어렵다.

[2]에서는 번역시스템 오류를 5가지로 분류하였는데, 이 5가지 분류는 오탈자(Missing Words), 어순(Word Order), 틀린 글자(Incorrect Words), 신조어(Unknown Words), 구두법(Punctuation) 등이다. 본 논문의 궁극적인 목적 역시 탐지된 오류를 위와 같이 분류하거나 좀 더 세분화하는 것이다.

번역시스템의 오류에 대하여 많은 연구가 이루어졌다. 그 중 대표적인 것이 원문으로부터 번역결과에 이르기까지 사람과의 상호작용을 통하여 오류 가능성을 제시하고 사용자로 하여금 교정하는 방법[3]과 기계번역시스템에 의해 번역된 결과를 사후편집(post-editing)하는 방법[4,5,6] 등을 예로 들 수 있다.

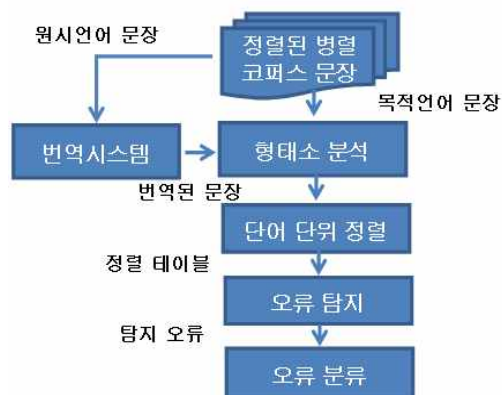
먼저, [3]에서는 번역시스템에 존재하는 오류 유형을 크게 원시 문장 오류, 번역시스템 오류, 도메인 특화 정보 부재 등으로 분류하고, 사람과의 상호작용을 통하여 가능한 오류를 탐지하고 사용자가 확인하여 수정해나가는 방법으로 보다 양질의 번역 결과를 얻으려 시도했다. 원시 문장 오류는 원시 문장에 존재하는 철자, 띄어쓰기 및 구문분석 오류로써, 코퍼스로부터 대량의 유사 문장을 이용하여 원시 문장의 오류 패턴을 자연어 처리 기법을 통하여 탐지하였으며, 이렇게 탐지된 가능한 오류는 사용자의 판단을 통하여 수정하도록 하였다. 또한 번역시스템 오류와 도메인 특화 오류는 사용자의 피드백을 통하여 해결하고자 하였다. 하지만, 이런 오류 탐지 방법 역시 사람의 개입의 없으면 이루어질 수 없으며, 대량의 오류를 탐지하기에는 역시 역부족이라는 단점을 갖고 있다.

다음으로, [4,5,6] 등에서는 번역시스템을 통해 번역된 번역결과를 사후 편집하는 방법으로써, 사용자의 교정 결과를 학습하거나 대량의 병렬 코퍼스를 이용하여 통계적인 기계번역(SMT: Statistical Machine Translation) [7,8]을 수행하여 종래의 번역시스템이 번역한 번역결과에 대해 사후 편집한다. 이때, 사용되는 SMT 번역결과는 번역시스템의 성능개선을 위한 것이 아니라 번역결과에 반복되어 나타나는 오류를 교정하여 보다 정확한 번역결과를 얻기 위한 목적으로 사용되고 있다. 이때 교정된 결과는 번역시스템이 갖고 있는 문제점을 개선시키지 못하는 문제점이 있으며, 또한 번역시스템을 새로운 도메인에 적용 시 이런 방법은 내부적으로 갖고 있는 문제점을 보완하기에 역부족이라는 문제점을 갖게 된다.

3. 오류 탐지

3.1 시스템 구성

제안하고자 하는 방법에 대한 전반적인 이해를 돕기 위해 다음 그림 1를 이용하여 시스템 각 모듈의 구성 및 각 모듈 간의 데이터 흐름을 설명하고자 한다.



[그림 1] 시스템 흐름도

시스템 구성에는 어떤 오류가 있는지 여부를 평가할 번역시스템, 목적어문어와 번역결과를 형태소 단위로 분절할 목적어문어 형태소 분절기, 단어 단위 정렬, 오류 탐지, 및 오류 분류 등 모듈로 구성된다.

전체 시스템 먼저, 번역할 원시언어 문장과 그에 대응되는 목적어문어 문장으로 정렬된 병렬코퍼스를 입력으로 받는다. 다음, 원시문장을 번역시스템으로 자동번역을 한다. 이렇게 번역된 문장을 목적어문어 문장과 비교하면 정답문장 대 테스트 문장으로 간주할 수 있다. 그 다음, 정답 문장과 테스트 문장을 단어 단위로 비교하기 위해 동일한 형태소 분절기로 분절을 수행한다. 분절된 두 문장은 단어 단위 정렬 모듈에 입력되어 정답 문장의 단어가 번역 결과 문장에도 나왔는지 여부를 판단하여 일련이 번호로 비교 결과를 디코딩 테이블에 표시한다. 오류 탐지 모듈에서는 위의 정렬된 비교 테이블을 이용하여 수평적, 수직적 비교를 한다. 마지막으로, 오류 분류에서는 지금까지 생성한 테이블에서 유형 별 오류를 추출하여 분류한다. 이렇게 분류된 오류는 엔지니어에게 제공하여 번역시스템의 오류를 개선할 수 있도록 한다.

3.2 단어 단위 정렬 방법

본 논문의 핵심은 단어 단위 정렬과 오류 탐지를 위한 디코딩이므로 한 예를 이용하여 구체적으로 설명하고자 한다. 제안한 방법을 한중 번역시스템에 적용하였기 때문에 설명의 편리를 위하여 본 예에서 사용한 원시언어 문장 역시 한국어 문장이고, 목적어문어 문장은 중국어 문장이다. 그 들로는 각각 “<S snum= 3>금융 및 실물시장 상황이 급변하는 가운데 14일(현지시간) 미국 시카고 상품거래소의 증개인들이 손짓 거래 주문을 내고 있다. </S>” 과 “<S snum=3>急剧变化的金融及实物市场状况中, 14日(当地时间), 美国芝加哥商品交易所的经纪人们通过首饰来下达交易订单 °</S>” 이다. 그리고, 이 한국어 문장을 자동 번역한 중국어 문장은 “<S snum=3>金融和实物市场情形急剧变化,14号(当地时间)美国芝加哥商交所的经纪人们提出手势交易预订 °</S>” 이다.

단어 단위의 정렬은 크게 2가지 부분으로 나누는데, 먼저 형태소 단위 분할된 병렬코퍼스 상의 중국어 문장과 그에 대응되는 자동 번역된 중국어 문장에 대하여 서로 비교하면서 디코딩하는 단계, 그 다음, 오류를 탐지 정확도를 높이기 위한 절(Clause) 단위로 분할하는 단계가 포함한다. 여기서 절 단위 분할은 정렬의 정확도를 향상시키기 위해서다. 한 문장에는 같은 단어, 같은 기호, 숫자 등이 중복되어 나오는 확률이 높기 때문에 절 단위로 분할하면, 정렬 정확도를 높일 수 있다.

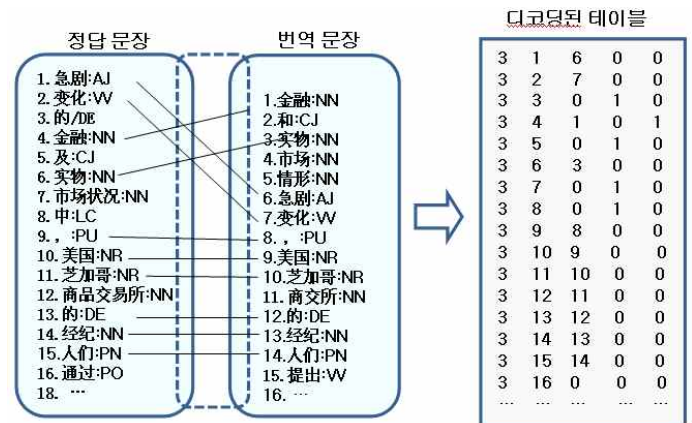
디코딩 정렬 테이블(decoding alignment table)을 구성하는 방법은 다음과 같은 과정을 거치게 된다.

1. 형태소 단위로 분할된 정렬된 두 문장을 입력으로 받아, 그 문장 인덱스 번호로 디코딩 테이블의 첫 번째 열을 구성한다. 만약 새로운 문장이 들어오면, 첫 번째 열은 그 문장의 인덱스 번호로 바뀌게 된다.
2. 입력된 문장 중 병렬 코퍼스 상의 중국어 문장의 각 토큰에 대해 1부터 시작되는 순차적인 번호를 부여한다. 이렇게 부여된 숫자로 디코딩 테이블의 두 번째 열을 구성한다.
3. 입력된 문장 중 자동 번역된 결과 문장의 각 토큰에 대하여 위와 같이 1부터 시작되는 순차적인 번호를 부여한다.
4. 입력된 문장 중 병렬 코퍼스 상의 중국어 문장을 순차적으로 방문하면서, 자동 번역된 결과 문장 중 같은 토큰이 있는지를 순차적으로 탐색한다. 만약 있으며, 자동 번역 결과 토큰의 번호, 없으면 0으로 하여 디코딩 테이블의 세 번째 열을 생성한다.
5. 디코딩 테이블을 순차적으로 방문하면서 두 번째 열과 세 번째 열을 붙이던 AND 연산을 하여 디코딩 테이블의 네 번째 테이블을 생성한다.
6. 디코딩 테이블의 세 번째 열을 첫 번째 열이 같은 조건 하에서 0을 제외한 숫자들을 중심으로 순차적으로 방문하면서, 보다 큰 숫자는 큐(Queue)에 넣고 0을, 보다 작은 숫자가 있으면 1을 생성하여 디코딩의 다섯 번째 열을 생성한다.

그림 2는 상기 예를 디코딩 테이블 구성하는 방법에 따른 문장 간 디코딩된 정렬 테이블이다. 그림 2을 참조하면, 중국어 두 문장을 형태소 분석한 결과인 각 토큰 및 각 토큰에 해당하는 품사 정보 쌍들을 일련의 순차적인 번호와 함께 각 라인에 표시한 것을 알 수 있다. 이 디코딩된 결과 테이블에 대하여 설명하면 다음과 같다.

1. 첫 번째 칼럼은 병렬코퍼스 상의 문장 단위로 정렬된 문장 인덱스 번호이다. 그림을 통해 사용한 예는 병렬 코퍼스 상의 3 번째 문장임을 알 수 있다. 이렇게 함으로써, 문장 단위로 정렬된 각 문장을 구별할 수 있다.
2. 두 번째 칼럼은 형태소 단위로 분할된 중국어 정답 문장의 각 토큰 리스트에 대해 1로부터 시작하는 순차적으로 번호를 부여한 결과이다. 이때 0으로부터 시작하지 않고, 번호 1로부터 시작하는 이유는 두 문장 정렬에서 정렬되지 않는 토큰에 대하여 값 0을 부여하기 위해서다. 또한, 이 칼럼은 정렬정보인 나머지 각 칼럼을 생성할 수 있도록 기준이 되는 역할을 한다.

3. 세 번째 칼럼은 자동 번역한 테스트 중국어 문장에 대한 형태소 분석 결과를 두 번째 칼럼을 기준으로 정렬한 단어 리스트 번호이다. 이때, 리스트 번호는 순차적이지 않다. 그 원인은 두 번째 칼럼인 중국어 용어를 기준으로 하였기 때문이다. 이때, 위와 마찬가지로 분할된 형태소는 번호 1부터 시작하지만, 정렬되지 않을 경우 값을 0으로 한다. 이 예에서 “市场状况:NN”에 대한 정렬정보가 없으므로 0으로 부여하였다. 이 칼럼의 역할은 목적언어에 대한 정렬을 통하여 각종 오류를 탐지하기 위한 것이다.
4. 네 번째 칼럼은 두 번째와 세 번째의 값을 AND연산을 하여 값이 1이면 0을, 0이면 1을 부여한다. 즉, 이 칼럼은 정렬된 문장에서 수평적 비교를 통하여 서로 다른 값을 가지는 부분을 찾아내기 위함이다. 이렇게 함으로써, “市场状况:NN”와 “市场:NN“, “情形:NN”이 서로 다를 수 있다. 즉, 대역어 사전 오류를 탐지한 것이다.
5. 다섯 번째 칼럼은 세 번째 칼럼에 대하여 수직적 비교를 통하여 순차적이지 않은 부분을 탐지하기 위함이다. 이 부분에서 탐지된 오류는 대부분 용언구 패턴 또는 구문분석 등 오류이다.



[그림 2] 단어 정렬 및 디코딩 테이블

본 논문에서의 오류 탐지는 패턴기반 기계번역시스템에서 흔히 발생할 수 있는 신조어 및 대역어 오류, 용언구 패턴 오류 등 지식 오류에 대한 탐지와 분석 오류, 구문 분석 오류, 등 엔진 오류에 대한 탐지를 말한다. 오류 탐지는 기본적으로 디코딩된 정렬 테이블 참조하여 알 수 있다. 다음은 각 오류에 대한 탐지 방법을 상세히 설명하고자 한다.

1. 신조어 및 대역어 오류: 이 오류는 수평적 비교를 통하여 이루어진다. 즉, 네 번째 칼럼의 각 행에는 값이 존재하지만 이에 반하여 “市场状况:NN”의 대응되는 행의 값은 “1”로 되어있고, 한국어 품사가 명사로 태깅되었을 경우 신조어 또는 대역어 오류로 추정하고 추출할 수 있다. 이는 기계번역 시 용어 사전에 대응되는 엔트리(entry)가 없거나 대응되는 대역어가 잘 못 구축되었을 수 있다.

2. 용언구 패턴 오류 탐지: 이 종류의 오류는 수직적

비교를 통하여 이루어진다. 예를 들어, 용언 “變化”의 정답 중국어 문장에서는 2번째 위치이지만, 번역 결과에서는 위치가 7로 교차되었다. 따라서, 용언구 패턴 오류로 탐지할 수 있다.

3. 분절 오류 탐지: 이 분류의 오류는 정답 문서 문절 위치를 통하여 알 수 있다. 예를 들어, 정답문서 “美國:NR”앞에서 분절되었으므로, 번역문장 9번째 위치에서도 분절 될 가능성이 높다. 분절 기호 콤마가 없다면, 분절 오류로 간주할 수 있다.

4. 실험

4.1 실험 데이터

본 논문에서 제안한 방법의 가능성을 증명하고 평가하기 위해, 한-중 병렬 문서를 대상으로 실험을 하였다. 대상 뉴스는 동아 경제 뉴스¹⁾이며, 실험에 사용한 데이터는 200 문장 쌍을 랜덤으로 추출하여 사용하였다. 200문장으로 한정지은 것은 실험 결과에 대한 판정 시간이 너무 오래 걸리기 때문에 그 수량을 한정할 수밖에 없었다.

4.2 실험 방법 및 결과

실험은 제안한 방법으로 추출한 오류에 대하여 정확성 여부 즉, 정확률을 평가하였으며, 또한, 추출된 오류가 전체 오류 중 얼마나 많은 비중을 차지하는지를 알기 위해 재현률도 측정하였다.

실험에서 오류의 정확률을 향상하기 위해, “的”와 인용부호와 같이 자주 쓰이면서, 번역시스템 성능에 덜 영향을 미치는 단어는 매칭이 되지 않아도 오류로 간주하지 않았다. 특히 인용부호는 한국어 뉴스에서는 거의 없지만, 중국어 정답문서에는 이것을 철저히 지킨 특징을 갖고 있었다.

실험결과를 오류 탐지 정확률은 76.3%로 저조했지만, 재현률은 83.4%로 괜찮은 성능을 보였다.

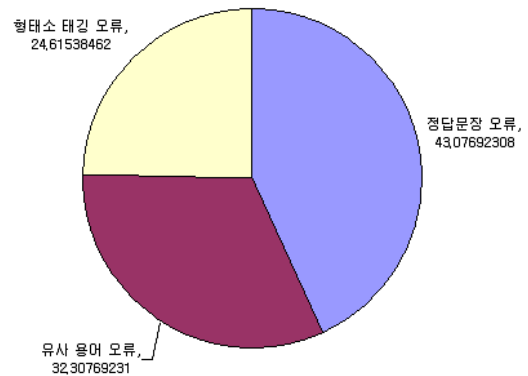
4.3 실험 결과 분석

실험에서 탐지했지만, 잘 못 탐지한 오류와 탐지 못한 오류에 대하여 분석한 결과 그림 3과 같은 몇 가지 문제점들을 발견하였다.

1) 정답문장 오류: 잘 못 탐지한 오류 중에는 정답문장 오류가 가장 많았다. 정답 문장에 나타나는 오류는 틀리게 번역한 오류, 원본 내용 외의 단어가 더 들어간 있는 경우, 의역을 한 경우 등 오류가 있었다. 대표적인 예로 원시언어 문장에는 “손짓(手勢)” 이란 단어가 있는데, 병렬 코퍼스 상의 목적언어 문장 즉, 정답문장에는 “首飾(장신구)” 로 되어 있어 번역 작업자가 타이핑을 잘못된 오류이다. 또 다른 예로 “도봉구 아파트값도 3.3m 당 1000만원 돌파” 이 문장에서 “도봉구” 를 목적언어 문장에서는 “서울 도봉구” 로 의미를 확장하여 번역하다 보

니, 오류로 탐지 된 사례도 있다. 이러한 정답 문장에 오류가 있어 오류를 잘 못 탐지되었다. 실험을 통하여 이 유형에 속하는 오류가 43.07%로 가장 많았다. 이 종류의 오류는 중-한 사전을 이용하여 그 뜻을 찾아봄으로써, 어느 정도 오류를 줄일 수 있다. 예를 들어, “장신구” 와 같은 단어는 의미가 단일 의미이고, 대역어도 다형성을 가지지 않으므로, 틀린 단어임을 알 수 있다. 이렇게 함으로써, 잘 못 탐지된 오류를 어느 정도 필터링할 수 있다.

잘 못 탐지된 오류



[그림 3] 잘 못 탐지한 오류 유형

2) 유사 용어 오류: 이 유형의 오류는 정답 문장에 나타난 용어와 번역된 용어가 서로 빈번하게 같이 쓰이는 경우를 말한다. “及”와 “和”는 서로 같은 접속사로 어느 걸 써도 무방하다. 또 다른 예로, “商品交易所(상품거래소)”, “商交所(축약형)”도 같은 의미이고, 글자의 형태가 다르다고 해서 다른 단어라고 할 수 없다. 이 종류의 오류는 대량의 병렬코퍼스만 확보한다면, SMT의 단어 정렬 도구를 이용하여 이와 같은 유사 용어를 확보함으로써, 이런 유사 용어 불일치는 오류로 간주하지 않음으로써 해결 가능하다.

3). 중국어 형태소 태깅 오류: 중국어 문장에 대하여 형태소 분할을 하는 잘 못 태깅 함으로써, 생기는 오류이다. 본 논문에서 사용한 중국어 형태소 분할기는 한국 전자통신연구원(ETRI)에서 개발한 것으로 분할 정확률이 97.48%로, 명사와 동사 모두 가능한 다품사 단어에 대한 태깅 오류와 분할 오류가 본 실험에 가장 많이 영향을 주었다.

5. 결론 및 향후 연구

본 논문에서는 병렬 코퍼스를 이용하여 번역시스템에 존재하는 오류를 탐지하고, 분석하여 분류하는 방법을 제안하였다. 번역시스템의 잠재적인 오류는 상대적으로 예측이 가능하지만, 오류를 탐지하여, 추출하려면 이중언어 가능한 학자의 도움이 필요하며, 많은 시간과 노력이 필요로 한다.

하지만, 본 논문에서는 잘 정렬된 병렬문장을 이용하여 단어 단위로 형태소 분석하여 정렬함으로써, 오류를

1) www.donga.com

자동으로 추출하는 방법을 제안하였다. 실험을 통하여 정확률은 다소 떨어지지만, 괜찮은 재현률을 얻을 수 있어 그 유용성을 증명할 수 있었다.

향후 계획으로는 위의 실험을 통하여 얻은 분석결과를 토대로 오류 탐지 정확률을 향상하여 번역시스템에 유용한 도구로 활용할 수 있도록 하는 것이다.

6. 참고 문헌

- [1] K. Papneni, S. Roukos, T. Ward, and W.J. Zhu, Bleu: A Method for Automatic Evaluation of Machine Translation, IBM Research Report, RC22176, 2001.
- [2] D. Vilar, J. Xu, L.F. D'Haro, H. Ney, "Error analysis of Statistical Machine Translation Output," In Proc. of the 5th Int. Conf. on Language Resource and Evaluation, 2006.
- [3] Y.A. Seo, C.H. Kim, S.I. Yang and Y.G. Kim, "Getting Professional Translation through User Interaction," MT Summit XI, PP.413-419, 10~14, Sep. 2007.
- [4] K. Knight and I. Chander, "Automated Postediting of Documents," Proc. of National Conf. on Artificial Intelligence(AAID), 1994.
- [5] G. Michael, J.F. Gao, B. Chris, K. Alexandre, B.D. William, B. Dmitriy, V. Lucy, "Using Contextual Techniques and Language Modeling for ESL Error Correction," In conf. of IJCNLP, 2008.
- [6] Jakob Elming, "Transformation-based correction of rule-based MT,"
- [7] W. Weaver(1955) and Translation(1949), In: Machine Translation of Languages, MIT Press, Cambridge, MA.
- [8] P. Brown, S. Della Pietra, V. Della Pietra, and R. Mercer, "The Machine of Statistical Machine Translation: Parameter Estimation," Computational Linguistics, Vol.19, No.2, PP.263-311, 1991.