

관용구 기계번역을 위한
한-영 데이터셋 구축 및 평가 방법

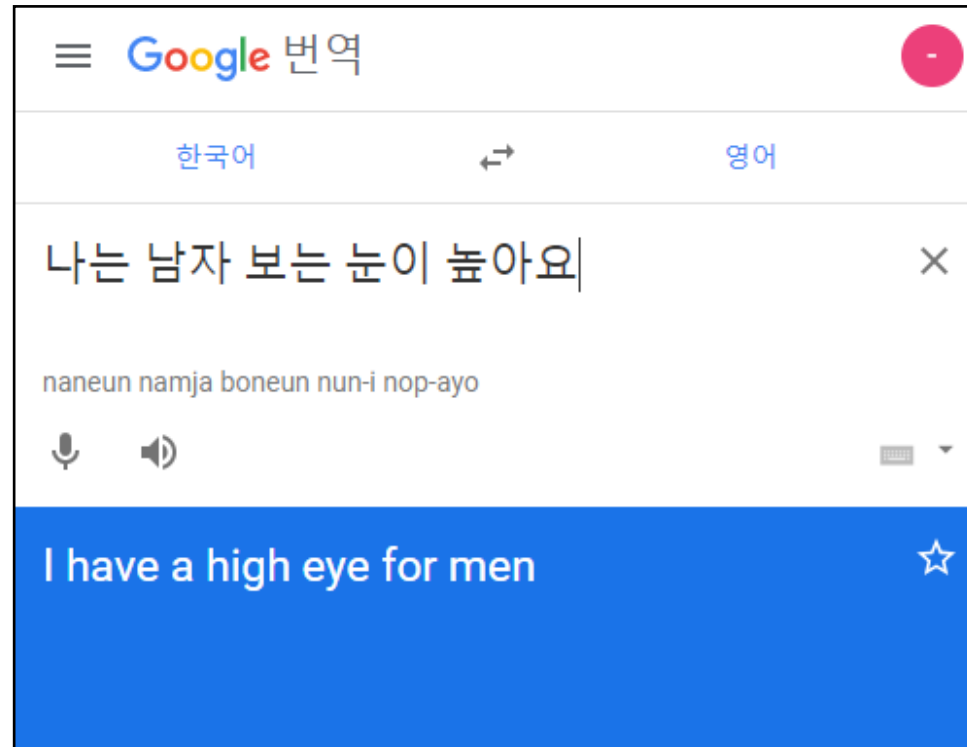
최민주

mjchoi0831@gmail.com

1. 서론

- NMT 기반 기계번역은 훌륭한 번역 성능을 보이나, 종종 오역이 발생함
 - 예) 관용구 번역
- 관용-구(慣用句)
 - 두 개 이상의 단어로 이루어져 있으면서 그 단어들의 의미만으로는 전체의 의미를 알 수 없는, 특수한 의미를 나타내는 어구(語句)
 - 예) '눈이 높다' = '눈' + '높다'
 - 뜻 : '좋은 것을 찾다', '안목이 높다'

예) 관용구 포함 문장 기계번역 결과



예) 관용구 포함 문장 기계번역 결과

The screenshot shows the Google Translate interface. At the top, it says "Google 번역" (Google Translate). Below that, the source language is "한국어" (Korean) and the target language is "영어" (English). The input text is "나는 남자 보는 눈이 높아요" (I have a high eye for men). The phonetic transcription below is "naneun namja boneun nun-i nop-ayo". The output text is "I have a high eye for men". A light blue highlight is placed over the output text, and a tooltip shows the correct translation: "I have high standards for man".

☰ Google 번역

한국어 ↔ 영어

나는 남자 보는 눈이 높아요

naneun namja boneun nun-i nop-ayo

🎤 🔊 📄 ▼

I have a high eye for men ☆

I have high standards for man

1. 서론

- 기계번역은 함축적인 의미를 지닌 관용구를 정확하게 번역할 수 없음
- 기계번역 모델이 관용구를 효과적으로 학습하고 번역 결과를 평가할 수 있도록 관용구 번역에 특화된 데이터셋과 평가방법이 필요

1. 서론

- 본 논문에서는...
 - 관용구 기계번역을 위한 한-영 번역 쌍 데이터셋 'KISS' 구축 방법 제안
 - 관용구 번역 결과의 품질을 평가하기 위해 블랙리스트를 이용하는 방법 소개
 - KISS를 이용하여 블랙리스트를 구축하는 방법 제안

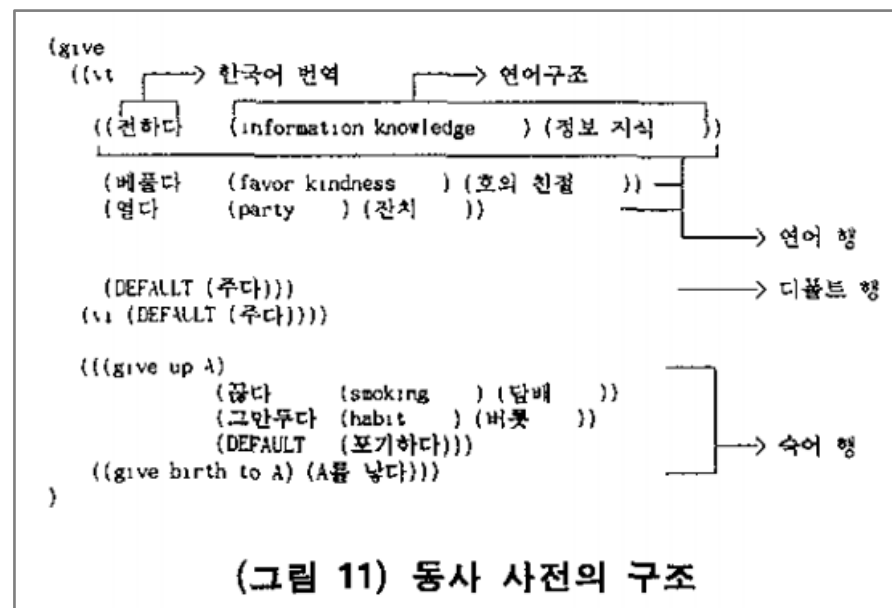
- KISS : **K**orean-english **I**dioms in **S**entences data**S**et
- <https://github.com/Judy-Choi/KISS>

2. 관련 연구

- NMT 이전에는...
 - 구-대-구 방식으로 관용구의 구조 정보를 이용하여 번역

- 사전 기반 방식

- 문장 내에 관용구 사전에 존재하는 관용구가 포함되어 있으면 관용구로 번역
- 단점
 - 사전 구축에 많은 시간과 노력 필요
 - 문장 단위로 학습하고 번역하는 NMT에는 적합하지 않음



[1] 이호석, 김영택. 영어-한국어 기계번역을 위한 언어와 숙어 트랜스퍼 사전. (구) 정보과학회논문지. 20.7: 976-987. 1993.

2. 관련 연구

- NMT에서 관용구를 학습/평가하려면 뭐가 필요할까
 - 관용구가 포함된 다량의 문장 데이터셋
 - 한-영 번역 쌍 말뭉치로부터 관용구 포함 번역 쌍 추출
⇒ 관용구 학습 데이터셋 구축
 - 관용구 번역에 특화된 번역 평가 지표
 - 관용구 학습 데이터셋을 이용
⇒ 블랙리스트(평가 지표) 생성

3. KISS : 관용구 포함 한-영 번역 쌍 데이터셋 구축

표준국어대사전
수록 관용어

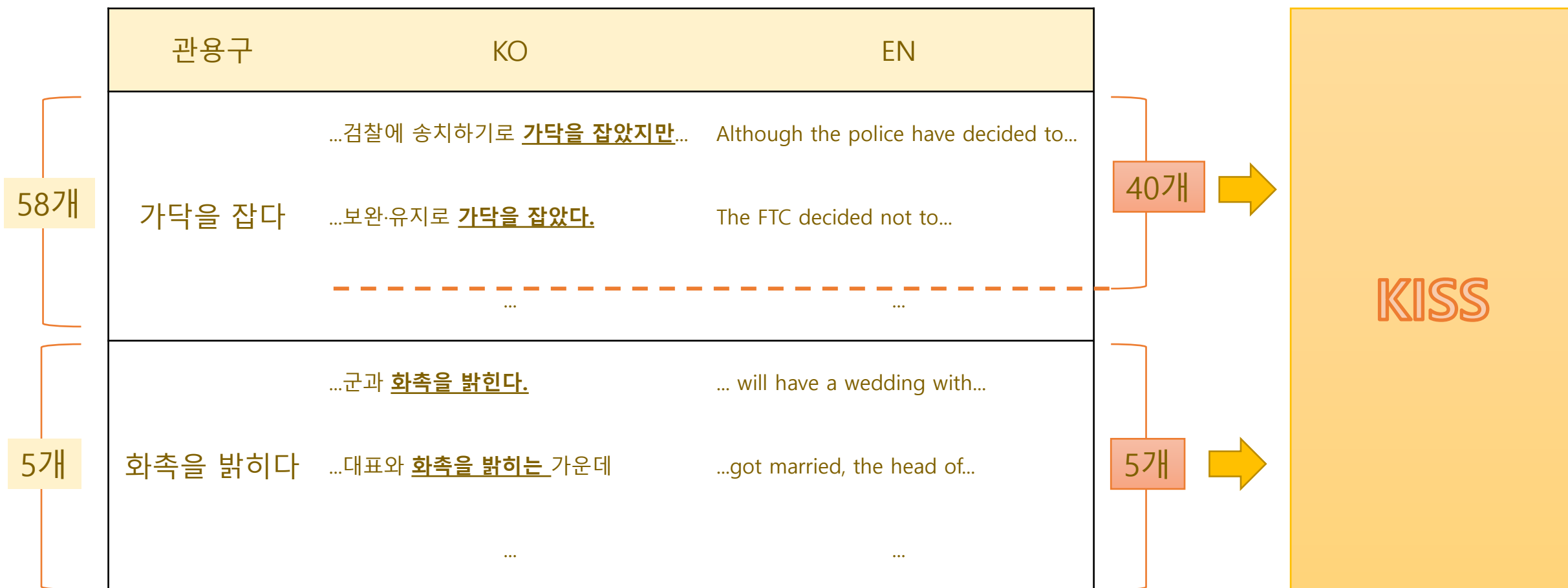
AI Hub 한국어-영어 번역(병렬) 말뭉치

...
눈이 높다
마침표를 찍다
총대를 메다
화촉을 밝히다
...



KO	EN
나는 여자 보는 <u>눈이 높아요</u>	I have high standards for woman.
아버지의 명예회복을 위한 김지훈의 기나긴 여정이 마침내 <u>마침표를 찍었다.</u>	Kim Ji-hoon's long journey ... has finally come to an end.
기획재정부가 혁신성장 관련 규제 완화에 <u>총대를 멘다.</u>	The Ministry of Economy and Finance takes charge of ...
중국 출신 할리우드 스타 장쯔이 ...와 내년 <u>화촉을 밝힌다.</u>	Chinese Hollywood star Zhang Ziyi (28) will marry ...
...	...

3. KISS : 관용구 포함 한-영 번역 쌍 데이터셋 구축



3. KISS : 관용구 포함 한-영 번역 쌍 데이터셋 구축

- 관용구 수집
 - 표준국어대사전 온라인 사이트 → 3,887개 관용구 목록 다운로드
- 한-영 문장 번역 쌍 추출
 - AI Hub 한국어-영어 번역(병렬) 말뭉치
 - 430개 한국어 관용구 포함 18,808개 번역 쌍 추출
 - 동일한 관용구 포함 문장 4~40개로 제한
 - 420개 한국어 관용구 포함 7500개 번역 쌍 추출 => KISS 구축
- KISS : Korean-english Idioms in Sentences dataSet



github.com/Judy-Choi/KISS

KISS : Korean-english Idioms in Sentences dataSet

관용구	한국어 원문	영어 번역 쌍
눈이 높다	나는 여자 보는 <u>눈이 높아요.</u>	I have <u>high standards</u> for woman.
마침표를 찍다	아버지의 명예회복을 위한 김지훈의 기나긴 여정이 마침내 <u>마침표를 찍었다.</u>	Kim Ji-hoon's long journey to restore his father's honor has finally <u>come to an end.</u>
총대를 메다	기획재정부가 혁신성장 관련 규제 완화에 <u>총대를 멘다.</u>	The Ministry of Economy and Finance <u>takes charge of</u> the deregulation related to innovative growth.
화촉을 밝히다	중국 출신 할리우드 스타 장쯔이(章子怡, 28)가 현 남자친구인 미국의 억만장자 비비 네보(42) 와 내년 <u>화촉을 밝힌다.</u>	Chinese Hollywood star Zhang Ziyi (28) will <u>marry</u> her current boyfriend, U.S. billionaire Aviv Nevo (42), next year.

Problem...

Correct

관용구	눈이 높다
한국어 원문	나는 여자 보는 <u>눈이 높아요.</u>
영어 번역 쌍	I have <u>high standards</u> for woman.
관용구	마침표를 찍다
한국어 원문	아버지의 명예회복을 위한 김지훈의 기나긴 여정이 마침내 <u>마침표를 찍었다.</u>
영어 번역 쌍	Kim Ji-hoon's long journey to restore his father's honor has finally <u>come to an end.</u>

Error

관용구	운을 떴다
한국어 원문	정부가 노인 연령 기준을 높이는 방안에 대해 <u>운을 뒀다.</u>
영어 번역 쌍	The government has been <u>lucky</u> about ways to raise the criteria of age for senior citizens.
관용구	유명을 달리하다
한국어 원문	오키나와 현지사인 오나가 다케시가 지난 8일 췌장암으로 <u>유명을 달리했다.</u>
영어 번역 쌍	Takeshi Onaga, an incumbent Governor of Okinawa, became <u>famous</u> for pancreatic cancer on the 8th.

Problem...

- 구축한 데이터셋에 다량의 오역 쌍 존재
 - 데이터셋에 포함된 다량의 오역을 제거할 수 있는 방법 필요
- 관용구 번역 시 다량의 오역 발생
 - 기계번역 오류를 탐지할 수 있는 품질 평가 지표 필요

Problem...

- 구축한 데이터셋에 다량의 오역 쌍 존재
 - 데이터셋에 포함된 다량의 오역을 제거할 수 있는 방법 필요
- 관용구 번역 시 다량의 오역 발생
 - 기계번역 오류를 탐지할 수 있는 품질 평가 지표 필요

1 shot 2 kill solution

Blacklist Method

4. 블랙리스트 평가 방법

- 직역으로 인한 번역 오류 탐지
- 단어-대-단어 구조의 관용구를 번역한 결과로부터 오역 여부를 판별
- 원리
 - 관용구를 한-영 번역한 결과에 블랙리스트 단어가 1개 이상 포함되어 있으면 오역으로 간주

관용구	블랙리스트
꼬집어 말하다	nip pinch twitch
눈 높다	eye
운을 떴다	lucky
유명을 달리하다	famous

4. 블랙리스트 평가 방법

Correct

관용구	눈이 높다
블랙리스트	eye
한국어 원문	나는 여자 보는 <u>눈이 높아요.</u>
영어 번역 쌍	I have <u>high standards</u> for woman.

Error

관용구	눈이 높다
블랙리스트	eye
한국어 원문	나는 여자 보는 <u>눈이 높아요.</u>
영어 번역 쌍	I have <u>high eye</u> for woman.

블랙리스트를 이용한 오역 탐지

4. 블랙리스트 평가 방법

- 블랙리스트 구축
- 단어의 의미를 이용한 블랙리스트 단어 추출

관용구	<u>꼬집어</u> 말하다
블랙리스트	nip pinch twitch
관용구	<u>눈</u> 높다
블랙리스트	eye

- 오역결과를 이용한 블랙리스트 단어 추출

관용구	운을 때다
한국어 원문	정부가 노인 연령 기준을 높이는 방안에 대해 <u>운을 땀다</u> .
영어 번역 쌍	The government has been lucky about ways to raise the criteria of age for senior citizens.
블랙리스트	lucky
관용구	유명을 달리하다
한국어 원문	오키나와 현지사인 오나가 다케시가 지난 8일 췌장암으로 <u>유명을 달리했다</u> .
영어 번역 쌍	Takeshi Onaga, an incumbent Governor of Okinawa, became famous for pancreatic cancer on the 8th.
블랙리스트	famous

5. 평가

- 블랙리스트 구축
 - 420개 관용구 중 275개 관용구에 대한 블랙리스트 구축
 - 오역이 거의 없는 관용구 제외
 - 예) '그건 그렇고' → 'By the way' 로 대부분 정확하게 번역됨
 - 영어로 직역되는 관용구 제외
 - 예) '(사람의 마음을) 가지고 놀다' → 'play with' 로 직역되므로 블랙리스트 구축 불가
- 블랙리스트를 이용한 오역 없는 번역 쌍 추출
 - KISS 로부터 275개 관용구 포함 3,461 개 번역 쌍 추출

5. 평가

- 오역 제외한 번역 쌍을 이용한 기계번역 서비스 품질 평가
 - 관용구 275개
 - 번역 쌍 3,461개

	Google 번역	Naver Papago	Kakao i 번역
블랙리스트 탐지	1,179	1,093	1,049
블랙리스트 미탐지	2,282	2,368	2,412
* 번역 정확도 (%)	65.93	68.41	69.69
평균 BLEU 점수	30.04	13.47	33.83

* 번역 정확도 : 전체 문장 중 블랙리스트 미탐지된 문장 수

5. 평가

- 번역 품질과 BLEU 점수가 비례하지 않는 예시

관용구	가닥을 잡다
한국어 원문	당초 강 위원장은 국정감사에서 이 문제가 불거지자 “감사원 감사를 받겠다” 며 버텼지만 즉각 사퇴로 <u>가닥을 잡았다</u> .
영어 번역 쌍	At the beginning, Representative Kang said, “I will be audited by the auditor,” when the matter was raised in the state audit.
Naver Papago 번역 결과	Initially, Kang endured the issue during a parliamentary audit, saying he would undergo an audit by the Board of Audit and Inspection, but he immediately <u>decided to</u> step down.
BLEU 점수	4.7930e-76

6. 결론

- 관용구가 포함된 다량의 문장 데이터셋 구축
 - KISS : 420개 관용구, 7,500개 한-영 번역 쌍 데이터셋
- 관용구 번역에 특화된 번역 평가 지표 생성
 - 275개 관용구, 3,461개 번역 쌍에 대한 블랙리스트
- 블랙리스트를 이용한 기계번역 서비스의 번역 정확도 측정 방법 제안

6. 결론

- 응용 가능한 연구 분야
 - NMT에서 관용구뿐만이 아니라 속어, 속담과 같이 함축적인 의미를 지닌 어휘를 학습하기 위한 데이터셋 구축
 - NMT 번역 품질 평가
 - 한-영 번역 쌍 뿐만이 아니라 다른 종류의 언어 쌍에도 적용 가능
- 향후 연구 방향
 - 관용구 데이터셋을 이용해 NMT 에 관용구를 효과적으로 학습시키는 방법
 - 번역 품질 결과를 비교하여 최적의 관용구 학습 방법 도출



감사합니다 – Thank you