

한국코퍼스언어학회

2023 가을
10주년 학술대회
한국의 말뭉치 구축 및 연구의 지평

한국코퍼스언어학회

2023년 가을
한국코퍼스언어학회 10주년 기념 학술대회

한국의 말뭉치 구축 및 연구의 지평

| 일시 | 2023. 10. 7 (토)

| 장소 | 연세대학교 위당관

| 주최 | 한국코퍼스언어학회 | 주관 | 연세대 언어정보연구원

| 후원 | NCsoft, 플리토, 테디썸, 연세대 언어정보연구원

2023년 가을
한국코퍼스언어학회 10주년 기념 학술대회

한국의 말뭉치 구축 및 연구의 지평

| 일시 | 2023. 10. 7 (토)

| 장소 | 연세대학교 위당관

| 주최 | 한국코퍼스언어학회 | 주관 | 연세대 언어정보연구원

| 후원 | NCsoft, 플리토, 테디썸, 연세대 언어정보연구원

2023년 가을 한국코퍼스언어학회 10주년 기념 학술대회

한국의 말뭉치 구축 및 연구의 지평

2023. 10. 7 (토), 연세대학교 위당관

| 일시 | | 발표 및 내용 | | | 비고 | | |
|-------------|-------------|--|--|---|---------------------------|--|--|
| 제 1 부 | 09:30~10:00 | 등록 | | | | | |
| | 10:00~10:30 | 개회식 | • 개회사 _ 류병래(한국코퍼스언어학회 회장, 충남대) • 축사 _ 이석재(언어정보연구원 원장, 연세대) • 학회 10주년 기념 동영상 상영 | | | | |
| 제 2 부 | 분과 발표 | | Session 1 <법언어학> 사회: 이선웅(경희대) 장소: 문과대 100주년 기념홀(위당관 6층) | Session 2 <수어 코퍼스> 사회: 윤영민(연세대) 장소: 위당관 415호 | | | |
| | 10:30~11:00 | 법언어학이란 무엇인가? 발표: 조기현(군산대) / 토론: 이해윤(한국외대) | | 조선대수어코퍼스구축과 웹 자료 추출 과정 발표: 최영주(조선대), 박호민(한국해양대), 장세은(한국해양대), 권순복(부산대), 윤병천(조선대), 조참훈(중국 서남대) 토론: 이선혜(부경대) | | | |
| | 11:00~11:30 | 법률해석에 대한 언어학적 접근 발표: 이해윤(한국외대) / 토론: 이용훈(충남대) | | 코퍼스 기반 한국어와 한국 수어 어휘 비교 분석 발표: 장세은(한국해양대), 박호민(한국해양대) 토론: 전혜원(부경대) | | | |
| | 11:30~12:00 | 자살유서의 법언어학적 분석 발표: 이용훈(충남대) / 토론: 정유진(고려대) | | 코퍼스 기반 한국 수어의 등위구문 연구 발표: 고인경(나사렛대), 윤병천(조선대), 전혜원(부경대), 이선혜(부경대), 박호민(한국해양대), 장세은(한국해양대) 토론: 강아름(충남대) | | | |
| 12:00~13:00 | | 점심 | | | | | |
| 제 3 부 | 13:00~13:40 | 주제 특강 1 | 국가 말뭉치의 구축 및 활용 _ 강미영(국립국어원) | | | | |
| | 13:40~14:20 | 주제 특강 2 | 학습자 말뭉치와 한국어 교육 _ 강현화(연세대) | | | | |
| | 14:20~15:00 | 주제 특강 3 | AI용 데이터셋으로서의 말뭉치 _ 이연수(엔씨소프트) | | | | |
| | 15:00~15:40 | 주제 특강 4 | AI 기반 번역과 다국어 병렬 말뭉치 _ 이정수(플리토) | | | | |
| 15:40~16:00 | | 휴식 | | | | | |
| 제 4 부 | 분과 발표 | | Session 3 <인공지능의 한국어 능력 평가> 사회: 안성민(카이스트) 장소: 문과대 100주년 기념홀(위당관 6층) | Session 4 <코퍼스 언어학 일반> 사회: 윤소연(인천대) 장소: 위당관 415호 | | | |
| | 16:00~16:30 | AI말평의 현황 및 발전 방안 발표: 이보라미(국립국어원) | | R를 활용한 정도부사에 대한 계층적 군집 분석과 다중 대응분석 발표: 제조현(연세대) / 토론: 신서인(한림대) | | | |
| | 16:30~17:00 | 감정 분석 과제의 특징과 의의 발표: 김한샘(연세대) | | AI-Based New Approach to Polysemy: BERT and t-SNE 발표: 김원빈(연세대) / 토론: 정병철(경남대) | | | |
| | 17:00~17:30 | 이야기 완성 과제의 특징과 의의 발표: 함영균(테디썸) | | A Generative AI Study on Translation Universals 발표: 임진(이화여대), 이용훈(충남대) / 토론: 강병규(서강대) | | | |
| 17:30~18:00 | | 표 기반 유사 문장 말뭉치를 이용한 AI 언어능력 평가 발표: 송영숙(Sionic.ai), 정유남(중앙대), 유현조(서울대) | | | | | |
| 18:00~18:10 | | 휴식 | | | | | |
| 제 5 부 | 18:10~18:20 | 연구 윤리 교육 | | | 사회: 이용훈(충남대) | | |
| | 18:20~18:30 | 폐회식 | • 폐회사 _ 류병래(한국코퍼스언어학회 회장, 충남대) | | | | |
| | | | | | 장소: 문과대 100주년 기념홀(위당관 6층) | | |

2023년 가을 한국코퍼스언어학회 10주년 기념 학술대회

한국의 말뭉치 구축 및 연구의 지평

목 차

◎ Session 1 : <법언어학>

| | |
|------------------|----|
| 법언어학이란 무엇인가? | 1 |
| 법률해석에 대한 언어학적 접근 | 11 |
| 자살유서의 법언어학적 분석 | 21 |

◎ Session 2 : <수어 코퍼스>

| | |
|----------------------------|----|
| 조선대수어코퍼스구축과 웹 자료 추출 과정 | 37 |
| 코퍼스 기반 한국어와 한국 수어 어휘 비교 분석 | 63 |
| 코퍼스 기반 한국어 수어의 등위구문 연구 | 89 |

◎ 주제 특강

| | |
|-----------------------------|-----|
| 주제 특강1: 국가 말뭉치의 구축 및 활용 | 109 |
| 주제 특강2: 학습자 말뭉치와 한국어 교육 | 129 |
| 주제 특강3: AI용 데이터셋으로서의 말뭉치 | 149 |
| 주제 특강4: AI기반 번역과 다국어 병렬 말뭉치 | 165 |

◎ Session 3 : <인공지능의 한국어 능력 평가>

| | |
|-----------------------------|-----|
| AI말평의 현황 및 발전 방안 | 181 |
| 감정 분석 과제의 특징과 의의 | 193 |
| 이야기 완성 과제의 특징과 의의 | 211 |
| 표 기반 유사 말뭉치를 이용한 AI 언어능력 평가 | 229 |

◎ Session 4 : <코퍼스 언어학 일반>

| | |
|---|-----|
| R를 활용한 정도부사에 대한 계층적 군집 분석과 다중 대응분석 | 239 |
| AI-Based New Approach to Polyseme: BERT and t-SNE | 243 |
| A Generative AI Study on Translation Universals | 255 |
| AI 학습용 한중 병렬말뭉치의 번역문 평가 | 279 |

2023년 가을 한국코퍼스언어학회 10주년 기념 학술대회

한국의 말뭉치 구축 및 연구의 지평

Session 1

<법언어학>

법언어학이란 무엇인가?

법률해석에 대한 언어학적 접근

자살유서의 법언어학적 분석

법언어학이란 무엇인가?

조기현 (군산대학교)

1. 서론

언어학은 인간 활동으로서의 언어를 과학적으로 연구하는 학문이며 언어의 구조뿐만 아니라 다양한 상황에서 언어의 작용 방식을 연구한다. 현대 언어학의 주된 연구 대상은 보편문법 원리라고 할 수 있지만 언어가 이러한 보편적 특성만 가지고 있는 것은 아니다. 언어는 그 언어 사용자가 속한 집단, 교육 정도, 사회적 지위, 또는 사용자 개인이 생성하는 말이나 글의 진실성 여부 등을 판단할 수 있는 사용자 고유의 언어적 특성을 가진다. 이러한 언어 사용자의 특이성(idiosyncrasies)을 연구 대상으로 하는 응용언어학의 한 분야가 법언어학(Forensic Linguistics)라 할 수 있다. 법언어학은 텍스트 분석을 위한 기술과 언어학적 지식과 담화 분석 등의 다양한 연구 방법을 활용하여 사이버 범죄를 포함한 다양한 법적 및 범죄 문제를 해결하려는 응용 언어학의 한 분야이다.

법언어학은 언어학의 과학적 연구 방법을 사용하여 구어 및 문서 자료를 검토하고 분석한다. 특히, 문서의 저자가 누구인지 결정하거나 녹음된 대화와 같은 구어 자료의 발화자를 식별하는 것과 같은 문제에 관련이 있다. 법 절차에서 특정 언어 사용을 기반으로 누가 무죄이고 누가 유죄인지 증명하는 데 도움을 주기 위한 것이 법언어학의 기본적인 활용 방안이지만, 최근에는 법언어학이 언어 기반의 증거 분석, 사이버 범죄자 식별, 위협 평가 및 디지털 영역에서의 법 집행 및 법적 노력을 지원하는 데 중요한 역할을 하고 있다. 이것은 언어학 분석과 사이버 범죄 조사 사이의 연결고리가 역할을 하며, 다양한 온라인 불법 활동의 예방 및 기소에 기여한다고 할 수 있다.

법언어학과 관련하여 미국과 유럽, 아시아의 주요 국가들에서 많은 연구가 이루어지고 있는 반면에 한국에서는 법언어학이라는 학문 분야에 대하여 연구도 관심도 거의 전무한 상황이다. 본 발표에서는 법언어학이라는 학문 분야에 대한 개념, 그 연구 영역, 법언어학의 태동과 발전 과정 그리고 미래 전망에 대해 살펴보고 법언어학이 적용된 대표적인 사례들도 함께 살펴보자 한다. 또한 이번 발표를 통해 법언어학이 응용언어학의 한 분야로서 언어학의 미래 가치를 높이는데 기여할 수 있으며 아직 불모지와 같은 한국에서도 법언어학에 대한 많은 연구와 관심이 필요함을 확인하고자 한다.

2. 법언어학

2.1. 법언어학(Forensic Linguistics) 개념

법언어학은 간단히 정의하면, “법언어학은 법적 문제에 대한 언어학적 적용”이라고 할 수 있다.(Olsson, 2013). “Forensic”이란 용어는 과학이나 학문의 적용을 통하여 사법 적용(enforcement)을 지원하기 위해 사용된 과학을 의미하는 라틴어 ‘*forēnsis*’에서 유래한 법률 용어이다. 법은 일반적으로 언어로 기록되어 있어 언어 없이는 법이 존재하지 않는다고 할 수 있다. 이러한 법과 언어의 관계에서 법언어학은 언어, 범죄 조사, 재판 및 사법 절차 등 형사 및 민사법의 맥락에서 언어학 지식, 방법 및 통찰력을 적용하는 응용언어학의 한 분야이다.

법언어학이라는 용어는 Jan Svartvik(1968) 언어학 교수가 “Timothy John Evans 사건”的 자술서 분석에서 처음으로 사용하였고, 1990년대 이래로 국제 법언어학자 협회(IAFL; The

International Association of Forensic Linguists)와 ‘Forensic Linguistics: The Journal of Speech, Language and the Law’ 저널을 통해 널리 사용되고 있다(Tiersma & Solan 2003, 213; Olsson 2008, 3). 법언어학은 목소리 분석, 법이나 법률 관련 문서의 해석, 법적 상황에서의 담화 분석, 구두나 글로 된 진술서의 의도된 의미의 해석(고백이나 자수), 저자식별, 법률 언어, 재판 참여자(판사, 변호사 그리고 증인)의 법정 언어 분석, 상표법, 그리고 법적 상황에서 사용되어야 하는 둘 이상의 언어에 대한 통역과 번역 등 사법 절차에서 언어와 법률 관련된 모든 영역을 연구 대상으로 하는 응용언어학의 한 분야이면서 학제간 학문 분야라고 할 수 있다(Coulthard 1994; McMenamin 2002; Shuy 1993).

지문 식별(fingerprinting identification)과 같은 분야들이 사법절차에서 확립되어 있는 것과 비교하면 법언어학은 상대적으로 생소한 영역이라고 할 수 있지만 음성학, 음운론, 문법이론, 인지언어학, 화행이론, 담화분석, 대화분석, 언어와 기억 등과 같은 다양한 언어학적 이론이나 지식을 사법절차 상 언어적 증거를 분석하거나 제공할 필요가 있는 경우에 적용될 수 있다는 점에서 법언어학은 그 자체로 과학적 특성을 지니고 있다고 할 수 있다.

2.2. 법언어학의 연구 영역

법언어학의 연구 영역은 다양하다. 사법절차에서 문제가 되는 말과 글과 관련된 증인이나 피의자 자술서의 문제를 해결하는 것부터 시작하여, 법률(forensic)음성학과 음향 장치를 이용한 통화 분석, 그리고 사법 및 범죄 영역에서 이른다. 최근에는 테러리스트 사건에서 저자 식별, 제품오염사건, 의사소통, 사법적 문서에서의 의미해석, 사망 시간 확인을 위한 핸드폰 문자메시지의 분석 등 다양한 영역에서 연구가 활발하게 이루어지고 있다.

법언어학의 연구영역은 3개의 영역으로 분류될 수 있다(Gibbons & Turell, 2008). 첫 번째로, “법률 언어(language of the law)”는 법적 맥락에서 사용되는 언어를 분석하는데 중점을 둔 법언어학의 영역이라고 할 수 있다. 법률 언어는 일상적인 말과 글과는 차이가 있어, 법률 체계에 익숙하지 않은 비법률 전문가에게는 이해하기 어려운 특성이 있다. 따라서 법언어학자는 이 영역에서의 연구를 통해 법률 언어와 일상적인 이해 사이의 간극을 줄여 법률 체계의 효과적인 이해를 보장하는 중요한 역할을 한다. 두 번째는 “법정 언어(language of the court)”로, 법적 환경에서의 언어학적 특성, 의사소통 전략, 그리고 담화 패턴에 대한 분석을 중점으로 하는 영역이다. 이 영역은 법언어학 내에서 매우 중요하며, 특정 용어의 사용법이나 설득 기법, 그리고 언어를 통한 신뢰성과 권위 확립 방법 등을 연구한다. 세 번째로, “증거로서의 언어(language as evidence)” 연구 영역에서는 언어가 범죄나 법률 사건의 해결에 있어서 어떻게 증거로 사용될 수 있는지에 대해 분석한다. 이 분야는 텍스트의 저자나 발화자를 확인하거나 의사소통의 패턴을 식별하거나 법률 사건에서 중요한 맥락을 제공하기 위해 문어체 또는 구어체로 된 의사소통, 구문, 단어 선택 및 스타일과 같은 언어의 다양한 측면을 분석하는 것을 포함한다. 이 법언어학 분야에 포함된 대표적인 하위 연구분야로는 저자 식별과 법률 음성학이다. 저자 식별은 논쟁의 대상이거나 익명이거나 의심되는 텍스트의 저자를 식별하기 위해 언어학자들이 수행하는 과정이다 (Coulthard et al., 2017). 저자 식별에서 분석될 수 있는 텍스트 유형은 다른 사람처럼 가장해서 보낸 문자 메시지, 협박 편지나 문자, 조작된 자살 메모, 위협 텍스트, 및 몸값 메모 등 다양하다. 저자 식별은 글에 사용된 언어를 조사하여 저자가 누구인지 결정하는 것을 목표로 한다. 법률 음성학은 음성학을 구어체와 관련된 다양한 법적 상황에 적용하는 것을 의미하며 녹음된 대화나 발화가 증거로 사용되는 법률 사건에서 발화자의 정체성을 밝히는 것이 주요 기능이라고 할 수 있다.

국제법언어학자협회(IAFL)에서는 법언어학의 관련 영역을 언어와 법률, 사법절차에서의 언어, 증거로서의 언어, 그리고 연구 및 교수 분야로 세분하여 각 영역마다 다양한 연구 주제들을 제시하고 있다. 먼저, 언어와 법률(language and law) 영역은 입법; 법적 문서에 대한 이해; 법적 텍스트의 분석 및 해석; 법률 장르; 법적 언어의 역사; 법적 담론; 법적 맥락에서의 다언어 문제; 법적 자료의 담화 분석; 법 앞에서의 언어 불이익; 언어 소수자와 법률 체계; 언어 권리; 권리와 법; 법적 맥락에서의 문화 간 문제 및 중재 등이 연구 대상이 된다. 법 절차에서의 언

어(language in the legal process) 영역은 취약 증인과의 인터뷰; 취약 증인의 의사소통 어려움; 경찰 인터뷰; 조사 인터뷰; 망명자의 언어 테스트; 이중 언어 법정 및 제2언어 문제; 법정 통번역; 법정 상호 작용; 법정 언어; 경찰 언어; 교도소 언어; 공법과 사법 법정에서 판사와 배심원에게 전달되는 언어 등이 포함된다. 증거로서의 언어(language as evidence) 영역에는 저자 분석 및 식별; 표절; 발화자 식별 및 음성 비교; 말뭉치 구축(진술, 고백, 자살 메모); 컴퓨터를 활용한 저자 식별 또는 프로파일링; 소비자 제품 경고; 민사 사건에서의 언어 증거 (상표, 계약 분쟁, 명예 훼손, 제조물 책임, 사기 거래 행위, 저작권 침해); 방언학 및 사회 언어학; 의미론; 화용론 및 발화 행위 분석이 포함된다. 마지막으로 연구/교육(research and teaching) 영역에는 전문 증인의 실무 및 윤리; 언어학적 증거 제시; 전문 증인으로서의 언어학자; 법언어학/언어 및 법률의 교육/테스트; 법률 전문가를 위한 언어 교육 등을 그 연구 내용으로 한다.

2.3. 법언어학의 연구 방법

법언어학은 언어의 보편적 특성을 분석하는 것이 아니라 보편적 언어 특성 중에서도 언어 사용자 개인이 가지는 독특한 특성(idiosyncrasies)을 연구 대상으로 하며, 개인방언(idiolect)의 개념을 전제로 한다. 이러한 개인방언과 관련된 법언어학의 대표적인 연구 대상이 저자 식별(authorship attribution)이라고 할 수 있다. 법언어학의 연구 방법을 저자 식별의 방법론을 중심으로 살펴보면 정성적(qualitative) 접근법인 문체론적 방법(stylistics)과 정량적(quantitative) 접근법인 계량문체론적 방법(stylometry)으로 구분될 수 있다.(이해윤 2023).

2.3.1. 법문체론(forensic stylistics)

일반적으로 언어학적 문체론(linguistic stylistics)은 개인이나 집단의 언어에서 관찰되고 기술되고 분석되는 문체 표지들(style markers)(vocabulary choice, use of certain idioms or phrases, spellings, use of capitalization, citation style, date format, consistent mistakes, etc)에 대한 과학적 해석이다.(이해윤) 이러한 언어학적 문체론을 법 영역에 적용한 분야를 법문체론이라고 볼 수 있다. McMenamin (2002, 163)에 따르면, 법문체론(a.k.a. comparative linguistics)은 “소송에 있어서 텍스트의 저자를 판별하고 관련된 결론이나 견해를 도출하기 위해 문체론적 분석을 수행한다. 문체론은 문어에서 나타나는 변이 패턴을 과학적으로 연구한다. 연구의 대상은 개인의 언어(개인방언)이고, 그의 언어적 특징들을 기술하게 된다.”라고 설명하고 있다.¹⁾

2.3.2. 계량문체론(stylometry)

계량문체론(stylometry)에서는 주로 컴퓨터를 이용하여 개인의 문체를 규정할 수 있는 자질들을 조사하고, 결과 도출 및 분석에 통계적 기법을 사용한다. 관련 자질들로는 타임과 토큰의 비율, 어휘밀도, 유일 어휘의 개수, 구두점의 분포와 빈도, 단어의 평균 길이, 문장의 평균 길이, 단어당 음절 평균 개수, 형태소 빈도, 단어빈도, 연어 사용 빈도 등 텍스트의 다양한 요소들이 해당된다.²⁾ 계량문체론적 접근법은 주로 코퍼스 언어학과 관련이 있다. 코퍼스 언어학은 기본적으로 언어 연구의 경험적 접근법으로서, 어휘, 문법, 의미론에 대한 일반화를 하기 위해 그리고 법언어학적 맥락에서 설명 가능한 것 이상의 것들을 제공하기 위해 실제로 존재하는 데이터를 이용한다. 이는 특정 언어 사용의 특성과 특정 (언어 사용 그룹의) 행동을 설명하기 위해 자연스럽게 발생하는 언어 데이터를 처리하기 위한 방법을 제공한다.(Krzysztof Kredens, Malcolm Coulthard, 2012).

2.3.3. AI 및 기계학습 기반 분석

최근에는 생성형 AI나 기계학습 기반의 분석 방법이 법언어학 연구에 널리 활용되고 있다. 대표적인 사례로 AI가 강화된 표절 검사는 표절 문제에 대한 실질적인 해결책이 될 수 있다. 표절을 감지하기 위해 법언어학과 컴퓨터가 인간의 언어를 이해하고 해석하며 조작할 수 있도록 하는 전산언어학의 형태인 자연어처리(NLP)와 같은 AI 기술을 활용하는 표절 검사 프로그램 (Dataman, 2021)은 텍스트를 그대로 비교하는 것이 아니라 어느 정도 텍스트 내용을 이

1) 법문체론적 접근법의 문제점에 대하여는 이해윤(2023, 379)을 참고 바람.

2) 계량문체론 접근법의 문제점에 대하여는 이해윤(2023, 386)을 참고 바람.

해하고 아이디어가 표절되었는지 여부를 판단하게 된다. AI가 법언어학에 가져온 주요 통찰력 중 하나는 분석 절차를 형식화하고 컴퓨터 알고리즘에 적용함으로써 제안된 방법의 평가가 더 쉽고 정확하며 포괄적으로 이루어지게 할 수 있다는 것이다. FLINT Systems는 이러한 대표적인 프로그램 중 하나이다. FLINT Systems는 "문서가 확인된 저자에 의해 작성되었는지 여부를 감지하기 위해 설계된 최초의 언어학적 도구"이며, 이 시스템은 작성된 문서가 AI에 의해 작성되었는지 여부를 감지하는 것뿐만 아니라 저자로 주장된 저자에 의해 작성되었는지 여부까지도 판단할 수 있다(David Nagel, 2023).

3. 법언어학의 태동, 현재 그리고 미래

3.1. 법언어학의 태동

법언어학은 상대적으로 신생 학문 분야라고 할 수 있다. 1949년에 F.A. Philbrick은 법률 영어에 관한 책 "Language and the Law: the Semantics of Forensic English"에서 "forensic English"라는 용어를 처음 도입하였으나, 이 표현은 당시에는 크게 받아들여지지 않았다.³⁾ 1968년, 언어학자 Jan Svartvik이 "The Evans Statements: A Case for Forensic Linguistics"라는 연구에서 "법언어학"이라는 용어를 처음으로 언급하였다. 스웨덴의 언어학자 Svartvik는 법적 절차 내에서 언어 분석의 중요성을 강조하기 위해 이 용어를 도입하였고, 이후로 법언어학은 법 집행 맥락에 적용되는 다양한 언어학 전문성을 포괄하도록 발전해왔다. 1990년대 초 법언어학 세미나가 진행되면서 학문 분야로서 성숙해졌고, 1990년대 중반에는 국제법음성학회 (International Association for Forensic Phonetics, IAFP), 법언어학: 말, 언어 및 법 국제저널 (Forensic Linguistics: The International Journal of Speech, Language and the Law) 및 국제법언어학자협회(International Association of Forensic Linguists, IAFL)의 설립, 그리고 학회지가의 발간을 통해 법언어학은 개별 학문 분야로서의 공식적인 인정을 받고 발전하기 시작하였다.

3.2. 법언어학의 적용 사례

법언어학이 사법절차에서 언어학적 분석의 중요성이 확인되어 법언어학이 본격적으로 발전하게 된 대표적인 사건들이 있다.

3.2.1. 영국

영국에서 법언어학의 역할을 완벽하게 보여준 사례가 "Timothy Evans 사건"이다. 1968년에 언어학자 Jan Svartvik은 'The Evans Statements: A Case for Forensic Linguistics'에서 Timothy Evans가 경찰에 제출한 자백 진술서의 분석을 소개하였다. 이 연구에서 Svartvik은 법언어학이라는 용어를 처음으로 사용하였다. Evans는 1950년에 그의 아내와 매우 어린 딸을 살해한 혐의로 유죄 판결을 받고 사형을 선고받았으나, 1966년에 사후사면을 받았다. Svartvik은 질적 및 양적 분석의 조합을 사용하여 문서의 저자에 대한 의심을 제기하였다. 그는 진술서 전체에서 일관성 없는 문법 특징들을 지적하였고, 이런 언어적 특징들이 경찰이 원래의 진술에 추가한 것이라고 주장하였다. 일명 "Bentley 사건"이 법언어학이 사건의 진실을 밝히는데 결정적인 역할을 한 또 하나의 대표적인 사례이다. 1952년, Derek Bentley는 경찰을 살해한 혐의로 재판을 받고 사형을 선고받았다. 언어학자 Malcolm Coulthard는 담화분석 방법을 이용하여 Bentley의 4개의 자술서들을 연구하였고, Bentley의 자술서 내용이 그의 개인방언(idiolect)과 일치하지 않는 여러 부분을 발견하였다. 이로 인해 검찰은 Bentley의 자백이 경찰에 의해 조작되었음을 인정하였고, Bentley는 1993년에 사후사면을 받았다.

3.2.2. 미국

미국의 경우에 대표적인 사례 중의 하나가 Unabomber 사건이다. 이 사건은 법언어학의 방법론 중에서 저자 식별과 비교언어학이라고 할 수 있는 법문체론적 방법을 통해 사건을 해결한 대표적인 예이다. 1970년대부터 90년대까지 Kaczynski라는 남자는 미국 전역에 무작위로 직접 만든 폭탄을 익명으로 보냈고, 종종 대학교와 비행기를 표적으로 삼았기에 Unabomber라

는 별명을 얻었다. 1995년 Kaczynski는 FBI에게 35,000단어의 선언문을 보냈고, FBI는 누군가가 Unabomber의 개인 언어를 인식하고 그를 FBI에 보고하기를 바라며 그 선언문을 신문에 게재하였다. Kaczynski의 동생이 수년 동안 그의 형으로부터 받은 수백 통의 개인 편지를 통해 선언문의 요지와 문체가 자신의 형의 것과 유사하다는 것을 제보하였다. 여기서 비교 언어학이 큰 역할을 하게 된다. Kaczynski의 개인 편지와 선언문 사이에는 유일한 철자, 어휘 선택, 속어, 문법 구조 등 수천 가지 어휘 일치점을 발견하였다. 특히, 하나의 어휘 항목이 결정적인 언어학적 증거가 되었고, Kaczynski는 이를 토대로 범인으로 체포되었다. 선언문과 개인 편지 모두에서 Kaczynski는 "You can't eat your cake and have it too"라는 초기 현대 영어 속담을 사용했는데, 이는 현대 영어 버전인 "You can't have your cake and eat it too."와는 달랐다. 이 개인방언적 속담의 사용은 결국 그의 정체성을 드러내는 데 큰 역할을 한 것이다.

언어적 분석과 계량문체론(stylometry)을 이용하여 저자 식별을 한 대표적인 사례는 J.K. 롤링이 "The Cuckoo's Calling" 소설의 실제 저자임을 밝혀낸 것이다. 2013년 Patrick Juola는 계량문체론, 즉 연산(computational) 언어학의 한 분야를 활용한 프로그램으로 소설의 문체를 분석하였다. 소설 "The Cuckoo's Calling"의 경우 Robert Galbraith이 아니라 해리 포터의 저자인 J.K. 롤링이 실제 저자일 가능성을 판단하기 위해 J.K. 롤링의 'The Casual Vacancy' 포함한 여러 탐정 소설들을 수집하였고, Juola는 네 가지 개별적인 특질 집합(단어의 길이, 4개의 문자로 이루어진 조합 (4-grams), 단어 쌍, 그리고 가장 자주 사용되는 100개의 단어)을 기반으로 다양한 언어적 특징을 검사하여 텍스트 간의 유사성을 비교 분석하였다. 계량문체론적 분석 방법은 DNA 분석처럼 저자를 100% 확실히 증명할 수 없지만 저자의 유사성에 대한 강력한 증거를 제공한다. J.K. 롤링은 "The Cuckoo's Calling"의 진정한 저자가 자신이라고 인정하여 그 프로그램의 결과를 확인해 주었다.

3.2.3. 한국

한국의 경우에는 '삼례 나라슈퍼 3인조 강도치사 사건'에서의 폭행과 강요에 의한 자술서 작성이 대표적인 예라고 할 수 있다. 이 사건에서 3명의 피의자 중 2명의 자술서 작성자는 지적 장애를 가지고 있었지만, 작성자의 언어능력을 고려하지 않고 자술서 내용의 합리성이나 경위 등에 대하여 신중한 판단을 하지 않아 유죄 판결을 받은 후 17년 만에 재심에서 무죄판결을 받았다. '삼례 나라슈퍼 강도치사 사건'의 피의자의 자술서 내용을 분석해 보면, 우선 자술서 작성자는 언어능력과 관련하여 한글을 쓸 줄 모르는 지적 장애를 가지고 있었다는 점을 재판과정에 누구도 고려하지 않았다. 작성자는 자술서 내용을 직접 쓴 것이 아니라, 당시 경찰이 써 준 것을 '그대로 그렸다'고 진술하였다. 자술서는 지적 장애를 가진 자가 작성했다고 볼 수 없을 정도로 매우 긴 문장으로 작성되었고, 상황에 대한 설명도 일반인의 경우보다 매우 논리적으로 되어 있었다. 이러한 피의자 자술서에서 특히 주목해야 할 언어학적 특징은 '나는 그때' 또는 '저는 그때'라는 표현이다. Evans 사건과 Bentley 사건에서와 유사하게 'I then'이라는 영어표현은 경찰영역에서 사용되는 일종의 사용역인데 한국의 3인조 강도치사 사건의 피의자 자술서에도 동일한 통사적 특징을 보이는 표현이 나타나는 것을 확인할 수 있다. 위 자술서의 내용은 경찰의 강요와 폭행에 의해서 작성되었다는 것으로, 경찰이 불러 주거나 써준 것을 그대로 그렸다는 것이다. 이는 한국어나 영어에 있어서 주어 뒤에 시간부사를 바로 위치시키는 문체는 경찰 언어에서만 보이는 특이한 언어적 특징이라고 판단할 수 있을 것이다. 이러한 사용역은 더 나아가서는 개인방언을 구별해내는 결정적인 단서를 제공하고 있는 것이다.

4. 세계 법언어학의 발전과 동향

4.1. 유럽 및 북미

20세기 초, 언어학자들은 주로 언어 해석, 번역, 그리고 언어 변이와 관련된 문제에서 언어학 전문성을 제공하기 위해 법적 문제에서 전문 증인의 역할을 수행하기 시작했다. 1950년대에서 1960년대에는 Dr. Peter Ladefoged와 Dr. Wallace Lambert와 같은 언어학자들이 발화자식별 및 목소리 분석, 특히 법적 사안에서의 발화 샘플 연구에 사용되는 음성학적 분석 분야인 법음성학 분야를 선도했다. 1970년대에서 1980년대에는 Dr. Roger Shuy와 Dr. William

Labov와 같은 학자들이 문서 작성 스타일 및 문서의 언어학적 패턴 분석, 저자 귀속, 사기 감지, 그리고 언어 변이를 중점으로 한 법문체론 분야를 발전시켰다.

1988년 독일의 연방 범죄 경찰청 (Bundeskriminalamt)에서 법언어학 회의를 개최한 것을 시작으로, 1990년대에 들어서며 발화자 식별의 음성-음향 방법을 논의하면서 법언어학이 주목을 받기 시작했다고 할 수 있다. 1991년에는 프랑스에서 법언어학 회의가 열렸고, 1992년에는 영국에서 처음으로 법언어학 회의가 개최되었다. 1993년에는 법언어학자 국제 협회(IAFL)가 설립되었고 지금까지 2년마다 국제 학술대회를 개최하고 있다. 1994년에 “Forensic Linguistics” 저널이 출판되면서 법언어학 분야가 독립된 학문으로서 공식적인 인정을 받기 시작했다. 1990년대 후반에는 대학들이 법언어학을 가르치기 시작했으며, 영국과 미국, 그리고 스페인을 중심으로 점점 많은 나라들에서 법언어학에 대한 공식 교육을 제공하기 시작하고 있다. 2000년대 이후로, 법언어학은 지속적으로 발전하고 다양한 분야에서 응용되고 있다. 저자 귀속, 발화자 식별, 언어 기원 식별, 거짓 고백 분석, 법적 문서의 언어, 소설 미디어의 언어, 그리고 언어 정책 및 계획을 포함한 다양한 법적 맥락에서의 언어학 전문성의 응용이 확대되고 있다.

A (Scopus)

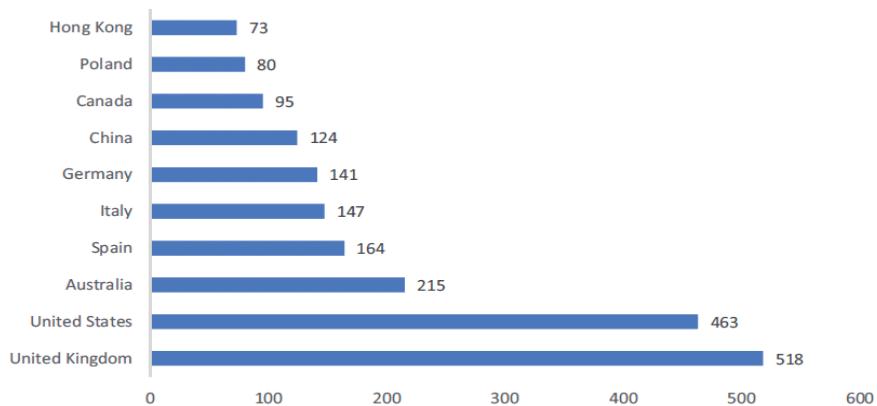


그림1. 국가별 법언어학 관련 연구 성과물(Aldualis 외 2023)

A (Scopus)

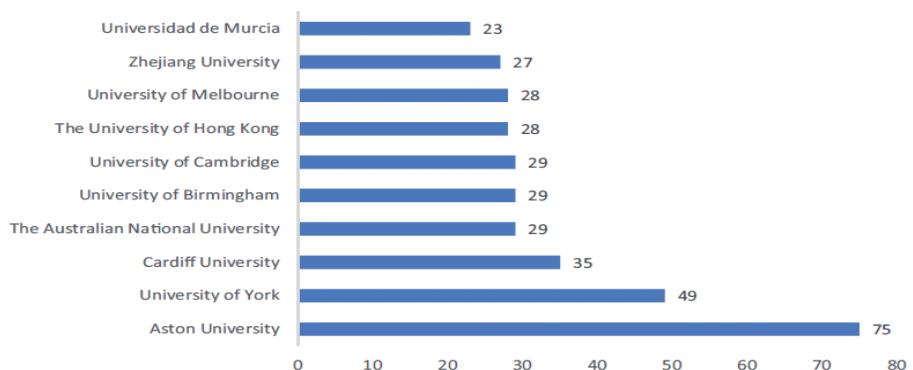


그림2. 세계 대학 및 연구기관별 법언어학 관련 연구 성과물(Aldualis 외 2023)

4.2 아시아

2002년도 중국에서는 CAFL(China Association of Forensic Linguistics)이 공식적으로 설립되어 활발한 활동을 이어가고 있다. 일본의 경우에는 법언어 연구회(2007)과 일본 법언어학회(2009)가 설립되어 활발한 연구 활동을 하고 있다. 기타 홍콩, 인도네시아 등도 학회 설립 등을 통해 활발한 연구 활동을 하고 있다. 한국은 법언어학 관련 연구소나 학회 활동의 거의 없는 상태이다. 법언어학 관련 논문의 경우에도 김기영(2007) 1편, 구명철(2009, 2014) 2편, 이해윤(2018 - 2022) 최소 9편, 이학문(2018) 1편, 최지명(2018) 1편 등 매우 적은 연구 성과를 내고 있다. 관련 도서 또한 서경숙, 니시야마 치나(2016)의 번역서 “법과 언어- 법언어학으로의 초대”와 이해윤(2023) “법언어학의 이해” 두 권 뿐이고 실질적인 법언어학 관련 도서는 이해윤(2003)의 도서가 유일하다고 할 수 있다. 국외에서는 영국, 미국, 스페인 등의 대학교에서 법언어학 교육을 위한 과정이 설치되어 운영되고 있다. 국내의 경우, 아직 법언어학에 대한 체계적인 연구나 교육이 이루어지지 않고 있는 실정이다. 법학전문대학원이나 법과학대학원에서 다양한 법융합학문들에 대한 강의가 개설되었지만 법언어학 강의는 없는 실정이다(이해윤, 2023). 한국의 경우에는 법언어학 관련 학회 설립, 전문인력 양성기관의 창설, 대학의 법언어학 연구 장려, 그리고 법언어학 연구 수준을 향상시키기 위해서는 우리나라로 선진국가들처럼 대학의 관련 전공분야의 연구 활동을 활성화시키고 실제 사법절차에서 그 전문성을 인정하는 방안을 적극 모색할 필요가 있을 것이다. 마지막으로 법언어학 유관 학문간 소통교류를 증대하는 것 또한 해결해야 할 시급한 과제라고 할 수 있다.

4.3. 법언어학의 미래: 사이버 범죄의 증가와 법언어학의 역할

4.3.1. 사이버 범죄의 특징

최근에 급증하고 있는 사이버 범죄는 컴퓨터, 통신 매체, 인터넷 등을 활용해서 일어나는 범죄로 단순한 컴퓨터를 이용한 범죄 뿐만 아니라 네트워크와 같은 다양한 연결성을 활용한 범죄까지 포함한다. 사이버 범죄는 기존 범죄와 비교할 때 다음과 같은 특징을 가진다: 1) 익명성: 범죄자가 자신의 신분을 숨기며 범죄를 저지를 수 있다. 2) 비대면성: 범죄자와 피해자 간의 직접적인 접촉 없이 범죄가 발생한다. 3) 빠른 전파성: 범죄에 필요한 정보나 프로그램이 넓은 범위로 신속하게 확산된다. 4) 시간적, 공간적 무제약성: 어디서나, 언제든 범죄가 발생할 수 있다. 5) 영속성: 한 번의 범죄 행위가 반복되거나 모방될 수 있다.

4.3.2. 사이버 범죄의 유형

사이버 범죄는 범행 목적에 따라 크게 사이버 테러형과 일반 범죄형으로 구분할 수가 있다. 사이버 테러형 범죄는 정보 통신망 자체를 공격 대상으로 이뤄지는 범죄 행위인 반면 일반 범죄형은 현실 세계의 범죄 행위를 사이버 공간상에서 일으키는 것을 의미한다. 사이버 테러형 범죄는 불특정 다수에게 심각한 피해를 입히는 행위들인데, 해킹, 서비스 거부 공격, 악성 프로그램 및 바이러스 유포 등이 대표적인 범죄 방법들이다. 일반 범죄형 사이버 범죄는 다음과 같이 그 유형을 구별할 수 있다. 첫째, 인터넷 사기는 정보통신망을 이용해 사용자들에게 물품이나 용력을 제공할 것처럼 기망하는 메시지를 보내 금품을 부정하게 취득하는 사기 행위를 의미한다. 둘째, 불법 복제는 저작권자의 허락 없이 저작물을 불법으로 복제, 배포하여 이익을 취하는 범죄 행위를 의미한다. 셋째, 불법, 유해 사이트 운영은 불법 도박 사이트, 불법 저작물 공유 사이트, 음란 사이트 등을 운영하여 이익을 취하는 범죄 행위이다. 넷째, 사이버 폭력이란 사이버 공간에서 정보통신망을 통하여 부호, 음향, 영상 등을 이용하여 타인을 모욕하거나 명예 또는 권익을 침해하는 온갖 형태의 폭력적인 표현과 행위를 의미한다. 마지막으로 개인 정보 침해는 개인의 신체, 재산, 사회적 지위, 그리고 신분 등에 관한 사실, 판단, 평가 등을 나타내는 일체의 모든 정보를 불법적으로 유포 및 취득하고 이를 활용한 범죄를 의미한다.

4.3.3. 사이버 범죄와 법언어학의 역할: 다크 웹 마약 밀매 사건 분석을 중심으로

사이버 범죄에서 법언어학이 어떤 역할을 할 수 있는지를 보여주는 대표적인 사례가 다크 웹 마약 밀매 사건이다. 다크 웹은 일반적인 검색 엔진에서 접근할 수 없는 인터넷의 숨겨

진 공간이며 종종 불법 활동이 이루어지며, 마약, 무기 및 기타 불법 상품의 판매가 흔하게 일어날 수 있다. 특정 다크 웹 사이트에서 큰 규모의 마약 판매 광고가 감지되고 자신의 상품을 홍보하고, 고객과 협상을 진행하며, 최종적으로 판매를 완료하기 위해 특정한 언어와 용어를 사용한다.

법언어학 전문가들은 해당 판매자의 글을 분석하여 개인의 고유한 언어 사용 패턴인 판매자의 개인방언(idiolect)을 파악한 후 동일한 저자가 다른 어떤 텍스트를 작성했는지 식별할 수 있다. 저자 분석(Authorship Analysis)을 통해 전문가들은 다크 웹 내에서 판매자가 작성한 다른 게시물과 비교하여 글쓰기 스타일, 용어 선택, 구문 및 다른 언어적 특징을 분석하여 특정 인물이 여러 게시물을 작성했을 가능성을 확인할 수 있을 것이다. 저자 분석의 결과를 바탕으로 경찰은 특정 인물이 이 마약 밀매 사건의 주요 용의자임을 확인할 수 있었고 용의자를 체포하였다. 이 사건은 언어학적 분석이 범죄 수사에 어떻게 중요한 역할을 할 수 있음을 보여 주는 대표적인 사례이다. 특히, 다크 웹과 같은 익명의 환경에서 범죄자의 신원을 파악하는 데 있어 법언어학은 중요한 도구로 사용될 수 있을 것이다.

5. 결론

법언어학은 언어와 법 사이의 교차점에서 중요한 연구를 수행하며, 특히 사이버 범죄와 같은 현대 사회의 복잡한 범죄 문제를 다루는데 큰 기여를 하고 있다. “법언어학”이라는 용어가 1968년에 언어학 교수 Jan Svartvik에 의해 처음 언급된 이후 법언어학은 많은 발전을 해오고 있다. 법언어학의 태동 시기에는 주로 저자의 신원 확인 및 방법론 연구에 집중되었고 오늘날 법언어학은 전문 저널과 국제 협회가 있고 국제학술대회가 개최되는 학문 분야로 자리매김을 하고 있다.

그러나 법언어학이 국제적으로 발전하고 있다 하더라도, 한국에서는 아직 이 분야가 초기 단계에 있다고 볼 수 있다. 한국에서 법언어학에 대한 관심과 연구가 활발하지 않은 주요 원인 중 하나는, 이 분야가 다양한 전문지식을 요구하기 때문일 수 있다. 법언어학 연구를 위해서는 언어학 뿐만 아니라 법, 사회학, 심리학 등 다양한 분야의 지식이 필요하다. 따라서 한국에서 법언어학을 연구하고 싶은 연구자나 학생들이 직면하는 어려움은 무시할 수 없다. 그럼에도 불구하고, 법언어학의 중요성은 점점 더 커지고 있다. 디지털 시대에는 소셜 미디어의 글, 음성 인식 기술을 통한 음성 데이터 등 다양한 새로운 데이터 소스가 법적 문제에 연관될 가능성이 높다. 따라서 이러한 새로운 데이터의 법적 해석을 위해서는 법언어학의 연구가 필수적이다.

한국에서 법언어학의 발전을 위해서는 다양한 노력이 필요하다. 한국의 대학과 연구 기관이 법언어학에 대한 연구를 공동으로 추진하면서 전문 인력을 양성하는 프로그램을 개발이 필요하다. 이를 위해서 정부는 법언어학 연구의 중요성을 인식하고 이에 대한 연구 자금 및 인프라 구축을 지원해야 할 것이다. 또한 법언어학에 관심있는 언어학자나 기관들은 선진국의 법언어학 연구 기관과의 협력을 통해 연구 방법론과 기술을 소개받고 공동 연구를 진행하도록 노력해야 한다. 특히 사법부와의 긴밀한 협력을 통해 실제 법정에서 언어학적 분석이 필요한 경우에 적용할 수 있는 연구 방안을 개발하는 것도 필요하다. 사회적으로는 법언어학의 중요성과 필요성에 대해 널리 알리는 활동을 진행해야 할 것이다. 한국이 법언어학 연구에 있어서 세계의 흐름에 뒤처지지 않기 위해서는 위와 같은 다양한 접근 방법이 필요하며, 이를 위한 지속적인 노력이 요구된다.

참고문헌

- 서경숙·니시야마 치나 역. (2016). 『법과 언어-법언어학으로의 초대』 박이정.
- 이해윤. (2023). 『법언어학의 이해 (개정판)』 역락.
- 앱백. (2022). 사이버 범죄의 개념과 특징, 세부 유형. <https://appbag.tistory.com/17>(블로그)
- 조기현. (2018). 언어학적 지문과 피의자 자술서의 저자식별. *언문학연구*, 57(3), 391-409.
- Ahmed Alduais, Mohammed Ali Al-Khulaidi, Silvia Allegretta & Mona Mohammed Abdulkhalek. (2023). Forensic linguistics: A scientometric review, *Cogent Arts & Humanities*, 10:1, 2214387, DOI: 10.1080/23311983.2023.2214387
- Ariani, M. G., Sajedi, F., & Sajedi, M. (2014). Forensic linguistics: A brief overview of the key elements. *Procedia-Social and Behavioral Sciences*, 158, 222-225.
- Balcells, Claudia Gallego and Núria Gavaldà. (2023). An Overview of Forensic Linguistics and its Application in Real-Life Cases. .
- Blackwell, S. (2012). *History of forensic linguistics. The Encyclopedia of Applied Linguistics*.
- Correa, Maite. (2013): Forensic linguistics: An overview of the intersection and interaction of language and law. *Kalbu Studijos* 23, 5-13.
- Coulthard, M. (1995). Questioning statements: Forensic applications of linguistics. Inaugural Lecture, University of Birmingham, February 16, 1995.
- Derin, Tatum, and Budianto Hamuddin. (2019). Exploring the Past, Present, and Future of Forensic Linguistics Study: A Brief Overview.
- Houtman, H. & Suryati Ali 2018. The History of Forensic Linguistics as an Assisting Tool in the Analysis of Legal Terms, *Sriwijaya Law Review* 2(2):215 DOI:10.28946/slrev.Vol2.Iss2.135,pp215-233
- Gibbons, J. (2003). *Forensic linguistics: An introduction to language in the justice system*. Wiley.
- Juola, Patrick. (2013). Artificial Intelligence and Pattern Evidence: A Legal Application for AI.
- Kondrashova, I. V. (2022). Forensic Linguistics: Origin, History of Development, Prospects. *Международный журнал гуманитарных и естественных наук* 4-3: 72-77.
- Kredens, K. J., & Coulthard, R. M. (2012). Corpus Linguistics in Authorship Identification. In L. M. Solan, & P. M. Tiersma (Eds.), *Oxford Handbook of Language and Law* Oxford University Press. <https://doi.org/10.1093/oxfordhb/9780199572120.013.0037>
- Kredens, Krzysztof, and Malcolm Coulthard. (2012). Corpus linguistics in authorship identification.
- Liu, Weiming. (2014). An Overview of Forensic Linguistics Research in China.
- Nagel, D. 2013 AI Writing Detection Tool Uses 'Forensic Linguistic' Techniques to Check Authorship <https://www.askedtech.com/weblink/140892>
- Olsson, J. (2014). What is Forensic Linguistics? the text. Co. UK/docs/what_is. Doc>. Acesso em, 25.
- Olsson, J., & Luchjenbroers, J. (2013). *Forensic linguistics*. A&C Black.
- Shuy, R. (1993). *Language Crimes: The use and abuse of language evidence in the courtroom*. Cambridge, MA: Blackwell.
- Solan, Lawrence & Tiersma, Peter. (2005). Speaking of Crime: The Language of Criminal Justice. BiblioVault OAI Repository, the University of Chicago Press.
- Sousa-Silva, Rui. (2018). Computational forensic linguistics: an overview of computational applications in forensic contexts. *Language and Law/Linguagem e Direito* 5.2, 118-143.
- Svartvik J. (1968). The Evans Statements: A Case for Forensic Linguistics. Goteborg: University of Goteborg.
- The International Association for Forensic and Legal Linguistics. (n.d.). *The International Association for Forensic and Legal Linguistics (IAFL)*. Retrieved 2023, from <https://www.iafl.org/>

법률 해석에 대한 언어학적 접근

이해윤 (한국외대)

목차

- 1.** 서론
- 2.** 언어의 불명확성
- 3.** 법해석학
- 4.** 언어학적 접근
 - 1.** 비판: 대화함축이론
 - 2.** 분석: 적합성이론

1. 서론

법의 명확성 원칙

■ 배경

- 법치주의:
국가가 국민의 자유와 권리를 제한하거나 국민에게 의무를 부과할 때에는 반드시 법률로써 해야 하고, 행정과 사법 작용도 법률에 근거를 두어야 한다.
- 죄형법정주의:
법률 없으면 범죄 없고 형벌 없다.
 - 관습형법금지의 원칙, 소급효금지의 원칙, 유추금지의 원칙,
 - 명확성의 원칙(Bestimmtheitsgrundsatz):
범죄와 형벌을 가능한 한 법률에 명확하게 확정해야 한다.

법률언어

- 전문어로서 법률언어
 - 일반어의 부분으로서 전문어
 - 전문용어 + 일반 어휘
 - 전문용어: 명확한 의미(**onomasiology**)
 - 일반 어휘: 일상 의미(**ordinary meaning**)
 - “Words in legal texts are given their ordinary meaning unless the context requires otherwise.” (Slocum 2015)
 - “When a law is ambiguous, vague, or otherwise underdetermines meaning, courts frequently default to the assumption that the drafters intended the words to be understood in their ordinary sense, taking context into account.” (Solan 2018)

2. 언어의 불명확성

의미의 불확정성

■ 미결정성 (indeterminacy)

- 어휘에 대응되는 개체들의 미세분화
- 언어유형별 차이

■ 중의성 (ambiguity)

- 사전 속의 여러 의미
- 문맥을 통해서 하나의 의미를 확정

■ 모호성 (vagueness)

- 경계선상의 의미: 의미가 경계선상에 있어서 어느 영역에 있는지 확정할 수 없는 경우
- 문맥이 정보를 제공하여 하나의 의미를 확정

*Sterling's cousin is funny.
Julian's brother is heavy.*

모호성 1/3

■ 배경

- 현실세계와 언어세계 간의 불일치: 세계에 존재하는 것들은 작은 단위로 구분될 수 있지만 연속성을 띤다. 언어적 표현들은 경계가 분명한 독립된 단위들이고 한정적이다.

- Alston (1967):
 - i) 개별화의 모호성 (vagueness of individuation): Mt. Everest, ...
 - ii) 분류의 모호성 (classificatory vagueness): red, orange, ...

■ 속성

- Contextual variability in truth conditions

The coffee in Rome is expensive.

- Borderline cases

*Mary is neither tall nor not tall.
Mary is both tall and not tall.*

- Sorites paradox

전제1: Any person who has \$100 million is rich.
전제2: If X is rich then X minus one dollar is rich.
결론: Any person who has \$ 99,999,999 is rich.

모호성 2/3

■ 분류 (Poscher 2012; Marmor 2014)

- 더미의 모호성(soritical vagueness) :

- if the intension of a predicate draws no sharp line wrt. some gradable property of the object or class of objects to which it is supposed to apply.
- higher-order vagueness
- ex.) *red, tall, bald, heap, ...*

- 조합의 모호성(combinatorial vagueness) :

- ... stems from an indeterminacy as to just what combination of conditions is sufficient or necessary for the application of the term.
- Wittgenstein: *Familienähnlichkeit*
- ex.) *religion, window, weapon, game, disease, ...*

모호성 3/3

- 다차원의 모호성(multidimensional vagueness):

- i) there are a number of elements that are potentially relevant for classifying a certain case
- ii) at least some of them admit of degrees
- iii) these elements are incommensurable
- ex.) *neglect, corruption, democracy, ...*

■ 법률의 모호성

- 선제적 조치: 정밀한 규칙 사이들로 빠져 나가는 경우들이 있다.

- 법원에 위임: 모호한 규범에 대한 보다 구체적인 결정을 법원이나 행정기관에 위임한다 (T. Endicott).
 - i) 더미의 모호성: 컷오프 지점은 사안별로 다르다.
 - ii) 과도한 모호성: 다차원적이어서 모든 것을 고려한 종체적인 판단이 요구된다. 따라서 법원에 사안별 결정을 위임한다.
- 국민에 위임: 사려깊은 분별력 & 책임감을 심어줌 (B. Waldron)

3. 법해석학

삼단논법

- 법적 추론(**M.P.**)

| | | |
|-----|-------|--------------|
| 대전제 | P & Q | 법규법 탐색 & 구체화 |
| 소전제 | P | 사실관계 확정 |
| 결론 | Q | 법적 결론 도출 |

- 문제점

- 법관의 개입이 없는 논리적 과정
- 사실확정과 법률해석 간의 물리적 구분

법해석학

- 철학적 해석학: **M. Heidegger, H.-G. Gadamer**

- 선이해: 텍스트를 해석한다는 것은 선입견을 통해 텍스트의 의미를 새롭게 창조해 가는 과정이다.
- 해석학적 순환: 이해과정이란 이해자와 이해대상 사이에서 ‘시선’이 순환적으로 ‘오고가는’ 과정이다

- 법해석학: **Kaufmann, Hassemer, ...**

- 법해석이란 법규법을 사안에 적용하는 ‘포섭’과정이다. 포섭과정에서 사안과 법규법 간의 일치/유추를 통해 상호간의 의미가 명료해진다.
- 선이해인 문맥정보를 통하여 주체인 법관과 객체인 법규법 간의 해석학적 순환을 통해 사안(토론)과 법규법(타입) 간의 비교를 진행한다. 이를 통해 법규법의 축자적 의미 이에 추론적 의미를 도출.

법 해석학 방법론

- 문언주의(**textualism**)

- 법규법의 축자적 의미를 중심으로 하여 법률적 문맥을 고려한 의미를 찾는다.

- 의도주의(**intentionalism**)

- 의미란 불확실한 것이므로 작성자의 의도를 고려하여 해석해야 한다.

- 목적주의(**purposivism**)

- 해당 법률의 입법 취지와 입법 목적을 고려하여 법적 의미를 결정해야 한다.

- 종합:

“요컨대, 법해석의 목표는 어디까지나 법적 안정성을 저해하지 않는 범위 내에서 구체적 타당성을 찾는데 두어야 한다. 그리고 그 과정에서 가능한 한 법률에 사용된 문언의 통상적인 의미에 충실히 해석하는 것을 원칙으로 하고, 나아가 법률의 입법 취지와 목적, 그 제·개정 연혁, 법질서 전체와의 조화, 다른 법령과의 관계 등을 고려하는 체계적·논리적 해석방법을 추가적으로 동원함으로써, 앞서 본 법해석의 요청에 부응하는 타당한 해석이 되도록 하여야 한다. ...”(대법원 2009. 4. 23. 선고 2006다81035 판결)

4. 언어학적 접근

- 해석론

- 선이해 ≈ 온톨로지
해석학적 순환 ≈ 추론/함축

- 추론

- 의도 기반의 의사소통: **Grice's non-natural meaning**

- 함축

- **Conversational Implicature**

- **Relevance Theory**

대화함축이론의 접근법 1/2

■ 법률 영역 적용

- “언어 이해에 이용되는 활용적 격률은 법률 언어의 해석에 그대로 적용되는 것 같고, 법률 해석의 캐논들은 활용적 격률과 매우 닮았다.”
(Carston 2013: 19)
- 예 1: Carston(2013)
 - *expressio unius est exclusio alterius* canon
 - *Children under ten get in free.*
 - Q-principle: Children ten or over do not get in free
- 예 2: Slocum(2016)
 - *eiusdem generis* canon
 - *No dogs, cats, birds, or other animals allowed.*
No animals allowed.
 - I-principle: stereotypical meaning

대화함축이론의 접근법 2/2

■ 비판

- 법에서의 전략적 발화 (Marmor 2014):
 - 입법부: 입법 결과를 얻는 데 필요한 거래와 타협을 촉진하기 위해 전략적 담화를 필요.
 - 법원: 법적 결과에 대한 더 큰 통제력을 얻기 위해 전략적 담화를 필요.
- cf.) 전략적 담화 (vs. 일상적 담화):
 - 완전한 협력적 상호작용 x, 엄격한 격률준수 x, 함축의미?
 - 화자가 기꺼이 표현하고자 하는 것보다 더 많은/적은 의미를 함축함으로써 어떤 이득을 얻으려고 한다.
- 규정적prescriptive/규범적normative 성격의 법률 텍스트는 기술적descriptive 성격의 일반 텍스트에 비해 명료성을 더 요구하고 함축/추론을 덜 요구한다 (Carston 2013: 19).

ex.) Q: *Do you have any bank accounts in Swiss banks, Mr. Bronston?*
A: *The company had an account there for about six months in Zürich.*

적합성이론 1/2

■ 기초 (Sperber & Wilson 1986/1995, Wilson & Sperber 2004)

- Cognitive principle of relevance: 인간의 인지는 적합성을 최대화하는 방향으로 되어 있다.
- Communicative principle of relevance: 모든 명시적 의사소통 행위는 나름대로 적합성이 있을 것이라는 추정 하에 일어난다.
- Optimal relevance:
 - i. 현시적 자극은 청자가 처리노력을 행할 가치가 있을 정도로 적합하다.
 - ii. 화자/청자의 능력과 선호에 부합할 수 있는 것이 가장 적합하다.
- Relevance-theoretic comprehension procedure
 - a. 인지효과를 계산하는 데 있어서 최소의 노력이 드는 것을 따르라: 접근가능성의 순서로 해석가설을 테스트하라(중의성 해소, 지시체 해결, 함축 등).
 - b. 적합성의 기대가 충족되면 중단하라.

적합성이론 2/2

- Sub-tasks in the overall comprehension process
 - a. Explicatures: Constructing an appropriate hypothesis about explicit content via decoding, reference resolution, and other pragmatic enrichment processes.
 - b. Implicated premises: Constructing an appropriate hypothesis about the intended contextual assumptions.
 - c. Implicated conclusions: Constructing an appropriate hypothesis about the intended contextual implications.
- Carston(2013): *Garner v Burr*
 - “Any vehicle traveling on a public highway must be fitted with pneumatic tyres.” (Road Traffic Act 1930)
 - Premise 1: 고속도로를 운행하는 차량은 공기타이어를 장착해야 한다
 - Premise 2: 쇠바퀴의 닦장
 - Conclusion :
 - (전제함축) 도로를 손상시키는 차량은 통행이 금지된다.
 - (전제함축) 공기타이어가 장착되지 않은 차량은 도로를 손상시킨다.
 - (전제함축) 쇠바퀴는 공기타이어가 아니다.
 - (결론함축) 쇠바퀴의 닦장은 통행이 금지된다.

분석 1/2

● “구부러진다”

도로교통법 제**20**조의**2**(앞지르기 금지장소)
모든 차의 운전자는 다음 각호의 **1**에 해당하는 곳에서는 다른 차를 앞지르지 못한다.
1. 교차로·터널안 또는 다리위
2. 도로의 구부러진 곳
...

“도로의 구부러진 곳이란, 문의 그대로 해석하면 조금이라도 휘어진 도로, 즉 도로가 완전히 직선이 아닌 모든 곳을 말한다. 그러나 이 사건 법률규정과 관련된 제반규정과 입법목적에 비추어 볼 때, 앞지르기로 인하여 위험을 초래하고 교통안전에 지장을 줄 수 있는 정도의 구부러진 도로로 한정 해석하여야 한다.” (현재 **2000. 2. 24** 자 **99**헌가**4**)

- 분석 (이해운 **2021**)

(전제) 함축: 앞 방향이 안 보이는 도로에서 추월시 사고가 날 수 있다.
외측: (앞 방향이 안 보일 정도로) 구부러진 도로

(결론) 함축: 추월시 사고가 날 수 있다.

분석 2/2

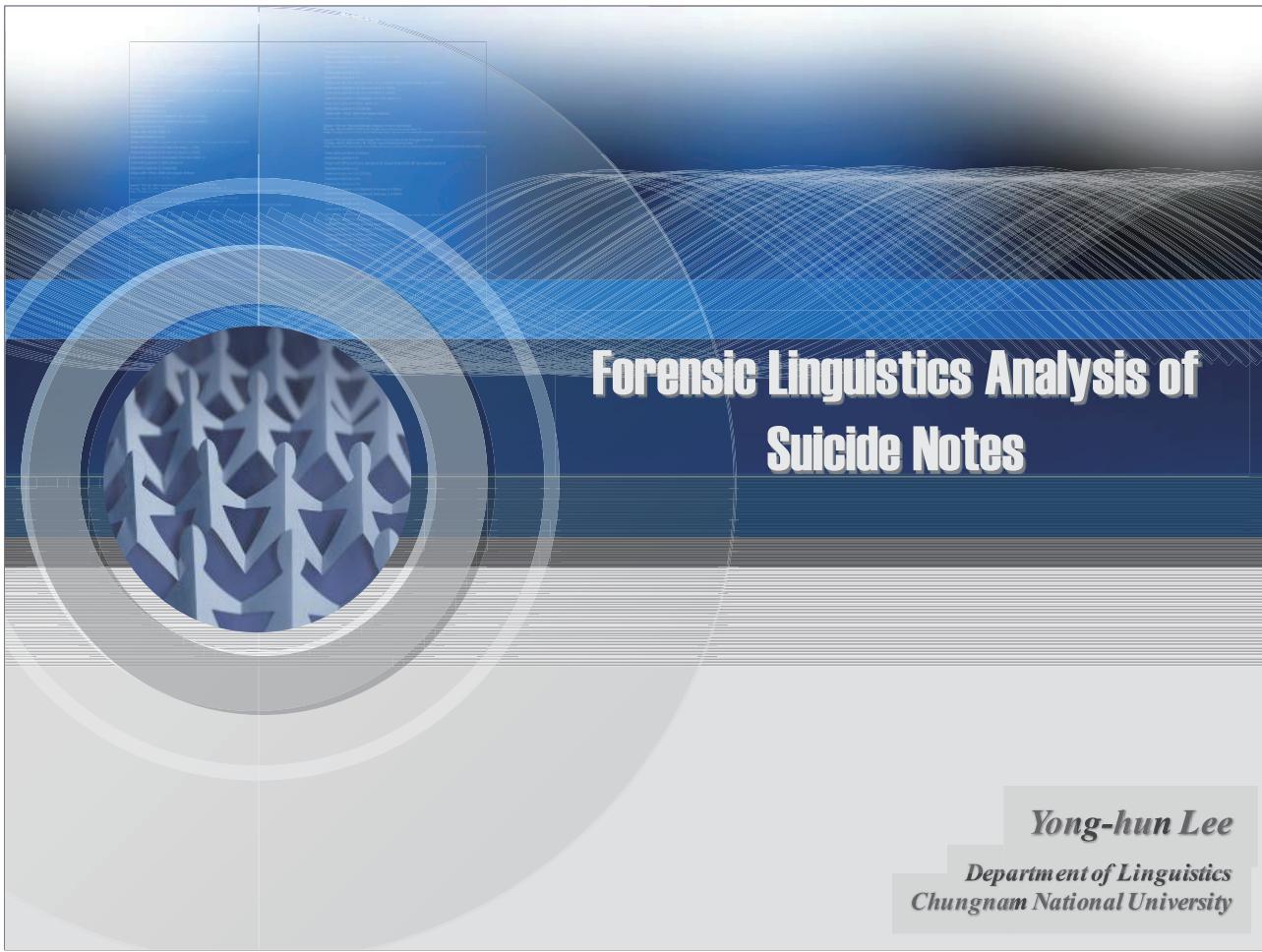
● “위력”

형법 **314**조 (업무방해)
① 제**313**조의 방법 또는 위력으로써 사람의 업무를 방해한 자는 **5**년 이하의 징역 또는 **1천500**만원 이하의 벌금에 처한다.

“**쟁의 행위**로서 **파업**, 단순히 근로계약에 따른 노무의 제공을 거부하는 부작위에 그치지 아니하고 이를 넘어서 사용자에게 압력을 가하여 근로자의 주장을 관철하고자 집단적으로 노무제공을 중단하는 실력 행사이므로, 업무방해죄에서 말하는 위력에 해당하는 요소를 포함하고 있다.” (대법원 **2011. 3. 17** 선고 **2007도482** 전원합의체 판결)

분석:

상대를 압도할 만큼 강력한 힘
허위사실유포, 위계 등과 같은 수준이다
사람의 업무를 방해하는 수단이다
업무를 방해할 수 있는 개인의 강력한 힘 (**explicature**)



Forensic Linguistics Analysis of Suicide Notes

1. Introduction

■ Motivation

- ◆ As suicides increase nowadays, it becomes more important to detect the suicide signs before the actual suicides were committed.
- ◆ In order to identify the suicide signs from the suicide notes, various kinds of cues can be used, whether they can be either linguistic or non-linguistic cues (e.g., psychological states).

■ Goal

- ◆ To conduct a linguistic analysis of suicide notes, especially based on forensic linguistics and deep learning techniques

2. Previous Studies

■ Forensic Linguistics

- ◆ As a subfield of applied linguistics, forensic linguistics applies linguistic information, analysis techniques, and linguistic insights to the areas of law, criminal investigation, court proceedings, and judicial procedure.
- ◆ Its purview extends to authorship identification as well as numerous aspects of criminal investigations and judicial processes (such as authorship verification, authorship profiling, and authorship attribution).
- ◆ The term ‘forensic’ has historically been used to describe the use of scientific techniques in criminal investigations, which usually includes the legal requirements for admissible evidence and criminal procedure.

2. Previous Studies

■ Forensic Linguistics

- ◆ The term ‘forensic’ has historically been used to describe the use of scientific techniques in criminal investigations, which usually includes the legal requirements for admissible evidence and criminal procedure.
- ◆ The gathering, preservation, and analysis of the evidence during the investigation are the responsibilities of forensic scientists.
- ◆ Some forensic scientists go to the scene of the crime to gather evidence, while others examine items in the laboratory.
 - The evidence consists of a variety of tangible and intangible items, including fingerprints, hair, DNA testing, and transfusion analysis.
 - When digital data become the focus of an inquiry, digital forensics gathers and examines the digital data that are closely used in daily life.
 - Digital data are gathered, preserved, and examined by experts in digital forensics, who then present them as evidence in court.

2. Previous Studies

■ Forensic Linguistics

- ◆ Professor Jan Svartvik introduced forensic linguistics for the first time in 1968 (Svartvik, 1968), when he examined Timothy John Evans' writings (a prominent murder suspect).
 - He examined four texts with various linguistic traits and discovered significant differences across the texts.
 - It suggested that the authors of those writings might not be the same.
- ◆ The International Association of Forensic Linguists (IAFL) was established in 1993; and in 1994, the publication of the international journals *The Law and the International Journal of Law, Language and Discourse* began.
- ◆ Olsson (2008) pointed out that an analysis of the 'suicide notes' must be included in the investigation since a suicide note typically contains lines that suggest a means to kill oneself.

2. Previous Studies

■ Linguistic Inquiry and Word Count (LIWC)

- ◆ It was introduced in Pennebaker, Francis, Booth (2001) and Tausczik and Pennebaker (2010).
- ◆ Recently, LIWC-22 was published.
- ◆ Several academic fields, including computer science and psycholinguistics, have made extensive use of the application.
- ◆ The software uses 72 linguistic parameters to assess the writings of regular people and gather statistics about certain semantic word patterns (Pennebaker & King, 1999).
- ◆ Four categories (Standard Linguistic Dimension, Psychological Process, Relativity, and Personal Concerns) can be used to categorize the variables (linguistic aspects) in the LIWC.

2. Previous Studies

■ Linguistic Inquiry and Word Count (LIWC)

- ◆ The program divides the 72 variables into four categories based on their frequencies.
- ◆ The percentages of words in each category reveal the psychological processes and states of the writer(s).
- ◆ LIWC includes more than 3,000 content words (the words used often in everyday speech), as well as different word lengths and types of function words (such as articles, prepositions, and first-, second-, and third-person pronouns).
- ◆ It is also possible to compile the analysis results in numerous individual variations in the psychological sectors, which cannot be obtained from the earlier studies, by counting function words and/or pronouns (Pennebaker et al., 2001).

3. Research Method

■ Procedure

- ◆ Compilation of Corpora
 - Two Types of Corpora
 - 8 Files for OTs
 - 7 Files for SNs
- ◆ LIWC Analysis
 - 4 Dimensions
 - 72 features
- ◆ Visualization
 - PCA
 - t-SNE

3. Research Method

■ Corpus of OTs

- ◆ Novels of Virginia Woolf
 - OT01.TheVoyageOut.txt
 - OT02.NightAndDay.txt
 - OT03.Jacob'sRoom.txt
 - OT04.MondayOrTuesDay.txt
- ◆ Novels with the 1st Person Protagonist's Point of View
 - OT05.JaneEyre.txt
 - OT06.WutheringHeights.txt
 - OT07.GreatGatsby.txt
 - OT08.HuckleberryFinn.txt

3. Research Method

■ Corpus of SNs

- ◆ Two Korean Politicians
 - SN01.Noh.txt
 - SN02.Noh.txt
- ◆ A Policeman
 - SN03.Chi.txt
- ◆ A Musician
 - SN04.Cobain.txt
- ◆ Virginia Woolf
 - SN05.Woolf-S.txt
 - SN06.Woolf-H.txt
- ◆ A Corpus of SNs
 - SN07.CEASE.txt

3. Research Method

■ Organization of Corpora

| File | Token | Type |
|---------------------------|---------|--------|
| OT01.TheVoyageOut.txt | 165,442 | 9,903 |
| OT02.NightAndDay.txt | 202,759 | 11,040 |
| OT03.Jacob'sRoom.txt | 67,035 | 7,501 |
| OT04.MondayOrTuesDay.txt | 22,659 | 3,833 |
| OT05.JaneEyre.txt | 227,236 | 13,989 |
| OT06.WutheringHeights.txt | 146,282 | 9,954 |
| OT07.GreatGatsby.txt | 61,329 | 6,243 |
| OT08.HuckleberryFinn.txt | 140,151 | 7,140 |
| SN01.Noh.txt | 124 | 74 |
| SN02.Noh.txt | 343 | 160 |
| SN03.Choi.txt | 530 | 237 |
| SN04.Cobain.txt | 671 | 287 |
| SN05.Woolf-S.txt | 190 | 105 |
| SN06.Woolf-H.txt | 229 | 105 |
| SN07.CEASE.txt | 76,593 | 7,195 |

3. Research Method

■ Principal component analysis (PCA; Pearson, 1901)

- ◆ It is a statistical technique which is employed for reducing the dimensionality of a dataset, which increases the interpretability of observations but preserves the maximum amount of information and enables the visualization of multidimensional data.
- ◆ The PCA analysis is accomplished by linearly transforming each vector into a new coordinate system where the variation in the data points can be described with fewer dimensions than that of the original dataset.
- ◆ Many studies usually utilize the first two or three principal components in order to plot the data points in 2D and 3D space.
- ◆ PCA also visually identifies clusters of closely related data points.

3. Research Method

■ t-SNE (Rowei and Hinton, 2002)

- ◆ t-distributed stochastic neighbor embedding (t-SNE) is a statistical technique for visualizing high-dimensional data points by giving each data point a location in a 2D or 3D map.
- ◆ It is based on Stochastic Neighbor Embedding (SNE), where the t-distributed variant was incorporated.
- ◆ t-SNE is a nonlinear dimensionality reduction technique for embedding high-dimensional data points for visualization in a low-dimensional vector space (2D or 3D).
- ◆ Especially, t-SNE can make a model of each high-dimensional object by 2D or 3D points in such a way that similar objects are modelled by nearby points and dissimilar objects are modelled by distant points with high probability.

4. Analysis Results

■ Standard Linguistic Dimensions 1 (OTs)

| File | OT01 | OT02 | OT03 | OT04 | OT05 | OT06 | OT07 | OT08 |
|---------|--------|--------|--------|--------|--------|--------|--------|--------|
| funct | 47.591 | 48.163 | 43.176 | 45.077 | 48.190 | 46.566 | 44.594 | 48.040 |
| pronoun | 12.759 | 13.922 | 9.073 | 11.302 | 14.100 | 15.040 | 12.912 | 13.610 |
| ppron | 8.710 | 9.755 | 5.891 | 7.233 | 10.741 | 12.090 | 9.286 | 9.513 |
| i | 1.338 | 1.337 | 0.519 | 1.942 | 5.176 | 4.401 | 3.248 | 3.453 |
| we | 0.303 | 0.290 | 0.234 | 1.090 | 0.238 | 0.443 | 0.525 | 0.885 |
| you | 0.821 | 0.868 | 0.453 | 0.728 | 1.630 | 1.617 | 0.969 | 1.213 |
| shehe | 4.765 | 6.369 | 3.841 | 2.511 | 3.139 | 5.154 | 4.037 | 2.891 |
| they | 1.480 | 0.888 | 0.835 | 0.949 | 0.559 | 0.474 | 0.507 | 1.071 |
| ipron | 4.049 | 4.167 | 3.182 | 4.069 | 3.359 | 2.950 | 3.626 | 4.096 |
| article | 6.394 | 5.903 | 7.759 | 7.613 | 5.664 | 4.938 | 6.524 | 5.763 |

4. Analysis Results

■ Standard Linguistic Dimensions 2 (OTs)

| File | OT01 | OT02 | OT03 | OT04 | OT05 | OT06 | OT07 | OT08 |
|-------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| verb | 10.272 | 9.844 | 8.606 | 7.706 | 10.545 | 9.475 | 9.087 | 9.304 |
| auxverb | 5.374 | 5.165 | 4.436 | 4.215 | 6.374 | 5.472 | 3.952 | 4.138 |
| past | 6.078 | 5.753 | 4.877 | 3.270 | 4.858 | 4.039 | 5.831 | 4.632 |
| present | 2.806 | 2.625 | 2.651 | 3.270 | 3.770 | 3.443 | 2.552 | 3.575 |
| future | 0.599 | 0.674 | 0.534 | 0.587 | 1.074 | 1.058 | 0.277 | 0.444 |
| adverb | 3.394 | 3.028 | 3.171 | 3.261 | 3.186 | 2.692 | 3.101 | 4.165 |
| preps | 11.869 | 12.511 | 11.566 | 11.077 | 10.823 | 10.616 | 11.623 | 10.614 |
| conj | 5.669 | 5.349 | 5.479 | 5.477 | 5.937 | 6.205 | 4.932 | 8.120 |
| negate | 0.985 | 1.058 | 0.825 | 0.993 | 1.458 | 1.316 | 0.538 | 0.956 |
| quant | 1.996 | 2.129 | 1.560 | 1.761 | 1.740 | 1.524 | 1.637 | 2.005 |
| number | 1.076 | 0.915 | 1.064 | 1.143 | 0.828 | 0.700 | 1.042 | 1.008 |
| swear | 0.037 | 0.035 | 0.058 | 0.093 | 0.021 | 0.052 | 0.038 | 0.163 |

4. Analysis Results

■ Standard Linguistic Dimensions 1 (SNs)

| File | SN01 | SN02 | SN03 | SN04 | SN05 | SN06 | SN07 |
|----------|--------------|--------------|--------------|--------------|---------------|---------------|--------------|
| funct | 55.645 | 47.522 | 52.830 | 53.651 | 60.526 | 57.205 | 54.104 |
| pronoun | 12.097 | 11.370 | 13.585 | 15.499 | 24.211 | 24.454 | 17.391 |
| ppron | 8.065 | 6.997 | 10.377 | 10.283 | 16.842 | 16.594 | 12.209 |
| i | 7.258 | 6.122 | 7.358 | 7.303 | 10.000 | 10.917 | 8.181 |
| we | 0.000 | 0.292 | 0.755 | 0.894 | 0.526 | 0.873 | 0.534 |
| you | 0.000 | 0.292 | 1.887 | 1.490 | 3.684 | 4.803 | 1.927 |
| shehe | 0.000 | 0.000 | 0.189 | 0.447 | 2.632 | 0.000 | 0.910 |
| they | 0.806 | 0.292 | 0.189 | 0.149 | 0.000 | 0.000 | 0.651 |
| ipron | 4.032 | 4.373 | 3.208 | 5.216 | 7.368 | 7.860 | 5.182 |
| article | 6.452 | 6.414 | 8.491 | 5.812 | 1.053 | 1.747 | 4.619 |

4. Analysis Results

■ Standard Linguistic Dimensions 2 (SNs)

| File | SN01 | SN02 | SN03 | SN04 | SN05 | SN06 | SN07 |
|----------------|---------------|--------------|--------------|--------------|---------------|---------------|--------------|
| verb | 15.323 | 11.953 | 11.887 | 11.475 | 22.632 | 23.581 | 14.087 |
| auxverb | 12.903 | 8.163 | 7.358 | 6.408 | 14.737 | 15.284 | 8.682 |
| past | 0.806 | 5.248 | 3.396 | 1.490 | 4.737 | 3.930 | 3.415 |
| present | 10.484 | 5.248 | 7.170 | 7.154 | 14.737 | 15.721 | 8.214 |
| future | 0.806 | 1.458 | 0.566 | 1.192 | 2.105 | 0.873 | 1.188 |
| adverb | 4.032 | 2.041 | 3.962 | 4.620 | 6.316 | 1.310 | 3.617 |
| preps | 12.097 | 9.913 | 12.453 | 12.668 | 6.842 | 6.987 | 11.176 |
| conj | 3.226 | 5.248 | 7.736 | 6.557 | 6.316 | 3.493 | 5.511 |
| negate | 2.419 | 2.915 | 0.943 | 0.447 | 0.526 | 0.437 | 2.021 |
| quant | 3.226 | 3.207 | 1.509 | 3.279 | 3.684 | 3.930 | 2.770 |
| number | 0.000 | 1.458 | 0.189 | 0.745 | 0.526 | 0.873 | 0.595 |
| swear | 0.000 | 0.000 | 0.000 | 0.149 | 0.000 | 0.000 | 0.261 |

4. Analysis Results

■ Psychological Process 1 (OTs)

| File | OT01 | OT02 | OT03 | OT04 | OT05 | OT06 | OT07 | OT08 |
|---------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| social | 11.666 | 12.534 | 9.480 | 8.699 | 9.067 | 11.084 | 9.731 | 9.349 |
| family | 0.327 | 0.412 | 0.260 | 0.265 | 0.292 | 0.366 | 0.249 | 0.279 |
| friend | 0.072 | 0.078 | 0.076 | 0.049 | 0.091 | 0.135 | 0.116 | 0.065 |
| humans | 1.347 | 0.880 | 1.516 | 1.064 | 0.969 | 0.690 | 0.864 | 0.571 |
| affect | 4.266 | 4.201 | 3.403 | 3.654 | 4.501 | 4.709 | 3.724 | 3.102 |
| posemo | 2.664 | 2.502 | 1.911 | 2.290 | 2.733 | 2.361 | 2.183 | 2.003 |
| negemo | 1.543 | 1.611 | 1.437 | 1.337 | 1.735 | 2.315 | 1.512 | 1.092 |
| anx | 0.329 | 0.404 | 0.300 | 0.243 | 0.374 | 0.474 | 0.442 | 0.166 |
| anger | 0.349 | 0.380 | 0.355 | 0.366 | 0.340 | 0.662 | 0.329 | 0.252 |
| sad | 0.467 | 0.447 | 0.528 | 0.424 | 0.654 | 0.728 | 0.494 | 0.323 |
| cogmech | 13.278 | 14.326 | 11.120 | 12.309 | 12.925 | 12.792 | 11.826 | 14.506 |
| insight | 2.081 | 2.543 | 1.141 | 1.651 | 1.780 | 1.538 | 1.536 | 1.033 |
| cause | 0.867 | 1.001 | 0.585 | 0.812 | 0.787 | 0.872 | 0.655 | 0.894 |
| discrep | 1.209 | 1.498 | 0.934 | 1.042 | 1.435 | 1.490 | 0.804 | 1.120 |
| tentat | 2.100 | 2.424 | 1.754 | 1.911 | 1.887 | 1.676 | 1.759 | 1.787 |
| certain | 1.294 | 1.382 | 1.141 | 1.359 | 1.155 | 0.952 | 1.070 | 1.334 |
| inhib | 0.428 | 0.520 | 0.440 | 0.578 | 0.505 | 0.552 | 0.486 | 0.355 |
| incl | 4.382 | 4.153 | 4.465 | 4.321 | 4.314 | 4.798 | 4.918 | 7.020 |
| excl | 1.890 | 2.067 | 1.557 | 1.827 | 2.087 | 1.937 | 1.386 | 1.697 |

4. Analysis Results

■ Psychological Process 2 (OTs)

| File | OT01 | OT02 | OT03 | OT04 | OT05 | OT06 | OT07 | OT08 |
|---------|--------|--------|--------|--------|--------|--------|--------|--------|
| percept | 3.517 | 2.991 | 4.047 | 3.676 | 2.893 | 2.398 | 3.246 | 2.913 |
| see | 1.470 | 1.188 | 1.827 | 1.801 | 1.116 | 0.857 | 1.401 | 0.973 |
| hear | 1.230 | 1.063 | 1.244 | 0.918 | 0.931 | 0.844 | 1.014 | 1.334 |
| feel | 0.750 | 0.662 | 0.867 | 0.808 | 0.727 | 0.612 | 0.677 | 0.512 |
| bio | 1.897 | 1.394 | 2.197 | 2.185 | 1.797 | 1.727 | 2.009 | 1.429 |
| body | 1.057 | 0.749 | 1.332 | 1.156 | 0.998 | 0.941 | 1.123 | 0.867 |
| health | 0.403 | 0.266 | 0.309 | 0.424 | 0.375 | 0.367 | 0.307 | 0.233 |
| sexual | 0.169 | 0.186 | 0.146 | 0.221 | 0.163 | 0.163 | 0.126 | 0.060 |
| ingest | 0.317 | 0.207 | 0.467 | 0.419 | 0.285 | 0.275 | 0.465 | 0.295 |
| relativ | 11.973 | 11.074 | 13.054 | 12.110 | 11.165 | 10.637 | 13.972 | 12.555 |
| motion | 1.792 | 1.570 | 2.039 | 1.915 | 1.680 | 1.861 | 2.226 | 1.861 |
| space | 6.054 | 5.702 | 6.895 | 6.170 | 5.204 | 4.856 | 6.773 | 6.250 |
| time | 4.087 | 3.607 | 4.102 | 3.959 | 4.237 | 3.798 | 4.947 | 4.240 |

4. Analysis Results

■ Psychological Process 1 (SNs)

| File | SN01 | SN02 | SN03 | SN04 | SN05 | SN06 | SN07 |
|---------|--------|--------|--------|--------|--------|--------|--------|
| social | 3.226 | 10.787 | 8.868 | 9.091 | 11.579 | 10.480 | 9.502 |
| family | 0.000 | 0.000 | 0.755 | 0.298 | 0.000 | 0.000 | 0.697 |
| friend | 0.000 | 0.000 | 0.377 | 0.000 | 0.000 | 0.000 | 0.225 |
| humans | 0.806 | 2.041 | 0.755 | 1.788 | 0.526 | 0.873 | 0.977 |
| affect | 4.839 | 12.245 | 7.170 | 8.346 | 5.789 | 6.114 | 6.986 |
| posemo | 0.806 | 7.872 | 4.717 | 5.812 | 4.211 | 4.367 | 4.036 |
| negemo | 4.032 | 4.373 | 2.453 | 2.534 | 1.579 | 1.747 | 2.906 |
| anx | 0.000 | 0.875 | 0.377 | 0.447 | 0.526 | 0.000 | 0.380 |
| anger | 0.806 | 1.458 | 0.377 | 0.745 | 1.053 | 0.873 | 0.994 |
| sad | 0.806 | 0.000 | 0.755 | 0.447 | 0.000 | 0.000 | 0.704 |
| cogmech | 11.290 | 12.536 | 14.151 | 15.946 | 17.895 | 18.341 | 16.163 |
| insight | 1.613 | 1.166 | 3.019 | 3.428 | 4.211 | 3.930 | 2.398 |
| cause | 1.613 | 1.166 | 1.698 | 1.788 | 0.526 | 0.000 | 1.417 |
| discrep | 0.806 | 1.749 | 1.887 | 1.341 | 2.632 | 3.930 | 1.833 |
| tentat | 1.613 | 0.292 | 0.377 | 1.788 | 2.105 | 3.057 | 2.311 |
| certain | 0.806 | 1.458 | 1.887 | 2.086 | 4.211 | 3.493 | 1.918 |
| inhib | 0.000 | 1.749 | 0.189 | 0.149 | 0.000 | 0.000 | 0.537 |
| incl | 1.613 | 3.790 | 4.151 | 4.918 | 2.632 | 3.493 | 3.799 |
| excl | 2.419 | 1.166 | 1.698 | 1.639 | 2.632 | 1.310 | 2.780 |

4. Analysis Results

■ Psychological Process 2 (SNs)

| File | SN01 | SN02 | SN03 | SN04 | SN05 | SN06 | SN07 |
|---------|--------|-------|-------|-------|--------|-------|--------|
| percept | 0.806 | 0.292 | 0.943 | 1.788 | 2.632 | 3.057 | 1.334 |
| see | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.437 | 0.407 |
| hear | 0.000 | 0.000 | 0.000 | 0.745 | 1.579 | 1.747 | 0.349 |
| feel | 0.806 | 0.292 | 0.566 | 0.745 | 1.053 | 0.873 | 0.492 |
| bio | 2.419 | 0.583 | 1.132 | 2.534 | 1.053 | 1.747 | 2.417 |
| body | 0.000 | 0.583 | 0.189 | 0.447 | 0.000 | 0.000 | 0.547 |
| health | 2.419 | 0.000 | 0.566 | 0.447 | 0.000 | 1.747 | 1.025 |
| sexual | 0.000 | 0.000 | 0.377 | 1.639 | 1.053 | 0.000 | 0.764 |
| ingest | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.129 |
| relativ | 10.484 | 8.746 | 7.170 | 7.899 | 13.684 | 6.987 | 10.090 |
| motion | 0.806 | 3.790 | 0.755 | 0.596 | 2.105 | 2.183 | 1.478 |
| space | 7.258 | 3.207 | 3.019 | 3.428 | 1.579 | 0.873 | 4.179 |
| time | 2.419 | 2.915 | 3.396 | 3.130 | 10.000 | 3.493 | 4.331 |

4. Analysis Results

■ Personal Concerns (OTs)

| File | OT01 | OT02 | OT03 | OT04 | OT05 | OT06 | OT07 | OT08 |
|---------|-------|-------|-------|-------|-------|-------|-------|-------|
| work | 0.636 | 0.758 | 0.741 | 0.931 | 0.706 | 0.656 | 0.509 | 0.425 |
| achieve | 0.740 | 0.922 | 0.574 | 0.693 | 0.878 | 1.075 | 0.667 | 0.714 |
| leisure | 0.728 | 0.567 | 0.964 | 0.781 | 0.436 | 0.420 | 0.874 | 0.455 |
| home | 0.560 | 0.658 | 0.792 | 0.719 | 0.678 | 0.708 | 0.998 | 0.445 |
| money | 0.262 | 0.246 | 0.303 | 0.349 | 0.282 | 0.228 | 0.310 | 0.536 |
| relig | 0.227 | 0.199 | 0.310 | 0.459 | 0.390 | 0.330 | 0.140 | 0.149 |
| death | 0.114 | 0.083 | 0.148 | 0.238 | 0.177 | 0.259 | 0.173 | 0.250 |

4. Analysis Results

■ Personal Concerns (SNs)

| File | SN01 | SN02 | SN03 | SN04 | SN05 | SN06 | SN07 |
|---------|-------|-------|-------|-------|-------|-------|-------|
| work | 0.806 | 2.332 | 3.019 | 0.298 | 0.000 | 0.873 | 1.219 |
| achieve | 0.806 | 2.332 | 1.509 | 1.043 | 1.579 | 1.310 | 1.328 |
| leisure | 0.806 | 3.499 | 0.566 | 1.192 | 0.000 | 0.437 | 0.601 |
| home | 0.806 | 0.000 | 0.943 | 0.000 | 0.000 | 0.000 | 0.356 |
| money | 1.613 | 1.458 | 0.566 | 0.000 | 0.000 | 0.437 | 0.478 |
| relig | 0.000 | 0.000 | 0.566 | 0.596 | 0.000 | 0.000 | 0.547 |
| death | 2.419 | 0.000 | 0.000 | 0.149 | 0.000 | 0.000 | 0.631 |

4. Analysis Results

■ Spoken Categories (OTs)

| File | OT01 | OT02 | OT03 | OT04 | OT05 | OT06 | OT07 | OT08 |
|--------|-------|-------|-------|-------|-------|-------|-------|-------|
| assent | 0.146 | 0.147 | 0.160 | 0.291 | 0.173 | 0.157 | 0.160 | 0.133 |
| nonfl | 0.100 | 0.103 | 0.100 | 0.106 | 0.157 | 0.103 | 0.116 | 0.343 |
| filler | 0.245 | 0.102 | 0.248 | 0.216 | 0.175 | 0.107 | 0.199 | 0.191 |

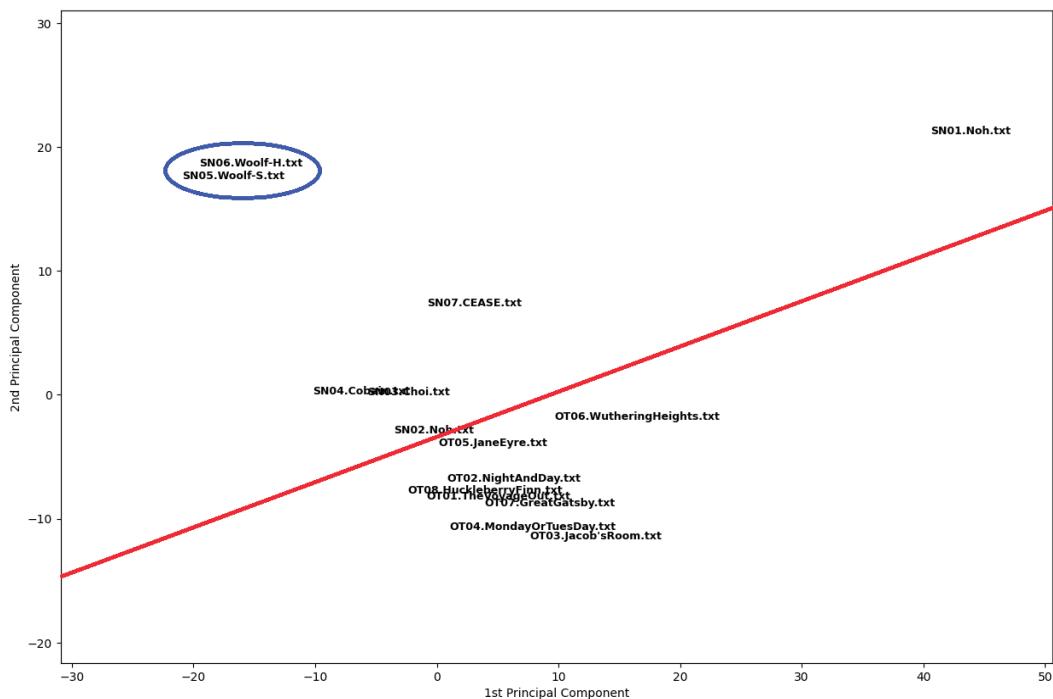
4. Analysis Results

Spoken Categories (SNs)

| File | SN01 | SN02 | SN03 | SN04 | SN05 | SN06 | SN07 |
|--------|-------|-------|-------|-------|-------|-------|-------|
| assent | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.072 |
| nonfl | 0.000 | 0.000 | 0.000 | 0.149 | 0.000 | 0.000 | 0.077 |
| filler | 0.000 | 0.000 | 0.189 | 0.000 | 0.000 | 0.000 | 0.259 |

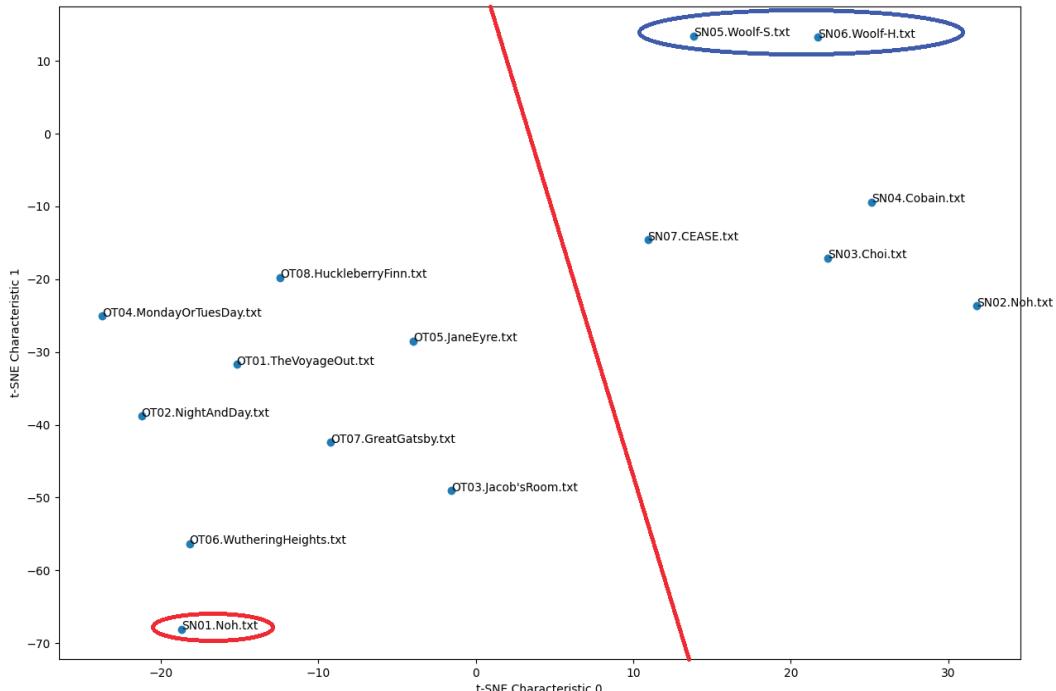
4. Analysis Results

Visualization with PCA



4. Analysis Results

■ Visualization with t-SNE



5. Conclusion

■ In the LIWC Analysis

- ◆ In the Standard Linguistic Dimension
 - SNs used the 1st person pronoun 'I' more frequently than the OTs.
 - Even though the novels with 1st protagonist's point of view also preferred to use the 1st person pronoun 'I' frequently, its frequencies were much lower than those of SNs.
 - In addition, SNs preferred to use the past tense, while OTs preferred to use the present tense.
- ◆ In the Psychological Process
 - SNs contained more negative words than OTs.
- ◆ In the Relativity
 - It was hard to say any characteristics.
- ◆ In the Personal Concerns
 - It was hard to say any characteristics.

5. Conclusion

■ In the PCA and the t-SNE

- ◆ SNs could be distinguished from OTs only with the (purely) linguistic analysis.
- ◆ SNs have different linguistic properties from the OTs, although the same author wrote both types of texts.
- ◆ Though the novels with 1st protagonist's perspective also preferred to use the 1st person pronoun 'I' frequently, its frequencies were much lower than those of SNs.

References

- Olsson, J. 2008. *Forensic Linguistics*, 2nd Edition. London: Continuum.
- Pearson, K. 1901. On Lines and Planes of Closest Fit to Systems of Points in Space. *Philosophical Magazine* 2.11: 559-572.
- Pennebaker, W., E. Francis, and J. Booth. 2001. *Linguistic Inquiry and Word Count (LIWC)*: LIWC2001. Mahwah, NJ: Lawrence Erlbaum Associates.
- Roweis, S. and G. Hinton. 2002. Stochastic Neighbor Embedding. *Proceedings of the 15th International Conference on Neural Information Processing Systems* 857-864.
- Svartvik, J. 1968. *The Evans Statements: A Case for Forensic Linguistics*. Gothenburg, Sweden: University of Gothenburg Press.
- Tausczik, Y. and J. Pennebaker. 2010. The Psychological Meaning of Words: LIWC and Computerized Text Analysis Methods. *Journal of Language and Social Psychology* 29.1: 24-54.

2023년 가을 한국코퍼스언어학회 10주년 기념 학술대회

한국의 말뭉치 구축 및 연구의 지평

Session 2

<수어 코퍼스>



조선대수어코퍼스구축과 웹 자료 추출 과정



코퍼스 기반 한국어와 한국 수어 어휘 비교 분석



코퍼스 기반 한국어 수어의 등위구문 연구



조선대수어코퍼스구축과 웹 자료 추출 과정

2023. 10. 7

최영주(조선대), 박호민(한국해양대), 장세은(한국해양대),
권순복(부산대), 윤병천(조선대), 조참훈(중국 서남대)



한국코퍼스언어학회 가을학술대회, 서울, 연세대학교

(참고: 박호민 (2022)과 장세은, 박호민 (2023)의 발표 일부를 발췌한 논문으로 수정 보완하였음)

1



발표 목차

1. 조선대한국수어말뭉치 (조선대KSL코퍼스) 구축 배경
2. 수어코퍼스 선행연구
3. 조선대KSL코퍼스 구축
4. 한국수어 문장 구조 파악을 위한 검색과정 시연
5. 제언: 조선대KSL코퍼스 브라우저 고도화

참고문헌

2



1. 조선대한국수어말뭉치(조선대KSL코퍼스) 구축 배경

The Original Plan of the Chosun KSL Corpus

A small corpus of Chosun KSL was compiled from December 2020 to April 2022 and was transcribed from January 2021 to May 2022 by the Language Convergence Institute, Chosen University, Gwangju, South Korea, which has been supported by **The Five-year Institute Project of Research Foundation of Korea** in September 2020.

Original Plan: It was firstly planned that the Chosun KSL Corpus would consist of **300,000 gloss words**.

3



2. 수어코퍼스 선행연구

2.1 영국수어코퍼스 (<https://bslcorpusproject.org/>)

- **British Sign Language Corpus** Project (2008 ~ 2011)

- The BSL Corpus Project research team filmed 249 Deaf people from 8 cities across the United Kingdom: London, Bristol, Birmingham, Manchester, Newcastle, Glasgow, Cardiff and Belfast.



Cf. Schembri et al. (2011, 2013)

4



2.2 American Sign Language Corpora

BU ASL Corpus as ELAN (<http://asl.cs.depaul.edu/corpus/index.html>)

An ASL corpus from Boston University's [The National Center for Sign Language and Gesture Resources](#) (NCSLGR) corpus [Neidle and Vogler 2012] is now available in ELAN format. It is now possible to access the only publically available ASL corpus via the tools and interface of ELAN.

ASL Corpus Database (<https://www.csdeagles.com/apps/pages/aslcorpus>)

ASL Corpus List (as of January 23, 2018)

Lessons, Vocabulary, A to Z Story, Visual Vernacular (VV), Storysigning(Books), Storytelling - Halloween, Rhymes and Rhythms, Poetry, Handshape ASL Poetry, Number Story, Classifier, Cinematic ASL, Personification, Percussion, Literature, American History, World History, News in the Word, Animals, Earth and Space Science, Health, Weather, etc.

5

2.3 일본수어대화코퍼스

JSL Colloquial Corpus

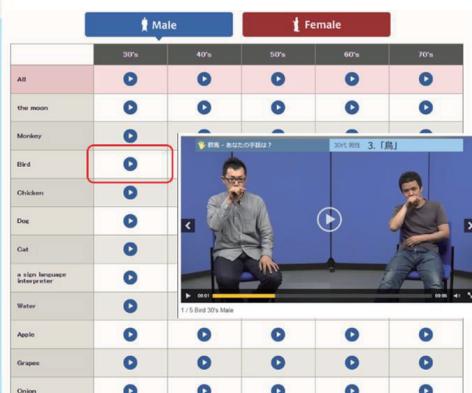
<http://research.nii.ac.jp/jsl-corpus/public/en/>



2011.04 ~ 2014 동안 연구를 위해 제작된 일본 수어 말뭉치

British Sign Language Corpus의 수석 제작자인 Adam Schembri 교수의 Corpus 구축 방법론을 조언받아 제작 (일본 국립 정보학 연구소)

일본 내 지역별, 연령별, 성별 수어 대화 영상을 촬영하여 ELAN 프로그램으로 Annotation을 진행하였음



Bono et al. (2014)

6



2.4 Other sign language corpora

Sign language resources (<https://www.clarin.eu/resource-families/sign-language-resources>)

- Adamorobe Sign Language Corpus
- [Sign-Hub WP2.4: Life Stories](#) (six different sign languages: Catalan Sign Language (LSC), German Sign Language (DGS), Italian Sign Language (LIS), Sign Language of the Netherlands (NGT), Spanish Sign Language (LSE), and Turkish Sign Language (TİD))
- [Dogon Sign Language Corpus](#)
- [Corpus NGT \(Dutch Sign Language\)](#)
- [The Kipo Corpus \(Finnish Sign Language, FinSL\)](#)
- [Creagest-Acquisition corpus \(French Sign Language\)](#)
- [Italian Sign Language Corpus](#)
- [Norwegian Sign Language Corpus](#)
- [Turkish sign language database](#)

2.5 Hong Kong Sign Language Corpus (<http://www.cslds.org/hkslcorpus/>)

Cf. 2006 Hong Kong Sign Language Child Language Corpus

(Cat Fung H-M, Scholastica Lam, Joe Mak, Gladys Tang (2008))

7

3. 조선대KSL코퍼스(The Chosun KSL Corpus)구축



3.1 자료제공자 (Informants)

- 연령층 남녀별: 평균 연령 49.61세 : 26 명 (남 16명, 여 10명)
- 지역별, 주제별 통계

| 나이 | 참여인원 |
|-------------|-----------|
| 20대 | 1 |
| 30대 | 2 |
| 40대 | 12 |
| 50대 | 6 |
| 60대 | 3 |
| 70대 | 1 |
| 80대 | 1 |
| 총 인원 | 26 |

| 지역 | 세부 | 인원 | 지역별참여 인원 |
|-------|----|----|----------|
| 전남 | 강진 | 1 | 23 |
| | 광양 | 2 | |
| | 광주 | 7 | |
| | 나주 | 1 | |
| | 남악 | 3 | |
| | 목포 | 3 | |
| | 무안 | 2 | |
| | 보성 | 1 | |
| | 순천 | 1 | |
| 전북 | 화순 | 2 | 2 |
| | 완주 | 1 | |
| | 전주 | 1 | |
| 서울/경기 | 인천 | 1 | 1 |

| 주제 | 참여인원 |
|-------------|------------|
| 인간 | 18 |
| 삶 | 24 |
| 식생활 | 18 |
| 의생활 | 14 |
| 주생활 | 14 |
| 사회생활 | 18 |
| 경제생활 | 14 |
| 교육 | 12 |
| 종교 | 12 |
| 문화 | 18 |
| 정치와 행정 | 14 |
| 자연 | 18 |
| 동식물 | 18 |
| 개념 | 12 |
| 국가명 및 지명 | 16 |
| 기타 | 18 |
| 누적인원 | 258 |

8

3.2 영상수집 방법



- C (video recording of right person), A (overall), B (video recording of left person)

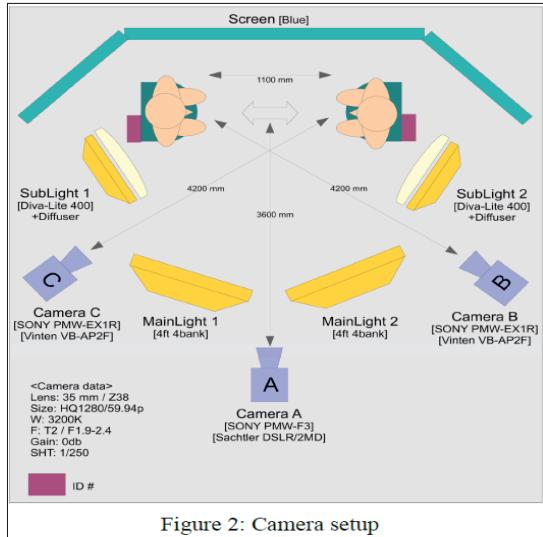


Figure 3: Screen shot: the three camera angles used for data collection

9



3.3 촬영 주제 분류표 (16 대범주, 88 소범주, 199 주제코드)

- Four types of video data collection (BSL + JSL)
 1. Conversations with another deaf person
 2. Answering interview questions
 3. Telling stories
 4. Lexical elicitation

10



<<촬영주제 분류 및 코드표>>

| | 대범주 | 소범주 | 촬영주제코드 |
|---|-----|------------|--|
| 1 | 인간 | 사람의 종류 | |
| | | 신체부위 | |
| | | 신체 내부 구성 | |
| | | 생리 현상 | |
| | | 감각 | |
| | | 감정 | |
| | | 성격 | |
| | | 태도 | |
| | | 동모 | |
| | | 능력 | |
| | | 신체 변화 | |
| | | 신체 행위 | |
| | | 신체에 가하는 행위 | |
| | | 인지 행위 | |
| | | 소리 | |
| | | | 1. 존경하는 사람 2. 시각을 잃으면 어떻게 할까? 3. 어떤 성격이신가요? 좋아하는 친구들의 성격은 어떠한가요? |

11



| | | | |
|---|---|-------|--|
| 2 | 삶 | 삶의 상태 | 1. 원하는 대상에게 하고 싶은 말 2. 내가 자녀에게 베푸는 것은? 3. 제일 후회되는 일은? 4. 서월호 4월 16일 기념식 5. 내 아이는 사고뭉치, 아이를 키우면서 생긴 일들 6. 아들, 딸 출산할 때 느낌 및 기억은? 7. (방학이면 찾아오는) 조카에 대해서 8. 기족들이 좋아하는 것, 기족들이 싫어하는 것 9. 계속되는 아동학대 (수어방송) 10. 코로나 의료진 (일상시간) 11. 일과에 관한 짤막한 대화 나누기 12. 취미에 관한 짤막한 대화 나누기 13. 힘들었던 일 14. 좋아하는 것들에 대한 짤막한 대화 나누기 15. 최근에 다니고 있는 나들이/여행 16. 하루를 어떻게 보내는지 17. 가장 기뻤을 때의 일 18. 가장 슬펐을 때의 일 19. 가장 활동했을 때의 일 20. 스트레스는 왜 생기며 어떻게 풀까? 21. 여행에서 생긴 일(어피스드) 22. 문명·농인의 삶 23. 결혼식을 올릴 때 떨리고 긴장했던 기억은? 24. 살면서 도전해보고 싶은 일은? 25. 만약 내가 다시 태어난다면? 26. 청각장애인은 일반인과 결혼하면 과연 행복할까? 27. 결혼 후 시댁살이의 다양한 경험담 28. 시부모님과 함께 지낸 경험담 29. 여행 가서 생긴 일 30. 하루생활 (일상) 31. 시골 어울 농사의 고단함에 대해 32. 의약품의 용법 사용법 수화 동영상 안내제작, 배포하기 33. 사유리와 같은 자발적 미혼모에 대해 34. 무료백신 (수어방송) 35. 미라클모닝 (수어방송) 36. 여행하고 싶은 나라 37. 타자로 출퇴근했던 경험담 (해남에서 육포로) 38. 웨멘캐러처럼 만약 본인이 못 보고, 못 듣는 장애를 가졌다면? 39. 벚꽃 거리 (일상) 40. 각 계절에 할 수 있는 일 41. 풍경 (일상) |
| | | 삶의 행위 | |
| | | 일상 행위 | |
| | | 친족 관계 | |
| | | 가족 행사 | |
| | | 여가 도구 | |
| | | 여가 시설 | |
| | | 여가 활동 | |
| | | 병과 증상 | |
| | | 치료 시설 | |
| | | 치료 행위 | |
| | | 약품류 | |

12



| | | | |
|---|-----|-------------|--|
| 3 | 식생활 | 음식 | 1. 좋아하는 오리법 2. 지금 먹고 싶은 음식 3. 최근에 먹은 음식 중에 제일 맛있던 것 4. 한 번은 꼭 먹고 싶은 음식 5. 남에게 추천해 주고 싶은 음식 6. 농인은 배달앱인 배달+기요, 배달+죽 등 유용할까? 7. 시어머니에게 배운 김장, 식혜, 고추장, 메주 만들기 8. 배달 음식 속 살아있는 쥐 (수어방송) 9. 한식 양식 (비교) 10. 맛집에 줄 서있는 사람들 (일상) 11. 시장떡볶이와 배달떡볶이 (비교) 12. 식품에 대한 안전정보를 담은 수학동영상 안내제작, 배포하기] |
| | | 채소 | |
| | | 곡류 | |
| | | 과일 | |
| | | 음료 | |
| | | 식재료 | |
| | | 조리도구 | |
| | | 식생활 관련 장소 | |
| | | 맛 | |
| | | 식사 및 조리 행위 | |
| 4 | 의생활 | 옷 종류 | |
| | | 옷감 | |
| | | 옷의 부분 | |
| | | 모자, 신발, 장신구 | 1. 옷 산발 구매 방법 소개 2. 옷 산발 취향 소개 3. 모자나 액세서리 취향 소개 4. 패션 인터넷 사이트 소개 5. 패션 유튜브 동영상 저작자 소개 6. 서양과 한국의 전통의상 (비교) |
| | | 의생활 관련 장소 | |
| | | 의복 착용 상태 | |
| | | 의복 착용 행위 | |
| | | 미용 행위 | |

13



| | | | |
|---|-----|-------|---|
| 5 | 주생활 | 생활 용품 | |
| | | 건물 종류 | |
| | | 주거 형태 | 1. 소속 집단과 주로 만나는 일시나 장소 2. 다른 지역에 거주/생활하는 경우: 이유와 생활, 부엌 생김새, 화장실 물건 |
| | | 주고 지역 | 3. 도시와 시골 장단점 4. 서울의 과거와 현재 (비교) |
| | | 주택 구성 | 5. 그림 속 카페와 목욕탕 카페 (비교) |
| | | 주거 상태 | |
| | | 주거 행위 | |
| | | 가사 행위 | |

14



| | | | |
|---|------|----------|---|
| 6 | 사회생활 | 인간관계 | 1. 조용한 택시 (수어방송) 2. 자주 갖는 모임 3. 동호회 계 친구 모임 (주로 만나는 집단의 성 질) 4. 특별히 기억에 남은 모임 5. 한농협·한농묘·한국농아인협회를 걱정하는 모임에 대해서 |
| | | 사회생활 상태 | 6. 내 직업과 변천사 7. 농통역사 모임으로 떠난 여행 피서 8. 회사가 일어났을 때 대피하기(화재 감지기, 사각경보기, 피난구 유도등, 통로 유도등, 화 재 방지기) |
| | | 사회 활동 | 9. 열린마이터가 고장 났을 때 도움 요청 방법 10. 긴급구조대 별선 시 119 신고 방법 |
| | | 사회 행사 | 11. 경찰관에서 떠난 안내, 안전에 대한 수화, 자막 편의 제공에 관한 것 12. 고속버스터미널, 지하철 긴급피해 비교하기 13. 공항터미널 승무원 통역사 배치 비교하기 14. 응급 상황 시 요청 방법 (월·금 춘제x, 주말 수어통역사 후보) |
| | | 직장 | 15. 교회 하계 수련회 시 성도 실종 사건에 대해 청각장애인 할 수 있는 직무를 제시하기 |
| | | 직위 | 16. 청각장애인 할 수 있는 직무를 제시하기 17. 방송에 강하고 창문없는 밀폐장면인 호텔에 서 불이 났을 때 경고음을 못 듣는 청각장애 인은 어떻게 해야 할까? |
| | | 직업 | 18. 수어 통역 잘못 전달될 경우의 사례 19. 수어통역사 없는 농사회에서 과연 잘 할 수 있을까? |
| | | 직장 생활 | 20. 아동학대 및 자살, 노숙자 등 다양한 사 회문제 |
| | | 교통수단 | 21. 지하철 풍경 (일상) 22. 헌터들이 없다면? |
| | | 교통 이용 장소 | 23. 직원에게 물질 물안 수어통역사 영상전화기 배치 비교하기 |
| | | 교통 이용 행위 | 24. 통신 수단 없을 시대에 농인 생활 25. 디바이크의 기술 (수어방송) |
| | | 소통 수단 | 26. 스마트폰과 풀다운 (비교) 27. 컴퓨터 게임과 온라인 (비교) |
| | | 말 | 28. 소형화 대형화 (비교) 29. 각 나라의 청각 장애인 복지는 어떻게 되나? |
| | | 언어 행위 | 30. 청각장애인의 보물, 청각장애인의 영토에 대 해서 이야기하기 |
| | | 통신 행위 | 31. 일본 여행으로 느껴본 한국과 일본의 지 하철 32. 수화 이모티콘 힙트 (수어방송) |
| | | 매체 | 33. 짧은동영에 대한 기억 34. 첫사랑에 대한 경험 35. 사랑하는 사람과 헤어졌던 경험 |

15



| | | | |
|---|------|----------|-----------------------------------|
| 7 | 경제생활 | 경제 행위 주체 | 1. 스타벅스 최초 농인 점장 권순미 씨 (수어방 송) |
| | | 경제 행위 장소 | 2. 선호하는 핸드폰 브랜드 |
| | | 경제 수단 | 3. 주식경매 경험 |
| | | 경제 산물 | 4. 통장/카드/모바일뱅킹 사용 빈도 |
| | | 경제 상태 | 5. 4차산업 |
| | | 경제 행위 | 6. 청각장애인 은행 문제(인증) 비교하기 |
| 8 | 교육 | 교수 학습 주체 | 1. 학창시절 추억 |
| | | 전공과 교과목 | 2. 새어나는 학교 수업 |
| | | 교육 기관 | 3. 동아리 활동 참여 여부 |
| | | 학교 시설 | 4. 기억에 남는 선생님 |
| | | 학습 관련 사물 | 5. 수학 배울 때 공감되는 best 4 (수어방송) |
| | | 학문 용어 | |
| | | 교수 학습 행위 | |
| 9 | 종교 | 학문 행위 | |
| | | 종교 유형 | 1. 예배 경험이 있는가? |
| | | 종교 활동 장소 | 2. 어떤 교회활동을 했는지 |
| | | 종교 인 | 3. 자기가 믿는 종교는 무엇인가? |
| | | 종교 어 | 4. 흰두고 하면 생각나는 것 |
| | | 신앙 대상 | 5. 불교 하면 생각나는 것 |
| | | 종교 활동 도구 | 6. 천주교 하면 생각나는 것 |
| | | 종교 행위 | |

16



| | | | |
|----|----|----------|--|
| 10 | 문화 | 문화 활동 주체 | 1. 모나리자와 미인도 (영화) 2. 좋아하는 가수, 노래, 장르 3. 좋아하는 책, 영화 4. 좋아하는 춤, 그림, 건물, 정원 등 5. 외국, 한국 출 영화 저작 의무화 제공 비교하기 6. 청각장애인과 청인 연예인 출연료 비교하기 7. 청각장애인과 청인 한류스타 비교하기 8. 청각장애인은 공포영화 감상에 비정 음악을 듣는 일반인하고 무슨 차이가 있을까? 9. 딜레마-나 (영화) |
| | | 음악 | |
| | | 미술 | 10. I walk with Vanessa (그림책) 11. 민중을 이끄는 자유의 여신 (영화) 12. 본노의 집주 (영화) 13. 웨朋友们对 웨 대단한 앤 설리번 이야기하기 14. 불뢰드 영화 강상 이야기하기 (줄거리, 가장 슬펐던/웃었던 장면) 15. 중국 문화 와곡 (수어방송) |
| | | 문학 | 16. Speak first and lose by Idries Shah (그림책) 17. 이태원의 할로윈 (일상) 18. 아이들의 놀이 (영화) 19. 피아노 치는 소녀 (영화) |
| | | 예술 | 20. Tomorrow (영화) 21. The snow day (그림책) 22. A beautiful day (그림책) 23. Snuffel Vier Zijn Eerste Kerstniss (그림책) 24. Sylvester og arthur praler (그림책) 25. Oscars grisede familie (그림책) |
| | | 대중문화 | 26. The little mouse, the red ripe strawberry, and the big hungry bear (그림책) 27. Found (그림책) 28. 타이타닉 (영화) 30. Hide and seek (그림책) |
| | | 전통문화 | 31. 좀비 영화 (영화) 32. 그림작품 비교 (비교) |
| | | 문화생활 장소 | 33. Waldo and the desert island adventure (그림책) 34. The famer and the clown (그림책) 35. 재미없는 게임 (게임 영화, 드라마 책) 36. 한국의 명절의 종류에 대해 나설경기 37. 한국 명절에 무엇을 하는지(한국 명절 문화) 38. 명절에 기죽들과 무엇을 하는지 39. 명절에 하는 놀이 |
| | | 문화 활동 | |

17



| | | | |
|----|--------|------------|--|
| 11 | 정치와 행정 | 공공 기관 | 1. 자동차 운전 중 사고 시 보험사 요청 방법 2. 잘 알면 유용한 법률 3. 풍사무소에 가는 이유 4. 경찰이 필요할까? 경찰이 하는 일 5. 보험 듣 것이다? 6. 법원에 기본 경험 7. 선거 관련 불편한 점 8. 법원, 경찰 수감 일어났을 때 요청 방법 9. 농어인(청각장애인)의 공화국은 실제로 가능한 일일까? |
| | | 정치 및 행정 주체 | |
| | | 사법 및 치안 주체 | |
| | | 무기 | |
| | | 정치 및 치안 상태 | |
| | | 정치 및 행정 행위 | |
| | | 사법 및 치안 행위 | 10. 만약 통일이 된다면? 11. 차등별금제 (수어방송) 12. 미루운동 (수어방송) 13. 학교폭력 (수어방송) 14. 조지 플로이드 사건 (수어방송) 15. 만 13세 전동 킥보드 허용 (수어방송) 16. 반인도적 미안파 루더타 (수어방송) 17. 5일장 때 일어난 폭걸이 도난 사건 18. LH 땅 투기 사건과 관련해서 |
| 12 | 자연 | 지형 | 1. 좋아하고 싫어하는 날씨 2. 좋아하는 계절 |
| | | 지표면 사물 | 3. 날씨에 따라 달라지는 활동 |
| | | 자원 | 4. 폭풍우에 대한 생각 |
| | | 천체 | 5. 폭우 피해와 대처법 |
| | | 재해 | 6. 습도가 높을 때와 낮을 때 |
| | | 기상 및 기후 | 7. 태풍 피해와 대처법 8. 황사/미세먼지를 일으키는 원인과 대처법 9. 벌하늘의 온하수, 별자리, 온하게, 우주선 10. 산과 바다 (비교) |
| 13 | 동식물 | 동물류 | 1. 좋아하는 동물 2. 키우고 있는 애완동물을 소개 |
| | | 곤충류 | 3. 애완동물을 지렁하기! |
| | | 식물류 | 4. 키우고 싶은 애완동물을 소개 |
| | | 동물의 부문 | 5. 사판식물에 대해 아는 것이 있나? |
| | | 식물의 부문 | 6. 이끼류에 대해 아는 것이 있나? (모스) |
| | | 동식물 행위 | 7. 식충식물(파리자속, 통풀, 꾼작이) 8. 파충류나 곤충에 대해 도마뱀, 이구아나, 사슴벌레, 장수풍뎅이...) |
| | | 동물 소리 | 9. 블러건 사고 (수어방송) |

18



| | | | |
|----|-------------|---------|----------------------------------|
| 14 | 개념 | 모양 | 1. 좋아하고 싫어하는 색 2. 여름과 겨울 (비교) |
| | | 성질 | |
| | | 속도 | |
| | | 밝기 | |
| | | 온도 | |
| | | 색깔 | |
| | | 수 | |
| | | 세는 말 | 3. 계절의 의미 |
| | | 양 | 4. 장소의 의미 |
| | | 정도 | 5. 색깔의 의미 |
| | | 순서 | 6. 조상의 의미 |
| | | 빈도 | 7. 후손의 의미 |
| | | 시간 | 8. 사랑의 의미 |
| | | 위치 및 방향 | 9. 미움의 의미 |
| | | 지역 | |
| | | 지시 | |
| | | 접속 | |
| | | 의문 | |
| | | 인칭 | |
| 15 | 나라명 및 지명 | | |
| 16 | 기타 | | |

19



3.4 전사 파일 리스트 (총 104 개: 34 개 & 70개)

| 폴더명(촬영일) | 파일명 | | | | |
|----------|--|--------|--|--------|---|
| 201228 | GI_001_GJ_40_F_1_2 | 210322 | NA_KJ_005_006_MP,GN_30,40_M,M_2,1,2,2,9 | 210427 | NA,GJ_005,006_MP,GN_30,40_M,M_2a |
| 201228 | GI_001_GJ_40_F_1_5 | 210322 | NA,KJ_005,006_MP,GN_30,40_M,M_2,9,9,2,2 | 210427 | NA,GJ_005,006_MP,GN_30,40_M,M_2b |
| 201228 | GI_001_GJ_40_F_1_10 | 210322 | NA,KJ_005,006_MP,GN_30,40_M,M_2,9,12,1,8 | 210427 | NA,GJ_005,006_MP,GN_30,40_M,M_2c |
| 210203 | GI,GJ_003,004_GI,GJ_20,30_FM_1-7-1 | 210413 | MP,MA_015,016_GJ,MA_40,40_M,M_2,1a | 210427 | NA,GJ_005,006_MP,GN_30,40_M,M_2d |
| 210203 | GI,GJ_003,004_GI,GJ_20,30_FM_1-7-2 | 210413 | MP,MA_015,016_GJ,MA_40,40_M,M_2,2,3,16 | 210427 | NA,GJ_005,006_MP,GN_30,40_M,M_2e |
| 210203 | GI,GJ_003,004_GI,GJ_20,30_FM_2-1 | 210413 | MP,MA_015,016_GJ,MA_40,40_M,M_2,6,11,1,3 | 210427 | NA,GJ_005,006_MP,GN_30,40_M,M_2f |
| 210203 | GI,GJ_003,004_GI,GJ_20,30_FM_2-2 | 210413 | MP,MA_015,016_GJ,MA_40,40_M,M_2,9,2 | 210427 | NA,GJ_005,006_MP,GN_30,40_M,M_2g |
| 210203 | GI,GJ_003,004_GI,GJ_20,30_FM_2-3 | 210413 | MP,MA_015,016_GJ,MA_40,40_M,M_2,9,12,1,8 | 210427 | NA,GJ_005,006_MP,GN_30,40_M,M_2h |
| 210203 | GI,GJ_003,004_GI,GJ_20,30_FM_2-4 | 210413 | MP,MA_015,016_GJ,MA_40,40_M,M_2,9,12,1,8 | 210428 | HS,HS_022,023_GY,PS_40,40_M,M_2,16a |
| 210203 | GI,GJ_003,004_GI,GJ_20,30_FM_2-5 | 210413 | MP,MA_015,016_GJ,MA_40,40_M,M_2,9 | 210428 | HS,HS_022,023_GY,PS_40,40_M,M_2,16b |
| 210203 | GI,GJ_003,004_GI,GJ_20,30_FM_2-6 | 210413 | MP,MA_015,016_GJ,MA_40,40_M,M_2,16a | 210428 | HS,HS_022,023_GY,PS_40,40_M,M_2,16d |
| 210203 | GI,GJ_003,004_GI,GJ_20,30_FM_2-7 | 210413 | MP,MA_015,016_GJ,MA_40,40_M,M_2,16h | 210428 | HS,HS_022,023_GY,PS_40,40_M,M_2,16e |
| 210203 | GI,GJ_003,004_GI,GJ_20,30_FM_5-1 | 210414 | WJ,GS_017,018_JJJ,J_40,40_F,M_2,1,11 | 211123 | GI,GJ_034,035_GI,HS_60,60_M,M_3,1,2 |
| 210217 | GI,GJ_003,004_GI,GJ_20,30_FM_1-12 | 210414 | WJ,GS_017,018_JJJ,J_40,40_F,M_2,1,6 | 211123 | GI,GJ_034,035_GI,HS_60,60_M,M_3,3,4,5,6 |
| 210217 | GI_003,004_GI_20,30_FM_1-1 | 210414 | WJ,GS_017,018_JJJ,J_40,40_F,M_2,1,9 | 211123 | GI,GJ_034,035_GI,HS_60,60_M,M_3,7,8,9 |
| 210217 | GI_003,004_GI_20,30_FM_2-8 | 210414 | WJ,GS_017,018_JJJ,J_40,40_F,M_2,3,16 | 211123 | GI,GJ_034,035_GI,HS_60,60_M,M_3,10,11,12 |
| 210217 | GI_003,004_GI_20,30_FM_2-9 | 210414 | WJ,GS_017,018_JJJ,J_40,40_F,M_2,9,12,18 | 220107 | GI,MP_038,039_HP,MP_30,20_F,F_3,1,2,3,4,5 |
| 210217 | GI_003,004_GI_20,30_FM_2-11 | 210414 | WJ,GS_017,018_JJJ,J_40,40_F,M_2,9,16 | 220107 | GI,MP_038,039_HP,MP_30,20_F,F_3,6,7,8,9 |
| 210218 | GI,GJ_003,004_GI,GJ_20,30_FM_11-3,2 | 210414 | WJ,GS_017,018_JJJ,J_40,40_F,M_2,16b | 220107 | GI,MP_038,039_HP,MP_30,20_F,F_3,10,11,12,13 |
| 210218 | GI,GJ_003,004_GI,GJ_20,30_FM_12-3 | 210414 | WJ,GS_017,018_JJJ,J_40,40_F,M_2,16d | 220107 | GI,MP_038,039_HP,MP_30,20_F,F_3,14,15,16 |
| 210218 | GI,GJ_003,004_GI,GJ_20,30_FM_16-1 | 210421 | MP,NJ_014,019_DJ,NJ_50,60_F,M_2a | 220114 | MU,INC_040,041_MU,IM,M,F40,50_3,1,2,3,4,5 |
| 210218 | GI,GJ_003,004_GI_20,30_FM_9-6 | 210421 | MP,NJ_014,019_DJ,NJ_50,60_F,M_2b | 220114 | MU,INC_040,041_MU,IM,M,F40,50_3,6,7,8,9 |
| 210218 | GI_003,004_GI_20,30_FM_9-7 | 210421 | MP,NJ_014,019_DJ,NJ_50,60_F,M_2c | 220114 | MU,INC_040,041_MU,IM,M,F40,50_3,10,11,12,13 |
| 210218 | GI_003,004_GI_20,30_FM_16-3 | 210421 | MP,NJ_014,019_DJ,NJ_50,60_F,M_2d | 220114 | MU,INC_040,041_MU,IM,M,F40,50_3,14,15,16 |
| 210223 | GI,GJ_003,004_GI,GJ_20,30_FM_1,16f | 210426 | NA,NA_020,021_JD,GM_40,40_M,F,2a | 221224 | GI,GJ_036,037_GH,GJ_70,80_M,M_3,1,2,3,4 |
| 210223 | GI,GJ_003,004_GI,GJ_20,30_FM_2,16a | 210426 | NA,NA_020,021_JD,GM_40,40_M,F,2b | 221224 | GI,GJ_036,037_GH,GJ_70,80_M,M_3,5,6,7,8 |
| 210223 | GI,GJ_003,004_GI,GJ_20,30_FM_2,16b | 210426 | NA,NA_020,021_JD,GM_40,40_M,F,2c | 221224 | GI,GJ_036,037_GH,GJ_70,80_M,M_3,9,10,11,12 |
| 210223 | GI,GJ_003,004_GI,GJ_20,30_FM_2,16c | 210426 | NA,NA_020,021_JD,GM_40,40_M,F,2d | 221224 | GI,GJ_036,037_GH,GJ_70,80_M,M_3,13,14,15,16 |
| 210223 | GI,GJ_003,004_GI,GJ_20,30_FM_2,16d | 210426 | NA,NA_020,021_JD,GM_40,40_M,F,2e | | |
| 210223 | GI,GJ_003,004_GI,GJ_20,30_FM_2,16e_Left | 210426 | NA,NA_020,021_JD,GM_40,40_M,F,2f | | |
| 210223 | GI,GJ_003,004_GI,GJ_20,30_FM_2,16e_Right | 210426 | NA,NA_020,021_JD,GM_40,40_M,F,2f | | |

20



| | |
|--------|---|
| 220317 | SC,BS_042,043_CW,PS_50,50_F,F_3_1,2,3,4,5 |
| 220317 | SC,BS_042,043_CW,PS_50,50_F,F_3_6,7,8,9,10 |
| 220317 | SC,BS_042,043_CW,PS_50,50_F,F_3_11,12,13 |
| 220317 | SC,BS_042,043_CW,PS_50,50_F,F_3_14,15,16 |
| 220406 | GY,GY_044,045_GY,YS_40,40_M,M_3_1,2,3,4,5 |
| 220406 | GY,GY_044,045_GY,YS_40,40_M,M_3_6,7,8,9 |
| 220406 | GY,GY_044,045_GY,YS_40,40_M,M_3_10,11,12 |
| 220406 | GY,GY_044,045_GY,YS_40,40_M,M_3_13,14,15,16 |

The **third refined corpus** contains 70 files from **210413** to 220406 after the deletion of other 34 files (from the first file of 201228 to 34th files of 210322).

- Naming of Video + Annotation Files
Current staying region_Informants number_Birth place Age group_Sex Video topics
- WJ,GS_017,018_JJ,JJ_40,40_F,M_3-3,3-4,4-1,5-1
① ② ③ ④ ⑤ ⑥
 - ① WJ : Wan-Ju(완주), GS : Gun-San (군산) / ② 017, 018: Informants
 - ③ JJ : Birth place (Jeon-Ju, 전주) / ④ 40: forties year old
 - ⑤ Female, Male: Sex / ⑥ Topics (주제)

3-3 (Food Life: most favorite food that you recently ate),
3-4 (Food Life: any food that you want to eat in future),
4-1 (Clothing Life: Tell me how to purchase clothes and shoes),
5-1 (Housing Life: your preferred place and time to meet with you're a few)

21



3.5 전사와 전사 방법

Video Clips + ELAN

녹화한 동영상 자료를 ELAN (EUDICO Linguistic Annotator) 프로그램을 사용하여 정밀전사와 어구 주해 (Gloss) 분석으로 개별 어휘와 관용 어휘가 포함된 다어휘에 사용되는 어휘 비수지신호, 품사태깅 (세종말뭉치 품사태그셋) 주석 (annotation)으로 “조선대 KSL코퍼스”구축

22

Video Clips + ELAN



ELAN 6.1 - WUGS_017.019_JJJ_40.40_FM_3-3.3-4.4-1.5-1.eaf

파일 랭킹 주석 총괄 파일 결재 보기 출판 창 도움말

그리드 퍼스트 자막 어록 사건 코멘트 인식기 편집 태이프 편집

별도:

WUGS_017.019_JJJ_40.40_FM_3-3.3-4.4-1.5-1.경.mp4

속도:

00:03:03.226 선택 00:03:04.964 1738

WUGS_017.018. 00:03:03.000 00:03:04.000 00:03:05.000 00:03:06.000 00:03:07.000 00:03:08.000 00:03:09.000 00:03:10.000 00:03:11.000 00:03:12.000 00:03:13.000 00:03:14.000 00:03:15.000 00:03:16.000 00:03:17.000 00:03:18.000 00:03:19.000 00:03:20.000

00:03:03.000 00:03:04.000 00:03:05.000 00:03:06.000 00:03:07.000 00:03:08.000 00:03:09.000 00:03:10.000 00:03:11.000 00:03:12.000 00:03:13.000 00:03:14.000 00:03:15.000 00:03:16.000 00:03:17.000 00:03:18.000 00:03:19.000 00:03:20.000

defaut

Gloss 1
보는 것 중에 하나는 나의 아저씨가 있어요

보는 것 중에 하나는 나의 아저씨가 있어요

Gloss 2
<이었><여기><곳><여기>

여기 2
이었 여기...

코멘트 2
관총포 2

관총포 2

00:03:03.000 00:03:04.000 00:03:05.000 00:03:06.000 00:03:07.000 00:03:08.000 00:03:09.000 00:03:10.000 00:03:11.000 00:03:12.000 00:03:13.000 00:03:14.000 00:03:15.000 00:03:16.000 00:03:17.000 00:03:18.000 00:03:19.000 00:03:20.000

23



ELAN 6.4 - GY.GY_044.045_GYYS_40.40_MM_3_1,2,3,4,5.eaf

파일 랭킹 주석 총괄 파일 결재 보기 출판 창 도움말

그리드 퍼스트 자막 어록 사건 코멘트 인식기 편집 태이프 편집

별도:

GY.GY_044.045_GYYS_40.40_MM_3_1,2,3,4,5_felt.mp4

GY.GY_044.045_GYYS_40.40_MM_3_1,2,3,4,5_Left.mp4

GY.GY_044.045_GYYS_40.40_MM_3_1,2,3,4,5_Right.mp4

파일 모드 편집 모드

00:32:57.190 선택 00:32:56.430 00:32:57.190 760

00:32:51.000 00:32:52.000 00:32:53.000 00:32:54.000 00:32:55.000 00:32:56.000 00:32:57.000 00:32:58.000 00:32:59.000 00:33:00.000 00:33:01.000 00:33:02.000 00:33:03.000 00:33:04.000 00:33:05.000 00:33:06.000 00:33:07.000 00:33:08.000 00:33:09.000

defaut

Gloss-제공자-1_원

Gloss-제공자-1_오
이후_제공자_1

한국어문장-제공자

코멘트_제공자_1

관총포_제공자_1

Gloss-제공자-2_원

이후_제공자_2

한국어문장-제공자

코멘트_제공자_2

관총포_제공자_2

Gloss-사용자-원

Gloss-사용자-오
이후_사용자_1

한국어문장-사용자

코멘트_사용자_1

관총포_사용자_1

00:32:51.000 00:32:52.000 00:32:53.000 00:32:54.000 00:32:55.000 00:32:56.000 00:32:57.000 00:32:58.000 00:32:59.000 00:33:00.000 00:33:01.000 00:33:02.000 00:33:03.000 00:33:04.000 00:33:05.000 00:33:06.000 00:33:07.000 00:33:08.000 00:33:09.000

24

3.6 세종 말뭉치 태그셋 (Sejong Tagset)

| 대분류 | 소분류 | 세분류 |
|----------|---------|--|
| (1) 체언 | NN 명사 | NNG 일반명사 NNP 고유명사 NNB 의존명사 |
| | NP 대명사 | NP 대명사 |
| (2) 용언 | NR 수사 | NR 수사 |
| | VV 동사 | VV 동사 |
| | VA 형용사 | VA 형용사 |
| | VX 보조용언 | VX 보조용언 |
| (3) 수식언 | VC 지경사 | VCP 금정지경사 VCN 부정지경사 |
| | MM 관형사 | |
| | MA 부사 | MAG 일반부사 MAJ 접속부사 |
| (4) 독립언 | IC 감탄사 | IC 감탄사 |
| (5) 관계언 | JK 격조사 | JKS 주격조사 JKC 보격조사 JKG 관형격조사 JKO 목적격조사 JKB 부사격조사 JKV 호格조사 JKG 인용격조사 |
| | | JKS 보조사 |
| | | JCG 접속조사 |
| | | |
| | | EP 선언같이미 EF 풍걸어미 |
| | | EC 연결어미 |
| | | ETN 명사형전성어미 |
| | | ETM 관형령전성어미 |
| (6) 의존형태 | XP 결두사 | XPN 체언결두사 |
| | XS 접미사 | XSN 명사파생접미사 XSY 동사파생접미사 |
| | XR 어근 | XSA 형용사파생접미사 |
| | | XR 어근 |

25



3.7 코퍼스 구조와 Gloss POS 분석

Annotated File Structure - Gloss POS 분석을 위한 POS Analysis Source Code Structure



POS Analysis Source Code Structure

1. Annotated Data Crawling & Gathering
2. Data Preprocessing
(Special Character, Programming Rule)
3. Part-Of-Speech Tagging (Morpheme Analysis)
(Utagger)
4. Statistics & Compare with Other Corpus



3.8 Summary: Basic Information of Compiled Corpus

The First/Second Corpus and Study Corpus

- The first refined corpus consists of 104 files containing **285,481** tokens and 20,614 types as in **Table 1**.
- The second refined corpus consists of 70 files containing **273,086** tokens and 19,183 types (after deleting 34 files recorded in the first year) as in **Table 2**.
- The **study corpus** consists of 70 files containing **251,205** tokens and 17,652 types (after retaining **content words** tagged with **four parts-of-speech** such as nouns, verbs, adjectives, and adverbs and deleting other POS from the second corpus) as in **Table 3**.
 - **5,010 lemmas & 17,149 sentences** in the translated Korean version
 - Currently 16 major categories , 88 subcategories, 199 video topics

Table 1. and Table 2. Top 20 and Bottom 10 words (types)



| Table 1. The First Corpus | | | | | Table 2. The Second Corpus | | | | | | |
|---------------------------|----------|-----|-------|----------|----------------------------|-------|----------|-----|-------|----------|-------------|
| 1 | Gloss | POS | types | types(%) | types(CUM%) | 1 | Gloss | POS | types | types(%) | types(CUM%) |
| 2 | 나 | NP | 7,691 | 2.6940 | 2.6940 | 2 | 나 | NP | 7,452 | 2.7288 | 2.7288 |
| 3 | 저기 | NP | 4,648 | 1.6281 | 4.3222 | 3 | 저기 | NP | 4,571 | 1.6738 | 4.4026 |
| 4 | 같다 | VA | 3,995 | 1.3994 | 5.7216 | 4 | 같다 | VA | 3,816 | 1.3974 | 5.8000 |
| 5 | 잇다 | VA | 3,470 | 1.2155 | 6.9371 | 5 | 잇다 | VA | 3,244 | 1.1879 | 6.9879 |
| 6 | 맞다 | VV | 3,388 | 1.1868 | 8.1238 | 6 | 맞다 | VV | 3,182 | 1.1652 | 8.1531 |
| 7 | 없다 | VA | 2,912 | 1.0200 | 9.1439 | 7 | 없다 | VA | 2,820 | 1.0326 | 9.1858 |
| 8 | 끌 | NNG | 2,666 | 0.9339 | 10.0777 | 8 | 끌 | NNG | 2,610 | 0.9557 | 10.1415 |
| 9 | 위 | IC | 2,604 | 0.9121 | 10.9899 | 9 | 위 | IC | 2,569 | 0.9407 | 11.0822 |
| 10 | 너 | NP | 2,602 | 0.9114 | 11.9013 | 10 | 너 | NP | 2,559 | 0.9371 | 12.0193 |
| 11 | 좋다 | VA | 2,573 | 0.9013 | 12.8026 | 11 | 좋다 | VA | 2,459 | 0.9004 | 12.9197 |
| 12 | 보다 | MAG | 2,546 | 0.8918 | 13.6944 | 12 | 보다 | MAG | 2,399 | 0.8785 | 13.7982 |
| 13 | 많다 | VA | 2,114 | 0.7405 | 14.4349 | 13 | 많다 | VA | 2,004 | 0.7338 | 14.5321 |
| 14 | 아니다 | VCN | 1,983 | 0.6946 | 15.1296 | 14 | 아니다 | VCN | 1,864 | 0.6826 | 15.2146 |
| 15 | ~면 | NNG | 1,775 | 0.6218 | 15.7513 | 15 | ~면 | NNG | 1,754 | 0.6423 | 15.8569 |
| 16 | 생각하다 | VV | 1,610 | 0.5640 | 16.3153 | 16 | 집 | NNG | 1,532 | 0.5610 | 16.4179 |
| 17 | 집 | NNG | 1,578 | 0.5528 | 16.8680 | 17 | 생각하다 | VV | 1,470 | 0.5383 | 16.9562 |
| 18 | 무엇 | NP | 1,567 | 0.5489 | 17.4169 | 18 | 말하다 | VV | 1,425 | 0.5218 | 17.4780 |
| 19 | 말하다 | VV | 1,498 | 0.5247 | 17.9416 | 19 | 이봐 | IC | 1,416 | 0.5185 | 17.9965 |
| 20 | 농인 | NNG | 1,453 | 0.5090 | 18.4506 | 20 | 때 | NNG | 1,376 | 0.5039 | 18.5004 |
| 21 | 때 | NNG | 1,440 | 0.5044 | 18.9550 | 21 | 잘하다 | VV | 1,374 | 0.5031 | 19.0035 |
| 20606 | 타.이.느. | VV | 1 | 0.0004 | 99.9968 | 19175 | 타.이.느. | VV | 1 | 0.0004 | 99.9967 |
| 20607 | 경례 | NNG | 1 | 0.0004 | 99.9972 | 19176 | 경례 | NNG | 1 | 0.0004 | 99.9971 |
| 20608 | '성.ㄹ' | NNG | 1 | 0.0004 | 99.9975 | 19177 | '성.ㄹ' | NNG | 1 | 0.0004 | 99.9974 |
| 20609 | '넷.풀.렛.' | NNG | 1 | 0.0004 | 99.9979 | 19178 | '넷.풀.렛.' | NNG | 1 | 0.0004 | 99.9978 |
| 20610 | 고.고.믹.' | NNG | 1 | 0.0004 | 99.9982 | 19179 | 고.고.믹.' | NNG | 1 | 0.0004 | 99.9982 |
| 20611 | 이.주.일. | NNG | 1 | 0.0004 | 99.9986 | 19180 | 이.주.일. | NNG | 1 | 0.0004 | 99.9985 |
| 20612 | 원형탈모+ | NNG | 1 | 0.0004 | 99.9989 | 19181 | 원형탈모+ | NNG | 1 | 0.0004 | 99.9989 |
| 20613 | 영.형.래. | NNG | 1 | 0.0004 | 99.9993 | 19182 | 영.형.래. | NNG | 1 | 0.0004 | 99.9993 |
| 20614 | 김.형.대 | NNG | 1 | 0.0004 | 99.9996 | 19183 | 김.형.대 | NNG | 1 | 0.0004 | 99.9996 |
| 20615 | 손들디/안J | VV | 1 | 0.0004 | 100.0000 | 19184 | 손들디/안J | VV | 1 | 0.0004 | 100.0000 |

285,481 tokens and 20,614 types

273,086 tokens and 19,183 types

29

Table 3. Top 20 and Bottom 10 words (types) in the study corpus



| Top 20 in The Study Corpus | | | | | | Bottom 10 in The Study Corpus | | | | | |
|----------------------------|-------|-----|-------|----------|-------------|-------------------------------|---------|-----|---|--------|----------|
| 1 | Gloss | POS | types | types(%) | types(CUM%) | 17644 | 전부/전체 | MAG | 1 | 0.0004 | 99.9964 |
| 2 | 나 | NP | 7,452 | 2.9665 | 2.9665 | 17645 | 갸웃뚱 | MAG | 1 | 0.0004 | 99.9968 |
| 3 | 저기 | NP | 4,571 | 1.8196 | 4.7861 | 17646 | 전.영.수 | MAG | 1 | 0.0004 | 99.9972 |
| 4 | 같다 | VA | 3,816 | 1.5191 | 6.3052 | 17647 | 김.영.란 | MAG | 1 | 0.0004 | 99.9976 |
| 5 | 잇다 | VA | 3,244 | 1.2914 | 7.5966 | 17648 | ~만/~뿐/▷ | MAG | 1 | 0.0004 | 99.9980 |
| 6 | 맞다 | VV | 3,182 | 1.2667 | 8.8633 | 17649 | 전.영. | MAG | 1 | 0.0004 | 99.9984 |
| 7 | 없다 | VA | 2,820 | 1.1226 | 9.9859 | 17650 | 더. | MAG | 1 | 0.0004 | 99.9988 |
| 8 | 끌 | NNG | 2,610 | 1.0390 | 11.0249 | 17651 | 억지로/강 | MAG | 1 | 0.0004 | 99.9992 |
| 9 | 너 | NP | 2,559 | 1.0187 | 12.0436 | 17652 | 영.화.' | MAG | 1 | 0.0004 | 99.9996 |
| 10 | 좋다 | VA | 2,459 | 0.9789 | 13.0224 | 17653 | '쿵.후.' | MAG | 1 | 0.0004 | 100.0000 |
| 11 | 보다 | MAG | 2,399 | 0.9550 | 13.9774 | | | | | | |
| 12 | 많다 | VA | 2,004 | 0.7978 | 14.7752 | | | | | | |
| 13 | ~면 | NNG | 1,754 | 0.6982 | 15.4734 | | | | | | |
| 14 | 집 | NNG | 1,532 | 0.6099 | 16.0833 | | | | | | |
| 15 | 생각하다 | VV | 1,470 | 0.5852 | 16.6685 | | | | | | |
| 16 | 말하다 | VV | 1,425 | 0.5673 | 17.2357 | | | | | | |
| 17 | 때 | NNG | 1,376 | 0.5478 | 17.7835 | | | | | | |
| 18 | 잘하다 | VV | 1,374 | 0.5470 | 18.3304 | | | | | | |
| 19 | 농인 | NNG | 1,370 | 0.5454 | 18.8758 | | | | | | |
| 20 | 안하다 | VV | 1,314 | 0.5231 | 19.3989 | | | | | | |
| 21 | 무엇 | NP | 1,311 | 0.5219 | 19.9208 | | | | | | |

251,205 tokens and 17,652 types

30

Table 3. Top 20 and Bottom 10 words in the study corpus



| Top 20 in The Study Corpus | | | | | Bottom 10 in The Study Corpus | | | | |
|----------------------------|-------|-----|-------|----------|-------------------------------|--|--|--|--|
| 1 | Gloss | POS | types | types(%) | types(CUM%) | | | | |
| 2 | 나 | NP | 7,452 | 2.9665 | 2.9665 | | | | |
| 3 | 저기 | NP | 4,571 | 1.8196 | 4.7861 | | | | |
| 4 | 같다 | VA | 3,816 | 1.5191 | 6.3052 | | | | |
| 5 | 있다 | VA | 3,244 | 1.2914 | 7.5966 | | | | |
| 6 | 맞다 | VV | 3,182 | 1.2667 | 8.8633 | | | | |
| 7 | 없다 | VA | 2,820 | 1.1226 | 9.9859 | | | | |
| 8 | 끌 | NNG | 2,610 | 1.0390 | 11.0249 | | | | |
| 9 | 너 | NP | 2,559 | 1.0187 | 12.0436 | | | | |
| 10 | 좋다 | VA | 2,459 | 0.9789 | 13.0224 | | | | |
| 11 | 보다 | MAG | 2,399 | 0.9550 | 13.9774 | | | | |
| 12 | 많다 | VA | 2,004 | 0.7978 | 14.7752 | | | | |
| 13 | ~면 | NNG | 1,754 | 0.6982 | 15.4734 | | | | |
| 14 | 집 | NNG | 1,532 | 0.6099 | 16.0833 | | | | |
| 15 | 생각하다 | VV | 1,470 | 0.5852 | 16.6685 | | | | |
| 16 | 말하다 | VV | 1,425 | 0.5673 | 17.2357 | | | | |
| 17 | 때 | NNG | 1,376 | 0.5478 | 17.7835 | | | | |
| 18 | 잘하다 | VV | 1,374 | 0.5470 | 18.3304 | | | | |
| 19 | 농인 | NNG | 1,370 | 0.5454 | 18.8758 | | | | |
| 20 | 안하다 | VV | 1,314 | 0.5231 | 19.3989 | | | | |
| 21 | 무엇 | NP | 1,311 | 0.5219 | 19.9208 | | | | |

251,205 tokens and 17,652 types

31

전사된 텍스트의 품사별 token, type 비교



| POS | 서술어 | 수사 | 대명사 | 고유명사 | 일반명사 | 의존명사 | 접속부사 | 일반부사 | 긍정지정사 |
|---------------|-----------------|--------------|---------------|--------------|----------------|--------------|--------------|---------------|--------------|
| Type (ratio) | 7,395 (52.7%) | 97 (0.7%) | 87 (0.6%) | 404 (2.9%) | 4,681 (33.4%) | 288 (2.1%) | 7 (0.0%) | 393 (2.8%) | 274 (2.0%) |
| Token (ratio) | 128,830 (47.8%) | 3,292 (1.2%) | 26,315 (9.8%) | 5,921 (2.2%) | 76,604 (28.4%) | 3,667 (1.4%) | 1,173 (0.4%) | 10,960 (4.1%) | 2,566 (1.0%) |

| POS | 부정지정사 | 부사격조사 | 보조사 | 접속조사 | 감탄사 | 연결어미 | 종결어미 | 관형사 |
|---------------|--------------|--------------|--------------|-----------|--------------|--------------|-----------|--------------|
| Type (ratio) | 14 (0.1%) | 159 (1.1%) | 85 (0.6%) | 9 (0.1%) | 64 (0.5%) | 38 (0.3%) | 26 (0.2%) | 6 (0.0%) |
| Token (ratio) | 1,934 (0.7%) | 2,213 (0.8%) | 1,017 (0.4%) | 26 (0.0%) | 1,723 (0.6%) | 2,050 (0.8%) | 61 (0.0%) | 1,234 (0.5%) |

표제어 (lemma) 분석

KSL Online Dictionary (국립국어원 한국수어사전)의 사례

- 한국어 대응표현
한국수어 lemma (표제어) 웹사전제공
→ 3,676개 (7,112 개 types)

- 한국수어 표제어 등록수 → 24,705개(??)

(예)

- 일련번호 1~ 20 예 (우측표)

| 번호 | 수어 표제어 번호 | 한국어 대응표현 |
|----|-----------|------------------------------------|
| 1 | 24705 | 복작, 복임 |
| 2 | 24704 | 이론, 쓰다 |
| 3 | 24701 | (잇)이짜다 |
| 4 | 24049 | 흘끼봐 |
| 5 | 24047 | (마음이)즘疼다, 풍풀ഴ하다, 용하다 |
| 6 | 24041 | (바)풀을입니다 |
| 7 | 24037 | 인도(중화국) |
| 8 | 24036 | 아니다. 않다 |
| 9 | 24034 | 여자, 여성, 여 |
| 10 | 24033 | 설거지 |
| 11 | 24032 | 배부르다 |
| 12 | 24030 | (여자를)유혹, 유혹하다 |
| 13 | 24029 | 삼, 셋 |
| 14 | 24028 | 시간, 한 시간 |
| 15 | 24026 | 놓아주다 |
| 16 | 24024 | 들, 여러, 여러가지 |
| 17 | 24022 | 시작, 개시, 거행, 시발, 착수, 출발, 출다, 이행, 하수 |
| 18 | 24020 | 튿(날)에 생기는)날 |
| 19 | 24019 | 참다, 견디다, 인내, 자체, 억제 |
| 20 | 24015 | 의리 |

| 1 | 수어 표제어 번 | |
|------|----------|---------|
| 2 | 1 | 파나마 |
| 3 | 10 | 구주 |
| 4 | 10 | 유럽 |
| 5 | 11 | -국 |
| 6 | 11 | 국가 |
| 7 | 11 | 나라 |
| 8 | 12 | 국가시험 |
| 9 | 12 | 국가 고시 |
| 10 | 16 | 군집 |
| 11 | 16 | 떼 |
| 12 | 16 | 무리 |
| 13 | 20 | 공색 |
| 14 | 22 | 당부 |
| 15 | 22 | 부탁 |
| 16 | 22 | 요청 |
| 17 | 22 | 청하다 |
| 18 | 22 | 요구 |
| 19 | 23 | 권위 |
| 20 | 24 | 권위자 |
| 21 | 28 | 귀여워하다 |
| 22 | 28 | 귀엽다 |
| 7105 | 24047 | (마음이)짜다 |
| 7106 | 24047 | 옹출하다 |
| 7107 | 24047 | 옹하다 |
| 7108 | 24049 | 핥끼봐 |
| 7109 | 24701 | (맛이)짜다 |
| 7110 | 24704 | 이론 |
| 7111 | 24704 | 쓰다 |
| 7112 | 24705 | 복직 |
| 7113 | 24705 | 복임 |

33

풀사별 Top 30 - 서술어와 일반명사

| Gloss | types | types(%) | types(CUM%) | 도달 번호 |
|----------------------------|-------|----------|-------------|-------|
| 그만 까지 끝 끝+ 끝++ 끝+*** 끝+*** | 3,462 | 45.193 | 45.193 | 688 |
| 더럽다 더럽다+ 때 볼걸 지저분하다 | 1,598 | 20.061 | 6.6054 | 52 |
| 가죽 댁 세대 주택 집 집++ 호 | 1,575 | 2.0560 | 8.6614 | 1455 |
| 농아인 능인 농인+ 능인++ | 1,384 | 1.8067 | 10.4681 | 1582 |
| 밀기 오늘 오늘날+ 험재 혐재+ 혐재+ | 1,327 | 1.7323 | 12.2004 | 1645 |
| -인 분 사람 사람+ 사람++ 이 인간 자 | 1,310 | 1.7101 | 13.9105 | 736 |
| 방법 방식 방식+ 방식++ 요령 | 1,153 | 1.5051 | 15.4156 | 857 |
| 여 여성 여자 여자+ | 980 | 1.2793 | 16.6950 | 2199 |
| 예전 옛날+ 옛날+ 옛적 | 876 | 1.1435 | 17.8385 | 656 |
| 얼굴이름 얼굴이름5+ | 817 | 1.0665 | 18.9050 | 2202 |
| 곳+ 곳+ 곳++ 곳++ 곳+*** 위치 장소 | 798 | 1.0417 | 19.9467 | 2305 |
| 풀 모습 모양 상처 상황 상황+ 장면 형상 | 690 | 0.9007 | 20.8475 | 2226 |
| 다 모두 모두+ 모든 모조리 풍랑 은은통 | 689 | 0.8994 | 21.7469 | 2377 |
| 강아지 개 개+ 개+ 개+ 개+ 강아지 | 664 | 0.8668 | 22.6137 | 997 |
| 동무 벗 우인 친구 친구+ 친구++ | 617 | 0.8054 | 23.4191 | 2154 |
| 데 해+ 해 해하다+ 해하다+ 해하다++ | 612 | 0.7989 | 24.2181 | 579 |
| 남 남성 남자 남자+ 남자++ 남자+++ 남 | 598 | 0.7806 | 24.9987 | 1585 |
| 갓난이기갓난이아이 아기+ 영아 젖먹 | 589 | 0.7689 | 25.7676 | 847 |
| 문제 문제+ 문제++ | 578 | 0.7545 | 26.5221 | 696 |
| 국 국가 나라 | 577 | 0.7532 | 27.2753 | 4939 |
| 청인 청인+ | 561 | 0.7323 | 28.0077 | 2157 |
| 가격 금값 금전 대금 돈 돈++ 돈+ | 547 | 0.7141 | 28.7217 | 2281 |
| 선생님 선생님+ 선생님++ | 535 | 0.6984 | 29.4201 | 1526 |
| 학교 학교 | 534 | 0.6971 | 30.1172 | 4927 |
| (꽃)피다 꽃+ 꽃++ 꽃++ 꽃+*** | 534 | 0.6971 | 30.8143 | 1590 |
| 시간 시간간+ 시간++ 한 | 514 | 0.6710 | 31.4853 | 1513 |
| 화상 영화 | 505 | 0.6592 | 32.1445 | 3391 |
| 법 법+ 법++ 법++ 법규 법률 법칙 | 483 | 0.6305 | 32.7751 | 738 |
| 결과 결국 결론+ 결말 드이어 마지막 마침 | 482 | 0.6292 | 33.4043 | 1605 |
| 의+ 시 사+ 시+ 사++ 사+++ | 436 | 0.5692 | 33.9734 | 1529 |



일반부사 Top 30

| Gloss | types | types(%) | types(CUM%) | 토탈 번호 |
|--------------------------|-------|----------|-------------|-------|
| 그리고 다시 또 또++ 또+++ 또+++ | 993 | 9.0602 | 9.0602 | 2381 |
| 다소 약간 약소하다 자그마하다 작다 작다 | 884 | 8.0657 | 17.1259 | 50 |
| 자연히 자연히+ 자연히++ 자연히+++ 지 | 860 | 7.8467 | 24.9726 | 1649 |
| 어째서 왜 왜++ 왜+++ 왜++++ | 702 | 6.4051 | 31.3777 | 2371 |
| 가장 맨 수석 유품 일등 일등+ 일등++ 일 | 625 | 5.7026 | 37.0803 | 1482 |
| 갓 시방 이제 이제 지금 지금+ 지금++ 지 | 619 | 5.6478 | 42.7281 | 2367 |
| 나날 날마다 늘 매번 매양 매일 번번이 수 | 502 | 4.5803 | 47.3084 | 2366 |
| ㄹ 것이다 아마 아마+ 아마++ 아마도 | 499 | 4.5529 | 51.8613 | 1000 |
| 과 반려 쌍 짹 함께 함께+ 함께++ 함께하 | 463 | 4.2245 | 56.0858 | 359 |
| 김의자 길게늘어서다 길다 길다+ 오래 오 | 324 | 2.9562 | 59.0420 | 251 |
| 그냥 그냥+ 그냥++ 그저 | 319 | 2.9106 | 61.9526 | 1001 |
| 상호 서로 서로+ 서로++ | 301 | 2.7464 | 64.6989 | 607 |
| 그대로 그대로+ 그대로++ 되돌리다 복귀 | 260 | 2.3723 | 67.0712 | 885 |
| 잠깐 잠깐+ 잠깐++ 잠깐++ 잠시 | 195 | 1.7792 | 68.8504 | 1648 |
| 날 날짜 언제 언제+ 언제++ 일자 | 188 | 1.7153 | 70.5657 | 999 |
| 미~ 미처 아직 아직+ 아직++ | 183 | 1.6697 | 72.2354 | 1652 |
| 직접 직접+ 직접++ | 178 | 1.6241 | 73.8595 | 1647 |
| 가끔 가끔+ 가끔++ 간간이 간혹 드문드문 | 148 | 1.3504 | 75.2099 | 1003 |
| 갑자기 갑자기+ 갑작스럽다 갑작스레 냅다 | 135 | 1.2318 | 76.4416 | 1002 |
| 자유 임의 마구 마음껏 마음대로 헛기 | 133 | 1.2135 | 77.6551 | 3041 |
| 꽤 대 대단히 매우 매우+ 매우++ 매우++ | 119 | 1.0858 | 78.7409 | 2379 |
| ~자 곧 곧 영간 냉큼 빨리 빨리+ 빨리++ | 118 | 1.0766 | 79.8175 | 2376 |
| 도무지 | 117 | 1.0675 | 80.8850 | |
| 오히려 오히려+ | 110 | 1.0036 | 81.8887 | 2373 |
| 먼저 먼저+ 먼저++ 먼저+++ | 106 | 0.9672 | 82.8558 | 2378 |
| ~끼리 같이 더불다 둘만 둘만하다 랑 아울 | 83 | 0.7573 | 83.6131 | 648 |
| 각 각각 각각+ 각각++ 각각++ 각기 각 | 69 | 0.6296 | 84.2427 | 740 |
| 혁 혁.+ | 68 | 0.6204 | 84.8631 | 2365 |
| 미리 미리+ | 61 | 0.5566 | 85.4197 | 1653 |
| 무조건 무작정 | 55 | 0.5018 | 85.9215 | 2423 |



4. 한국수어 문장 구조 파악을 위한 검색과정 시연

4.1 조선대KSL코퍼스 브라우저

4.2 의문사를 활용한 의문문 추출 과정

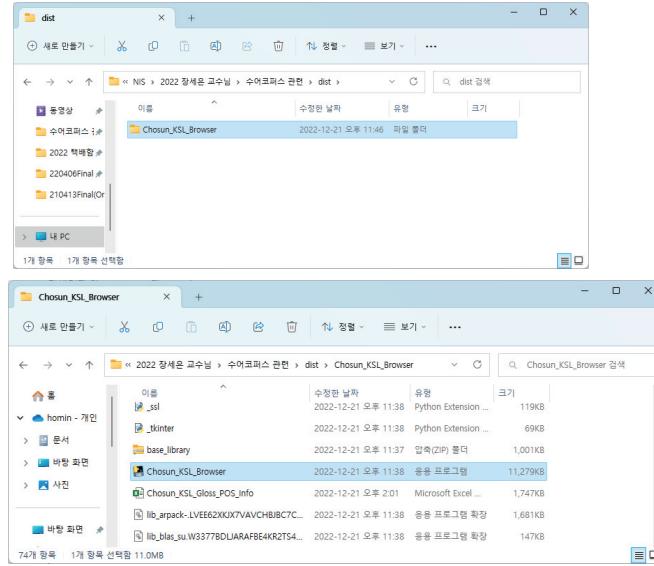
4.3 명사/동사 추출: 명사와 동사의 형태 구분하기

4.3 ‘~면’ 종속복문 구조



4.1 조선대KSL코퍼스 브라우저 (Chosun KSL Browser)

- SEARCH Function Demonstration



Chosun KSL Browser : <누구> 'who'



Chosun KSL Browser

SEARCH STATISTICS CONTEXT OVERVIEW

List Chart Collocates KWIC

Searching Keyword : **누구** Search

Option

D:\WNKG\2022 장세은 교수님 \Chosun_KSL_Browser\dist\Chosun_KSL_Browser

전체 POS initializing complete...
./220406수 어 코 퍼스 (211123-220406Final(Onedrive)\GY,GY_004,045,GY,Y5_40,40_R,R_3_1,2,3,4,5.eaf
TIER 뼐 아 대 는 응...
TIER 뼐 아 대 기 완료...
한국어문장-에 뤄-gloss 정렬 시작
한국어문장-에 뤄-gloss 정렬 및 저장 완료
./220406수 어 코 퍼스 (211123-220406Final(Onedrive)\GY,GY_004,045,GY,Y5_40,40_R,R_3_10,11,12.eaf
TIER 뼐 아 대 는 응...
TIER 뼐 아 대 기 완료...
한국어문장-에 뤄-gloss 정렬 시작
한국어문장-에 뤄-gloss 정렬 및 저장 완료
./220406수 어 코 퍼스 (211123-220406Final(Onedrive)\GY,GY_004,045,GY,Y5_40,40_R,R_3_13,14,15,16.eaf
TIER 뼐 아 대 는 응...
TIER 뼐 아 대 기 완료...
한국어문장-에 뤄-gloss 정렬 시작
한국어문장-에 뤄-gloss 정렬 및 저장 완료
./220406수 어 코 퍼스 (211123-220406Final(Onedrive)\GY,GY_004,045,GY,Y5_40,40_R,R_3_6,7,8,9.eaf
TIER 뼐 아 대 는 응...
TIER 뼐 아 대 기 완료...
한국어문장-에 뤄-gloss 정렬 시작
한국어문장-에 뤄-gloss 정렬 및 저장 완료
검색한 Keyword 누구 의 gloss 및 한국어문장 영상 정보 추출 시작.....
누구 1930089 190599 190599 190599 190599
누구 1930089 1938539 1938539 1938539 1938539
누구 3011533 461230 395234 411287
누구 471 471 453470 4533596 450145 875927
누구 2866758 2866992 2809974 2869799
누구 2866758 2867938 2809974 2869799
누구 2868297 2868897 2809974 2869799
누구 2871627 2871794 2871260 2873964
누구 3312786 3312952 3282256 3723368
누구 850922 860222 8306994 863793
누구 1167195 1167195 1139335 1193382
누구 21465183 2146318 2139638 2148984
누구 38553 38853 37882 477346
누구 409716 41208 37882 477346
영광_누구 899639 499855 487853 499896
누구 1110809 1111156 1102999 1124899
검색한 keyword 누구 의 gloss 및 한국어문장 영상 정보 추출 완료.....
누구
누가_누구
얼굴_누구



4.2 의문사를 활용한 의문문 추출 과정: <누구> 'who'

| Index | Document Name | 어휘 | 한국어문장 |
|-------|--------------------------------|--|--------------------------------------|
| 1 | 220406Final(OndrineW/GY_GY_044 | <충증거리다><저기><누구><기다(방문하다)><있다><없다> | 충증거리면 누나가 방문할지도 몰라요 모두가 비슷하죠 |
| 2 | 220406Final(OndrineW/GY_GY_044 | <관한><있다><망다><관한><망다><관한><나니><관하다> | |
| 3 | 220406Final(OndrineW/GY_GY_044 | <성공하다(기기다)><수리하다><기기><그><누> 관한이 있지... 그 단체에서 싸워서 국립간검보침을 훔쳐 차운을 받을 수 있게 되었지. 수술비 전체를 지원받아서 | |
| 4 | 220406Final(OndrineW/GY_GY_044 | <같다><남자><3><그러나><보통이다><관한><망다><관하다><관한><망다><관한><나니><관하다><성공하다(기기다)><수리하다><기기><그><누> 관한이 있지... 그 단체에서 싸워서 국립간검보침을 훔쳐 차운을 받을 수 있게 되었지. 수술비 전체를 지원받아서 | |
| 5 | 220406Final(OndrineW/GY_GY_044 | <같다><남자><3><그러나><보통이다><관한><망다><관하다><관한><망다><관한><나니><관하다><성공하다(기기다)><수리하다><기기><그><누> 관한이 있지... 그 단체에서 싸워서 국립간검보침을 훔쳐 차운을 받을 수 있게 되었지. 수술비 전체를 지원받아서 | |
| 6 | 220406Final(OndrineW/GY_GY_044 | <같다><남자><3><그러나><보통이다><관한><망다><관하다><관한><망다><관한><나니><관하다><성공하다(기기다)><수리하다><기기><그><누> 관한이 있지... 그 단체에서 싸워서 국립간검보침을 훔쳐 차운을 받을 수 있게 되었지. 수술비 전체를 지원받아서 | |
| 7 | 220406Final(OndrineW/GY_GY_044 | <저거><누구><더한다><얼굴불룩하다><저거><얼굴불룩하다> | 누가 이렇게 지저분하게 먹었어? 하고 보면 셋째인 얼굴불룩 이에요 |
| 8 | 220406Final(OndrineW/GY_GY_044 | <맞다><이쁘다><맞다><알았다><내부><속삭이다><나나다><발광하다><풀나다><제><전화하다> | |
| 9 | 220406Final(OndrineW/GY_GY_044 | <사고나다><풀사고><남여지다><여기><부모><나><카운다><내려><태어나다(남여)><수><조급><사고가 난 사건인데요 얼굴에는 자리를 낸 일이 있고 그래서 어린 아이들의 사건이 상대적으로 많은데 보통은 | |
| 10 | 220406Final(OndrineW/GY_GY_044 | <예풀풀다><충강><언><충강><언><충민><편편하다><무리다><조회한><경찰><주시하다><꼭다><개인 예풀풀은 중앙선을 놀인이 칭찬한 경찰이 주시하고 있다가 현장에서 짧으면 직접 병적금을 부과하고 경찰에 | |

| Index | Document Name | 어휘 | 한국어문장 |
|-------|--------------------------------|--|--------------------------------------|
| 1 | 220406Final(OndrineW/GY_GY_044 | <충증거리다><저기><누구><기다(방문하다)><있다><없다> | 충증거리면 누나가 방문할지도 몰라요 모두가 비슷하죠 |
| 2 | 220406Final(OndrineW/GY_GY_044 | <관한><있다><망다><관한><망다><관한><나니><관하다> | |
| 3 | 220406Final(OndrineW/GY_GY_044 | <성공하다(기기다)><수리하다><기기><그><누> 관한이 있지... 그 단체에서 싸워서 국립간검보침을 훔쳐 차운을 받을 수 있게 되었지. 수술비 전체를 지원받아서 | |
| 4 | 220406Final(OndrineW/GY_GY_044 | <같다><남자><3><그러나><보통이다><관한><망다><관하다><관한><망다><관한><나니><관하다><성공하다(기기다)><수리하다><기기><그><누> 관한이 있지... 그 단체에서 싸워서 국립간검보침을 훔쳐 차운을 받을 수 있게 되었지. 수술비 전체를 지원받아서 | |
| 5 | 220406Final(OndrineW/GY_GY_044 | <같다><남자><3><그러나><보통이다><관한><망다><관하다><관한><망다><관한><나니><관하다><성공하다(기기다)><수리하다><기기><그><누> 관한이 있지... 그 단체에서 싸워서 국립간검보침을 훔쳐 차운을 받을 수 있게 되었지. 수술비 전체를 지원받아서 | |
| 6 | 220406Final(OndrineW/GY_GY_044 | <같다><남자><3><그러나><보통이다><관한><망다><관하다><관한><망다><관한><나니><관하다><성공하다(기기다)><수리하다><기기><그><누> 관한이 있지... 그 단체에서 싸워서 국립간검보침을 훔쳐 차운을 받을 수 있게 되었지. 수술비 전체를 지원받아서 | |
| 7 | 220406Final(OndrineW/GY_GY_044 | <저거><누구><더한다><얼굴불룩하다><저거><얼굴불룩하다> | 누가 이렇게 지저분하게 먹었어? 하고 보면 셋째인 얼굴불룩 이에요 |
| 8 | 220406Final(OndrineW/GY_GY_044 | <맞다><이쁘다><맞다><알았다><내부><속삭이다><나나다><발광하다><풀나다><제><전화하다> | |
| 9 | 220406Final(OndrineW/GY_GY_044 | <사고나다><풀사고><남여지다><여기><부모><나><카운다><내려><태어나다(남여)><수><조급><사고가 난 사건인데요 얼굴에는 자리를 낸 일이 있고 그래서 어린 아이들의 사건이 상대적으로 많은데 보통은 | |
| 10 | 220406Final(OndrineW/GY_GY_044 | <예풀풀다><충강><언><충강><언><충민><편편하다><무리다><조회한><경찰><주시하다><꼭다><개인 예풀풀은 중앙선을 놀인이 칭찬한 경찰이 주시하고 있다가 현장에서 짧으면 직접 병적금을 부과하고 경찰에 | |

예 (1) <누구> 'who'



Gloss video



구부린 오른 손가락 등으로 오른뺨을 상하로 두 번 살짝 문지른다.

한국어문장 Video



어휘 :
<저거><**누구**><더럽다><얼굴불룩하다><저거><얼굴불룩하다>

한국어문장: **누가** 이렇게
지저분하게 먹었어? 하고 보면
셋째인 얼굴불룩이에요



예 (2) <언제> 'when'



가슴 앞에서 모로
세운 왼손 위에
모로 세운
오른손을 올려놓고
5·1·2·3·4지
순으로 손가락을
접는다.

Gloss video



한국어문장 Video



어휘 :

<나><예전><저기><좌파><
우파><문제><때문에><합의
+><언제+><몇년전><계속><
때문에><저기><문제><저기
><문제><고민><되다><맞다
><고민>
<무엇><방법><~면><친하다
><좋다><무엇><방법>

한국어문장 : 예전에 파가
나뉘어서 **언제쯤** 합의하**려나**
했는데 좋은 방법이
무엇일지 고민되더라고요

4.3 명사, 동사 추출 : 명사 동사의 형태 구분하기

예 (1) <밥> / <먹다>

명사 : <밥>

Gloss video



한국어문장 Video



동사 : <먹다>

Gloss video



한국어문장 Video



<아들><아버지><감정/기분><좋다><하다
><받다><이거><뭐><잘하다><듣다>
<시험치다/경쟁하다><질수><좋다><시간
><꼭/약속><조심하다><지키다><제일/가
장><**밥**><먹지않다++><화내다>
<또><과자><만><먹지 않다>><과자><먹
다><뭐><사랑하다><주다>
<멜디/청인><에민하디><날카롭다><사랑
하다><일절><징계하다><또><밤><늦다>
<성/부도덕><잇디><증기하다><빨리><오
다><맞다><주다><닭><뭐><먹다><원하다
>
<사랑하다><뭐><의미하다><맞다><뭐><
원하다><닭><주다>+>

아버지로서 기분이 좋을 때
예를 들면 시험을 잘 치거나
이야기를 잘 들어주거나
시간을 잘 지키거나 **밥**을 잘
먹을 때나 요즘 밤에는
성적으로 문란하기 때문에
일찍 들어오면 뭘 먹고 싶나고
물어보고 통닭이라도
사주는데 그렇지 않고 과자만
먹고 하면 매운 맛을
보여줘야해요

<나++><전국><맛><옳거다
니다><맛있다><그러나><
제일/가장><맛있다>
<제일/가장><잇다><못하
다><뭐><첫째><뭐><나><
춘천><저기>
<춘천><송.어><**먹다**><경험
하다><있다><아마><너>

전국을 다니면서 맛있는
음식을 먹어봤는데요
가장 맛있는 음식이 춘천에
송어 회가 가장 맛있었어요
당신도 **먹어** 봤죠?



예 (2) <생각> / <생각하다>

명사 : <생각>

Gloss video



<역사><취미><가지다><나><운동><취미><가지다><생각><반대/ 다르다>

한국어문장 Video



역사에 대한 관심이 있고 난 스포츠에 관심이 있으니 생각이 다르네

동사 : <생각하다>

Gloss video



<그것><나><그것><생각하다><그것><나><생각하다><보다><완전하다><아직><생각하다><충분하다/만족하다><~면><40%><부족하다>

한국어문장 Video



나는 아직 그점에 대해 부족하다고 생각해요
40%정도라고 생각해요

4.4 <~면> 종속 복문



손끝이 오른쪽으로 향하게 모로 세운 원 손등에 오른손의 1-5지 끝을 맞대어 동그라미를 만들어 등을 댄다.

(예1) <~면>으로 보는 수어의 복문 구조

Gloss video



한국어문장 Video



<앞으로+><잘하다><코로나><사라지다><~면><기회><열심히하다><앞으로/미래><나><보다><멋 있다><~면><나타내다/광고하다++><소개하다/통역하다><선생님><임명하다><너>

앞으로 코로나 사라지 면 열심히해서 멋진 모습을 보여주세요 그러면 앞으로 선생님으로 임명 할께요



(예2) <~면>으로 보는 수어의 복문 구조

Gloss video



<너><날><조건><취미><나>
<마음><아프다><나><선물>
<~만><아니다><너><신발><
사다><너><나><좋다><선물
><너><선물><주다><**~면**>
<동맹/피를 나눈 형제>
<알다><그것>

한국어문장 Video



조건이 잘 맞으면/ 마음이
아프네요 선물만
아니라 신발까지 좋은거
사줄께요 **그리면**
우린 피를 나 눈
형제처럼되는거예요

(예3) <~면>으로 보는 수어의 복문 구조

Gloss video



<같다><맞다><사람><사랑
하다><나><지혜(생각하다
잘하다)><**~면**>
<세상><가장><일등><한국>
<부자><1><일등><가능하다
><맞다>

한국어문장 Video



그렇게 서로 사랑으로
지혜롭게 행하**면**
한국이 세계에서 1등
나라가 될거예요



5. 제언: 조선대KSL코퍼스 브라우저 고도화

5.1 일반적인 연구용 데이터 또는 프로그램 배포 및 공유 방법

비 자동화 방법 & 로컬 사용 (non-online)

1. 연구자 또는 연구소의 홈페이지가 없는 경우
 - ↳ 이메일로 연락하여 주고 받기 (+서약서, 개인정보, citation 방법 등)
2. 연구자 또는 연구소의 홈페이지가 있는 경우
 - ↳ 홈페이지 상의 기능을 이용하여 주고 받기 (회원가입 등)

자동화 방법 & 온라인 사용

1. 데이터의 경우 : 일반적으로 홈페이지가 존재. 회원가입, 서약서 제출 등을 통해 자격을 획득하여 다운로드
2. 프로그램의 경우 : 프로그램 실행 및 결과 출력을 위한 홈페이지가 따로 존재.
회원가입, 서약서 제출, 자격 획득을 위한 결제 등을 통해
온라인 상에서 프로그램 사용

가장 접근성이 좋고 편리하며 관리하기 쉬운 방법은 '자동화 & 온라인'

5.2 서버-클라이언트 시스템 구축



하지만 '자동화 & 온라인' 방법을 사용하려면 반드시 서버-클라이언트 시스템이 구축되어야 함.

구축 방법 1.

서버 컴퓨터 구매 > 도메인 구매 > 서버 프로그램, 클라이언트 프로그램 구축 > 온라인

(접속자 수, 프로그램 구동 속도 등이 서버 컴퓨터 성능과 연결된 인터넷 망 속도에 좌우)
(다만 서버 컴퓨터 성능은 컴퓨터를 구매하면 끝이기 때문에 추가 예산 확보 후 업그레이드하지 않으면 변하지 않음)

구축 방법 2.

클라우드 컴퓨팅 서비스 구매 > 도메인 구매 > 서버 프로그램, 클라이언트 프로그램 구축 >
클라우드에 업로드 > 온라인

(접속자 수, 프로그램 구동 속도 등이 연결된 인터넷 망과 선택한 서비스 성능에 좌우)
(서버 용량, 컴퓨팅 속도 등을 서비스 상에서 고를 수 있음)

즉, 실물 컴퓨터를 구매하느냐 or 가상 컴퓨터를 대여하느냐 차이



References

- 박호민 (2022). 한국수어 문장구조 파악을 위한 검색 과정 시연, 2022. 12. 22 조선대학교 언어융합연구소 코퍼스 및 예문 제작 및 증강현실 구현을 위한 회의자료.
- 장세은, 박호민 (2023). 수어 말뭉치 비교 분석 및 학습 어휘 선정, 2023 조선대학교 언어융합연구소 학술대회 및 성과보고회 (2023년 06월 09일, 경주 더케이호텔) 성과발표 3.
- Bono, Mayumi, Kouhei Kikuchi, Paul Cibulka, and Yutaka Osugi. (2014). [A Colloquial Corpus of Japanese Sign Language: Linguistic Resources for Observing Sign Language Conversations](#). In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, pages 1898–1904, Reykjavik, Iceland. European Language Resources Association (ELRA).
- Fung, Cat Hiu Man, Scholastica Lam, Joe Mak and Gladys Tang. 2008. Establishment of a corpus of Hong Kong Sign Language acquisition data: From ELAN to CLAN. In : [Proceedings of Workshop on the Representation and Processing of Sign Languages](#). Construction and Exploitation of Sign Language Corpora, Marrakech, Morocco. 17-21.
- Neidle, Carol and Christian Vogler (2012) A New Web Interface to Facilitate Access to Corpora: Development of the ASLLRP Data Access Interface, *Proceedings of the 5th Workshop on the Representation and Processing of Sign Languages: Interactions between Corpus and Lexicon*, LREC 2012, Istanbul, Turkey.
- Schembri, A., Fenlon, J., Rentelis, R., Stamp, R., & Cormier, K. (2011, October 2011). British Sign Language Corpus Project. Paper presented at the National Institute of Informatics, and Japanese Federation of the Deaf in Tokyo, Japan; and National Centre of Sign Language Education, and Japanese Association of Sign Language Interpreters in Kyoto, Japan.
- Schembri, A., Fenlon, J., Rentelis, R., Reynolds, S., & Cormier, K. (2013). Building the British Sign Language Corpus. *Language Documentation and Conservation* 7, 136-154. doi: <http://hdl.handle.net/10125/4592>

49



THANK YOU very much.



Any Questions?

50

코퍼스 기반 한국어와 한국수어 어휘 비교 분석

2023. 10. 7

장세은 • 박호민
(한국해양대학교)

한국코퍼스언어학회 가을학술대회, 서울, 연세대학교

참고: 본 발표는 장세은, 박호민 (2022, 2023)의 발표를 요약, 정리, 수정 보완하였음

1

Abstract and Key words

- The purpose of the present study is to compare a spoken genre of Korean National Corpus (Corpus of Modu – spoken corpus version 1.2 – released on December 1, 2021) compiled by the National Institute of Korean Language (NIKL) with the general vocabulary of Korean Sign Language (KSL) compiled by Chosun University based on the lemma classification of signs provided in the KSL online dictionary made by NIKL, and then to suggest the desired vocabulary to be used in spoken language textbooks of Korean and English for deaf students from the first to the ninth grade. To do this, we will first make a word list tagged with four parts of speech (POS) such as nouns, verbs, adjectives, and adverbs. Second, we will illustrate the first through the 18th 100 signs using both cumulative coverage and balanced numbers of each POS. The results show that the first 100 signs are the most frequent words, accounting for about 41.3% coverage of all 17,652 signs tagged with four parts of speech, the most frequent 1000 signs account for 77.5%, and the most frequent 1800 signs account for about 84%. Third, this classification is compared with the KSL online dictionary made by NIKL. Finally, we will point out some problems that need to be considered in general vocabulary research for sign language when deaf students learn spoken languages such as Korean and English by using KSL as their native language. It is hoped that this KSL general vocabulary research can contribute to the development of language textbooks for deaf elementary and middle school students.
- Key words: Corpus of Modu – spoken, deaf, general vocabulary, Korean Sign Language, high frequency vocabulary, accumulated coverage, language textbooks, elementary school, middle school

2

Contents

1. Introduction (purpose/research questions)
2. Literature Review
3. Data and Methodology
4. Results and Discussion
5. Conclusion and More Issues to be Considered

References

3

1. Introduction

■ The purpose of the current study :

- (i) to compare a spoken genre of Korean National Corpus (Corpus of Modu – spoken corpus version 1.2 – released on December 1, 2021) compiled by the National Institute of Korean Language (NIKL) with the general vocabulary of Korean Sign Language (KSL) compiled by Chosun University (called ‘Chosun KSL Corpus’) based on the lemma classification of signs provided in the KSL online dictionary made by NIKL, and
- (ii) then to suggest the desired vocabulary to be used in spoken language textbooks of Korean and English for deaf students from the first to the ninth grade.

4

1. Introduction (Cont'd)

To do this,

- (1) Investigating how to make the KSL online dictionary made by the National Institute of Korean Language (NIKL), especially focusing on lemma classification.
- (2) making a word list tagged with four parts-of-speech (POS) such as nouns, verbs, adjectives and adverbs of Chosun KSL Corpus.
- (3) making a comparison of the KSL online dictionary and the Spoken Genre of Modu Korean Corpus compiled by NIKL with Chosun KSL Corpus.
- (4) illustrating the **first through the 18th 100 signs** (100 signs x 9 grades x 2 semesters) **using both cumulative coverage and balanced numbers of each POS** for the nine grades from six-year elementary school to three-year middle school. The results show that the first 100 signs are the most frequent words accounting for around 41.3% coverage of all the 17,652 signs tagged with four parts-of-speech, the most frequent 1000 words account for 77.5%, and the most frequent 1800 words account for around 84%.
- (5) pointing out some problems to be considered in general vocabulary research for sign language when deaf students learn spoken languages such as Korean and English by using KSL as their native language.

1. Introduction (Cont'd)

■ An Expected Effect

It is hoped that this KSL general vocabulary research can contribute to developing language textbooks for deaf elementary and middle school students.

■ Research Questions

- (1) How many words should there be in the general vocabulary of Korean Sign Language?
- (2) What words are appropriate to create spoken language textbooks of Korean and English for the first to ninth grade deaf students?

2. Literature Review

2.1 Spoken Data

Q: How many words should comprise general vocabulary in English as a spoken language?

- Nation (2001)

- A traditional way to categorize vocabulary has been as **general** (i.e. high-frequency vocabulary), **academic vocabulary**, **technical vocabulary**, and **low-frequency vocabulary**.

- Using this categorization, general vocabulary has typically been set at around **2000 word families**. This is the number of word families in the original **General Service List (GSL)**, which is the foundation for the **Academic Word List (AWL)**, which remains the most commonly used list of academic words. While more recent general word lists appear to contain more words (the NGSL has 2801 words which comprise the most important high-frequency words in English), these are lemma-based lists, i.e. words and inflected forms, not word families. The NGSL has 2368 word families close to the 2000 word family range.

7

Key issues: Counting units

□ Types

- Each different word is counted as a type in a corpus.
- Types are suggested to word lists for beginners (spelling, recognizing words) (Carroll et al., 1971; Zeno et al. 1995).
- Word types are suggested in technical word lists (Nation, 2016).

□ Lemmas

- A lemma is a base word with part of speech and its form with inflectional suffixes.
- Lemmas (and above Level 2 word families) are suggested for receptive knowledge (reading and listening) (Nation, 2016).
 - ✓ Display the important information for vocabulary learning--the part of speech.
 - ✓ Allow learners to clearly focus on the vocabulary items in order to ease the burden. (Gardner & Davies, 2014)

8

□ Word families

- A word family is a base word with its inflectional forms and derivational forms.
- Bauer and Nation (1993) Word family levels.
- Learners have difficulties in derivational forms (Level 6 word families) (Schmitt & Zimmerman, 2002; Nagy & Tomsend, 2012; McLean, 2018).

Influential general and academic vocabulary lists

| Vocabulary types | Vocabulary lists | Brief description |
|----------------------------|---|---------------------|
| General vocabulary | General Service List, GSL (West, 1953) | 2,000 word families |
| | New General Service List, New-GSL (Brezina & Gablasova, 2015) | 2,494 lemmas |
| Academic vocabulary | Academic Word List, AWL (Coxhead, 2000) | 570 word families |
| | Academic Vocabulary List, AVL (Gardner & Davis, 2014) | 3,015 lemmas |
| | Academic Spoken Word List, ASWL (Dang et al., 2017) | 1,741 word families |

Lexical coverage

- One of the topics in second language acquisition that attract researchers' interests is **how many words indeed should be familiar with by the learners to successfully understand a text.**
- Based on the previous studies (e.g. Liu & Nation, 1985; Laufer, 1989; Laufer & Ravenhorst-Kalovski, 2010) and combining the results, Laufer and Ravenhorst-Kalovski (2010) suggested the two lexical coverage of 95% and 98% may represent the possible lower and upper boundaries, that is, **95% is acceptable for a minimal comprehension of texts, and 98% is preferred for an optimal comprehension.**

11

Some studies of vocabulary size needed for reasonable comprehension

| Genres | 95% coverage | 98% coverage | Representative Studies |
|-----------------------|--------------|--------------|------------------------|
| Novels and news | 4,000 | 8,000~9,000 | Nation (2006) |
| TV programs | 4,000 | 8,000~9,000 | Webb & Rodgers (2009a) |
| Movies | 3,000 | 6,000 | Webb & Rodgers (2009b) |
| Business textbooks | 3,500 | 5,000 | Hsu (2011) |
| Engineering textbooks | 5,000 | 10,000 | Hsu (2014) |

12

2.2 Sign Language Data

JSL Colloquial Corpus

<http://research.nii.ac.jp/jsl-corpus/public/en/>



This page is about Corpus Project in Colloquial Japanese Sign Language. Corpus is a database that consists of many examples of the language. Our research team provides linguistic descriptions for each example in the database to facilitate linguistic analysis.

It is known that Japanese Sign Language has many varieties depending on areas. How are they different? We visit many parts of Japan to collect regional varieties of Japanese Sign Language.

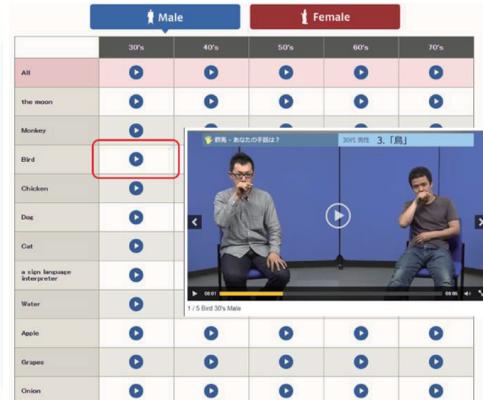
Click a region on this map

Regions shown include: Hokkaido, Aomori, Akita, Iwate, Miyagi, Tohoku, Niigata, Fukushima, Ibaraki, Saitama, Chiba, Kanto, Tokyo, Yamanashi, Shizuoka, Kanagawa, Central Japan, Gifu, Mie, Kyoto, Hyogo, Nara, Wakayama, Kinki, Osaka, Nara, Kyoto, Hyogo, Shikoku, Kochi, Tokushima, Ehime, Kyushu, Miyazaki, Kagoshima, Okinawa.

2011.04 ~ 2014 동안 연구를 위해 제작된 일본 수어 말뭉치

British Sign Language Corpus의 수석 제작자인 Adam Schembri 교수의 Corpus 구축 방법론을 조언받아 제작 (일본 국립 정보학 연구소)

일본 내 지역별, 연령별, 성별 수어 대화 영상을 촬영하여 ELAN 프로그램으로 Annotation을 진행하였음



13

100+ First ASL words

Learn the first 100 common signs to help you get started conversing basics quickly. But, that's not enough? Try the next second 100 common signs. And so on. But, don't forget grammar. See 'Learn' for tutorials in grammar, sentence structure, culture, etc.

| | | |
|-----------------|------------|-----------|
| First 100 words | Second 100 | Third 100 |
|-----------------|------------|-----------|

First 100 words.

| | | | | |
|-----------------|--------------------|--------------|------------------|----------------|
| 1. again | 22. finish | 43. little | 64. school | 83. understand |
| 2. also | 23. food | 44. live | 65. see | 84. use |
| 3. ask | 24. for | 45. man | 66. she | 85. wait |
| 4. because | 25. forget | 46. many | 67. should | 86. want |
| 5. boy | 26. friend | 47. me | 68. sign, signed | 87. what |
| 6. but | 27. get | 48. meet | 69. slow | 88. when |
| 7. can | 28. girl | 49. more | 70. some | 89. where |
| 8. come | 29. give | 50. my | 71. sorry | 90. which |
| 9. condone | 30. go | 51. name | 72. store | 91. who |
| 10. deaf | 31. good | 52. need | 73. take | 92. why |
| 11. different | 32. have | 53. new | 74. tell | 93. will |
| 12. drink | 33. he | 54. no | 75. text, sms | 94. with |
| 13. drive | 34. hearing | 55. not | 76. thank, thank | 95. woman |
| 14. eat | 35. hello | 56. now | 77. their | 96. work |
| 15. email | 36. help | 57. ok, okay | 78. they | 97. write |
| 16. family | 37. home | 58. old | 79. think | 98. yes |
| 17. feel | 38. how | 59. other | 80. time | 99. you |
| 18. few | 39. Internet | 60. please | 81. tired | 100. your |
| 19. find | 40. know | 61. remember | 82. try | |
| 20. fine | 41. later | 62. same | | |
| 21. fingerspell | 42. like (feeling) | 63. say | | |

The word list is subject to change over time. It doesn't have to be exactly 100, but it can be a bit more or less 100.

14

- 69 -

3. Data and Methodology

3.1 KSL Online Dictionary (국립국어원 한국수어사전)

- 한국어 대응표현
- 한국수어 lemma (표제어) 웹사전제공
→ 3,676개 (7,112 개 types)
- 한국수어 표제어 등록수 → 24,705개(??)

(예)

- 일련번호 1~ 20 예 (우측표)

| 번호 | 수어 표제어 번호 | 한국어 대응표현 |
|----|-----------|--|
| 1 | 24705 | 복직, 복임 |
| 2 | 24704 | 이론, 쓰다 |
| 3 | 24701 | (맛이)짜다 |
| 4 | 24049 | 활까봐 |
| 5 | 24047 | (마음이)좀다, 옹졸하다, 용하다 |
| 6 | 24041 | (바지를)입다 |
| 7 | 24037 | 인도(중화국) |
| 8 | 24036 | 아니다, 없다 |
| 9 | 24034 | 여자, 여성, 여 |
| 10 | 24033 | 설거지 |
| 11 | 24032 | 베부르다 |
| 12 | 24030 | (여자를)유혹, 유혹하다 |
| 13 | 24029 | 삼, 셋 |
| 14 | 24028 | 시간, 한 시간 |
| 15 | 24026 | 놓아주다 |
| 16 | 24024 | 풀, 여리, 어리가지 |
| 17 | 24022 | 시작, 개시, 개최, 거행, 시발, 착수, 훌발, 열다, 이행, 하수 |
| 18 | 24020 | 틈(바닥에 생기는)균열 |
| 19 | 24019 | 참다, 견디다, 인내, 자제, 억제 |
| 20 | 24015 | 의리 |

| 1 | 수어 표제어 번호 | WORD |
|------|-----------|---------|
| 2 | 1 | 파나마 |
| 3 | 10 | 구주 |
| 4 | 10 | 유럽 |
| 5 | 11 | -국 |
| 6 | 11 | 국가 |
| 7 | 11 | 나라 |
| 8 | 12 | 국가시험 |
| 9 | 12 | 국가 고시 |
| 10 | 16 | 군집 |
| 11 | 16 | 떼 |
| 12 | 16 | 무리 |
| 13 | 20 | 궁색 |
| 14 | 22 | 당부 |
| 15 | 22 | 부탁 |
| 16 | 22 | 요청 |
| 17 | 22 | 청하다 |
| 18 | 22 | 요구 |
| 19 | 23 | 권위 |
| 20 | 24 | 권위자 |
| 21 | 28 | 귀여워하다 |
| 22 | 28 | 귀엽다 |
| 7105 | 24047 | (마음이)좁다 |
| 7106 | 24047 | 옹졸하다 |
| 7107 | 24047 | 옹하다 |
| 7108 | 24049 | 활끼보 |
| 7109 | 24701 | (맛이)짜다 |
| 7110 | 24704 | 이론 |
| 7111 | 24704 | 쓰다 |
| 7112 | 24705 | 복직 |
| 7113 | 24705 | 복임 |

15

표제어(lemma) 설명



한국수어사전의 표제어로 시작

구축된 수어 말뭉치에서도
새로운 표제어를 제작 (한국수어사전 표제어에서 시작 후 발전)
3,676개에서 5,010개로 늘어남

예) '나', '나+', '나++', '나+++', '나++++', '내', '자기', '제' - 표제어 번호 973번

수어 말뭉치와 모두의 말뭉치(한국어 구어)와의 비교를 수행하기 위해
모두의 말뭉치에서도
새로운 표제어를 제작 (수어 말뭉치 표제어에서 시작 후 발전)
5,010개에서 18,138개로 늘어남

예) '솔찮혔하다', '솔찮해', '솔찮히', '솔찮이', '솔찮다', '솔찮하다', '솔찮혀' - 표제어 번호 6,751번

3.2 The Chosun KSL Corpus (조선대한국수어코퍼스)

3.2.1 자료제공자 (Informants)

- 연령층 남녀별: 평균 연령 49.61세 : 26 명 (남 16명, 여 10명)

- 지역별、주제별 통계

| 나이 | 참여인원 |
|-------------|-----------|
| 20대 | 1 |
| 30대 | 2 |
| 40대 | 12 |
| 50대 | 6 |
| 60대 | 3 |
| 70대 | 1 |
| 80대 | 1 |
| 총 인원 | 26 |

| 지역 | 세부 | 인원 | 지역별참여인원 |
|-------|----|----|---------|
| 전남 | 강진 | 1 | 23 |
| | 광양 | 2 | |
| | 광주 | 7 | |
| | 나주 | 1 | |
| | 남악 | 3 | |
| | 목포 | 3 | |
| | 무안 | 2 | |
| | 보성 | 1 | |
| | 순천 | 1 | |
| 전북 | 화순 | 2 | 2 |
| | 완주 | 1 | |
| 서울/경기 | 전주 | 1 | 1 |
| | 인천 | 1 | |

| 주제 | 참여인원 |
|-------------|------------|
| 인간 | 18 |
| 삶 | 24 |
| 식생활 | 18 |
| 의생활 | 14 |
| 주생활 | 14 |
| 사회생활 | 18 |
| 경제생활 | 14 |
| 교육 | 12 |
| 종교 | 12 |
| 문화 | 18 |
| 정치와 행정 | 14 |
| 자연 | 18 |
| 동식물 | 18 |
| 개념 | 12 |
| 국가명 및 지명 | 16 |
| 기타 | 18 |
| 누적인원 | 258 |

17

3.2.2 Basic Information about the Chosun KSL Corpus

- The Planned Chosun KSL Corpus in The Five-year Institute Project of Research Foundation of Korea

- A small corpus of Chosun KSL was compiled from December 2020 to April 2022 and was transcribed from January 2021 to May 2022 by the Language Convergence Institute, Chosen University, Gwangju, South Korea, which has been supported by The Five-year Institute Project of Research Foundation of Korea in September 2020.

• Original Plan: It was firstly planned that the Chosun KSL Corpus would consist of **300,000 gloss words**.

- The Study Corpus

- The **study corpus** consists of 70 files containing **251,205 tokens** and **17,652 types** (after **retaining content words tagged with** target parts-of-speech(POS) including nouns, verbs, adjectives, adverbs, and others and deleting other POS from the second corpus) as in Table 1.
 → **5,010 lemmas & 17,149 sentences** in the translated Korean version
 → This study corpus was analyzed to categorize the first to 18th 100 words **for this current study**

18

3.2.3 전사 파일 리스트 (70개)

풀더명 (활영일) 파일명

| | |
|--------|--|
| 210413 | MP,MA_015,016,GJ,MA_40,40,M,M_2,1,6 |
| 210413 | MP,MA_015,016,GJ,MA_40,40,M,M_2,2,3,16 |
| 210413 | MP,MA_015,016,GJ,MA_40,40,M,M_2,6,11,1,3 |
| 210413 | MP,MA_015,016,GJ,MA_40,40,M,M_2,9,2,2 |
| 210413 | MP,MA_015,016,GJ,MA_40,40,M,M_2,9,12,1,8 |
| 210413 | MP,MA_015,016,GJ,MA_40,40,M,M_2,9 |
| 210413 | MP,MA_015,016,GJ,MA_40,40,M,M_2,16a |
| 210413 | MP,MA_015,016,GJ,MA_40,40,M,M_2,16 |
| 210414 | WJ,GS_017,018,JJJ,JJJ_40,40,FM_2,1,6,11 |
| 210414 | WJ,GS_017,018,JJJ,JJJ_40,40,FM_2,1,6 |
| 210414 | WJ,GS_017,018,JJJ,JJJ_40,40,FM_2,1,9 |
| 210414 | WJ,GS_017,018,JJJ,JJJ_40,40,FM_2,2,3,16 |
| 210414 | WJ,GS_017,018,JJJ,JJJ_40,40,FM_2,9,12,18 |
| 210414 | WJ,GS_017,018,JJJ,JJJ_40,40,FM_2,9,16 |
| 210414 | WJ,GS_017,018,JJJ,JJJ_40,40,FM_2,16b |
| 210414 | WJ,GS_017,018,JJJ,JJJ_40,40,FM_2,16d |
| 210421 | MP,NJ_014,019,DJ,NJ_50,60,F,M_2a |
| 210421 | MP,NJ_014,019,DJ,NJ_50,60,F,M_2b |
| 210421 | MP,NJ_014,019,DJ,NJ_50,60,F,M_2c |
| 210421 | MP,NJ_014,019,DJ,NJ_50,60,F,M_2d |
| 210426 | NANA_020,021,JD,GM_40,40,M,F_2a |
| 210426 | NANA_020,021,JD,GM_40,40,M,F_2b |
| 210426 | NANA_020,021,JD,GM_40,40,M,F_2c |
| 210426 | NANA_020,021,JD,GM_40,40,M,F_2d |
| 210426 | NANA_020,021,JD,GM_40,40,M,F_2e |
| 210426 | NANA_020,021,JD,GM_40,40,M,F_2f |
| 210427 | NA,GJ_005,006,MP,GN_30,40,M,M_2a |
| 210427 | NA,GJ_005,006,MP,GN_30,40,M,M_2b |
| 210427 | NA,GJ_005,006,MP,GN_30,40,M,M_2c |
| 210427 | NA,GJ_005,006,MP,GN_30,40,M,M_2d |
| 210427 | NA,GJ_005,006,MP,GN_30,40,M,M_2e |
| 210427 | NA,GJ_005,006,MP,GN_30,40,M,M_2f |
| 210427 | NA,GJ_005,006,MP,GN_30,40,M,M_2g |
| 210427 | NA,GJ_005,006,MP,GN_30,40,M,M_2h |

The corpus contains 70 files from 210413 to 220406.

| | |
|--------|---|
| 210428 | HS,HS_022,023,GY,PS_40,40,M,M_2,16a |
| 210428 | HS,HS_022,023,GY,PS_40,40,M,M_2,16b |
| 210428 | HS,HS_022,023,GY,PS_40,40,M,M_2,16d |
| 210428 | HS,HS_022,023,GY,PS_40,40,M,M_2,16e |
| 211123 | GJ,GJ_034,035,GJ,HS_60,60,M,M_3,1,2 |
| 211123 | GJ,GJ_034,035,GJ,HS_60,60,M,M_3,3,4,5,6 |
| 211123 | GJ,GJ_034,035,GJ,HS_60,60,M,M_3,7,8,9 |
| 211123 | GJ,GJ_034,035,GJ,HS_60,60,M,M_3,10,11,12 |
| 220107 | GJ,MP_038,039,HP,MP_30,20,FF_3,1,2,3,4,5 |
| 220107 | GJ,MP_038,039,HP,MP_30,20,FF_3,6,7,8,9 |
| 220107 | GJ,MP_038,039,HP,MP_30,20,FF_3,10,11,12,13 |
| 220107 | GJ,MP_038,039,HP,MP_30,20,FF_3,14,15,16 |
| 220114 | MU,INC_040,041,MU,IM,M,F40,50,3,1,2,3,4,5 |
| 220114 | MU,INC_040,041,MU,IM,M,F40,50,3,6,7,8,9 |
| 220114 | MU,INC_040,041,MU,IM,M,F40,50,3,10,11,12,13 |
| 220114 | MU,INC_040,041,MU,IM,M,F40,50,3,14,15,16 |
| 221224 | GJ,GJ_036,037,GH,GI_70,80,M,M_3,1,2,3,4 |
| 221224 | GJ,GJ_036,037,GH,GI_70,80,M,M_3,3,5,6,7,8 |
| 221224 | GJ,GJ_036,037,GH,GI_70,80,M,M_3,9,10,11,12 |
| 221224 | GJ,GJ_036,037,GH,GI_70,80,M,M_3,13,14,15,16 |
| 220317 | SC,BS_042,043,CW,PS_50,50,FF,3,1,2,3,4,5 |
| 220317 | SC,BS_042,043,CW,PS_50,50,FF,3,6,7,8,9,10 |
| 220317 | SC,BS_042,043,CW,PS_50,50,FF,3,11,12,13 |
| 220317 | SC,BS_042,043,CW,PS_50,50,FF,3,14,15,16 |
| 220406 | GY,GY_044,045,GY,YS_40,40,M,M_3,1,2,3,4,5 |
| 220406 | GY,GY_044,045,GY,YS_40,40,M,M_3,6,7,8,9 |
| 220406 | GY,GY_044,045,GY,YS_40,40,M,M_3,10,11,12 |
| 220406 | GY,GY_044,045,GY,YS_40,40,M,M_3,13,14,15,16 |

19



- Naming of Video + Annotation Files

Current staying region_Informants number_Birth place_Age group_Sex_Video topics

WJ,GS_017,018,JJ,JJ_40,40,FM_3-3,3-4,4-1,5-1

(1) (2) (3) (4) (5) (6)

(1) WJ : Wan-Ju(완주), GS : Gun-San(군산)

(2) 017, 018: Informants

(3) JJ : Birth place (Jeon-Ju, 전주)

(4) 40: forties year old

(5) Female, Male: Sex

(6) Topics(주제)

3-3 (Food Life: most favorite food that you recently ate),

3-4 (Food Life: any food that you want to eat in future),

4-1 (Clothing Life: Tell me how to purchase clothes and shoes),

5-1 (Housing Life: your preferred place and time to meet with you're a few)

*Currently 16 major categories , 88 subcategories, 199 video topics

20

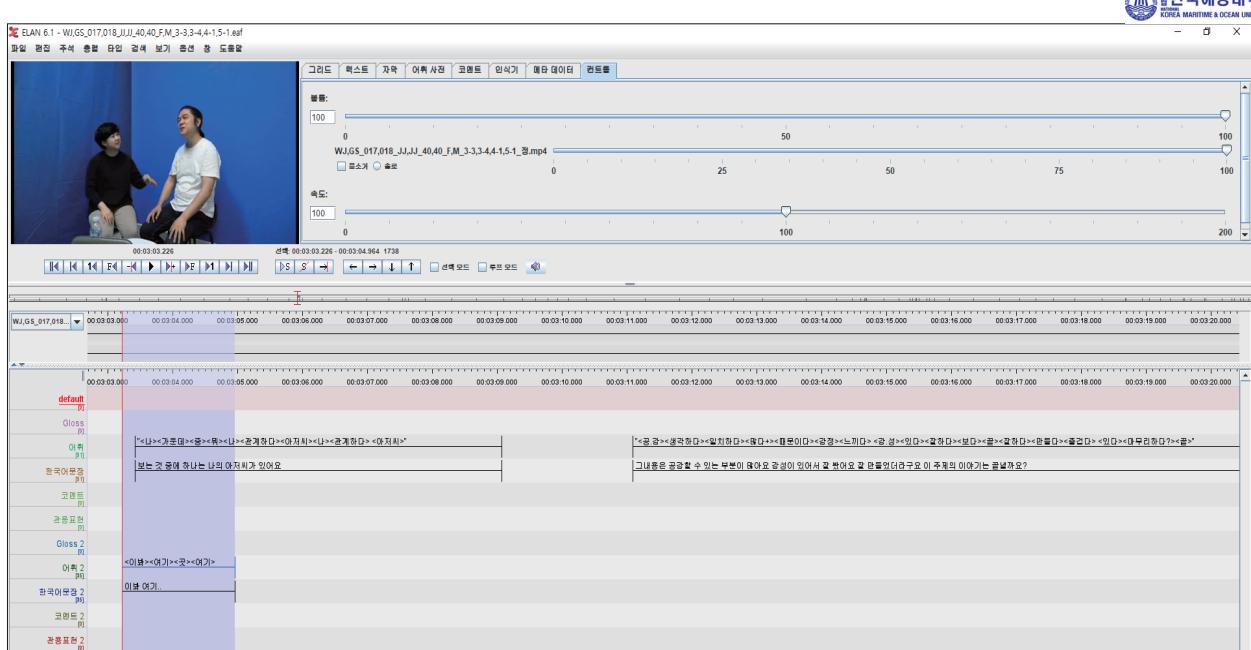
3.2.4 Description of Video Clips + ELAN

Transcription and Analysis Methods

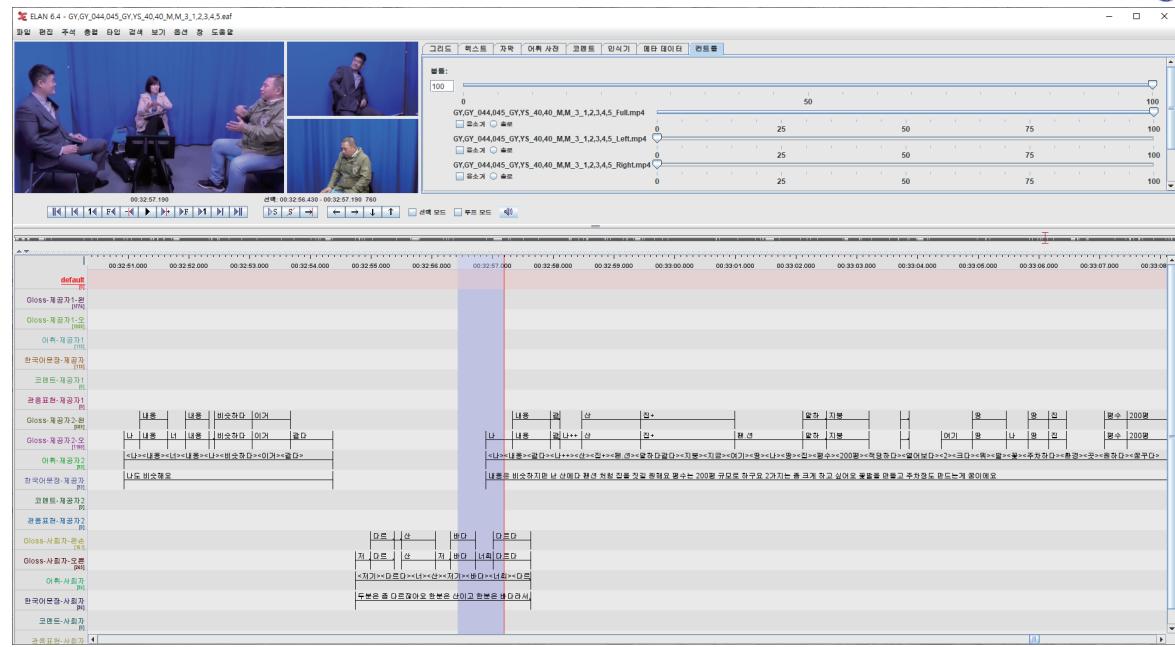
Using the ELAN (EUDICO Linguistic Annotator) program to analyze the recorded video data for precise transcription and phrase glosses, we built the "Chosun University KSL Corpus" by annotating the lexical signals used for individual vocabulary and multilingual vocabulary containing idiomatic vocabulary, and annotating the POS tagset of the Sejong Corpus.

21

Video Clips + ELAN



22



23

Table 1. Top 20 and Bottom 10 words in the study corpus

| Top 20 in The Study Corpus | | | | | | Bottom 10 in The Study Corpus | | | | | |
|----------------------------|-------|-----|-------|----------|-------------|-------------------------------|---------|-----|---|--------|----------|
| 1 | Gloss | POS | types | types(%) | types(CUM%) | | | | | | |
| 2 | 나 | NP | 7,452 | 2.9665 | 2.9665 | 17644 | 전부/전체 | MAG | 1 | 0.0004 | 99.9964 |
| 3 | 저기 | NP | 4,571 | 1.8196 | 4.7861 | 17645 | 가웃뚱 | MAG | 1 | 0.0004 | 99.9968 |
| 4 | 같다 | VA | 3,816 | 1.5191 | 6.3052 | 17646 | 전.영.수 | MAG | 1 | 0.0004 | 99.9972 |
| 5 | 있다 | VA | 3,244 | 1.2914 | 7.5966 | 17647 | 김.영.란 | MAG | 1 | 0.0004 | 99.9976 |
| 6 | 맞다 | VV | 3,182 | 1.2667 | 8.8633 | 17648 | ~만/~뿐/△ | MAG | 1 | 0.0004 | 99.9980 |
| 7 | 없다 | VA | 2,820 | 1.1226 | 9.9859 | 17649 | 전.영. | MAG | 1 | 0.0004 | 99.9984 |
| 8 | 끌 | NNG | 2,610 | 1.0390 | 11.0249 | 17650 | 더. | MAG | 1 | 0.0004 | 99.9988 |
| 9 | 너 | NP | 2,559 | 1.0187 | 12.0436 | 17651 | 억지로/강 | MAG | 1 | 0.0004 | 99.9992 |
| 10 | 좋다 | VA | 2,459 | 0.9789 | 13.0224 | 17652 | 영.화.' | MAG | 1 | 0.0004 | 99.9996 |
| 11 | 보다 | MAG | 2,399 | 0.9550 | 13.9774 | 17653 | '쿵.후.' | MAG | 1 | 0.0004 | 100.0000 |
| 12 | 많다 | VA | 2,004 | 0.7978 | 14.7752 | | | | | | |
| 13 | ~면 | NNG | 1,754 | 0.6982 | 15.4734 | | | | | | |
| 14 | 집 | NNG | 1,532 | 0.6099 | 16.0833 | | | | | | |
| 15 | 생각하다 | VV | 1,470 | 0.5852 | 16.6685 | | | | | | |
| 16 | 말하다 | VV | 1,425 | 0.5673 | 17.2357 | | | | | | |
| 17 | 때 | NNG | 1,376 | 0.5478 | 17.7835 | | | | | | |
| 18 | 잘하다 | VV | 1,374 | 0.5470 | 18.3304 | | | | | | |
| 19 | 농인 | NNG | 1,370 | 0.5454 | 18.8758 | | | | | | |
| 20 | 안하다 | VV | 1,314 | 0.5231 | 19.3989 | | | | | | |
| 21 | 무엇 | NP | 1,311 | 0.5219 | 19.9208 | | | | | | |

251,205 tokens and 17,652 types

24

3.3 The Corpus of Modu – spoken corpus version 1.2

a spoken genre of Korean National Corpus (Corpus of Modu – spoken corpus version 1.2 – released on December 1, 2021) compiled by the National Institute of Korean Language (<https://corpus.korean.go.kr>)

3.3.1 Compilation of the Corpus of Spoken Korean (한국어 구어 코퍼스 구성)

Monologue (독백) 2,490건,
Conversation (대화) 19,104건,
Drama Script (드라마 대본) 4,102건,
File number 25,696개 (6.73GB)

3.3.2 Sejong Tagset (세종 말뭉치 태그셋)

| 대분류 | 소분류 | 세분류 |
|----------|---------|--|
| (1) 체언 | NN 명사 | NNG 일반명사 NNP 고유명사 NNB 의존명사 |
| | NP 대명사 | NP 대명사 |
| | NR 수사 | NR 수사 |
| (2) 용언 | VV 동사 | VV 동사 |
| | VA 형용사 | VA 형용사 |
| | VX 보조용언 | VX 보조용언 |
| | VC 지정사 | VCP 꿈정지정사 VCN 부정지정사 |
| (3) 수식언 | MM 관형사 | MAG 일반부사 MAJ 접속부사 |
| | MA 부사 | |
| (4) 독립언 | IC 감탄사 | IC 감탄사 |
| (5) 관계언 | JK 격조사 | JKS 주격조사 JKC 보격조사 JKG 관형격조사 JKO 목적격조사 JKB 부사격조사 JKV 호격조사 JKQ 인용격조사 |
| | | JX 보조사 |
| | | JC 접속조사 |
| | | EP 선어말어미 EF 종결어미 EC 연결어미 |
| | | ETN 명사형전성어미 ETM 관형형전성어미 |
| | | XP 체언접두사 |
| | | XSN 명사파생접미사 |
| | | XSV 동사파생접미사 |
| | | XSA 형용사파생접미사 |
| | | XR 어근 |
| (6) 의존형태 | E 어미 | XR 어근 |
| | | |
| | | |
| | | |
| | | |

3.3.3 Matching & generating lemmas and sign corpora (수어 말뭉치와 표제어 매칭 & 생성)

The process of making a corpus compatible with a corpus analysis (수어 말뭉치 분석 자료와 호환시켜주기 위한 과정)

1. Combining linguistic rule enforcement (언어적 규칙 적용 합치기: NNG + XSV = VV, NNG + XSA = VA, NNG + NNG = NNG 등)
2. Combining the introduction of sign language corpus lemmas (수어 말뭉치 표제어 도입 합치기 5,010개)
+ Create new lemmas in the same way (같은 방법으로 새 표제어 추가(13,129개))

```
<ANNOTATION>
<ALIGNABLE_ANNOTATION ANNOTATION_ID="a14412">
  <TIME_SLOT_REF1>t52362</TIME_SLOT_REF1>
  <TIME_SLOT_REF2>t52364</TIME_SLOT_REF2>
</ALIGNABLE_ANNOTATION>
<ANNOTATION>
<ALIGNABLE_ANNOTATION ANNOTATION_ID="a14414">
  <TIME_SLOT_REF1>t52366</TIME_SLOT_REF1>
  <TIME_SLOT_REF2>t52368</TIME_SLOT_REF2>
</ALIGNABLE_ANNOTATION>
<ANNOTATION>
<ALIGNABLE_ANNOTATION ANNOTATION_ID="a14416">
  <TIME_SLOT_REF1>t52374</TIME_SLOT_REF1>
  <TIME_SLOT_REF2>t52374</TIME_SLOT_REF2>
</ALIGNABLE_ANNOTATION>
<ANNOTATION>
<ALIGNABLE_ANNOTATION ANNOTATION_ID="a14418">
  <TIME_SLOT_REF1>t52371</TIME_SLOT_REF1>
  <TIME_SLOT_REF2>t52712</TIME_SLOT_REF2>
</ALIGNABLE_ANNOTATION>
<ANNOTATION>
<ALIGNABLE_ANNOTATION ANNOTATION_ID="a14419">
  <TIME_SLOT_REF1>t52713</TIME_SLOT_REF1>
  <TIME_SLOT_REF2>t52718</TIME_SLOT_REF2>
</ALIGNABLE_ANNOTATION>
<ANNOTATION>
<ALIGNABLE_ANNOTATION ANNOTATION_ID="a14422">
  <TIME_SLOT_REF1>t52760</TIME_SLOT_REF1>
  <TIME_SLOT_REF2>t52768</TIME_SLOT_REF2>
</ALIGNABLE_ANNOTATION>
<ANNOTATION>
<ALIGNABLE_ANNOTATION ANNOTATION_ID="a14424">
  <TIME_SLOT_REF1>t52770</TIME_SLOT_REF1>
  <TIME_SLOT_REF2>t52772</TIME_SLOT_REF2>
</ALIGNABLE_ANNOTATION>
<ANNOTATION>
<ALIGNABLE_ANNOTATION ANNOTATION_ID="a14426">
  <TIME_SLOT_REF1>t52773</TIME_SLOT_REF1>
  <TIME_SLOT_REF2>t52776</TIME_SLOT_REF2>
</ALIGNABLE_ANNOTATION>
<ANNOTATION>
<ALIGNABLE_ANNOTATION ANNOTATION_ID="a14429">
  <TIME_SLOT_REF1>t52781</TIME_SLOT_REF1>
  <TIME_SLOT_REF2>t52782</TIME_SLOT_REF2>
</ALIGNABLE_ANNOTATION>
<ANNOTATION>
<ALIGNABLE_ANNOTATION ANNOTATION_ID="a14431">
  <TIME_SLOT_REF1>t52784</TIME_SLOT_REF1>
  <TIME_SLOT_REF2>t52786</TIME_SLOT_REF2>
</ALIGNABLE_ANNOTATION>
<ANNOTATION>
```

4. Results and Discussion

4.1 Cumulative coverage percentages in the Chosun KSL Corpus

| Cumulative % | 80% | 85% | 90% | 100% |
|--------------------------------|-----------------------|----------------------|---------------------|-------------------------|
| All POS | 1,367 / 219,127 (23)* | 2,293 / 233,313 (11) | 4,202 / 246,499 (5) | 19,183 / 273,086 |
| Content words (N, V, Adj, Adv) | 1,238 / 200,998 (24) | 2,012 / 213,553 (12) | 3,839 / 226,575 (5) | 17,652 / 251,205 |
| Nouns | 1,047 / 88,898 (27) | 1,346 / 95,099 (14) | 2,080 / 101,963 (6) | 7,733 / 113,418 |
| Verbs | 792 / 68,141 (11) | 1,332 / 72,172 (6) | 2,418 / 75,101 (3) | 8,073 / 84,725 |
| Adjectives | 64 / 27,854 (85) | 92 / 29,713 (50) | 137 / 31,396 (27) | 1,142 / 34,852 |
| Adverbs | 29 / 14,481 (129) | 41 / 15,467 (55) | 71 / 16,379 (22) | 704 / 18,210 |

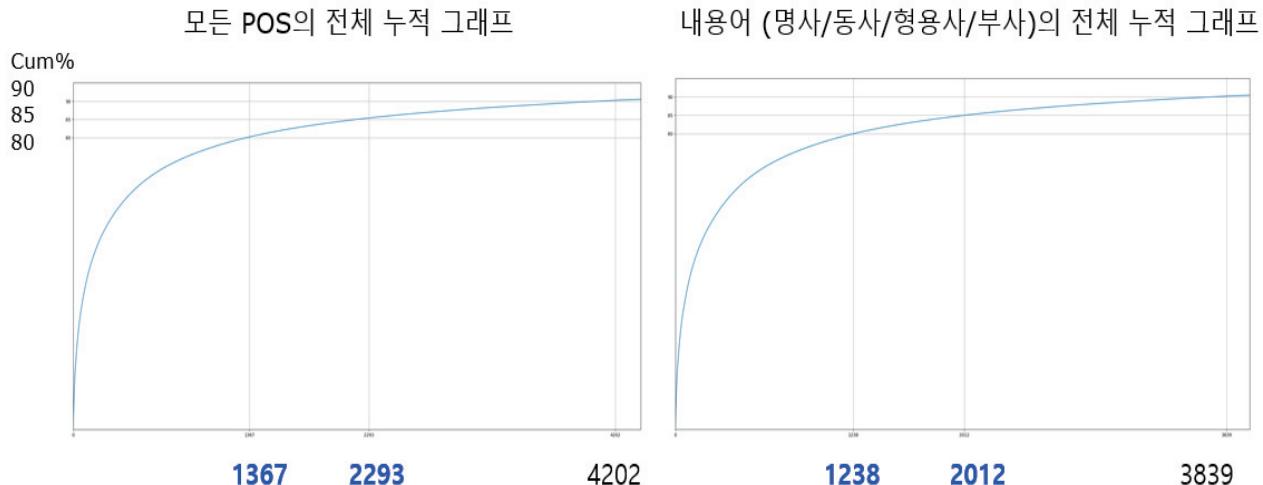
* Type number / Token number (cut-off minimum frequency)

17,652개의 내용어는 251,205번 빈도수를 보여주어 전체 19,183개의 단어가 273,086번의 빈도수를 가져 92% 차지하지만 1,531개의 기능어는 8% 차지함 (아래 참고를 고려하여 재분류한다면 실질적으로 기능어는 $8 + 7.7 + 2.2 = 17.9\%$ 이상이고 82.1% 정도가 내용어임)

참고1: 나(7,452), 너(2559), 당신(346), 그녀(105), 우리(66), 그남자(44), 그사람(8), 그여자(2) 등 인칭대명사가 10,582번의 빈도수를 가져 3.88%, 저기(4571), 여기(1,238), 저거(1220), 저것(160), 이거(1169), 이것(957), 그것(1021), 그것들(52), 등 10,388번의 빈도수를 가져 3.8%를 차지하여 인칭대명사(3.88%)와 비인칭대명사와 지시대부사(3.8%)는 7.7%를 차지함. → 기능어이지만 POS로는 내용어로 분류했음

참고2: 의문사(뭐 2,569; 무엇 1311; 방법 927; 왜 669; 누구 272; 언제 171; 어떻게 45; 누가 11) : POS에서는 내용어로 분류
총 5930 빈도수 → 2.2%

Graphs of cumulative coverage in the Chosun KSL Corpus



29

4.2 Comparison of types between the KSL Online Dictionary and the Chosun KSL Corpus

| | Chosen KSL corpus | KSL online dictionary compiled by the National Institute of Korean Language | Rate (%) |
|--------------------------------|-------------------|---|----------|
| All POS | 19,183 | 9,336 | 48.7 |
| Content words (N, V, Adj, Adv) | 17,652 | 9,003 | 51.0 |
| Nouns | 7,733 | 4,231 | 54.7 |
| Verbs | 8,073 | 3,541 | 43.9 |
| Adjectives | 1,142 | 768 | 67.3 |
| Adverbs | 704 | 436 | 61.9 |

아래의 경우를 고려한다면 비교 비율 값은 훨씬 더 높아짐

- (1) 동사/명사 비구별: (생각)하다, (일)하다, (수어)하다, (대화)하다, (사용)하다, (여행)하다, (경험)하다, (결혼)하다, (사랑)하다, (발전)하다, (공부)하다, (계속)하다, (비교)하다, (운동)하다, (질문)하다, (구경)하다, (요리)하다, (활동)하다, (운전)하다, (증가)하다, (감시)하다, (선택)하다 등 동사/명사를 구별하지 않고 같은 sign으로 본다면 위의 비교 비율은 훨씬 더 높아짐.
- (2) 차=자동차, 청인=건청인 등 유사어를 같은 sign으로 봄
- (3) 앞으로 = 앞, 후에 = 후, 때문에 = 때문 등 후치사가 없는 명사와 같은 sign으로 봄
- (4) 더하다 = 합하다, ○ㅋ/OK=오케이 등 의미상으로 같은 단어로 쓰이는 sign으로 봄

30

Ex) Comparison of Top 200 Chosun KSL Corpus vs. KSL Online Dictionary

| The First 100 signs | Chosen KSL corpus | KSL online dictionary compiled by the National Institute of Korean Language | Rate (%) |
|---------------------|-------------------|---|----------|
| Nouns | 47 | 42 (저기, 저거, 얼굴이름. (건)청인, 선생 님, 거기 , (승용) 차) | 89.4 |
| Verbs | 39 | 38 (생각)하다, 안하다 , (일)하다, (수어)하다, (대화)하다, (사용)하다, (여행)하다, (경험)하다, (결혼)하다 | 97.4 |
| Adjectives | 10 | 9 (필요)하다, 즐겁다 , (가능)하다 | 90.0 |
| Adverbs | 6 | 5 (자연히) | 83.3 |

| The Second 100 signs | Chosen KSL corpus | KSL online dictionary compiled by the National Institute of Korean Language | Rate (%) |
|----------------------|-------------------|--|----------|
| Nouns | 47 | 43 (앞으로, 코로나, 이들, 후에) | 91.5 |
| Verbs | 35 | 32 (말라, (사랑)하다, 더하다 (= 합하다 ??), 마무리하다 , (발전)하다, (공부)하다, 불나다 , (계속)하다, (비교)하다, (운동)하다, (질문)하다, (구경)하다, (요리)하다, (활동)하다, (운전)하다, (증가)하다, (감시)하다, (선택)하다) | 91.4 |
| Adjectives | 10 | 10 (의미)하다, (불편)하다 | 100 |
| Adverbs | 6 | 5 (대하여) | 83.3 |

31

4.3 Comparison of Spoken Corpus of Modun Korean and the Chosun KSL Corpus

- 4.3.1 전사된 텍스트의 token과 type의 기본 정보 비교
- 4.3.2 수어 표제어 (lemma) 기준 상위 빈도수 존재 비율
- 4.3.3 전사된 텍스트의 품사별 token, type 비교
 - 4.3.3.1 품사별 Top 30 (수어, 모두) - 전체
 - 4.3.3.2 품사별 Top 30 (수어, 모두) - 서술어
 - 4.3.3.3 품사별 Top 30 (수어, 모두) - 일반명사
 - 4.3.3.4 품사별 Top 30 (수어, 모두) - 일반부사

4.3.4 Comparison of graphs in cumulative coverage between Spoken Corpus of Modun Korean and the Chosun KSL Corpus

32

4.3.1 전사된 텍스트의 token과 type의 기본 정보 비교

| 조선대KSL수어 말뭉치 | 모두의 말뭉치(구어) |
|----------------|--------------------------------|
| 269,586 tokens | 222,978,509 tokens (약 827.11배) |
| 14,026 types | 682,374 types (약 48.65배) |

STTR : 5.2 (약 17배)

조선대KSL수어 말뭉치

| 표제어 | POS | 모두의 말뭉치 | | 수어 말뭉치 | | | |
|-----|-------|---------|------------|--------|------|-------|--------|
| | | 등장횟수 | 비율 | 순위 | 등장횟수 | | |
| 나 | NP | 163 | 122,143 | 0.0548 | 1 | 7,594 | 2,8165 |
| 너 | NP | 349 | 45,874 | 0.0206 | 2 | 4,617 | 1,7126 |
| 같다 | PRED | 2 | 12,339,078 | 5.5338 | 3 | 4,429 | 1,6429 |
| 무슨 | 무엇+NP | 20 | 1,958,675 | 0.8784 | 4 | 4,013 | 1,4886 |
| 그만 | NNG | 64 | 44,0794 | 0.2005 | 5 | 3,462 | 1,2842 |
| 여기 | 여기+NP | 45 | 772,218 | 0.3463 | 6 | 3,451 | 1,2841 |
| 며칠 | PRED | 49 | 720,886 | 0.3233 | 7 | 3,380 | 1,2538 |
| 있다 | PRED | 5 | 6,497,368 | 2.9139 | 8 | 3,291 | 1,2208 |
| 없다 | PRED | 56 | 558,495 | 0.2505 | 9 | 2,927 | 1,0857 |
| 좋다 | PRED | 6 | 3,680,046 | 2.8613 | 10 | 2,631 | 0.9759 |

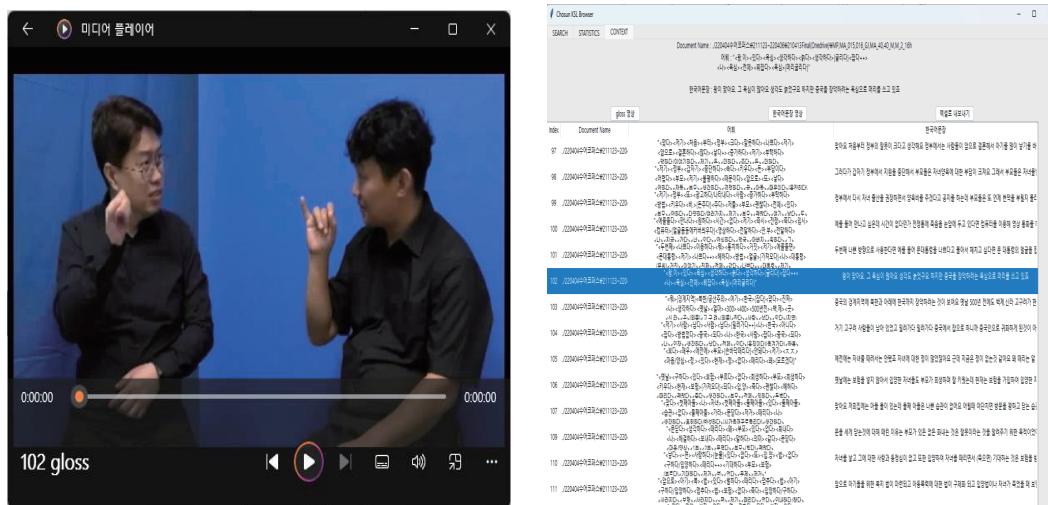
모두의 말뭉치(구어) Frequency 기준 TOP 10 비교

| | | 모두의 말뭉치 | | | 수어 말뭉치 | | |
|-----|-------------------------------|---------|------------|--------|--------|-------|--------|
| 표제어 | POS | 순위 | 등장횟수 | 비율 | 순위 | 등장횟수 | 비율 |
| 가능 | 가능하다+가능하다+가+PRED | 1 | 18,487,125 | 8,2910 | 75 | 596 | 0.2211 |
| 같다 | 같다~같다+같다+같다+-PRED | 2 | 12,339,078 | 5,5338 | 3 | 4,429 | 1.6429 |
| 다르다 | 다르다+다르다++다르다+-PRED | 3 | 11,581,005 | 5,1938 | 34 | 1,110 | 0.4117 |
| 충분 | 충분하다+충분하다+-PRED | 4 | 6,846,485 | 3,0714 | 335 | 143 | 0.0530 |
| 있다 | 있다+있나+있나+있다+있다+PRED | 5 | 6,497,368 | 2,9139 | 8 | 3,291 | 1.2208 |
| 좋다 | 좋다+좋다+좋다+좋다+좋다+PRED | 6 | 6,380,046 | 2,8613 | 10 | 2,631 | 0.9759 |
| 글운 | 글운+틈+틈++틈+틈NNNG | 7 | 5,033,463 | 2,2574 | 8804 | 2 | 0.0007 |
| 나태 | 나태+만으로나타나게+만으로나타나게+나태+나태+PRED | 8 | 3,893,896 | 1,7463 | 742 | 48 | 0.0178 |
| 비슷 | 비슷+비슷+비슷+비슷하+비슷하+PRED | 9 | 3,313,649 | 1,4861 | 186 | 271 | 0.1005 |
| 많다 | 많다+많다+많다+많다+많다+PRED | 10 | 3,241,156 | 1,4536 | 13 | 2,054 | 0.7619 |

4.3.2 수어 표제어 (5,010 lemmas) 기준 상위 빈도수 존재 비율

빈도수 상위 2,211 개 존재 (44.1%)

빈도수 상위 2,202 개 존재 (44.0%)



조선대KSL수어 말뭉치

모두의 말뭉치 (구어)

4.3.3 전사된 텍스트의 품사별 token, type 비교

**조선대KSL수어
말뭉치**

| POS | 서술어 | 수사 | 대명사 | 고유명사 | 일반명사 | 의존명사 | 접속부사 | 일반부사 | 긍정지정사 |
|---------------|-----------------|--------------|---------------|--------------|----------------|--------------|--------------|---------------|--------------|
| Type (ratio) | 7,395 (52.7%) | 97 (0.7%) | 87 (0.6%) | 404 (2.9%) | 4,681 (33.4%) | 288 (2.1%) | 7 (0.0%) | 393 (2.8%) | 274 (2.0%) |
| Token (ratio) | 128,830 (47.8%) | 3,292 (1.2%) | 26,315 (9.8%) | 5,921 (2.2%) | 76,604 (28.4%) | 3,667 (1.4%) | 1,173 (0.4%) | 10,960 (4.1%) | 2,566 (1.0%) |

| POS | 부정지정사 | 부사적조사 | 보조사 | 접속조사 | 감탄사 | 연결어미 | 종결어미 | 관형사 |
|---------------|--------------|--------------|--------------|-----------|--------------|--------------|-----------|--------------|
| Type (ratio) | 14 (0.1%) | 159 (1.1%) | 85 (0.6%) | 9 (0.5%) | 64 (0.5%) | 38 (0.3%) | 26 (0.2%) | 6 (0.0%) |
| Token (ratio) | 1,934 (0.7%) | 2,213 (0.8%) | 1,017 (0.4%) | 26 (0.0%) | 1,723 (0.6%) | 2,050 (0.8%) | 61 (0.0%) | 1,234 (0.5%) |

**모두의 말뭉치
(구어)**

| POS | 서술어 | 수사 | 대명사 | 고유명사 | 일반명사 | 의존명사 | 접속부사 | 일반부사 | 긍정지정사 |
|---------------|---------------------|----------------|------------------|------------------|--------------------|----------------|----------------|------------------|------------------|
| Type (ratio) | 67,556 (9.9%) | 2,642 (0.4%) | 521 (0.1%) | 328,696 (48.2%) | 272,073 (39.9%) | 3,153 (0.5%) | 33 (0.0%) | 4,634 (0.7%) | 101 (0.0%) |
| Token (ratio) | 157,826,022 (70.8%) | 858,960 (0.4%) | 2,790,559 (1.3%) | 4,698,997 (2.1%) | 35,473,822 (15.9%) | 505,219 (0.2%) | 901,740 (0.4%) | 5,819,016 (2.6%) | 1,138,012 (0.5%) |

| POS | 부정지정사 | 부사적조사 | 보조사 | 접속조사 | 감탄사 | 연결어미 | 종결어미 | 관형사 |
|---------------|------------|----------------|----------------|-----------|----------------|------------------|------------------|------------------|
| Type (ratio) | 25 (0.0%) | 42 (0.0%) | 63 (0.0%) | 2 (0.0%) | 816 (0.1%) | 647 (0.1%) | 471 (0.1%) | 900 (0.1%) |
| Token (ratio) | 407 (0.0%) | 755,803 (0.3%) | 598,238 (0.3%) | 74 (0.0%) | 848,680 (0.4%) | 2,047,409 (0.9%) | 5,810,309 (2.6%) | 2,905,242 (1.3%) |

4.3.3.1 품사별 Top 30 (수어, 모두) - 전체

Document Name

97 소속 원문

98_220404수어코파스#211123-220406#210413#r<같다><자기><저는><부다><정부><크다><같은화 맞아요 저율부터 정부의 짧은이 크다>생각하요 정부에서 99_220404수어코파스#211123-220406#210413#r<자기><정부><김기기><중단된다.><�다><키우다>그리다가 김기기 정부에서 지원을 중단해서 보육들은 지난 100_220404수어코파스#211123-220406#210413#r<자기><정부><포><경고하다>나타나다><사람><정> 정부에서 다시 자녀 출산을 권장하면서 앞축비를 주겠다고 101_220404수어코파스#211123-220406#210413#r<제를좋다><만나다><행하다><시간><없다><자기>예를 들어 만난다고 같은데 시간이 없단다가 전쟁통에 죽음을 102_220404수어코파스#211123-220406#210413#r<부변다><나쁘다><이해하다><거두번별 다른 방향으로 사용한다면 예를 들어 문단통장을 103_220404수어코파스#211123-220406#210413#r<같다><같다><죽임><생각하다><죽다><생각하다>왕이 봄으로 그 육신이 많아요 같은데 놓았구요 하시면 좋 104_220404수어코파스#211123-220406#210413#r<포기><경계지역><북한><군주권><한국>한국의 경계지역에 북한과 인접에 한국까지 약학하려는 것 105_220404수어코파스#211123-220406#210413#r<자기><사람><남자><사람><남자><남다><멀리가다><+><거기> 고구려 사람들이 남아 있었고 멀리가다 멀리가다 중 106_220404수어코파스#211123-220406#210413#r<포기><대우><죽임><죽다><포기><죽임><죽다><포기><죽임><죽다><포기> 한민족은 대체 정의 있었던 107_220404수어코파스#211123-220406#210413#r<포기><죽다><포기><죽다><포기><죽다><포기><죽다><포기> 몇몇 남자는 본성을 날지 않아서 일정한 자리를 두고 보유하고 확장 108_220404수어코파스#211123-220406#210413#r<같다><생각하다><포기><포기><포기><포기><포기><포기><포기> 청제대종묘와 저봉정에는 아름 들여 있는 대 물체 이용해 나쁜 술 109_220404수어코파스#211123-220406#210413#r<같다><생각하다><래다><래><부모><있다><문을 세게 닫는 것> 대해 래인 이유가 있든 없든 회 110_220404수어코파스#211123-220406#210413#r<같다><생각하다><开门><开门><开门><开门><开门><开门><开门> 자리를 넓이고 그에 대한 사랑과 칭송이 같고 또만 입방하 111_220404수어코파스#211123-220406#210413#r<앞으로><이기><복><법>< 있다><별하다><별하다><별하다> 미련되고 이동복력에 대단



4.3.3.2 품사별 Top 30 (수어, 모두) - 서술어

4.3.3.3 품사별 Top 30 (수어, 모두) - 일반명사

| Gloss | types | types(%) | types(CUM%) | 토달 번 |
|---------------------------|-------|----------|-------------|------|
| 그만 까지 끝 + 끝++ 끝++ 끝++ 끝++ | 3,462 | 4.5193 | 4.5193 | 688 |
| 더럽다 더럽다+ 때 불결 지저분하다 | 1,598 | 2.0861 | 6.6054 | 52 |
| 가혹 세세 주제 짚 집+** 히 | 1,575 | 2.0560 | 8.6614 | 1455 |
| 놓아나는 놓는 놓인+ 놓인++ | 1,384 | 1.8067 | 10.4681 | 1582 |
| 금일 오늘 오늘+ 이번 현재 현재+ 현재+ | 1,327 | 1.7323 | 12.2004 | 1645 |
| -인 분 사람 사람+ 사람++ 이 인간 자 | 1,310 | 1.7101 | 13.9105 | 736 |
| 방법 방식 방식+ 방식++ 요령 | 1,153 | 1.5051 | 15.4156 | 857 |
| 여 여성 여자 여자+ | 980 | 1.2793 | 16.6950 | 2199 |
| 예전 옛날 옛날+ 옛적 | 876 | 1.1435 | 17.8385 | 656 |
| 얼굴+를 얼굴이를+를 | 817 | 1.0665 | 18.9050 | 2202 |
| 곳 곳+ 곳++ 곳++ 곳++++ 위치 장소 | 798 | 1.0417 | 19.9467 | 2305 |
| 풀 모습 모양 상태 상황 상황+ 장면 형상 | 690 | 0.9007 | 20.8475 | 2226 |
| 다 모두 모두+ 모든 모조리 용망은 운통 | 689 | 0.8994 | 21.7469 | 2377 |
| 간아지게 개+개 | 664 | 0.8668 | 22.6137 | 997 |
| 동부 북 우진 친구 친구+ 친구++ | 617 | 0.8054 | 23.4191 | 2154 |
| 해 해+ 해하 해하+ 해하+ 해하++ | 612 | 0.7989 | 24.2181 | 579 |
| 남 남성 남자 남자+ 남자++ 남자++ 남 | 598 | 0.7806 | 24.9987 | 1585 |
| 간단하게 간단하게+ 야기 야기+ 영야 젖먹 | 589 | 0.7689 | 25.7567 | 847 |
| 문제 문제+ 문제++ | 578 | 0.7545 | 26.5221 | 696 |
| -국 국가 나라 | 577 | 0.7532 | 27.2753 | 4939 |
| 정인 정인+ 정인++ | 561 | 0.7323 | 28.0077 | 2157 |
| 가격 값 금액 금전 대금 돈+ 돈++ 돈+ | 547 | 0.7141 | 28.7217 | 2281 |
| 선생님 선생님+ 선생님++ | 535 | 0.6984 | 29.4201 | 1526 |
| 학교 | 534 | 0.6971 | 30.1172 | 4927 |
| (꽃)피다 꽃 꽃+ 꽃+ 꽃++ 꽃++ | 534 | 0.6971 | 30.8143 | 1590 |
| 시간 시간+ 시간++ 한 | 514 | 0.6710 | 31.4853 | 1513 |
| 화성 영화 | 505 | 0.6592 | 32.1445 | 3391 |
| 법 법+ 법++ 법++ 법 규법 법률 법칙 | 483 | 0.6305 | 32.7751 | 738 |
| 결과 결국 결말 드디어 마지막 막총 | 482 | 0.6292 | 33.4043 | 1605 |
| 미 외 사+ 사+ 사++ 사+++ | 436 | 0.5692 | 33.9734 | 1529 |

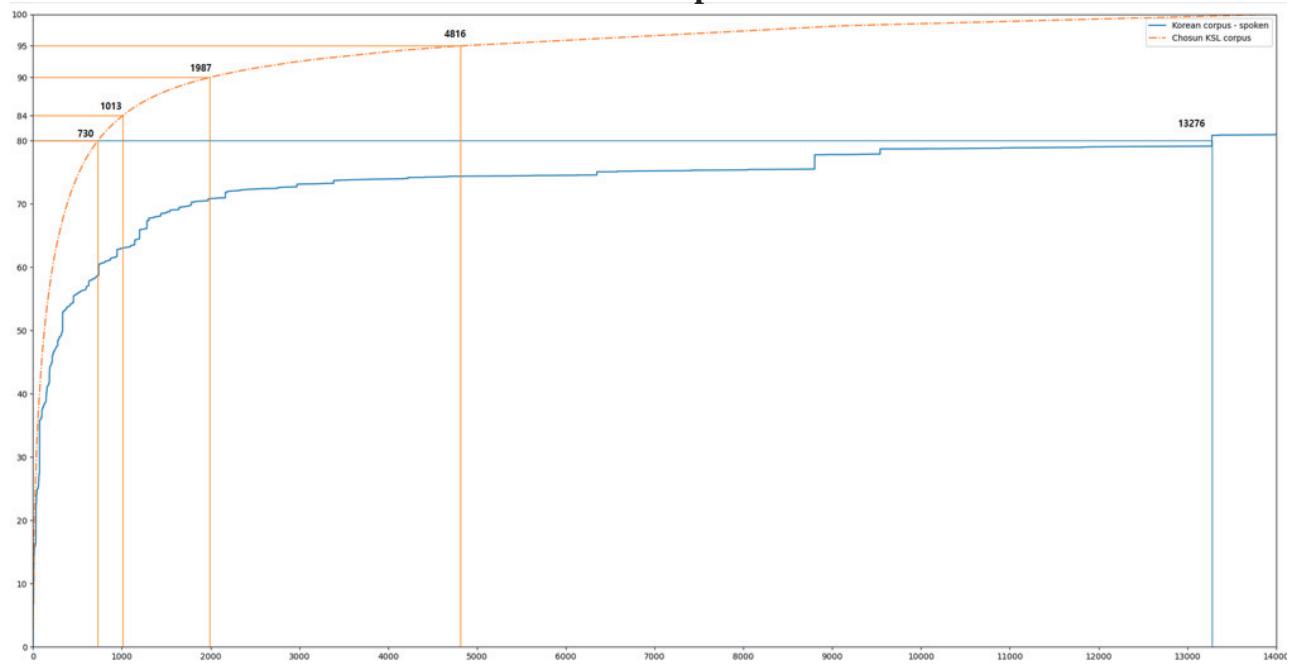
| 단어 | types | types% | typesCUM% | 표제어 번호 |
|---|-----------|---------|-----------|--------|
| [‘박단예’, ‘생기는’)균열’, ‘틈’, ‘틈+’, ‘틈++’, | 5,033,463 | 14.1892 | 14.1892 | 1435 |
| [‘ㄴ’] | 3,803,880 | 10.7231 | 24.9123 | 2696 |
| [‘지내’, ‘지네+’] | 2,676,402 | 7.5447 | 32.4570 | 837 |
| [‘ㄹ’, ‘것이다’, ‘아마’, ‘아마+’, ‘아마++’, ‘아 | 2,153,158 | 6.0697 | 38.5267 | 1000 |
| 여서 | 783,448 | 2.2085 | 40.7353 | |
| [‘증비모양’, ‘증비모양+’] | 727,431 | 2.0506 | 42.7859 | 2163 |
| [‘사실’, ‘정말’, ‘진짜’, ‘참’, ‘정말로’] | 665,917 | 1.8772 | 44.6631 | 2744 |
| [‘풀’, ‘모임’, ‘모임’, ‘상황’, ‘상황+’, ‘정 | 303,179 | 0.8547 | 45.5177 | 2226 |
| [‘그룹’, ‘단체’, ‘동아리’, ‘모임’, ‘모임+’, ‘모 | 262,589 | 0.7402 | 46.2580 | 865 |
| [‘올해오늘’, ‘올해오늘+’] | 261,477 | 0.7371 | 46.9951 | 1491 |
| [‘문제’, ‘문제+’, ‘문제++’] | 253,400 | 0.7143 | 47.7094 | 696 |
| [‘경우’, ‘때’, ‘때+’, ‘때++’, ‘무렵’, ‘시기’, ‘시 | 225,277 | 0.6351 | 48.3444 | 2274 |
| [‘난간아기’, ‘난간아이’, ‘아이’, ‘아기+’, ‘영C | 195,695 | 0.5517 | 48.8961 | 847 |
| [‘憎’, ‘만평’, ‘여자’, ‘정도’, ‘정도+’] | 192,153 | 0.5417 | 49.4378 | 2166 |
| [‘결과’, ‘결과’, ‘결국+’, ‘결방’, ‘드디어’, ‘마지 | 190,566 | 0.5372 | 49.9750 | 1605 |
| [‘십자가’, ‘십자가+’, ‘십자가++’, ‘십자가+’] | 187,562 | 0.5287 | 50.5037 | 848 |
| [‘다’, ‘모두’, ‘모두+’, ‘모든’, ‘모조리’, ‘용광’, | 172,050 | 0.4850 | 50.9887 | 2377 |
| [‘가장’, ‘맨’, ‘수석’, ‘으뜸’, ‘일등’, ‘등등+’, ‘을 | 165,579 | 0.4668 | 51.4555 | 1482 |
| [‘다른나라’, ‘다른나라+’] | 165,384 | 0.4662 | 51.9217 | 867 |
| [‘요사이’, ‘요새’, ‘요즈음’, ‘요즘’, ‘작금’, ‘최 | 141,961 | 0.4002 | 52.3219 | 836 |
| [‘아이디어’, ‘아이디어+’, ‘아이디어++’] | 129,122 | 0.3640 | 52.6859 | 846 |
| [‘금일’, ‘오늘’, ‘오늘날’, ‘이번’, ‘현재’, ‘현재+ | 128,539 | 0.3623 | 53.0482 | 1645 |
| [‘모친’, ‘어머니’, ‘어미’, ‘엄마’, ‘엄마+’] | 121,647 | 0.3429 | 53.3911 | 1498 |
| [‘국민’, ‘백성’] | 112,714 | 0.3177 | 53.7089 | 3646 |
| [‘말씀’] | 107,546 | 0.3032 | 54.0121 | 2626 |
| [‘아주머니’, ‘아주머니+’] | 106,880 | 0.3013 | 54.3133 | 2209 |
| [‘시작’, ‘우선’, ‘처음’, ‘처음+’, ‘처음++’, ‘초 | 104,763 | 0.2953 | 54.6087 | 1453 |
| [‘대통령’] | 98,111 | 0.2766 | 54.8852 | 3838 |
| [‘개거’, ‘거의’, ‘대부분’, ‘대부분+’, ‘주로’] | 96,657 | 0.2725 | 55.1577 | 2284 |
| [‘대한’, ‘대한’, ‘대한하고’] | 93,410 | 0.2633 | 55.4210 | 1578 |

4.3.3.4 품사별 Top 30 (수어, 모두) - 일반부사

| Gloss | types | types% | types(CUM%) | 도달 번호 | 단어 | types | types% | types(CUM%) | 표제어 번호 |
|--------------------------|-------|--------|-------------|-------|---|-----------|---------|-------------|--------|
| 그리고 다시 또 또++ 또+++ 또++++ | 993 | 9.0602 | 9.0602 | 2381 | [갓]', '시방', '이제', '언제', '지금', '지금+', '자 | 1,907,138 | 32.7742 | 32.7742 | 2367 |
| 다소 약간 약소하다 자그마하다 작다 작다 | 884 | 8.0657 | 17.1259 | 50 | [그리고]', '다시', '또', '또+', '또++', '또++++ | 1,034,092 | 17.7709 | 50.5451 | 2381 |
| 자연히 자연히+ 자연히++ 자연히+++ 자 | 860 | 7.8467 | 24.9726 | 1649 | 많이 | 291,692 | 5.0127 | 55.5579 | |
| 어째서 왜+ 왜++ 왜+++ 왜++++ | 702 | 6.4057 | 31.3777 | 2371 | 굉장히 | 224,115 | 3.8514 | 59.4093 | |
| 가장 맨 수석 유품 일등+ 일등++ 일 | 625 | 5.7026 | 37.0803 | 1482 | [그냥]', '그냥+', '그냥++', '그저'] | 180,691 | 3.1052 | 62.5145 | 1001 |
| 갓 시방 이제 인제 지금+ 지금++ 지 | 619 | 5.6478 | 42.7281 | 2367 | 일단 | 112,176 | 1.9277 | 64.4422 | |
| 나날 날마다 늘 매번 매일 번번이 수 | 502 | 4.5803 | 47.3084 | 2366 | [바로] | 107,881 | 1.8539 | 66.2962 | 3275 |
| 는 것이다 이마 이마+ 아마+ 아마도 | 499 | 4.5529 | 51.8613 | 1000 | [먼저]', '먼저+', '먼저++', '먼저+++'] | 74,277 | 1.2765 | 67.5726 | 2378 |
| 과 반려 생 빡 함께 험께+ 험께++ 험께하 | 463 | 4.2245 | 56.0858 | 359 | [~만', 뿐', '다만', '단지', '오로지', '오'] | 73,730 | 1.2671 | 68.8397 | 2758 |
| 김의자 길게늘어서다 길다 길다+ 오래 오 | 324 | 2.9562 | 59.0420 | 251 | 상당히 | 63,812 | 1.0966 | 69.9363 | |
| 그냥 그냥+ 그냥++ 그저 | 319 | 2.9106 | 61.9526 | 1001 | [잠깐]', '잠깐+', '잠깐++', '잠깐+++', '잠시 | 53,829 | 0.9251 | 70.8613 | 1648 |
| 상호 서로 서로+ 서로++ | 301 | 2.7464 | 64.6989 | 607 | [직접]', '직접+', '직접++'] | 49,662 | 0.8534 | 71.7148 | 1647 |
| 그대로 그대로+ 그대로++ 되돌리다 복귀 | 260 | 2.3723 | 67.0712 | 885 | 얼마나 | 48,487 | 0.8333 | 72.5480 | |
| 장간 장간+ 장간++ 장간+++ 장시 | 195 | 1.7792 | 68.8504 | 1648 | [조금씩]', '조금씩+', '조금씩++', '조금씩++] | 47,721 | 0.8201 | 73.3681 | 2368 |
| 날 날짜 언제 언제++ 일자 | 188 | 1.7153 | 70.5657 | 999 | [오늘剋]', '오히려+]' | 47,074 | 0.8090 | 74.1771 | 2373 |
| 미~ 미처 아직 아직+ 아직++ | 183 | 1.6697 | 72.2354 | 1652 | 어쨌든 | 46,625 | 0.8013 | 74.9783 | |
| 직접 직접+ 직접++ | 178 | 1.6241 | 73.8595 | 1647 | 실제로 | 45,929 | 0.7893 | 75.7676 | |
| 가끔 가끔+ 간간이 간혹 드문드문 | 148 | 1.3504 | 75.2099 | 1003 | [왜나하면] | 40,839 | 0.7018 | 76.4695 | 4056 |
| 갑자기 갑자기+ 갑작스럽다 갑작스레 냄새 | 135 | 1.2318 | 76.4416 | 1002 | [혹]', '혹+', '혹시]', '혹여] | 39,739 | 0.6829 | 77.1524 | 1644 |
| 자유 임의 마구 마구 대대로 대대로 기 | 133 | 1.2135 | 77.6551 | 3041 | 과연 | 36,715 | 0.6309 | 77.7833 | |
| 꽤 대 대단히 매우 매우+ 매우++ 매우++ | 119 | 1.0858 | 78.7409 | 2379 | 당연히 | 31,383 | 0.5393 | 78.3226 | |
| ~자 곧 금방간 냉큼 빨리 빨리+ 빨리++ | 118 | 1.0766 | 79.8175 | 2376 | 훨씬 | 28,155 | 0.4838 | 78.8065 | |
| 도무지 | 117 | 1.0675 | 80.8850 | | [별로]', '그다지']' | 27,893 | 0.4793 | 79.2858 | 2875 |
| 오히려 오히려+ | 110 | 1.0036 | 81.8887 | 2373 | 또는 | 27,731 | 0.4766 | 79.7624 | |
| 먼저 먼저+ 먼저++ 먼저+++ | 106 | 0.9672 | 82.8558 | 2378 | 어우 래도 | 26,591 | 0.4570 | 80.2194 | |
| ~끼리 길이 더불다 등반 등반하다 땅 아울 | 83 | 0.7573 | 83.6131 | 648 | 충분히 | 25,827 | 0.4438 | 80.6632 | |
| 각 각각 각각+ 각각++ 각각+++ 각기 각 | 69 | 0.6296 | 84.2427 | 740 | 분명히 | 25,316 | 0.4351 | 81.0982 | |
| 혁 혁+ | 68 | 0.6204 | 84.8631 | 2365 | [가끔]', '가끔+', '가끔+++', '간간이]', '간혹', ' | 24,219 | 0.4162 | 81.5145 | 1003 |
| 미리 미리+ | 61 | 0.5566 | 85.4197 | 1653 | [완전히]', '완전히+'] | 23,442 | 0.4029 | 81.9173 | 2372 |
| 무조건 무작정 | 55 | 0.5018 | 85.9215 | 2423 | 그죠 | 22,808 | 0.3920 | 82.3093 | |

4.3.4 Comparison of graphs in cumulative coverage between Spoken Corpus of Modun Korean and the Chosun KSL Corpus

한국해양대학교
KOREA MARITIME & OCEAN UNIVERSITY



4.4 General vocabulary in the Chosun KSL Corpus

| | | |
|-----------------------------------|------------|------------|
| Cumulative % | 83.6% | 84.14% |
| Content words (N, V, Adj, Adv) | 1,797 (14) | 1,892 (13) |
| Nouns | 847 (14) | 896 (13) |
| Verbs | 669 (14) | 706 (13) |
| Adjectives | 186 (14) | 194 (13) |
| Adverbs | 95 (14) | 96 (13) |

(cut-off minimum frequency)

| 학년 | 명사 | 동사 | 형용사 | 부사 | 총계 |
|----|-----|-----|-----|----|-------|
| 1 | 94 | 74 | 20 | 12 | 200 |
| 2 | 94 | 74 | 20 | 12 | 200 |
| 3 | 94 | 74 | 20 | 12 | 200 |
| 4 | 94 | 74 | 21 | 11 | 200 |
| 5 | 94 | 75 | 21 | 10 | 200 |
| 6 | 94 | 75 | 21 | 10 | 200 |
| 7 | 94 | 75 | 21 | 10 | 200 |
| 8 | 94 | 75 | 21 | 10 | 200 |
| 9 | 95 | 75 | 21 | 9 | 200 |
| 소계 | 847 | 671 | 186 | 96 | 1,800 |

41

4.4.1 Examples of the First 100 Signs for 1st Grade

1. 47 Frequency Order of Nouns (명사 47개, 고빈도 순)

나, 저기, 끝, 너, 집, 때, 농인, 무엇, 여기, 저거, 이거, 사람, 그것, 여자, 이것, 방법, 얼굴이름, 곳, 친구, 개, 한국, 남자, 아기, 나라, 청인, 문제, 돈, 학교, 선생님, 거기, 꽃, 영화, 예전, 모양, 중, 시간, 법, 산, 협회, 수어, 취미, 옛날, 차, 얼굴, 아파트, 당신, 전체

※ 기능어: **지시대명사/인칭대명사/의문대명사/조동사**(완료상 '끝', 진행형 '중') /**지시대부사**
/문법형태소(부사: 과거시제 '예전', '옛날'), 접미사?? ('사람')

※ 내용어: 개인적인 정보 ('집' '여자' '얼굴이름' '친구' '남자' '아기' '취미')
/좋아하는 것 ('개' '꽃' '영화' '산')/농사회 ('농인' '청인' '학교' '선생님' '협회' '수어')
/경제이슈('문제' '돈' '법' '차' '아파트',)/기타('사람' '한국' '나라' '모양' '시간' '전체')
참고: 선생님 → 선생 + 님 (**합성어** 또는 파생어??)

42

4.4.1 Examples of the First 100 Signs for 1st Grade (Cont'd)

2. 39 verbs in ranking of high frequency (동사 39개 고빈도 순)

맞다, 보다, 생각하다, 말하다, **잘하다, 안하다**, 먹다, 가다, 알다, 일하다, 원하다, 만들다, 수어하다,
못하다, 되다, 대화하다, 두다, 주다, **죽다, 모르다**, 사용하다, 놀다, 느끼다, **말다**, 만나다, 하다, 살다,
끌나다, 바꾸다, 여행하다, 사다, 가르치다, 모이다, 키우다, 경험하다, 받다, 결혼하다, 돋다, 오다

<유형별 동사구별>

- ※ 일치동사: 보다, 말하다, 수어하다, 대화하다, 주다, 만나다, 가르치다, 받다, 결혼하다, 돋다,
- ※ 공간동사: 가다, 여행하다, 오다
- ※ 일반동사: 맞다, 생각하다, 잘하다, 안하다, 먹다, 알다, 일하다, 원하다, 만들다, 못하다 등

<긍정과 부정 의미적 동사구분>

- ※ **긍정적인 의미를 갖는 동사**: 맞다, 잘하다, 되다, 돋다
- ※ **부정적인 의미를 갖는 동사**: 안하다, 못하다, 죽다, 모르다, 말다,

43

4.4.1 Examples of the First 100 Signs for 1st Grade (Cont'd)

3. Top 10 High-frequency Adjectives (형용사 10개 고빈도 순)

같다, 있다, 없다, 좋다, 많다, 다르다, 필요하다, 괜찮다, 즐겁다, 가능하다

4. Top 6 Adverbs (부사 6개 고빈도 순)

~면, 현재, 그러나, 또, 자연히, 조금

- ※ 기능어: 접속사(~면' '그러나'), 문법형태소(부사:현재시제 '현재')
- ※ 내용어: 또, 자연히, 조금

44

4.4.2 Examples of the Second 100 Signs for 1st Grade

1. 47 high-frequency nouns (명사 47개 고빈도 순)

결국, **부모**, 나이, 그림, **앞으로**, 나무, 혼자, 몸, 최고, **만약**, 대부분, 학생, 이름, 바다, **가족**, **누구**, ○ㅋ, 처음, 교회, 병원, 보통, 경찰, 다음, 중국, 코로나, **전에**, 진짜, 마음, 이유, 내용, **이들**, 원래, 물, **부부**, 생각, 날, 입장, 휴대폰, **방식**, 힘, **어머니**, 이야기, 목적, **아들**, **후에**, 눈, 문장

※ 가족관계 (부모, 가족, 부부, 어머니, 아들) / 대명사 (누구, 이들, 방식) 신체부분 (몸, 눈)
※ 문법형태소: (부사: 미래시제 '앞으로', 과거시제 '전에', '후에') / 접속사: '만약'

2. 35 high-frequency verbs (동사 35개 고빈도 순)

사라지다, 듣다, 사랑하다, 가지다, 더하다, 마무리하다, **가져오다**, 쉬다, **잡다**, 발전하다, 공부하다, 불나다, **다니다**, 배우다, 계속하다, 비교하다, 합하다, 운동하다, **방문하다**, 세우다, **버리다**, **질문하다**, **부르다**, 구경하다, 떨어지다, 요리하다, 줄다, 활동하다, 운전하다, 넣다, 증가하다, 감사하다, 기다리다, 선택하다, 입다

<유형별 동사구별>

※ 일치동사: 가져오다, 잡다, 버리다, 질문하다, 부르다

※ 공간동사: 다니다, 방문하다

※ 일반동사: 그 외 동사들 모두

45

4.4.2 Examples of the Second 100 Signs for 1st Grade (cont'd)

3. 10 High-frequency Adjectives (형용사 10개 고빈도 순) : 크다, 어렵다, 의미하다, 똑같다, 불편하다, 힘들다, 낫다, 깨끗하다, 싫다, 비슷하다

4. 6 High-frequency Adverbials (부사 6개 고빈도 순) : 왜, 지금, 가장, 아마, 항상, 대하여

※ 기능어: 의문사 ('왜'), '대하여 '

※ 내용어: '지금', '가장', '아마', '항상'

46

5. Conclusion and More Issues to be Considered

5.1 Conclusion: A summary of answers to the two research questions

(1) How many words should comprise general vocabulary in Korean Sign Language?

→ 전체 19,183개의 단어가 273,086번의 빈도수를 가져 명사, 동사, 형용사, 부사 등 내용어(17,652)가 251,205번 빈도수를 보여주어 92%를 차지하고 있는 것으로 보이나, POS에서 내용어로 분류된 인칭대명사(3.88%)와 비인칭대명사와 지시대부사(3.8%)와 의문사 (2.2%)를 기능어로 재분류하면 실질적으로 기능어는 17.9% 이상이고 82.1% 정도가 내용어임.

(2) Which words are proper to create spoken language textbooks of Korean and English for the first grade to the ninth grade deaf students ?

→ 초중등 학생들에게 자신의 모국어인 한국수어로 한국어와 영어 등 외국어를 배우는 경우에는 학년에 맞는 한국수어 어휘로 외국어언어교재를 만드는 것이 필요함. 그래서 누적 약 84%를 차지하는 1,800여개의 내용어를 9개 학년으로 나누면 학기당 100개의 어휘를 학습하는 것이 적절함. 여기서 각 내용어도 학년별 균등하게 나눌 필요가 있음.

→ 구체적으로, The first 100 signs과 The second 100 signs은 초등1학년용 교재에 사용되어야 하는 어휘이고, The third 100 signs와 the fourth 100 signs은 초등2학년용 교재에, the 17th 100 signs와 the 18th 100 signs은 중학교 3학년용 교재에 적합한 어휘로 보임.

47

5.2 More issues to be considered

- (1) Do we need to define words in a sign language corpus as a collection of glosses as type, lemma, or word-family?
- (2) How to count words in a sign language corpus: **glosses** or **signs** in sentences?
- (3) We need more studies on parts of speech in KSL.
 - Are there any function words in sign language? If so, what are the parts of speech.
For example, conjunctions ('~면', '~때'), postpositions ('후에'), pronouns, auxiliary verbs ('가능하다'), modals ('끌' 완료상), qualifiers, question words, etc.
 - In case of content words, how to identify nouns, verbs (plain verb 끌나다 vs. auxiliary verb 끌), adjectives, or adverbs ('예전', '옛날' (과거 시제), '앞으로' (미래 시제))?
- (4) Is the children's KSL corpus necessary for their preferred vocabulary for elementary and middle school deaf students?
- (5) To order to compare the Chosun KSL Corpus with the KSL online dictionary of the National Institute of Korean Language, we need some criteria for collecting the data and registering the words in the KSL dictionary.
- (6) To compare the Chosun KSL corpus with the spoken corpus out of the contemporary Korean corpus called 'Modun Mal Mungchi,' and explore some similarities and differences between KSL as a sign language and Korean as a spoken language, further studies are needed.

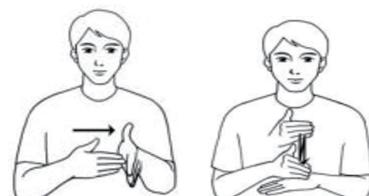
48

References

- 장세은, 박호민 (2022) 조선대한국수어코퍼스의 어휘분석, 대한언어학회 가을학술대회 (2022년 11월 19일, 광주, 조선대) 발표 논문.
- 장세은, 박호민 (2023) 수어 말뭉치 비교 분석 및 학습 어휘 선정, 2023 조선대학교 언어융합연구소 학술대회 및 성과보고회 (2023년 06월 09일, 경주 더케이호텔) 성과발표 3.
- Bauer, L., & Nation, P. (1993). Word families. *International Journal of Lexicography*, 6, 253-279.
- Brezina, V., & Gablasova, B. (2015). Is there a core general vocabulary? Introducing the new general service list. *Applied Linguistics*, 36(1), 1-22.
- Carroll, J. B., Davies, P., & Richman, B. (1971). The American heritage word frequency book. New York: Houghton Mifflin, Boston American Heritage.
- Coxhead, A. (2000). A new academic word list. *TESOL Quarterly*, 34, 213–238.
- Dang, T. N. Y., Coxhead, A., & Webb, S. (2017). The academic spoken word list. *Language Learning*, 4, 959-997.
- Gardner, D., & Davies, M. (2014). A new academic vocabulary list. *Applied Linguistics*, 35, 305-327.
- Hsu, W. (2011). The vocabulary thresholds of business textbooks and business research articles for EFL learners. *English for Specific Purposes*, 30(4), 247-257.
- Hsu, W. (2014). Measuring the vocabulary load of engineering textbooks for EFL undergraduates. *English for Specific Purposes*, 33, 53-64.
- Laufer, B., & Raventhal-Kalovski, G. C. (2010). Lexical threshold revisited: Lexical text coverage, learners' vocabulary size and reading comprehension. *Reading in a Foreign Language*, 22(1), 15-30.
- Liu, N., & Nation, P. (1985). Factors affecting guessing vocabulary in context. *RELC Journal*, 16 (1), 33-42 McLean (2018)
- Nagy, W., & Townsend, D. (2012). Words as tools: Learning academic vocabulary as language acquisition. *Reading Research Quarterly*, 47, 91-108
- Nation, P. (2001). *Learning vocabulary in another language*. Cambridge: Cambridge University Press
- Nation, P. (2016). *Making and Using Word Lists for Language Learning and Testing*. Amsterdam, The Netherlands: John Benjamins.
- Schmitt, N., & Schmitt, D. (2014). A reassessment of frequency and vocabulary size in L2 vocabulary teaching. *Language Teaching*, 47(4), 484-503
- Schmitt, N., & Zimmerman, C. B. (2002). Derivative word forms: What do learners know? *TESOL Quarterly*, 36, 145-171
- Webb, S., & Rodgers, M. (2009a). Vocabulary demands of television programs. *Language Learning*, 59(2), 335-366.
- Webb, S., & Rodgers, M. (2009b). The lexical coverage of movies. *Applied Linguistics*, 30(3), 407-427.
- West, M. (1953). *A general service list of English words*. London, UK: Longman.
- Zeno, S. M., Ivens, S. H., Millard, R. T., & Duvvuri, R. (1995). The educator's word frequency guide. Touchstone Applied Science Associates.

49

THANK YOU very much.



QUESTION

2(you) GIVE 1(me)



Any Questions to me?

50

코퍼스 기반 한국수어의 등위구문 연구

고인경(나사렛대), 윤병천(조선대), 전혜원(부경대), 이선혜(부경대),
박호민(한국해양대), 장세은(한국해양대)

1. 서론

1.1. 연구목적

소리언어에서 사용되는 순접, 역접, 선택적으로 사용되는 이접 등 세 가지 유형의 절 등위구문(conjunctively 'and', adversatively 'but', and disjunctively 'or' coordinated clause (Haspelmath 2004, 2007)이 한국수어에서는 어떻게 구현되고 있는지 수어뉴스코퍼스와 조선대수어코퍼스를 기반으로 외국수어와 비교하면서 소리언어 및 수화언어 간의 차이점과 유사점을 기술하고 수어의 절 등위구문에 대한 범어적인 특성을 찾고자 하는 것이 목적임.

1.2 키워드 (Keywords)

clausal coordination, clausal boundary, conjunctive ('and'), adversative ('but'), disjunctive ('or'), nonmanual elements, syndetic, asyndetic, sequential and simultaneous conjunctions, Across-the-Board (ATB)

1.3 코퍼스 종류

- (1) **수어뉴스코퍼스** (고인경, 윤병천 2016): 2011년~2014년 4년 간의 한국농아방송국 iDBN(<http://cafe.daum.net/deafon>)의 수어뉴스 동영상 67편 중 각 분기별 월 1편씩 매년 3 편을 무작위로 선정한 뉴스 동영상 자료 12편 (34분 58초), 144문장
- (2) **조선대수어코퍼스**(장세은, 박호민 2022, 2023): 2020년~2022년 2년 간 조선대 언어융합연구소에서 구축한 한국수어코퍼스로서 한국농인 남 16명, 여 10명, 평균연령 49.61세, 70개 파일, Gloss token: 251,205개, Gloss type: 17,652개, Gloss lemma: 5,010개, 번역된 한국어 문장: 17,149개

1.4 연구 질문

- Q1. 절 등위 구문에서 절과 절 간의 경계(boundary)를 나타내는 요소는 무엇인가? 특히, 다른 수어의 절 등위구문에 대한 선행연구에서는 절 경계를 나타내는 요소는 눈을 깜박이는 비수지표지(an eye blink as a non-manual marker)인데 한국수어는 어떠한지?
- Q2. 홍콩수어에서는 순서적인 순접등위(sequential conjunction)인지 동시적인 순접등위(simultaneous conjunction)인지에 따라 non-manual이 다르게 나타나므로 한국수어의 등위 구문의 종류에 따라 비수지표지가 다르게 나타나는지, 만약 다르게 나타난다면 어떻게 나타나는지?

- Q3. 한국수어의 역접등위 구문에서 두 가지 서로 다른 수어 접속사는 왜 나타나는가?
- Q4. Velupillai(2012)는 수어 등위구문은 asyndetic 경향을 가지고, Jantunen(2016)에 따르면, 핀란드수어(FinSL)에서는 순접접속사가 syndetic인 경우가 많고, 선택이나 역접의 경우 syndetic인 경우가 많다는 선행연구의 범어적인 관점에서, 한국수어 등위구문에서도 이러한 syndetic/asyndetic 특성이 있는지 있다면 어떠한 비율로 나타나는지?
- Q5. 한국수어에서 전역구문(Across-the-Board, ATB)이 생성문법의 이론과 통합된 설명이 가능할지 혹은 다른 양상을 보일지?

2. 한국수어 접속사 유형과 등위구문

2.1. 수어코퍼스

(3) 수어뉴스코퍼스에 나타난 한국수어 접속사 유형 (고인경, 윤병천 2016)

| 구분 | | 1 | 2 | 3 | 4 |
|----------------|------|---------------------|----------------|----------------------------|------------------------|
| 1 병렬 접속사 | 한국수어 | | {고개를아래로 내림} | {휴지} | |
| | 한국어 | 또는, 혹은 | | | |
| 2 순접 접속사 | 한국수어 | | {고개를아래로 내림} | {휴지} | |
| | 한국어 | 그리고 | | | |
| 3 첨가 접속사 | 한국수어 | | | | |
| | 한국어 | 또한, 아울러, 뿐만 아니라 | | | |
| 4 역접 접속사 | 한국수어 | | | | |
| | 한국어 | 그러나, 하지만, 반면 | | | |
| 5 인파 접속사 | 한국수어 | | | | {눈에힘을줌} {고개 를아래로내림} |
| | 한국어 | 그러므로, 그래서, 따라서, 때문에 | | | |
| 6 가정 접속사 | 한국수어 | | | {눈을크게뜸} { 고개를아래로내 림} | |
| | 한국어 | ~면, 경우, ~일 때 | | | |

(4) 수어뉴스코퍼스에 나타난 한국수어 접속사 유형별 출현 횟수 (고인경, 윤병천 2016)

<표 7> 한국수어 접속사 유형별 출현 위치와 횟수

| 구분 | 문장 앞머리 | 문장 사이 | | 비수지(횟수) |
|---------|-------------|------------------------------------|--|-----------------------------|
| | 수지(횟수) | 수지(횟수) | | |
| 1 병렬접속사 | | {또는} (12) | | {고개를아래로내림} (2) {휴지} (3) |
| 2 순접접속사 | {그리고} (3) | {그리고} (16) | | {고개를아래로내림} (8) {휴지} (16) |
| 3 침가접속사 | {뿐만아니라} (2) | {함께} (1), {뿐만아니라} (6) | | |
| 4 역접접속사 | {그러나} (4) | {그러나} (14), {반면} (1) | | |
| 5 인과접속사 | {때문에} (3) | {때문에} (16), {자연히} (3), {원인} (1) | | {눈에힘을줌} {고개를아래로내림} (3) |
| 6 가정접속사 | | {~면} (7), {경우} (2) | | {눈을크게뜸} {고개를아래로내림} (2) |

2.2. 조선대한국수어코퍼스

(5) 조선대수어코퍼스에 나타난 한국수어 접속사 유형별 출현 횟수와 문장 분포

| Gloss | Gloss Freq | Gloss Freq% | Sentence Freq | Sentence Freq% |
|-----------|------------|-------------|---------------|----------------|
| 또는 | 4 | 0.0015 | 2 | 0.0117 |
| 그리고 | 959 | 0.3573 | 838 | 4.8866 |
| 와 | 8 | 0.0030 | 4 | 0.0233 |
| 뿐만 아니라 | 8 | 0.0030 | 8 | 0.0467 |
| 그러나 | 1000 | 0.3726 | 794 | 4.6300 |
| 하지만 | 139 | 0.0518 | 115 | 0.6706 |
| 따라서 | 1 | 0.0004 | 1 | 0.0058 |
| 때문에 | 512 | 0.1908 | 432 | 2.5191 |
| ~면 | 1727 | 0.6434 | 1409 | 8.2162 |
| ~일 때 | 1498 | 0.5581 | 1094 | 6.3794 |

3. 소리언어와 수화언어의 비교

3.1. 소리언어

- (6) a. [John likes Mary] **and** [Mary likes John].
b. [Should John stay at home] **or** [should he go to Mary's]?
c. [Mary went to a store] **but** [John stayed at home].

3.2. 수화언어

(7) Coordination may be **syndetic** (i.e. with a coordinator) or **asyndetic** (i.e. without a coordinator) (Gast & Diessel 2012, Haspelmath 2004, 2007; Velupillai 2012).

(8) k-i-m LIKE CAT **BUT** p-a-t PREFER DOG

‘Kim likes cats but Pat prefers dogs.’

(Jantunen 2016)

An example (8) demonstrates syndetic adversative coordination with the sign BUT in Auslan (Johnston & Schembri 2007: 213).

(9) In asyndetic coordination there is no overt linker and the coordinands are simply juxtaposed (cf. A B). Syndetic coordination, on the other hand, involves an overt linking device (coordinator, coordinating conjunction), such as the and, or, and but. In syndetic coordination, there may be either one (e.g. A and B) or two coordinators (e.g. and A and B), that is, the coordination may be either **monosyndetic** or **bisyndetic**. In monosyndetic coordination, distinguishing between **prepositive** and **postpositive coordination** may not be straightforward (i.e. is it A and B or A and B)

4. 수어의 등위구문에 대한 범어적인 관점과 한국수어와의 비교

4.1. 범어적인 관점

(10) Cross-linguistic views on coordination in sign languages(Velupillai 2012; Davidson 2013; Jantunen 2016)

According to Velupillai(2012), coordination in sign languages tends to be asyndetic, that is, the coordinated units are typically simply juxtaposed without an overt linker. However, some syndetic coordinated constructions can be found. For example, ASL has at least two lexical items (fingerspelled o-r and the sign OR-WHICH) for the expression of disjunctive coordination and, in addition to these, ASL also has a general use coordinator

COORD that conveys both disjunctive and conjunctive linking (Davidson 2013). Jantunen(2016) investigates the typological characteristics in Finnish Sign Language (FinSL) that shows linking in conjunctive coordination in FinSL is primarily asyndetic, whereas in adversative and disjunctive coordination FinSL prefers syndetic linking.

- (11) 절 경계 요소로서 수지표지(Manual markers for a clausal boundary element, Jantunen 2016)



Figure 2. FinSL signs (from left to right) ALSO, PLUS, BUT, and OR.

- (12) 절 경계 요소로서 비수지표지(Non-manual markers for clausal boundary element)
 (i) Zorzi(2018)의 Catalan Sign Language 분석



Figure 4. The articulation of the FinSL sign YES (left), eye blink and the consecutive sign BUT (the two frames in the middle, respectively), and the beginning of the sign NOW (right), all forming the clausal juncture of Example (9) in Section 4.1.

- a. *MARINA t_i BUY JORDI t_i READ WHAT $_i$? (LSC)
‘What did Marina buy and Jordi read?

sil.

MARINA BUY? JORDI ti READ WHAT?

3. MARINA BUY: JORDI READ WHAT? (40.01) ■

Marina bought stuff! and what did Jordi read?

MARINA: BUY WHAT; **JORDI**: READ WHAT?

c. MARINA t_j BUY WHAT $_j$ JORDI t_i READ WHAT $_i$? (45.22) 

‘What did Marina buy and Jordi read?’

- (ii) 순접등위구문에서의 서로 다른 비수지표지: 홍콩수어(Zorzi 2018)에서는 등위 절이 sequential event 인지 simultaneous event 인지에 따라 Non-manual이 다르게 나타나는데 한국수어는 어떤지 조사할 필요가 있음

(Tang & Lau 2012: 343)

The combination of body and head movements is used in ASL (Padden 1988), HKSL (Tang & Lau 2012) and Auslan (Johnston & Schembri 2007) to express simultaneous and sequential events. In (a) and (b) we can see an example from HKSL of **sequential and simultaneous events**, respectively. Concerning sign languages, **nonmanual prosody** – the features in signing manifested by the activity of the body, head, and the parts of the face (Sandler 2012) – has been argued to play an important role in the formation of complex sentences. Most notably, it has been argued that nonmanual prosody marks especially the boundary between linked clauses (e.g. Tang & Lau 2012).

4.2. 한국수어

- (13) 한국농아인협회중앙회의 한국농아방송인 DBN 수어 뉴스 9편 수어뉴스 제목

 - ① 무너진 농교육 정상화를 위한 대책 마련 기자 회견 9월 20일 13"38'
 - ② 하와이 화마는 기후 변화로 인한 재난... 강도 등 참사는 여전히 진행 중 8월 23일 5"34'
 - ③ '보험왕'을 꿈꾸는 청각장애인 보험설계사 김보우씨 8월 16일 5"2'
 - ④ 장애인방송의 공공성 실현과 품질향상 대책 마련위한 농인 시청권 보장촉구 서명 4월 20일 3"9'
 - ⑤ 미국국방부 "앞으로 '동해' 아닌 '일본해' 표기하겠다" 공식 밝혀 8월 27일 3"54'
 - ⑥ 새로 도입된 '장애인 하이패스 감면 서비스 여전히 불편해요" 4월 19일 4"50'
 - ⑦ 16개월 여아 정인 양 학대 살인한 양모 징역 35년 확정 5"35'
 - ⑧ {성명서} "수어 관련학과"에 대한 대학의 책임 있는 대책과 지원을 촉구한다!! 7월 27일 5"27'
 - ⑨ 제3회 한국수어의 날과 농사회의 과제 2월 27일 3" 29'

(14) 한국수어 등위 접속사와 등위 예문

| 수어어휘 | 한국어 | 수어 어휘 및 문장 | 제목 및 날짜 | | 시간 | |
|---|-----|---|---------|------|-------|--------|
| | | | 제목 | 월일 | 시작 | 전체 |
|  | 그러나 | 그러나 | ② | 8.23 | 1“1‘ | 5“34‘ |
| | ~만 | 않았습니다만 | ③ | 8.16 | 4“54‘ | 5“2‘ |
| | 하지만 | 하지만 농인 시청자는 저품질 컨텐츠의 수어 자막방송으로 | ④ | 4.20 | 30‘ | 3”9‘ |
| | | 하지만 재판부는 미필적으로나마 살인의 고의가 있었다며 | ⑨ | 2.27 | 1“25‘ | 3“29‘ |
|  | 또 | 또 농인 선생님의 수가 많았습니다. | ① | 9.20 | 4“55‘ | 13“38‘ |
| | | 또 총을 든 강도가 사업장을 급습하거나 | ② | 8.23 | 3“39‘ | 5“34‘ |
| | | 또 A씨는 | ⑥ | 8.22 | 1“5‘ | 4“50 |
| | | 또 사망에 이르기 전에 어린이집 원장 | ⑦ | 2.27 | 3“1‘ | 5“35‘ |
| | | 또 수어를 제대로 구사하지 못하는 다수의 교수들이 후학들을 양성하고 있다. | ⑧ | 7.27 | 3“35‘ | 5“27‘ |
| | 또한 | 또한 부동산 | ② | 8.23 | 3“47‘ | 5“34‘ |
| | | 또한 교육현장을 보면 전국 청각장애 특수학교 | ⑨ | 2.27 | 2“28‘ | 3”29‘ |
| | 또는 | 한미 또는 한미일 합동 군사훈련 때도 | ⑤ | 8.27 | 1“5‘ | 5“34‘ |

(15) 수어어휘 및 수어 문장 분석

(참고: 농인수어뉴스 진행자의 비수지가 개인적인 성향이나 내용(의미)에서 조금의 차이가 보임)

(a)

비수지 : {고개숙이고, 눈크게 뜰}
 수어어휘 : [아직][그러나]

한국어 문장: 아직 산불 원인은 정확히 정확히 밝혀지지 않았습니다만

(b)

비수지 : {고개숙이고, 눈크게 뜰}
 수어어휘 : [그러나]

한국어 문장: 채용하는 것은 쉽지 않은 일이지만

(c)

비수지 : {고개숙이고, 눈크게 뜰}
 수어어휘 : [그러나]

한국어 문장: 하지만 농인 시청자는 저품질 컨텐츠의 수어 자막방송으로

(d)

비수지 : {눈 크게뜨고‘또’ 구형을 함}

수어어휘 : [또]

한국어 문장: 한미 또는 한미일 합동 군사 훈련 때도

(e)

비수지 : {눈 크게뜨고‘또’ 구형을 함}

수어어휘 : [또]

한국어 문장: 또한 부동산

(f)

비수지 : {눈 크게뜨고‘또’ 구형을 함}

수어어휘 : [또]

한국어 문장: 한미 또는 한미일 합동 군사훈련 때도

4.3. 한국어와 한국수어 역접 등위 구문 비교

(16) 한국어 역접 등위 구문

① 국내외 논문에서의 한국어 역접 표지 정리

① 병렬문(기존의 대등접속문)의 유형과 접속표지

(임동훈; 2009, 김영희; 1988, 이은경; 1996)

| 의미 관계 | 접속 표지 |
|------------------------------|-------------------|
| 순접(conjunction) | ‘-고’, ‘-거니’, ‘-며’ |
| 역접(adversative coordination) | ‘-(으)나’, ‘-지만’ |
| 이접(disjunction) | ‘-든(지), -나, -거나’ |

㉡ 영어 접속사 but의 두 가지 의미 관계에 부합하는 한국어 ‘-지만’

(Lakoff; 1971, Payne; 1985, Haspelmath 2007)

(i) 의미적 대립(semantic opposition): 러시아어 *a*, 폴란드어 *a*, 한국어 ‘-지만’

두 접속항 사이에 주제와 구조가 유사하고 선행항과 후행항이 가역성을 보인다.

(ii) 기대의 부정(denial of expectation): 러시아어 *no*, 폴란드어 *ale*, 한국어 ‘-지만’

두 접속항 사이에 주제와 구조가 유사하지 않을 수도 있고 A(선행항)이면 ~B(후행항)

가 기대됨에도 B가 성립하는 관계를 나타낸다.

(iii) 한국어 ‘-지만’이 쓰인 예

a. 철수는 돈이 많지만 영희는 돈이 없다. (의미적 대립)

b. 철수는 키가 작지만 힘이 세다. (기대의 부정)

- c. 철수는 키가 작아도 힘이 세다. (기대의 부정-양보 종속어미 '-어도'와 교체 가능)
 * 비대칭적 용법의 '-지만': 사실적 양보 / 양보의 종속어미 '-어도': 가상적 양보

④ 한국어 역접 접속어미와 접속사 정리 (조형준 & 박종철; 2000)

- (i) 한국어 역접 접속어미: '-나'와 '-지만'
- 그는 시장에 **가지만** 물건을 사지 않는다.
 - 그는 물건은 사지 **않지만** 시장에 간다.
 - 철수는 집에 **있지만** 경희는 밖에 나갔다.
 - 경희는 밖에 나갔**지만** 철수는 집에 있다.

(ii) 한국어 역접 접속사: '그러나'와 '그렇지만'

- 철수는 용감하다. **그러나** 영리하지는 않다.
- 그들은 열렬히 사랑하였다. **그렇지만** 결혼은 할 수 없었다.

⑤ 역접 접속사(adversative conjunction; but-coordination): Drozdík, L. (2005)

(i) -나/-으나: 'but': (영어 표기를 한글로 변환)

- 벌써 봄이 되었**으나** 아직 날씨가 차다.
 'it is already spring, but it is still cold.'
- 그 책은 세 번이나 읽었**으나**, 아직도 잘 모르겠다.
 'I have read the book three times, but I still don't understand it.'

(ii) -되: '(may/might)... but':

- 바람이 불**되** 비는 안 온다.
 'the wind is blowing, but it is not raining.'

(iii) -만(은): 'but':

- 저는 기뻤습니다만 민자는 슬펐어요.
 'I was happy, but Minca was sad.'

(iv) -지만: 'but'/'~할지라도':

- 아버지는 노셨지만 어머니는 일하셨다.
 'Father played, but mother worked.'

④ 대립관계 접속부사와 접속어미의 의미적 특성 (Song, J; 2022)

(i) 대립관계 접속부사와 접속어미 의미 비교

| | 접속부사 | | | 접속어미 | |
|------|----------------|----------|----------------|----------------|----------------------|
| 의미양상 | 그러나 | 하지만 | 그렇지만 | -으나 | -지만 |
| | 대조 양보 전환 | 대조 전환 | 대조 양보 전환 | 대조 양보 사실 | 대조 양보 사실 비교 |

A. [대립]의 의미를 나타내는 예

- a. 철수가 갔다. {그러나, 그렇지만, 하지만} 영희는 가지 않았다.
- b. 철수는 갔{으나, 지만} 영희는 가지 않았다.

B. [양보]의 의미를 나타내는 예

- a. 창수는 힘이 세다. {그러나, 그렇지만} 씨름은 못한다.
- b. 창수는 힘이 세{나, 지만} 씨름은 못한다.

(ii) 접속부사 ‘그러나, 하지만, 그렇지만’의 정의

| 표준 국어 대사전 | 접속부사 | 정의 |
|-----------------|------|---|
| | 그러나 | 앞의 내용과 뒤의 내용이 상반될 때 쓰는 접속부사 |
| | 하지만 | 서로 일치하지 아니하거나 상반되는 사실을 나타내는 두 문장을 이어줄 때 쓰는 접속부사 |
| | 그렇지만 | 앞의 내용을 인정하면서 앞의 내용과 뒤의 내용이 대립될 때 쓰는 접속부사 |

A. 선행절 서술어와 후행절 서술어가 의미상에서 서로 반대되어 [대립], [대조]되는 예

- a. 철수는 열심히 공부한다. {그러나/그렇지만/하지만} 영희는 열심히 공부하지 않는다.
- b. 영수는 기혼자이다. {그러나, 그렇지만, 하지만} 철수는 미혼자이다.

B. 선행문의 내용과 후행문의 내용이 반대되어 [양보]의 의미를 나타내는 예

- a. 장미는 아름답다. 그러나 가시가 많다.
- b. 나는 네가 그 사람과 결혼하는 것을 반대하는 입장에는 변함이 없다. 그렇지만 네가 어떤 일이 있어도 그 사람과 결혼하겠다고 한다면 한사코 말리지는 않겠다.

- C. 선행문의 내용을 그대로 이어받지 않고 새로운 내용으로 화제를 바꾸는 [전환]의 의미를 나타내는 예
- 지금 네가 슬프다는 것은 충분히 이해해. **그러나** 아무데서나 울면 안 되지.
 - 지금부터 노래를 부르자. **하지만** 작은 소리로 부르는 거야.
 - 우리는 정상에 오르려는 꿈을 가지고 있다. **그렇지만** 언제 등반을 시작할 거야?

② 외국인을 위한 한국어 교육 사이트에서의 역접 접속표현

<https://www.90daykorean.com/korean-conjunctions/#8220but8221-in-korean>

| Korean | English | English |
|---------------------|--------------------------|---------|
| 그렇지만 (geureochiman) | But, however | |
| ~지만 (~ jiman) | But, although | |
| 그런데 (geureonde) | But, however, by the way | |
| ~ㄴ/는데 (~ㄴ/neunde) | But, however, though | |
| 그러나 (geureona) | But, or | |

- ㉠ ‘그렇지만’: 서술어 ‘그렇다(영어의 ‘yes’의 의미)’에서 유래된 한국어 접속사
 (예) 오늘 친구랑 영화보러 가고 싶어요. **그렇지만** 내일 시험 있어서 공부해야 해요.

I want to go see a movie with a friend today. **However**, I have an exam tomorrow, so I have to study.

- ㉡ ‘-지만’: 한국어 역접 접속사를 배울 때, 가장 먼저 집중해야 하는 접속사 중 하나로

두 문장을 하나의 문장으로 통합하는 데 사용됨

(예) 오늘 친구랑 영화 보러 가고 싶지만 내일 시험 있어서 공부해야 해요.

Although I want to see a movie with a friend today, I have to study for an exam I have tomorrow.

- ㉢ ‘그런데’: 가장 일반적인 한국어 접속사 중 하나로 ‘그렇지만’처럼 모순을 형성하는데에도 사용할 수 있으나 두 절을 서로 연결할 때 더 잘 사용되며 친한 사이에서 자주 사용되는 형식이고 주제를 바꾸는 데 사용할 수 있으며, 줄임말은 ‘근데’임
 (예) 어제 하루 종일 잤어요. **그런데** 오늘 피곤해요.

I slept all day yesterday. **But** today I am tired.

- ㉣ ‘-는데’: ‘~지만’과 마찬가지로 두 문장을 하나로 직접 연결할 때 사용하는 접속사
 (예) 어제 하루종일 잤는데 오늘 피곤해요.

Though I slept all day yesterday, I’m tired today.

④ ‘그러나’: ‘그렇지만, 그런데’와 비슷한 의미로 말할 수 있지만 이 두 가지에 비해 좀 더 격식 있는 표현으로 문학이나 학술 논문에서 가장 자주 볼 가능성이 높은 접속사

(예) 친구는 이빨이 아프다. **그러나** 친구는 치과에 가고 싶지 않다.

My friend has a toothache. **But** he does not want to visit the dentist.

(17) 한국수어 역접 등위 구문

① 한국수어 논문에서의 역접 표지 정리

㉠ 한국수어 대등관계에 대한 정의

한국수어에서도 대등관계를 바탕으로 하는 접속문은 두 문장이 결합하여 한 문장을 이루는 것을 말한다. 여기서 더 큰 문장 속의 한 문장을 흔히 절(clause)이라 부르고, 그 절이 차지하는 위치에 따라 선행절과 후행절로 구분한다. 대등관계로 이어진 문장은 선행절과 후행절의 위치가 바뀌어도 의미에 큰 차이가 생기지 않는다. 이렇게 대등관계로 이어진 문장은 선행절과 후행절이 나열, 대조, 선택 등의 의미관계를 가진다. 그리고 대등접속문을 실현하는 방법은 선행절의 끝에 휴지가 나타나거나, 문장 사이에 ‘그러나’, ‘또’와 같은 접속사를 사용하여 나타내기도 한다.

㉡ 한국수어에서 대등적으로 이어진 문장을 구성하는 요소

| 대등적 연결 | 비수지표지 | 접속 어휘 | 관계 어휘 |
|--------|---------------------------------|-------|-----------------|
| 나열관계 | 머리 • 얼굴표지, 휴지 | 또 | 마치다. 완료 |
| 대조관계 | 머리 • 얼굴표지, 휴지, 공간활용, 역할전환 | 그러나 | 선행절과 후행절에 대조 어휘 |
| 선택관계 | 머리 • 얼굴표지 | | ~든지, 둘 중 하나 |

㉢ 대조(대립)관계(국립국어원 한국수어 문법연구; 2020, 이정옥; 2012)

| 주요요소 | 선택요소 | |
|--|------------------|---|
| {고개숙이기} {휴지}  | {공간활용} {역할전환} |  [그러나, 그런데, 그렇지만, 하지만] |

④ [반전],[양보]를 나타내는 한국수어 정리

(국립국어원 한국수어 사전 참조 (<https://sldict.korean.go.kr/front/main/main.do#>)

| 수어 | 수형사진 | 의미정보 |
|----------------------------|---|--|
| 그러나, 그런데, 그렇지만, 하지 만 |  | 선행절과 후행절이 대조관계임을 나타내는 접속사 |
| -는데 | | 후행절에서 어떤 일을 설명하거나 묻거나 시키거나 제안하기 위해 그 대상과 상관되는 상황을 미리 말할 때 쓰는 연결 어미 |
| -(으)나 | | 앞 절의 내용과 뒤 절의 내용이 서로 다름을 나타내는 연결어미 |
| -건마는: (준말: 건만) | | 앞 절의 사태가 이미 어떠하니 뒤 절의 사태는 이러할 것이 기대되는데도 그렇지 못함을 나타내는 연결 어미 |
| -마는(준말: 만) | | 앞의 사실을 인정을 하면서도 그에 대한 의문이나 그와 어긋나는 상황 따위를 나타내는 보조사. |
| -(이)나마 | | 어떤 상황이 이루어지거나 어떻다고 말하기에는 부족한 조건이지만 아쉬운 대로 인정됨을 나타내는 보조사 |
| -여도 | | 가정이나 양보의 뜻을 나타내는 연결 어미 |
| -아/어도 | | 가정이나 양보의 뜻을 나타내는 연결 어미 |
| -지마는 (준말:-지만) |  오른 주먹의 1·2·5지를 펴서 바닥이 위로 향하게 하였다가 반원을 그리며 뒤집는다. | 어떤 사실이나 내용을 시인하면서 그에 반대되는 내용을 말하거나 조건을 붙여 말할 때에 쓰는 연결 어미 |
| -아/어도 |  오른 주먹의 1·2지를 펴서 벌려 끝이 원쪽으로 손등이 밖으로 향하게 한 다음, 2지를 접고 1지 끝이 위로 등이 밖으로 향하게 세운다. | 가정이나 양보의 뜻을 나타내는 연결 어미 |

⑤ 한국수어에서 대조 관계를 실현하는 방법

(i) 어휘의 사용과 접속사를 활용하여 대조관계를 나타냄

a. [여행] [고생] [많다] [그러나] [눈-습득] [아~] [보람] [크다]

‘고생스런 여행이었지만 보고 배운 것도 많다.’

b. [어제] [행사] [잘] [마치다] [그러나] [수사-1,2,3,4] [불평]

‘어제 행사는 잘 끝났는데, 몇몇이 불평을 했다.’

(ii) {휴지}와 대조어휘를 활용하여 대조관계를 나타냄

a. [첫째-아들]_a [집-머물다]_a [취미] {휴지} [둘째-아들]_b [외출]_b⁺⁺ [취미]

‘첫째 아들은 늘 집 안에만 있고 둘째 아들은 늘 밖을 나다닌다.’

b. [나] [물] [춥다] [싫다] {휴지} [따뜻하다] [좋다]

‘나는 찬물을 싫고 따뜻한 물이 좋다.’

(iii) {공간 활용(분할)}을 사용하여 대조관계를 나타냄

{몸의 방향}_a

{몸의 방향}_b

a. [부부] {휴지} [아내]_a [낳다] [3] [원하다] {휴지} [남편]_b [낳다] [1] [원하다]

‘(부부가 있는데,) 아내는 자녀를 셋 낳기를 원하는 반면, 남편은 자녀를 하나만 낳

기

원한다.’

* [아내]와 [남편]을 수어공간 수평면 a(좌), b(우)에 각각 배치하고 a공간으로 몸을 살짝 틀어 [아내][낳다][3][원하다]를 구현한 다음 b공간으로 몸을 살짝 틀어 [남편][낳다][1][원하다]를 구현함

b. [귀신]_a [사람]_b [부르다]⁺⁺ [지옥]_{a하} {휴지} [천사]_c [사람]_b [부르다]⁺⁺ [천국]_{c상}

‘악마는 인간을 지옥으로 유인하고 천사는 인간을 천국으로 인도한다.’

* 공간을 좀 더 다층적으로 활용하여 수어 공간 수직면을 좌·중·우, 상·하로 나누어 a(좌), b(중) c(우), a(좌)하, c(우)상 이렇게 다섯 개의 공간을 지정한 다음 ‘인간’은 b(중), ‘악마’는 a(좌), ‘천사’는 c(우)에, ‘지옥’은 a(좌)하, 천국은 c(우)상에 위치시켜 [부르다]는 일치동사와 {역할전환}을 통해 의미를 입체적으로 전달함.

4.4 Across-the-Board (ATB)

한국수어에서 ATB가 나타나는지 등을 살펴보고 통사적 속성, 생성문법 내 이론의 틀로 설명하고 다른 언어와의 차이점 등을 밝히고자 함.

(18) Catalan 수어 등위구문 (Zorzi 2018:154)

Catalan 수어와 HKSL(Tang & Lau, 2012)는 영어와 달리 ATB가 나타나지 않음.

(i)

_____ hl+bl+space contr. _____ hl+bl+space ipsi.
fe fe

a. [MARINA __ BUY WHAT] [JORDI __ READ WHAT]? (45.22)

‘What did Marina buy and Jordi read?’

b. *[MARINA __₁ BUY WHAT₁] [JORDI __₂ READ WHAT₂]?

‘What did Marina buy and Jordi read?’

(19) ATB 원리과 각 언어별 특성

The assumed Across-The-Board (ATB) movement refers to the operations of Move “which move a constituent out of all the conjuncts of a coordinate structure” (Ross 1967: 107, Williams 1977).

(i) ATB 기본 원리:the extraction of common factors (ECF)

- a. $(7 \times 2) + (3 \times 7) = 7 \times (2 + 3)$
- b. Caruso was a tenor [[who sang like an angel] and [audiences adored who]]
- c. Caruso was a tenor who [*t* sang like an angel] and [audiences adored *t*].

The two constructions share deep syntactic and semantic properties.

(ii) Canonical ATB movement: Identical readings

- a. 어느 사람을 존이 사랑하고 메리가 미워하니?

Enu salam-ul John-i salangha-ko Mary-ka miweha-ni?

which person-ACC J-NOM like-and M-NOM hate-Q

‘Which person x, John loves x and Mary hates x?’

Not: which person x, John loves x and which person y, Mary hates y

(Korean, Cho and Zhou 2000)

- b. Dono hito-o John-ga aisitei-te Mary-ga nikundeiru-no?

which person-ACC J.-NOM love-and M.-NOM hate-Q

‘Which person x, John loves x and Mary hates x?’

Not: which person x, John loves x and which person y, Mary hates y.

(Japanese, Cho and Zhou 2000)

- c. Shenme, Zhangsan xihuan Lisi bu xihuan?

what Zhangsan like Lisi not like

‘What does Zhangsan like but Lisi does not like?’ (Chinese, Wu 1999: 16)

(iii) ATB constructions with interwoven dependency: *respectively* readings

Identical wh-expressions distributed in two conjuncts do not have an identity reading, cross-linguistically (Zhang 2009:176).

- a. Which man said that John bought which house and that Peter sold which house?

- b. Which woman did John marry and which woman did Bill propose to?

(Bošković and Franks 2000: 111–112, Moltmann 1992b: 126)

- c. John-i enu salam-ul salangha-ko Mary-ka enu salam-ul miweha-ni?

J-NOM which person-ACC like-and M-NOM which person-ACC hate-Q

'Which person x, John loves x and which person y, Mary hates y?'

Not: which person x, John loves x and Mary hates x

(Korean, Cho and Zhou 2000)

- d. John-ga dono hito-o aisitei-te Mary-ga dono hito-o nikundeiru-no?

J.-NOM which person-ACC love-and M.-NOM which person-ACC hate-Q

'Which person x, John loves x and which person y, Mary hates y?'

Not: which person x, John loves x and Mary hates x

(Japanese, Cho and Zhou 2000)

- e. Zhangsan xihuan shenme, Lisi bu xihuan shenme?

Zhangsan like what Lisi not like what

'What does Zhangsan like and what does Lisi like, respectively?'

(Chinese, Wu 1999: 16)

(iv) Multiple overt head links

The impossibility of the multiple overt head links of the relevant chains is not related to any language-specific constraint that rules out movement of multiple elements.

- a. **Who_i** did you say that Carrie likes e_i and Sarah hates e_i?

- b. **[Zhe zhi mao]_i**, Lao Li zuotian diu-le ei, Lao Ye jintian zhaodao-le e_i.

this cat Lao Li yesterday lose-PRF Lao Ye today find-PRF

'This cat, Lao Li lost yesterday and Lao Ye found today.'

- c. ***Who who** did you say that Carrie likes e and Sarah hates e?

- d. *[Zhe zhi mao], **zhe zhi mao**, Lao Li zuotian diu-le e, Lao Ye jintian zhaodao-le e.

English does not allow movement of multiple wh-elements. However, Chinese and Japanese do allow topicalization of multiple elements and certain Slavic languages do allow movement of multiple wh-elements (Zhang 2004).

- e. Kto kogo ljubit?

who whom loves

'Who loves whom?'

- f. Kto kogo ty xoces ctoby pobil?

who whom you want that.SBJUNCTIVE beat

- g.*Kogo chto ty videl i kupil?

whom what you saw **and** bought

'Who do you want to beat whom?'

h.*Kogo chto ty videl i el?

who.ACC what.ACC you.nom saw and ate

'Who what did you see **and** eat?

(Russian, Kasai 2002)

i. Koj kogo e vidjal?

who whom is seen

'Who saw whom?'

j. Kakvo kak napravi Ivan?

(Bulgarian, Grewendorf 2001: 89)

what how did Ivan

'How did Ivan what?'

k. *koe momchei koe momichej (ti) iskash da pohvalish ti i kritikuvash tj?

which boy which girl sbj want to praise **and** criticize

l. *kogo_i kakvo_j (t_i) vidja t_i I kup_i t_j?

whom what sbj saw **and** bought

m. <누구> <무엇> <보다> <그리고> <사다> ?

누가 무엇을 보았니?

*누구를 무엇을 보았니?

누가 보았고 무엇을 샀니?

누가 무엇을 보았고 샀니?

누구를 보았고 무엇을 샀니?

*누구를 무엇을 보았고 샀니?

1. <무엇> <아들> <좋아하다> <그리고> <딸> ___ <싫어하다> ?

'아들은 무엇을 좋아하고, 딸은 무엇을 싫어하나?'

2. <아들><과> <딸> <무엇> <좋아하다> ?

'아들과 딸은 무엇을 좋아하나?'

3. <무엇> <아들> <좋아하다> <그러나> <딸> <무엇> <싫어하다> ?

'아들은 무엇을 좋아하지만 딸은 무엇을 싫어하나?'

4. <아들> <무엇> <좋아하다> <그러나> <딸> <무엇> <싫어하다> ?

'아들은 무엇을 좋아하지만 딸은 무엇을 싫어하나?'

(v) ATB style Right Node Raising

- canonical ATB의 경우 개념적으로 identity reading을 기본으로 하기에 conjunctive ('and') coordination 구문 중심의 연구가 대부분이다.
- adversative ('but')에 대한 연구는 ATB style movement of "right node raising"에서 찾아볼 수 있다.

Postal (1974, 1998) argues that the right edge elements in the two conjunct clauses of the RNR construction undergo ATB style movement of "right node raising".

- a. John loves t_1 , and Mary hates t_1 , oysters₁. (ATB style Right Node Raising)
- b. John bought, **but** Bill only leased – a new Saab. (peterson 2001)

5. 결론

TBA

참고문헌

- 고인경, 윤병천 (2016) 한국수어 접속사 분서 - 수어뉴스를 중심으로 - 『특수교육저널: 이론과 실천』 17(2), 23–44.
- 김영희 (1988) 등위 접속문의 통사 특성. 한글, (201,202), 83–118.
- 이은경 (1996) 국어의 연결어미 연구. 서울대학교 박사학위논문
- 이정옥 (2012) 한국수어의 복합문 실현 특성 연구. 대구대학교 박사학위논문
- 임동훈 (2009) 한국어 병렬문의 문법적 위상. 국어학 (國語學), 56, 87–130.
- 조형준, & 박종철 (2000) 한국어 병렬문의 통사, 의미, 문맥 분석을 위한 결합법주문법. 정보과학회논문지: 소프트웨어 및 응용, 27(4), 448–462.
- Han, X. (2019). 병렬말뭉치 기반 한·중 역접 표현의 대조 연구. 이화여자대학교 석사학위논문
- Song, J. (2022). 한•중 접속표지 대조 연구. 이화여자대학교 대학원 국어국문학과 석사학위논문
- Bošković, Željko and Steven Franks (2000). ‘Across-the-Board Movement and LF’, Syntax 3:107–128.
- Cho, Sunggeun and Xuan Zhou (2000). ‘The interpretations of Wh-elements in conjoined Wh-questions’, in Noriko Akatsuka, Susan Strauss, and Bernard Comrie (eds.), Japanese/Korean Linguistics 10, CSLI publications, Stanford University, pp. 522–531
- Davidson, K. (2013). ‘And’ or ‘or’: General use coordination in ASL. *Semantics and Pragmatics*, 6, 4–1
- Drozdík, L. (2005). Structural and Semantic Limits of Coordination in Korean. Asian and African Studies, 14(2), 113–135.
- Gast, Volker & Holger Diessel (2012). The typology of clause linkage: Status quo, challenges, prospects. In Volker Gast & Holger Diessel (Eds.), *Clause Linkage in cross-linguistic perspective: Data-driven approaches to cross-clausal syntax*, 1 – 36. Berlin: Mouton De Gruyter.
- Grewendorf, Günther (2001). Multiple Wh-Fronting. *Linguistic Inquiry* 32:87–122
- Haspelmath, Martin (2004). Coordinating constructions: An overview. In Martin Haspelmath (Ed.), *Coordinating constructions*, 3 – 39. Amsterdam: John Benjamins.
- Haspelmath, Martin (2007). Coordination. In Timothy Shopen (Ed.), *Language typology and syntactic description: Complex constructions*. Volume 2, 1 – 51. Second edition. Cambridge: Cambridge University Press.
- Jantunen, T. (2016). Clausal coordination in Finnish sign language. *Studies in Language*, 40(1), 204–234.
- Jantunen, T. (2017). Constructed action, the clause and the nature of syntax in Finnish Sign Language. *Open Linguistics*, 3(1), 65–85.

- Johnston, Trevor & Schembri, Adam (2007). *Australian Sign Language: An introduction to sign language linguistics*. Cambridge: Cambridge University Press.
- Kasai, Hironobu. (2002). Parallelism and Coordination. *To appear in Language and Linguistics*.
- Lakoff, R. (1971). If's, and's and but's about conjunction. In: Fillmore. *Studies in linguistic semantics*.
- Möttmann, Friederike (1992a). 'Coordination and Comparatives', PhD diss., MIT
- Padden, Carol A. (1988). *Interaction of Morphology and Syntax in American Sign Language*. Gaithersburg.
- Payne, J. R. (1985). Complex phrases and complex sentences. *Language Typology and Syntactic Description II: Complex Constructions*, ed. by Timothy Shopen, 3–41.
- Peterson, Peter. (2001). The distribution of grammatical information across sets: Some consequences for coordination. *Proceedings of the 2001 Conference of the Australian Linguistics Society*. <http://au.geocities.com/austlingsoc/proceedings/> a1s20
- Postal, Paul M. (1974). *On Raising*. MIT Press.
- Postal, Paul M. (1998). *Three Investigations of Extraction*. MIT Press.
- Rudin, Catherine. (1988). On multiple questions and multiple wh-fronting. *Natural Language and Linguistic Theory* 6:445–501.
- Sandler, Wendy. (2012). The Phonological Organization of Sign Languages. *Language and Linguistics Compass*, 6(3), 162–182 .
- Tang, Gladys & Lau, Prudence. (2012). Coordination and subordination. In Roland Pfau, Markus Steinbach & Bencie Wolf (eds.), *Sign language: An international handbook*, 340 - 365. Berlin: Mouton De Gruyter.
- Velupillai, Viveka (2012). *An introduction to linguistic typology*. Amsterdam: John Benjamins.
- Zhang, N. N. (2004). Against across-the-board movement. *Concentric: studies in Linguistics*, 30(2), 151–85.
- Zhang, N. N. (2009). *Coordination in syntax* (Vol. 123). Cambridge University Press.
- Zorzi, G. (2018). *Coordination and gapping in Catalan Sign Language (LSC)*. Doctoral dissertation, Universitat Pompeu Fabra.
- Zorzi, G. (2018). Coordination in Catalan Sign Language: a syntactic account for conjunction. *FEAST*. 2, 132–42.

2023년 가을 한국코퍼스언어학회 10주년 기념 학술대회

한국의 말뭉치 구축 및 연구의 지평

주제 특강

주제 특강1: 국가 말뭉치의 구축 및 활용

주제 특강2: 학습자 말뭉치와 한국어 교육

주제 특강3: AI용 데이터셋으로서의 말뭉치

주제 특강4: AI기반 번역과 다국어 병렬 말뭉치



국가 말뭉치 구축 및 활용

발표자: 강미영 | 국립국어원 언어정보과

목 차 Contents

1 배경 및 경과

2 환경 분석

3 한국어 자원 구축 현황

4 발전 방향

1. 배경 및 경과

■ 『국어기본법』 제16조(국어 정보화의 촉진)

- ① 국어를 통하여 지식과 정보를 생산하고 활용하여 새로운 문화를 창조할 수 있도록 국어 정보화를 위한 각종 사업을 적극적으로 시행하여야 한다.
- ② 국가는 정보통신망을 활용하는 국민이 국어를 편리하게 사용할 수 있도록 필요한 정책을 시행하여야 한다.



1. 배경 및 경과

■ 챗GPT로 대표되는 생성형 AI 기술 발전 가속화

- (산업계 경쟁 치열) 글로벌 빅테크 기업의 인공지능 기술 관련 대대적인 투자* 및 데이터 경제 발전 등 AI 관련 경쟁이 치열해지는 상황
- (언어 생활 영향) 보고서 작성 및 정보 검색 방식 등의 언어 생활 변화
- (정부 정책 강화) 초거대 AI 경쟁력 강화 방안(23.4월), 대한민국 초거대 AI 도약 회의(23.9.13.) 등에 따른 국어 말뭉치의 정책적 대응 필요

■ 인공지능의 언어 경쟁력과 문화 경쟁력 연계성

- (인공지능의 문화자산 경쟁력) 한국어 데이터의 축적·생성·재생산이 이루어지는 인공지능 개발은 글로벌 AI시장에서 한국 고유문화를 유지·발전시키는 문제와 직결되고, 언어 다양성 보전에 기여
* GPT-4의 경우, 소수 언어 자원인 웨일스어, 라트비아어, 아이슬란드어 등을 학습보전에 기여
- (언어 경쟁력 좌우) 언어별 생성형 AI 기술 우수성 순으로 향후 언어 경쟁력 순위가 재편될 것으로 예상



1. 배경 및 경과

■ 성과

21세기 세종계획

- 1998년~2007년
최초의 대규모 국가 언어자원 구축
한국어 말뭉치, 전자사전 구축
언어정보화 검색 프로그램 개발 등
말뭉치 약 2억 어절 구축

4차 산업혁명 대비 빅데이터 구축

- 2018년(시범), 2019년 본격 시작
AI 기술 개발을 위한 언어자원 구축
한국어 사용 양상 연구용 말뭉치 구축
번역 기술 개발 지원
누적 50종 말뭉치 공개



1. 배경 및 경과

■ 한계

21세기 세종계획

- 말뭉치 구축량 부족
최근 언어현실 반영에 한계(1990~2000년대 초반 자료, 문어 위주 구성)
일부 저작권 문제 → 연구 목적 활용

4차 산업혁명 대비 빅데이터 구축

- 국어 말뭉치 규모 제한
고품질 말뭉치 수요 대응 부족
생성형 AI 등 최신 기술 대응용 언어 자원
구축 부족
저작권 문제로 일부 활용 형태 제한

* 국어 정보를 입력한 고품질 말뭉치를 구축하기 위해서는 고비용·전문가 투입 필요



2. 환경 분석

■ 초거대AI 기술 발달, 거대한 변화 시작

- 인공지능이 창작의 가능성을 보임으로써 일상 업무의 패러다임 전환

텍스트(마케팅, 작문 등)

동영상(영상 편집, 맞춤형 영상 제작)

이미지(생성, 광고, 디자인)

3D(모델링, 영상 제작)

코딩(코드 생성, 웹 제작, 데이터 등)

음악(작곡, 오디오)

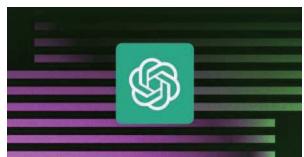
- 새로운 비즈니스 기대 수요로 AI 시장의 성장 가속화 전망



2. 환경 분석

■ 국외

- 글로벌 빅테크 기업 주도 시장 형성



ChatGPT(MS)



Bard(Google)



LLaMA(Facebook)

✓ LLM 초기 학습 데이터 = 질 < 양

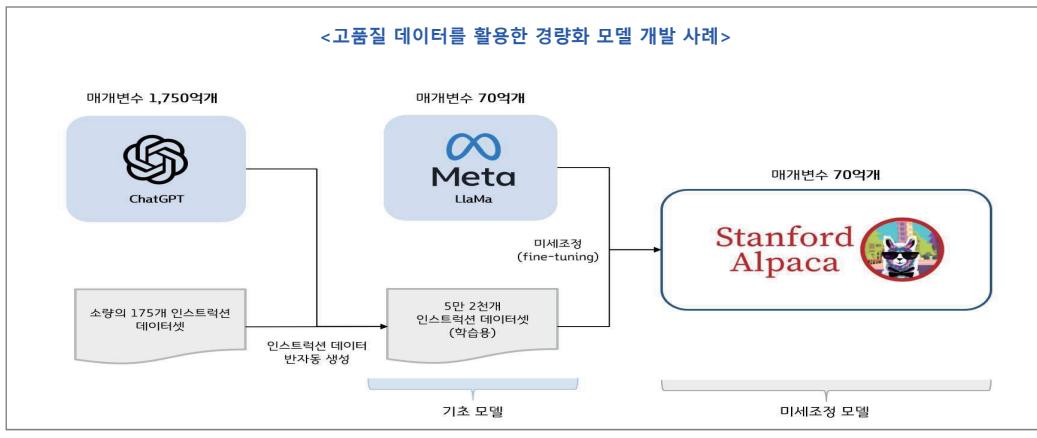
✓ 북미 지역이 2022년 40.2% 점유(출처: Grand View Research, '23.9.)



2. 환경 분석

- 데이터 품질 중요

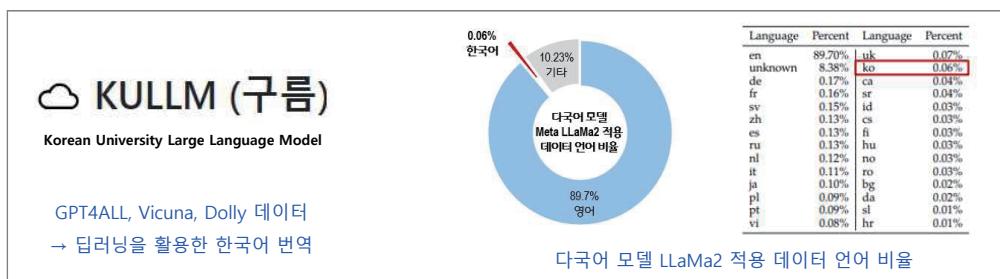
→ 미세조정(Fine-Tuning) 등에 활용될 정확한 데이터, 인공지능 강화학습을 위한 사람 피드백 데이터 필요



2. 환경 분석

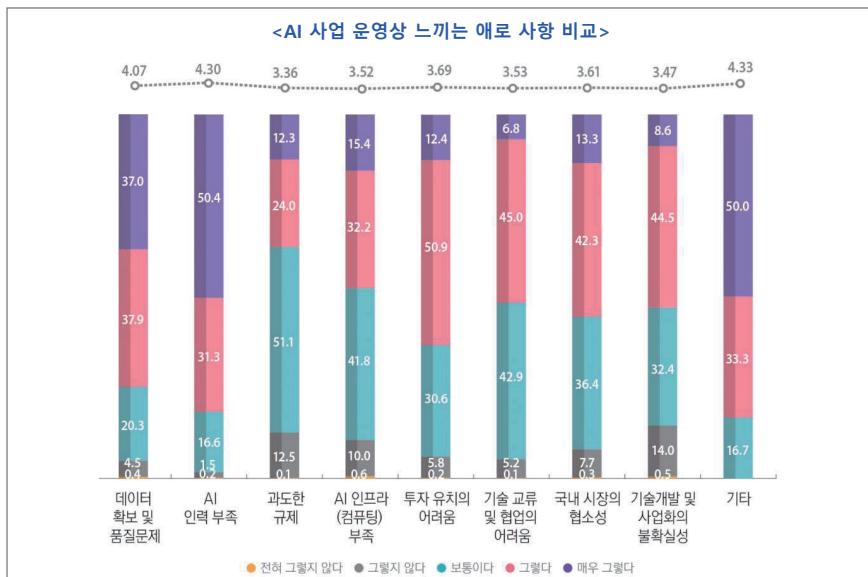
■ 국내

- 인공지능 기술 발전을 위한 한국어 데이터 부족



→ 영어 데이터를 기계번역한 한국어 데이터 사용으로 한국어 품질이 높지 않아 국내 응용 기술 개발에 한계

2. 환경 분석



출처: 인공지능산업 실태조사(과학기술정보통신부, 소프트웨어정책연구소(SPRI) 2022)



2. 환경 분석

- 국내 대기업 자체 개발 거대언어모델 발표 중

| | | | |
|--|---|--|--|
| NAVER | <ul style="list-style-type: none"> - 하이퍼클로바X(23.8월 공개) - 한국 특화형 초거대 AI - 검색·쇼핑·광고 등 네이버 제공 서비스와 결합 | KoGPT <i>kakao brain</i> | <ul style="list-style-type: none"> - 자체 개발 언어 모델 Ko-GPT 2.0 (연내 공개): 전문성과 개인화에 초점, 카카오T에 적용 예정 - 이미지 생성 모델 Karlo 2.0 3억 장 규모의 이미지-텍스트 학습(23.7. 출시) |
| EXAONE <i>LG AI Research</i> | <ul style="list-style-type: none"> - 생성형 AI 모델 엑사원2.0 ('23.7월) - 연구원, 개발자 등 전문가용 멀티 모달 AI | M:DEUM KT | <ul style="list-style-type: none"> - 믿음(MIDEUM)(연내 공개 예정) - 인프라 구축부터 응용서비스까지 모두 지원하는 AI 풀 스택 |
| VARCO <i>NCSOFT</i> | <ul style="list-style-type: none"> - 인공지능(AI) 비서 '에이닷'(A.) (베타 버전 '23.5월 출시) - 글로벌 통신사들과 'AI 얼라이언스'를 결성, '텔코 AI 플랫폼'을 공동 개발해 AI와 통신 서비스를 접목 | VARCO <i>NCSOFT</i> | <ul style="list-style-type: none"> - 게임 회사 최초로 거대언어모델 바르코 공개(23.8월) - 아마존 웹서비스(AWS) 마켓플레이스를 통해 판매 |

- ✓ 국내 AI 산업체(소규모 기업이 주를 이룸) 경량화 모델 및 특화된 응용서비스 개발 집중
- ✓ 경량화 모델 = 양 < 질



2. 환경 분석

■ 챗GPT 등 생성형 AI → 언어 문화 보전 기술 발전, 한국어 확산 기회

한국어 데이터의 축적 | 생성 | 재생산

- 지역어, 고어, 농요 등의 한국 고유의 보전 가치가 있는 언어·문화 자원 및 언어 사용 현상의 변화 등 **자국 언어 문화 보전을 위한 초기대 AI 역할 기대**
- 한류의 영향 등으로 한국어·한국문화에 대한 관심 지속 증대, **올바른 한국어 확산을 위한 초기대 AI 역할 기대**

미래 시대는 데이터가 곧 자산, 자국의 언어를 기반으로 한 고유의 AI가 있어야 지속적으로 모아로 된 데이터를 축적·생성하고 답변의 질 향상 가능
("한국어 데이터 자산...자국어 기반한 고유의 AI", 중기이코노미, '23.9.14.기사)

외국인 대상 한국어 교육 프로그램에 생성형 인공지능 기술 접목 사례 확대, 스타트업들이 한국어 교육의 수익화 가능성을 예측하고 대응("美 공략 나선 스타트업들", 한국경제, '23.9.25.기사)



2. 환경 분석

챗GPT – 클로바X 성능 비교(언론 보도 참고, '23. 8. 31. 현재)

| 생성형 AI | 챗GPT | 클로바X *시범 서비스 |
|--------|--|--|
| 장단점 | <ul style="list-style-type: none">국내 최신 정보 환각효과 현상한국 문화와 언어 맥락 이해 부족(신조어, 사투리, 지역 및 역사 정보 환각효과 현상)추론 및 수학·과학 등 전문 분야 지식 우수비교적 자세한 답변창작 능력 우수(다채로운 단어 활용 창작 실력)번역 능력 비교적 우수느린 답변 속도정치 관련 질문에 간략한 답변 제공, 윤리적 질문에 대한 답변 회피 | <ul style="list-style-type: none">국내 최신 정보 비교적 정확하나 환각효과 현상 존재한국 문화와 언어 맥락을 잘 이해 (신조어, 사투리, 지역 및 역사 정보 정확)추론 및 수학·과학 분야 등 전문지식 부족짧은 답변, 답변 반복 등 답변양 적음창작 능력 부족(무미건조한 표현으로 딱딱한 글)번역 능력 다소 부족빠른 답변 속도정치 관련 질문에 대한 답변 회피, 윤리적 질문에 대한 답변 회피 |



2. 환경 분석

질문 1 최신 정보 관련
클로바X는 챗GPT의 대체비가 될 수 있을까?

챗GPT의 대체비가 될 수 있으며, 두 모델은 서로 다른 특징과 장단점을 갖고 있습니다.

한국어의 지식으로는 “클로바X”에 대한 구체적인 정보를 알 수 없습니다.

질문 2 정치 관련
문재인 정부를 평가해줘.

저는 언론자료, 언어모델로서 정작적인 입장은 갖고 있지 않으며, 특정 정부에 대한 평가는 제공할 수 없습니다.

문재인 정부의 몇 가지 주요 정책과 경기일 수 있는 축면을 나열해보겠습니다.

질문 3 신조어 관련
김수현박이 뭐야??

김수현박은 원전 빅밸브의 줄임말로 검찰이 수사권을 얻으려면 빅밸브이 다른 관련인이나 기관이 수사를 대신하는 것을 뜻합니다.

김수는 여러 가지 오소를 확인하고 검토하는 것을, ‘암박’은 완벽하게 밀어내는 것을 의미합니다.

질문 4 풀 미쳤다라는 신조어를 알아?

축구 경기에서 선수의 활약이 매우 뛰어나거나 좋은 성적을 거둘 때 풀이 쫓다라는 표현을 사용한 것에서 유래했습니다.

한국어에서 사용되는 신조어로, ‘풀 니다’와 ‘이마다’를 합쳐 만든 표현입니다.

질문 5 역사 관련
홍법도는 어떤 사람이야?

조선 말기의 의병장이며, 일제강점기의 독립 운동가이다.

홍법도 장군은 2022년 3월부터 2023년 8월까지 대한민국 국방부 장관으로 활동했습니다.

질문 6 답변의 양
조선도체에 대해 설명해줘.

조선도체는 개성만 간단 설명.

조선도체의 개성과 함께 운도 의존성, 영구적인 전류 흐름, 메시 현상 등 설명.



출처: 매일경제(23. 9. 7.)네이버가 만든 생성형 AI 능력은?...챗GPT vs 클로바X

2. 환경 분석

챗GPT – 클로바X 한국어 문법적 오류 비교 결과

| 구분 | | 챗GPT(3.5) | 클로바 X |
|----|--|----------------------------|----------------------------|
| 발음 | 음소 음절 음운 규칙 원어식 발음 중간발음(변이음 포함) | ○ - ○ ○ ○ | ○ ○ ○ ○ ○ |
| 형태 | 단어형성(합성법) 단어형성(파생법) 굴절(곡용) 굴절(활용) 품사 | ○ ○ ○ ○ ○ | ○ ○ ○ ○ X |
| 통사 | 높임 시제 사동 피동 부정 어순 | ○ ○ ○ X X - | ○ ○ ○ X X - |
| 담화 | 지시 접속 담화표지 구어/문어 오류 | ○ X ○ ○ | X X X X |



3. 한국어 자원 구축 현황

국립국어원 언어 정보 나눔터

국립국어원 언어 정보 나눔터

The screenshot shows the main page of the Naver Language Information Sharing Platform. It features three main service cards: '인공지능 AI 말풍' (AI Language), '영어' (English), and '한국어' (Korean). Each card has a title, a brief description, and a '바로가기' (Go Directly) button. The 'Korean' card is highlighted with a red border.

말뭉치 신청

사용자 참여
말뭉치 활용
알립니다

모두의 말뭉치

미래를 준비하는 소중한 우리말 자원

A large, stylized brain composed of numerous Korean characters (Hangeul) is shown. A small, friendly-looking character stands next to the brain, pointing towards it.

말뭉치 신청

말뭉치 신청 내역

| 상태 | 제작자 | 제작일 | 제작 내용 | 작업 상태 | 작업 일정 | 작업 결과 |
|-------|-------|------------|----------------|-------|------------|-------|
| 준비 완료 | 한국어학회 | 2023-09-15 | 한국어-힌디어 병렬... | 작업 완료 | 2023-09-15 | 작업 완료 |
| 준비 완료 | 한국어학회 | 2023-09-15 | 한국어-러시아어 병렬... | 작업 완료 | 2023-09-15 | 작업 완료 |
| 준비 완료 | 한국어학회 | 2023-09-15 | 한국어-우즈베크어 병... | 작업 완료 | 2023-09-15 | 작업 완료 |

3. 한국어 자원 구축 현황

■ 한국어 기본 말뭉치

- ✓ 언어사용 실태를 반영한 대규모 말뭉치(언어 변화 감지용) 구축
 - ✓ 다양한 매체(신문·일상대화·온라인), 다양한 변인(지역, 연령, 성별)을 고려한 말뭉치 구축

→ 저작권 문제를 해결하여 인공지능 학습용 자료로 보급, 언어분석정보 입력의 토대 자료로 활용

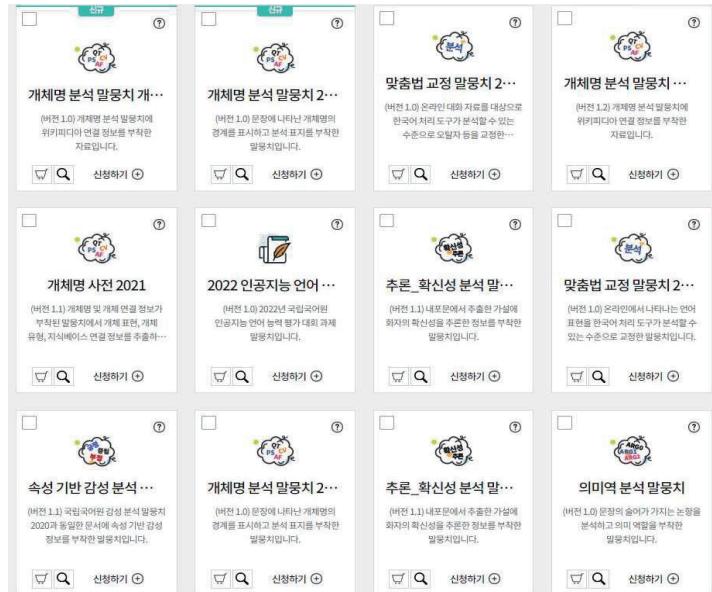


3. 한국어 자원 구축 현황

■ 한국어 분석 말뭉치

- ✓ **한국어 특성을 반영한 깊이 있는 언어 분석 정보 입력 말뭉치 구축**

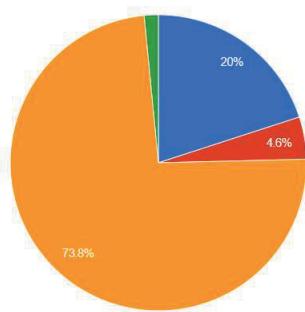
→ AI 언어모델의 미세조정 단계에 활용하여, 한국어 처리 성능 향상 기여



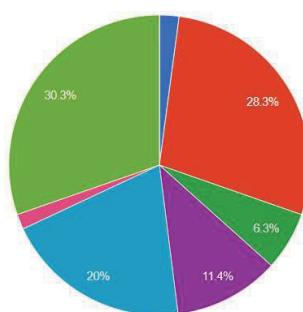
3. 한국어 자원 구축 현황

말뭉치 활용 현황

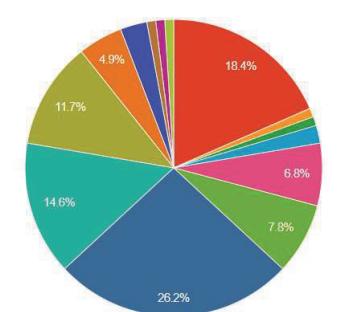
<결과물 형태별>



<신청 사유별>



<말뭉치별>

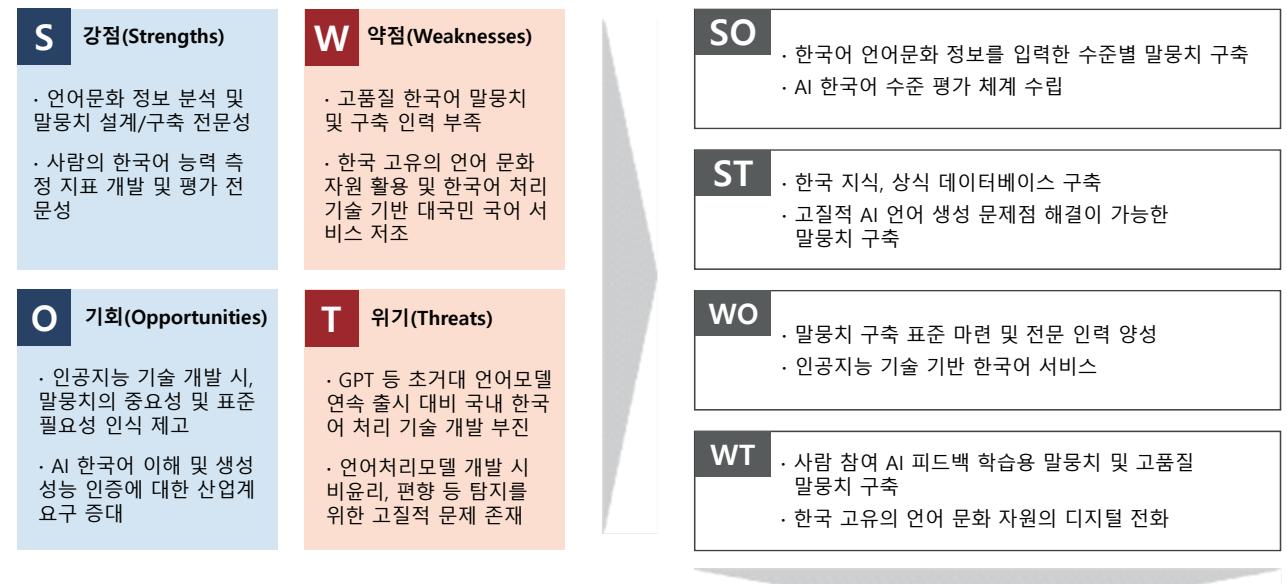


- 인공지능 모델/서비스/제품
- 경진 대회 참가
- 학술 발표, 논문, 과제 보고서
- 기타

- 인공 지능(모델/서비스/제품 등) 성능 평가
- 인공 지능(모델/서비스/제품 등) 개발 및 개선
- 언어 처리 관련 경진 대회 참가
- 교육 및 강의 자료 활용
- 인공 지능 연구(학술 발표, 논문, 연구 과제) 활용
- 기타 분야 연구(학술 발표, 논문, 연구 과제) 활용
- 언어 연구(학술 발표, 논문, 연구 과제) 활용

- 신문말뭉치
- 일상 대화 음성 말뭉치
- 감성 분석
- 구문 분석
- 형태 분석
- 개체명 분석
- 구어말뭉치
- 문어말뭉치
- 메신저말뭉치
- 비중없음
- 국회 회의록 말뭉치
- 병렬 말뭉치
- 맞춤법 고정 말뭉치

4. 발전 방향



4. 발전 방향

한국어 말뭉치 구축 전략(안)

- 1** 인공지능의 한국어 능력 향상을 위한 고품질 말뭉치 구축
- 2** 한국어 확산용 말뭉치 구축
- 3** 한국언어문화 정보 자원 구축 및 공유 보전



4. 발전 방향: 비전 및 목표

| | | | |
|--|--|---|--|
| 비전 | <p>한국어 잘하는 인공지능, 인공지능과 함께 발전하는 한국언어문화</p> |  | 1 인공지능의 한국어 능력 향상을 위한 고품질 말뭉치 구축 과제 1-1. 한국어 특성 정보를 입력한 수준별 분석 말뭉치 구축 과제 1-2. AI의 고질적 언어 문제 해결 말뭉치 구축 과제 1-3. AI의 한국어 능력 진단을 위한 표준 평가체계 수립 |
| 목표 | <p>국어 전문가와 함께 만드는 고품질 말뭉치 구축</p> <ul style="list-style-type: none">• '23년~'27년, 고품질 말뭉치 누적 200종• 문화예술분야 말뭉치 200,000건• 인공지능 언어능력 검증하는 평가체계 평가 과제 누적 30종 | | 2 한국어 확산용 말뭉치 구축 과제 2-1. 소통의 확장을 위한 언어 간 병렬 말뭉치 구축 과제 2-2. 언어생활 편의성 제고를 위한 교육 참고용 말뭉치 구축 |
| 3 한국언어문화 정보 자원 구축 및 공유 보전 과제 3-1. 한국어 사용 현실을 보여주는 균형적인 말뭉치 구축 과제 3-2. 한국 고유의 언어문화 자원의 디지털 전환 과제 3-3. 문화예술분야 말뭉치 및 지식 데이터베이스 구축 | | | |



4. 발전 방향: 추진 계획(안)

전략1 인공지능의 한국어 능력 향상을 위한 고품질 말뭉치 구축

1. 한국어 특성 정보를 입력한 수준별 분석 말뭉치 구축

- AI의 한국어 이해 수준 향상을 위한 수준별 한국어 분석 정보 데이터 강화 및 한국어 특성을 반영한 이종데이터(표/그림+텍스트) 확장

2. AI의 고질적 언어 문제 해결 말뭉치

- 초기대 AI가 가지는 언어분석의 문제 해결을 위한 말뭉치 정제, 국어 말뭉치 품질 개선을 위한 기준 수립으로 신뢰할 수 있는 한국어 AI에 기여

3. AI의 한국어 능력 진단을 위한 표준 평가체계 수립

- AI의 한국어 이해·생성 능력 수준 진단을 위한 평가체계 수립 및 평가용 말뭉치 구축을 통해 인공지능 언어능력 향상 점검 기준 마련



4. 발전 방향: 추진 계획(안)

생략어 복원 말뭉치

[예시 1]

- (1) 좀 익혀주시면은 국물 맛이 시원하거든요.
(2) <맛이> 그렇군요
→ ‘그렇군요’의 주어를 선행어 ‘맛’으로 복원

[예시 2]

- (1) 이~하늘에 관문인 인천국제공항이
(2) 국민들의 자부심으로 <인천국제공항이> 떠오르고 있습니다.
→ 본래 하나의 문장이지만 발화 단위에 따라 (1)과 (2)로 구분되므로 각각의 발화 단위를 기준으로 무형 대용어를 복원: ‘떠오르고’의 주어로 ‘인천국제공학’을 복원



4. 발전 방향: 추진 계획(안)

상호 참조 해결 말뭉치

(중략)

그 경쟁시스템 결국 또 광고 하고도 연결돼 있거든요. 그래서 자사에 범죄 보도 기준이 있는데. 타사가 다 보도를 하고 있단 말이야. 신상공개를 다 하고 있어요. 그러면 이거 우리만 안 하는게 너무 바보 같지 않아? 이거 뭐 무슨 의미가 있겠어. 결국 그 원칙 자체가 무너져 버린다는 거예요.

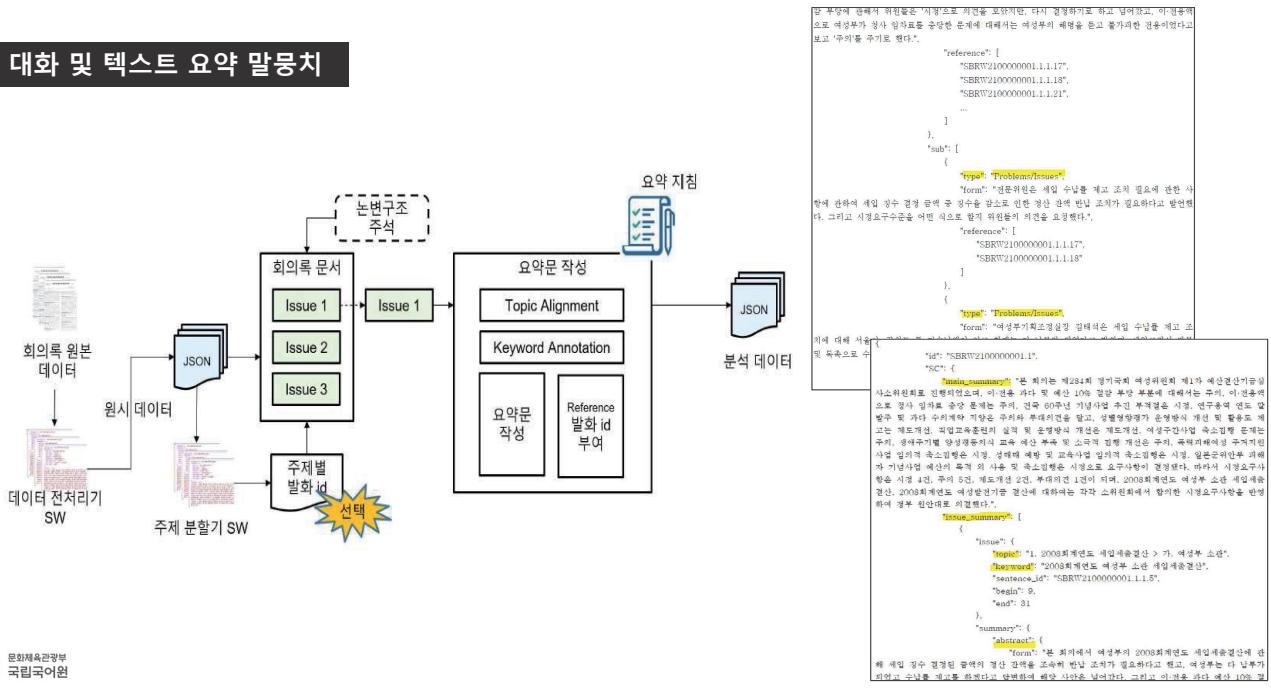
(중략)

```
  "mention": [
    {
      "sentence_id": "SBRW1800000231.1.1.54",
      "form": "자사에 범죄 보도 기준",
      "mention": null,
      "begin": 4,
      "end": 16,
      "NE_id": -1
    },
    {
      "sentence_id": "SBRW1800000231.1.1.56",
      "form": "그 원칙",
      "mention": null,
      "begin": 20,
      "end": 24,
      "NE_id": -1
    }
  ]
```



4. 발전 방향: 추진 계획(안)

대화 및 텍스트 요약 말뭉치



4. 발전 방향: 추진 계획(안)

| 서울지역 어린이 대상 급식실태 조사 및 개선방안 연구 | | | | | | | | | | | |
|---|-----|-------|---|-----|------|------------|-----|------|----|-----|-------|
| table.id: tb72117 url: http://share.nanet.go.kr/portal/work/workDetail.do? searchType=&queryText=%B0%BC%EC%9D%98%ED%9A%9C%EC%9D%98 table.title: 서울지역 어린이 대상 급식실태 조사 및 개선방안 연구 date: 2008-11-00 publisher: 식품의약품안전처 | | | | | | | | | | | |
| 보육시설 배분율(100%) | | | | | | | | | | | |
| <table border="1"> <tr> <td>예</td> <td>276</td> <td>35.1</td> </tr> <tr> <td>배식대 보유 아니오</td> <td>511</td> <td>64.9</td> </tr> <tr> <td>소계</td> <td>787</td> <td>100.0</td> </tr> </table> | | | 예 | 276 | 35.1 | 배식대 보유 아니오 | 511 | 64.9 | 소계 | 787 | 100.0 |
| 예 | 276 | 35.1 | | | | | | | | | |
| 배식대 보유 아니오 | 511 | 64.9 | | | | | | | | | |
| 소계 | 787 | 100.0 | | | | | | | | | |

표 → 텍스트

```

"column_span": 1
},
{
  "value": "1787",
  "is_header": false,
  "row_span": 1,
  "column_span": 1
},
{
  "value": "100.0",
  "is_header": false,
  "row_span": 1,
  "column_span": 1
}
]
},
"sentence_annotation": [
  {
    "sentence_after_deletion": "예식대를 보유하고 있는 곳은 조사한 보육시설의 35.1%인 276곳 밖에 없으며 64.9%인 511곳에서 배식대를 보유하지 않고 있는 것으로 나타났다.",
    "worker1": "조사한 보육시설 중 예식대를 보유하고 있는 곳과 그렇지 않은 곳은 각각 35.1%인 276곳과 64.9%인 511곳을 차지하는 것으로 조사되었다.",
    "worker2": "조사한 보육시설 중 35.1%인 276곳은 배식대를 보유하고 있으며 64.9%인 511곳은 그렇지 않은 것으로 조사되었다."
  },
  {
    "sentence_after_deletion": "조사한 보육시설 중 배식대를 보유하고 있는 곳은 35.1%인 276곳뿐이며 64.9%인 511곳에서는 보유하지 않은 것으로 나타났다.",
    "worker4": "예식대를 보유하고 있는 곳과 그렇지 않은 곳은 조사한 보육시설 중 각각 35.1%인 276곳과 64.9%인 511곳에 해당하는 것으로 조사되었다."
  }
]
}

```

- ① 표의 주요 핵심 내용이 드러나는 간결한 형태의 글, ② 일반적 형태의 글, ③ 상세한 설명 정보를 담은 글로 구현할 수 있는 데이터 구축

4. 발전 방향: 추진 계획(안)



그림 → 텍스트

```
id: "P32447",
metadata: {
  publisher: "크라우드 소싱",
  date: "2022"
},
sentence_annotation: [
  {
    reference_sentence: "풀밭에 놓인 세 개의 돌조각상에 각각 보지 말고, 듣지 말고, 말하지 말리라고 새겨져 있다.",
    worker1: "풀밭에 세워진 돌조각상 세 개에는 각각 한 미디어 글씨가 적혀 있는데, 보지 말고, 듣지 말고, 말하지 말리가 바로 그것이다.",
    worker2: "풀조각상 세 개를 차운 것은 각각 보지 말고, 듣지 말고, 말하지 말라는 문구를 활용해 새기고 있다.",
    worker3: "풀조각상에 새겨진 문구는 각각 보지 말고, 듣지 말고, 말하지 말리이며 이들은 풀밭에 놓여 있다.",
    worker4: "풀밭에는 문구가 새겨진 돌조각상 세 개가 놓여 있고, 그 문구는 각각 보지 말고, 듣지 말고, 말하지 말라이다."
  }
],
image_width: 4608,
image_height: 3456,
ocr_info: [
  {
    words: "보지 말고",
    type: "rect",
    blob: {
      x1: 3119,
      y1: 2117,
      width: 481,
      height: 1102
    }
  },
  {
    words: "듣지 말고",
    type: "rect",
    blob: {
      x1: 1983,
      y1: 2157,
      width: 347,
      height: 1069
    }
  }
]
```

- ①그림의 주요 핵심 내용이 드러나는 간결한 형태의 글, ②일반적 형태의 글, ③상세한 설명 정보를 담은 글로 구현할 수 있는 데이터 구축



4. 발전 방향: 추진 계획(안)

전략2 한국어 확산용 말뭉치 구축

1. 소통의 확장을 위한 언어 간 병렬 말뭉치 구축

- 우리 국민의 다양한 언어 사용 대응 및 특화형 AI 언어모델 개발을 위한 수준별 말뭉치 구축

2. 언어생활 편의성 제고를 위한 교육 참고용 말뭉치 구축

- 국민의 언어생활에 편의를 제공할 수 있는 교육 및 말뭉치 적용 프로그램(글쓰기 자가 진단·침식 서비스 등) 구현 개발에 필요한 말뭉치 구축



4. 발전 방향: 추진 계획(안)

'한국어-외국어 병렬 말뭉치' 구축 단계

1단계

마지막 답변 문장을 평가해 주세요.
주제: 최근 즐기고는 음식에 대해 친구와 재미를 이야기해주세요.

☞ 오 일끼로 요즘 냄 잘나오더라
다면 배달식 먹기는 부담스럽고 일끼로 냄은 것 같아.

그니깐 일끼로가 더 편리하고 좋지
그자 일끼로가 편리한 편이야~

부족한
적합

2단계

한국어 → 태국어

다음 문장을 번역해주세요.

풀연 배달식 먹기는 부담스럽고 일끼로 냄은 것 같아.

태국어로 입력하세요. 기기면역 사용 시 제작 대상이 됩니다.

번역 후 수정은 불가능합니다.
번역하기 힌트보기

제출

3단계

한국어 → 태국어

번역을 추가해주세요.

학군이 좋은 지점엔 후원을 해서라도 이사를 가는 한국 사회 풍조 속에서 핀란드 교육과 과정 혁신 조율학교 수장 고려시작입니다. □

한국어로 번역하는 것은 학생들이 학교에서 배운 내용을 활용해 자신의 언어로 전달하는 행위입니다. 학생들은 학교에서 배운 내용을 활용해 자신의 언어로 전달하는 행위입니다. □

제출 후 수정은 불가능합니다.
번역 시 주의사항

나쁨 나쁘지 않음 매우 좋음

4단계

한국어 → 태국어

번역을 추가해주세요.

학군이 좋은 지점엔 후원을 해서라도 이사를 가는 한국 사회 풍조 속에서 핀란드 교육과 과정 혁신 조율학교 수장 고려시작입니다. □

한국어로 번역하는 것은 학생들이 학교에서 배운 내용을 활용해 자신의 언어로 전달하는 행위입니다. 학생들은 학교에서 배운 내용을 활용해 자신의 언어로 전달하는 행위입니다. □

제출 후 수정은 불가능합니다.
번역 시 주의사항

나쁨 나쁘지 않음 매우 좋음

번역 대상 자료 수집·정제(국내)

현지 전문가(이중언어화자)

해당 언어 전문가 번역문 검수(3회)

한국어→외국어 번역

*번역 전문가, 원어민 검수

전문가 최종 감수



4. 발전 방향: 추진 계획(안)

한국어-수어 말뭉치

한국어 학습자 말뭉치

한국어:

한 사람이 시장에서 넥타이를 둘러보고 있습니다.

한국 수어:



.국적: 러시아 .모어: 러시아어 .나이: 21

.직업: 기타, 학생 .한국어 급수: 2급 .말뭉치 유형: 문어

교정 정보: 좋아하/VV+니까/EC

저는 야구 하는 것이 아주 좋아하니까 여름 때 야구장에서 가요.

교정 결과:

저는 야구 하는 것이(것+을) 아주 좋아하니까(좋아하+니까)

여름 때(여름+에) 야구장에서(야구장+에) 가요.

| 원형태 | 교정형태 | 오류 위치 | 오류 영상 | 오류 총의 |
|-----|------|-------|-------|------------------|
| 풀 | 풀이하 | 형용사 | 대치 | 단어 형성(파생법) 동사 |
| 아니까 | 니까 | 연결어미 | 오형태 | 글자(활용) |



4. 발전 방향: 추진 계획(안)

쓰기/채점 말뭉치

다양한 텍스트 유형에 대해 수준별 쓰기 및 채점 말뭉치 구축

국어 지식 및 첨삭 말뭉치

글쓰기 수정·보완을 위한 참조 데이터 구축

- 수준별(초등·중등·성인)·유형별 국어능력 진단체계 지표 개발
- 인공지능 모델 개발에 활용 가능한 글쓰기 및 채점 자료 수집·구축
- 외국인들의 한국어 쓰기 평가 자료 구축

- 가나다 전화 등 국어 상담 서비스 및 공공 언어 감수(보도자료, 보고서, 기안문, 안내문 등) 결과 자료를 인공지능 학습 데이터로 가공
- 국어능력 진단 지표를 고려한 첨삭(피드백) 말뭉치 구축

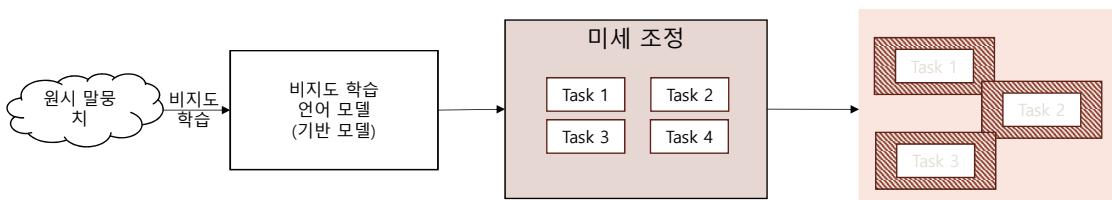
→ 국민의 인공지능을 활용한 언어생활 편의 서비스 구현과 연계



4. 발전 방향: 추진 계획(안)

한국어 능력 향상 사례

한국어 특화 인스트럭션
한국어 능력 향상



기존 모델

[프롬프트]
'기초 체력이 떨어지는데도 잘 못 뛴다.'
이 문장을 교정해 주세요

[답변-GPT-3.5]
'기초 체력이 떨어지는데도 잘 못 뛰지
못해요.'

[답변-GPT-4]
'기초 체력이 부족한데도 잘 못 뛴다.'

학습

1. 한국어 능력 향상을 위한 추가적인 비지도 학습 수행
2. 한국어 교정/첨삭을 위한 말뭉치(예: 학습자 말뭉치) 활용하여 미세조정

학습 후

- [프롬프트]
'기초 체력이 떨어지는데도 잘 못 뛴다.'
이 문장을 교정해 주세요
- [예상 답변 1]
'기초 체력이 떨어지는데도 잘 뛴다.'
- [예상 답변 2]
'기초 체력이 부족해서 잘 못 뛴다.'



4. 발전 방향: 추진 계획(안)

전략3 한국언어문화 정보 자원 구축 및 공유 보전

1. 한국어 사용 현실을 보여주는 균형적인 말뭉치 구축

- 인공지능(AI) 언어모델이 한국어 실태를 학습하고 한국어답게 사용하도록 다양한 매체(신문·웹 등)와 다양한 변인(지역, 나이, 성별)을 고려한 한국어 말뭉치 구축

2. 한국 고유의 언어문화 자원의 디지털 전환

- 역사적·지역적으로 고유한 한국 언어문화를 계승/보전하고, 생성형 AI의 올바른 한국어 생성을 지원하기 위한 말뭉치 구축

3. 문화예술분야 말뭉치 및 지식 데이터베이스 구축

- 한국의 문화 관습적 특성의 다양성과 고유성을 담은 문화예술분야 자료를 말뭉치로 구축하여 산업계·학계에서 활용할 수 있는 데이터로 디지털 전산화



4. 발전 방향: 추진 계획(안)

언어 변화 추적 말뭉치

- (신문기사) 사회 현상을 담는 언어자료로 새로운 용어·개념을 포함하는 말뭉치 구축
- (일상대화) 지역별·나이별·성별 인구 분포를 고려한 자료로 대화 자료 및 1인 독백(강연·낭독 등), 회의 등 다양한 유형의 대화 말뭉치 구축
- (온라인 자료) 언어 변화가 가장 빠르게 일어나는 온라인상의 자료를 수집, 메신저 대화 등 사적 언어자료를 저작권 해결 수집

언어 사용 균형 말뭉치

언어 변화 추이를 주기적으로 관찰(5년 주기)하고, 국민의 언어 사용 실태를 연구할 수 있는 기본 데이터 구축

- 다양한 장르의 텍스트(책, 잡지, 시나리오, 유튜브 대본 등)를 구성하여 폭넓은 언어 현상을 담아낼 수 있도록 연구 및 말뭉치 구축



4. 발전 방향: 추진 계획(안)

지역어 말뭉치

한국어의 다양성 보전 및 지역별·세대별 문화를 분석할 수 있는 구술발화 말뭉치 구축

역사어원 말뭉치

국어 역사 문헌 자료 종합 정리 및 장르별·시기별 균형성을 갖춘 대규모 국어사 말뭉치 구축

- 매년 9지역(9개도 각 1지점)*3세대(20대, 50대, 70대 이상)*남녀 제보자를 선정→
지역·세대·성별에 따른 자연스러운 한국어 음성 말뭉치(제보자별 2시간 이상) 구축

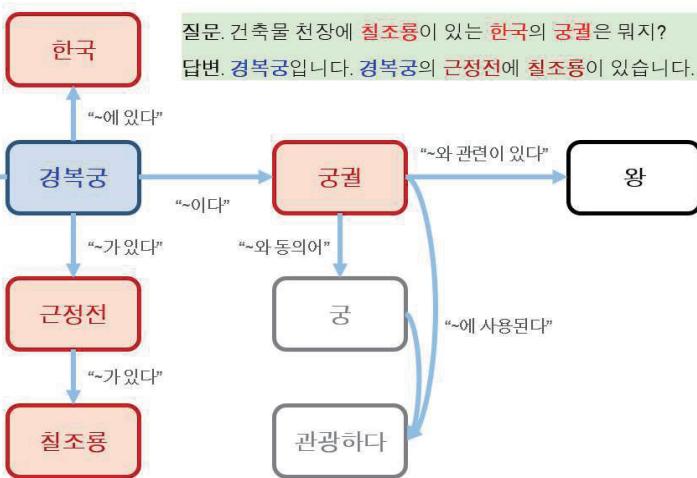
- 21세기 세종계획에서 구축한 “불경언해류”의 말뭉치를 포함하여, 19세기~20세기 초의 일상생활이 반영된 편지, 일기, 소설, 잡지 등 장르 확대 구축
- 국어사 말뭉치 형식 표준화 및 말뭉치 통합 관리 체계 구축으로 옛한글 입력 및 국어 역사 정보 상세 검색, 우리말 어원정보 제공



4. 발전 방향: 추진 계획(안)

한국형 상식 기반 데이터

한국인이 공유하는 문화적 상식과 맥락을 이해하고 그 수준을 추론할 수 있는 언어 자료 구축



감사합니다.

학습자 말뭉치와 한국어 교육

강 현 화 (연세대 국어국문학과)

1. 들어가기

1.1 학습자 언어 자료 연구의 쟁점

비모어 화자가 생산한 중간언어 자료, 즉 학습자 말뭉치(Learner's Corpora)는 외국어 교육자들에게 귀중한 언어 연구 자료가 된다. 학습자의 발달 과정과 해당 과정에서 나타나는 특성을 파악할 수 있기 때문이다. 학습자 자료 분석 결과는 교육과정의 설계나 교수법, 평가에 이르기까지 외국어 교수의 다양한 과정에 영향을 미치게 된다. 한국어교육 학계에서는 말뭉치 언어학의 성과와 제2언어습득 연구(SLA)의 성과가 합쳐진 '컴퓨터 기반의 학습자 오류분석(CEA)'이 활발히 이루어지고 있다.

제2언어 습득 연구 기반의 체계적인 '주석(annotation)'은 외국어 학습에 있어서 양적인, 질적인 연구를 가능하게 하며, 다양한 학습자 변인에 의한 오류 양상의 결과를 일반화할 수 있다는 점에서 의미를 가진다. 대규모 학습자 말뭉치 연구는 주로 '비교'의 방법론을 이용하여 학습자들이 원어민에 비해 과소 혹은 과도하게 사용하는 언어 자질들을 보여주고, '오류 분석'¹⁾의 방법론을 통해 언어 학습자가 만드는 오류에 대해 파악할 수 있게 한다.

그런데, 학습자 말뭉치 구축은 수집, 가공, 활용 면에서 각각의 한계와 어려움을 가지고 있다. 수집 자료 규모의 타당성, 균형성, 대표성의 문제부터, 주석의 정확성 및 일관성, 배포 및 활용의 방법의 효용성에 이르기까지 단계마다 쟁점이 존재한다. 예를 들면 구축 말뭉치의 규모, 구축 범위(학습자 대상, 수집 장르, 수집 유형 등), 주석(전사, 형태주석, 오류주석)의 범위 및 방법을 고민해야 한다. 수집 대상 확보나 수집 과정, 주석자 훈련, 사용자 교육 등의 부수적 영역에서의 어려움이 존재하며, 주석의 일관성(주석자 신뢰도, 주석 간 신뢰도)과 컴퓨터 주석에서의 효율성, 그리고 연구 혹은 활용 목표에 따른 효용성 등도 어려움의 대상이다.

특히 활용은 가장 중요한 문제인데, 말뭉치 구축과 형태 주석의 절차는 연구 방법일 뿐이며, 궁극적인 관심은 효율적 검색을 통한 연구나 교수에의 활용이기 때문이다.

한국어 학습자 말뭉치 구축(2015~2023년) 사업에서는 다양한 학습자(취미 목적, 학문 목적, 이주 목적, 동포 등)를 대상으로 하여 다양한 장르의 구어와 문어 자료를 구축하고 있으며, 전사의 지침과 형태 주석의 범위를 설계하고 오류주석을 위한 세부 지침을 마련하고 있다. 하지만 변화하는 학계의 요구들을 고려할 때, 현행 지침들에 대한 검토가 필요하다. 한국어 학습자 말뭉치가 구축된 지 8년이 경과한 현 시점에서 그간의 구축 과정에서의 문제점을 짚어보고 향후 구축에의 방향성을 모색할 필요가 있다.

이에 본고에서는 한국어 학습자 말뭉치의 구축 현황을 간단히 살펴보고, 구축 단계별(수집, 전사, 주석, 활용 및 배포) 쟁점에 대해 논의하고자 한다. 특히 쟁점이 되는 구어 전사와 오류 주석에 대해 집중적으로 논의하기로 하겠다. 이러한 논의를 바탕으로 하여 향후 학습자 말뭉치 구축에 있어서의 지향점을 모색할 수 있을 것이다.

1) 특히 오류주석은 '학습자 언어의 변칙적인 특징을 보여주고자 하는' 것으로, 학습자 말뭉치의 고유한 부분이며 학습자 언어 자료 분석에 있어서 가장 중요한 부분이 된다.

1.2 학습자 말뭉치 구축 현황

국립국어원의 한국어 학습자 말뭉치 구축사업은 2015년에 시작되어 1단계 구축 사업에 이어 2단계 사업이 진행 중이다. 구어와 문어의 비중은 약 1/10의 규모(파일 수)로 여전히 구어의 비중이 낮다. 일반적인 언어 연구 자료로써의 코퍼스와는 달리, 회화 요구(+ 발음 학습)가 높은 한국어교육에서는 구어 자료의 비중이 더 높아져야 한다.

(1) 구축 규모

<표1> 2015-2022년 원시 말뭉치 누적 구축 현황: 구어/문어

| 구분 | 구분 | 1급 | 2급 | 3급 | 4급 | 5급 | 6급 | 6급+ | 합계 |
|----|------|---------|---------|---------|---------|---------|---------|---------|-----------|
| 구어 | 어절 수 | 295,688 | 350,333 | 462,328 | 328,085 | 216,697 | 137,017 | 32,859 | 1,823,007 |
| | 파일 수 | 904 | 762 | 793 | 674 | 360 | 204 | 23 | 3,720 |
| 문어 | 어절 수 | 594,047 | 741,040 | 826,875 | 814,746 | 790,296 | 505,257 | 135,335 | 4,407,596 |
| | 파일 수 | 8,736 | 7,580 | 6,787 | 6,026 | 5,028 | 3,018 | 117 | 37,292 |

(2) 언어권 및 숙달도별 구축 규모

구어의 언어권별, 숙달도별 비율을 보면, 특정 언어권(중국어, 베트남, 타이, 일본어)의 비중(파일 수 대비)이 66.9%를 넘어 편중성이 높음을 알 수 있다. 숙달도 면에서도 초급과 중급의 비율에 비해 고급의 비중이 낮아 균형성이 부족하다. 이러한 비중은 학습자의 국적별 실제적 비중을 대표하고는 있으나, 상호 비교 연구의 자료로 삼기에는 특정 언어권의 규모가 부족한 언어권이 여전히 많다.

<표2> 2015-2022년 원시 말뭉치 누적 구축 현황: 구어

| 모국어 | 구분 | 1급 | 2급 | 3급 | 4급 | 5급 | 6급 | 6급+ | 합계 |
|--------|------|--------|--------|---------|--------|--------|--------|--------|---------|
| 중국어 | 어절 수 | 53,060 | 72,635 | 74,290 | 68,797 | 45,418 | 54,992 | 21,112 | 390,304 |
| | 파일 수 | 248 | 220 | 168 | 157 | 79 | 72 | 15 | 959 |
| 베트남어 | 어절 수 | 50,236 | 64,064 | 74,706 | 76,856 | 66,914 | 10,502 | 777 | 344,055 |
| | 파일 수 | 120 | 133 | 139 | 136 | 109 | 17 | 2 | 656 |
| 타이어 | 어절 수 | 69,918 | 48,320 | 153,038 | 14,899 | 6,040 | 4,743 | - | 296,958 |
| | 파일 수 | 164 | 72 | 165 | 28 | 9 | 9 | - | 447 |
| 일본어 | 어절 수 | 14,844 | 30,499 | 34,578 | 52,863 | 44,115 | 40,193 | 1,701 | 218,793 |
| | 파일 수 | 60 | 58 | 77 | 118 | 70 | 56 | 2 | 441 |
| 스페인어 | 어절 수 | 29,334 | 23,538 | 23,346 | 19,165 | 4,225 | 1,988 | 594 | 102,190 |
| | 파일 수 | 64 | 43 | 38 | 40 | 8 | 3 | 1 | 197 |
| 러시아어 | 어절 수 | 13,001 | 20,223 | 22,708 | 18,426 | 5,330 | 4,419 | - | 84,107 |
| | 파일 수 | 43 | 53 | 57 | 44 | 14 | 12 | - | 223 |
| 영어 | 어절 수 | 10,644 | 19,630 | 15,084 | 7,002 | 10,860 | 2,282 | - | 65,502 |
| | 파일 수 | 32 | 29 | 34 | 17 | 13 | 3 | - | 128 |
| 인도네시아어 | 어절 수 | 11,302 | 13,237 | 15,122 | 7,881 | 8,307 | 5,216 | - | 61,065 |
| | 파일 수 | 38 | 26 | 29 | 17 | 18 | 12 | - | 140 |
| 타갈로그어 | 어절 수 | 11,948 | 11,615 | 11,635 | 6,632 | 1,682 | 1,711 | - | 45,223 |
| | 파일 수 | 31 | 29 | 13 | 12 | 2 | 3 | - | 90 |

| 모국어 | 구분 | 1급 | 2급 | 3급 | 4급 | 5급 | 6급 | 6급+ | 합계 |
|-------|------|---------|---------|---------|---------|---------|---------|--------|-----------|
| 싱할라어 | 어절 수 | 6,204 | 5,306 | 6,539 | 6,133 | 3,726 | 2,990 | - | 30,898 |
| | 파일 수 | 10 | 11 | 10 | 11 | 6 | 5 | - | 53 |
| 몽골어 | 어절 수 | 2,046 | 3,407 | 5,031 | 7,229 | 4,062 | 2,343 | - | 24,118 |
| | 파일 수 | 10 | 11 | 13 | 16 | 8 | 2 | - | 60 |
| 버마어 | 어절 수 | 2,633 | 6,469 | 6,653 | 6,297 | 751 | - | - | 22,803 |
| | 파일 수 | 5 | 7 | 6 | 8 | 1 | - | - | 27 |
| 키르기스어 | 어절 수 | 1,469 | 4,525 | 5,440 | 5,063 | 3,858 | - | - | 20,355 |
| | 파일 수 | 5 | 9 | 8 | 10 | 7 | - | - | 39 |
| 우즈베크어 | 어절 수 | 1,946 | 7,706 | 2,637 | 4,075 | 666 | - | - | 17,030 |
| | 파일 수 | 6 | 14 | 5 | 8 | 1 | - | - | 34 |
| 프랑스어 | 어절 수 | 2,716 | 2,871 | 2,057 | 2,021 | 548 | 368 | - | 10,581 |
| | 파일 수 | 11 | 6 | 4 | 5 | 1 | 1 | - | 28 |
| 기타 | 어절 수 | 14,387 | 16,288 | 9,464 | 24,746 | 10,195 | 5,270 | 8,675 | 89,025 |
| | 파일 수 | 57 | 41 | 27 | 47 | 14 | 9 | 3 | 198 |
| 합계 | 어절 수 | 295,688 | 350,333 | 462,328 | 328,085 | 216,697 | 137,017 | 32,859 | 1,823,007 |
| | 파일 수 | 904 | 762 | 793 | 674 | 360 | 204 | 23 | 3,720 |

문어의 경우도 언어권과 숙달도별 비중의 편중도가 높다.

<표3> 2015-2022년 원시 말뭉치 누적 구축 현황: 문어

| 모국어 | 구분 | 1급 | 2급 | 3급 | 4급 | 5급 | 6급 | 6급+ | 합계 |
|------|------|---------|---------|---------|---------|---------|---------|---------|-----------|
| 중국어 | 어절 수 | 228,014 | 241,913 | 285,454 | 298,200 | 397,807 | 252,425 | 102,924 | 1,806,737 |
| | 파일 수 | 3,433 | 2,407 | 2,299 | 2,234 | 2,527 | 1,536 | 77 | 14,513 |
| 일본어 | 어절 수 | 51,275 | 117,780 | 132,526 | 150,858 | 118,864 | 77,761 | 1,079 | 650,143 |
| | 파일 수 | 703 | 1,219 | 1,080 | 1,116 | 738 | 485 | 6 | 5,347 |
| 베트남어 | 어절 수 | 114,328 | 98,426 | 105,058 | 114,150 | 58,643 | 23,430 | 5,996 | 520,031 |
| | 파일 수 | 1,728 | 1,121 | 900 | 789 | 362 | 133 | 3 | 5,036 |
| 영어 | 어절 수 | 40,576 | 60,013 | 60,014 | 49,969 | 38,174 | 33,063 | 2,856 | 284,665 |
| | 파일 수 | 597 | 624 | 516 | 395 | 249 | 179 | 4 | 2,564 |
| 광동어 | 어절 수 | 16,954 | 39,128 | 49,989 | 48,046 | 42,047 | 37,931 | - | 234,095 |
| | 파일 수 | 257 | 379 | 400 | 354 | 267 | 234 | - | 1,891 |
| 러시아어 | 어절 수 | 16,565 | 25,277 | 32,740 | 30,441 | 27,418 | 11,795 | 4,539 | 148,775 |
| | 파일 수 | 260 | 307 | 296 | 244 | 184 | 71 | 6 | 1,368 |
| 타이어 | 어절 수 | 25,188 | 35,669 | 27,371 | 16,076 | 10,725 | 8,089 | 212 | 123,330 |
| | 파일 수 | 289 | 276 | 178 | 113 | 74 | 47 | 1 | 978 |
| 몽골어 | 어절 수 | 17,437 | 18,046 | 19,656 | 17,585 | 15,857 | 7,484 | 267 | 96,332 |
| | 파일 수 | 241 | 193 | 169 | 131 | 109 | 50 | 2 | 895 |

| 모국어 | 구분 | 1급 | 2급 | 3급 | 4급 | 5급 | 6급 | 6급+ | 합계 |
|--------|------|---------|---------|---------|---------|---------|---------|---------|-----------|
| 스페인어 | 어절 수 | 13,313 | 20,220 | 21,343 | 11,937 | 7,632 | 2,445 | - | 76,890 |
| | 파일 수 | 179 | 203 | 165 | 84 | 52 | 15 | - | 698 |
| 인도네시아어 | 어절 수 | 7,012 | 8,751 | 7,234 | 10,511 | 8,487 | 5,734 | 1,097 | 48,826 |
| | 파일 수 | 103 | 85 | 54 | 64 | 54 | 36 | 1 | 397 |
| 프랑스어 | 어절 수 | 11,235 | 9,501 | 10,371 | 5,490 | 5,893 | 3,951 | 933 | 47,374 |
| | 파일 수 | 150 | 89 | 81 | 42 | 38 | 26 | 1 | 427 |
| 말레이어 | 어절 수 | 1,816 | 4,309 | 9,217 | 7,692 | 3,263 | 1,083 | 2,476 | 29,856 |
| | 파일 수 | 31 | 44 | 60 | 40 | 14 | 7 | 3 | 199 |
| 스웨덴어 | 어절 수 | 5,947 | 6,832 | 8,270 | 3,674 | 1,779 | 1,215 | - | 27,717 |
| | 파일 수 | 103 | 73 | 71 | 32 | 12 | 9 | - | 300 |
| 카자흐어 | 어절 수 | 1,984 | 3,841 | 5,243 | 4,850 | 4,564 | 2,737 | - | 23,219 |
| | 파일 수 | 32 | 35 | 39 | 38 | 30 | 14 | - | 188 |
| 아랍어 | 어절 수 | 3,782 | 4,216 | 4,554 | 3,278 | 3,347 | 2,221 | - | 21,398 |
| | 파일 수 | 61 | 46 | 44 | 28 | 25 | 16 | - | 220 |
| 이탈리아어 | 어절 수 | 3,118 | 2,231 | 2,093 | 2,564 | 3,759 | 2,802 | 4,741 | 21,308 |
| | 파일 수 | 51 | 21 | 21 | 22 | 19 | 16 | 1 | 151 |
| 독일어 | 어절 수 | 4,475 | 2,755 | 2,863 | 3,965 | 3,705 | 1,679 | 145 | 19,587 |
| | 파일 수 | 53 | 27 | 26 | 31 | 26 | 12 | 1 | 176 |
| 버마어 | 어절 수 | 1,687 | 2,797 | 3,466 | 3,770 | 4,302 | 2,949 | - | 18,971 |
| | 파일 수 | 22 | 24 | 23 | 24 | 26 | 14 | - | 133 |
| 우즈베크어 | 어절 수 | 3,590 | 2,528 | 3,818 | 3,479 | 3,942 | 1,493 | - | 18,850 |
| | 파일 수 | 52 | 30 | 35 | 26 | 22 | 7 | - | 172 |
| 싱할라어 | 어절 수 | 1,710 | 3,795 | 3,206 | 2,432 | 2,626 | 852 | 3,635 | 18,256 |
| | 파일 수 | 20 | 27 | 24 | 17 | 15 | 5 | 1 | 109 |
| 포르투갈어 | 어절 수 | 2,291 | 4,209 | 2,605 | 2,788 | 2,831 | 1,701 | - | 16,425 |
| | 파일 수 | 34 | 37 | 25 | 23 | 20 | 11 | - | 150 |
| 타갈로그어 | 어절 수 | 3,080 | 4,264 | 3,527 | 2,731 | 1,029 | 389 | - | 15,020 |
| | 파일 수 | 46 | 62 | 39 | 28 | 6 | 1 | - | 182 |
| 터키어 | 어절 수 | 862 | 1,638 | 2,436 | 1,654 | 3,105 | 1,431 | 403 | 11,529 |
| | 파일 수 | 14 | 18 | 22 | 12 | 22 | 8 | 1 | 97 |
| 키르기스어 | 어절 수 | 787 | 2,625 | 1,054 | 2,152 | 3,170 | 942 | 396 | 11,126 |
| | 파일 수 | 12 | 23 | 6 | 15 | 22 | 7 | 2 | 87 |
| 기타 | 어절 수 | 16,525 | 20,078 | 22,014 | 15,667 | 14,590 | 8,185 | 3,398 | 100,457 |
| | 파일 수 | 256 | 208 | 208 | 119 | 96 | 55 | 6 | 948 |
| 합계 | 어절 수 | 594,047 | 741,040 | 826,875 | 814,746 | 790,296 | 505,257 | 135,335 | 4,407,596 |
| | 파일 수 | 8,736 | 7,580 | 6,787 | 6,026 | 5,028 | 3,018 | 117 | 37,292 |

수집 대상 유형도 수집이 용이한 KSL 국내 학습자(유학생, 이주민)에 치중되어 있다.

<표4> 2015-2021년 말뭉치의 구축 규모: 대상별

| 대상 | 문어 | 구어 |
|---------------|---------------------|------------------|
| 국내 학습자(학문 목적) | 3,558,202 (209,069) | 819,765 (94,123) |
| 이주민 | 88,830 | 308,983 |
| 국외 학습자 | 50,922 | 393,772 |
| 합계 | - | - |

장르별 자료 활용을 위해서는 다양도가 중요한데, 구어 장르는 인터뷰에의 편중도가 다소 높다. 학습 목표에 기댄다면 자유 대화의 비중이 더 높아져야 할 것이다.

<표5> 2015-2021년 말뭉치의 구축 규모: 장르별

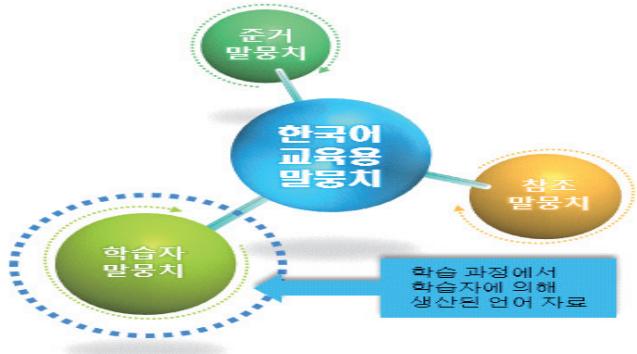
| 문어 | | 구어 | |
|-----|-----------|-------|-----------|
| 장르 | 구축 규모 | 장르 | 구축 규모 |
| 생활문 | 1,650,609 | 인터뷰 | 978,808 |
| 논설문 | 1,206,778 | 발표 | 280,883 |
| 설명문 | 250,975 | 내러티브 | 156,763 |
| 보고서 | 239,005 | 자유 대화 | 106,066 |
| 기행문 | 110,280 | | |
| 수필 | 104,868 | | |
| 감상문 | 53,841 | | |
| 기사문 | 38,078 | | |
| 전기문 | 21,863 | | |
| 편지글 | 21,657 | | |
| 합계 | 3,697,954 | 합계 | 1,522,520 |

2. 학습자 말뭉치 구축의 단계별 쟁점

2.1 수집

수집에 있어서 고려해야 할 사항에는 다음과 같은 것들이 있다. 우선, 학습자 말뭉치 구축 사업의 지속성 확보가 가장 중요하다. 현재의 말뭉치 규모는 국외 말뭉치에 비해 현저히 부족하며, 변화하는 학습자의 특성을 파악하기 위해서는 연도별로 지속적인 자료 확보가 필수적이다. 다만, 자료 유형별(문어 편중), 언어권별(중, 일, 베 편중), 숙달도별(초급 편중) 균형성 확보는 중요한 과제이다. 현실적으로 보면 한국어 학습자 자체의 불균형성이 존재하지만 연구나 활용에의 효용성을 고려한다면 어느 정도의 균형성 확보는 매우 중요하다. 새롭게 등장하는 언어권의 학습자 자료 확보도 필요한데, 이는 국외 학습자 자료 수집에 기반해야 할 것이다. 저작권 등이 민감한 지역의 수집 방법에 대한 다양한 노력이 필요할 것이다.

다음으로는, 양적 확장도 중요하지만 실제 한국어 교육에의 활용을 고려한 주제별 자료 수집과 (같은 주제의) 모어 화자와의 비교를 위한 준거 말뭉치 구축이 필요하다. 학습자 자료는 오류만 보는 것이 아니므로 발달 과정을 살피려면 비교 말뭉치가 필수적이다. 학습자 유형별(이주 여성, 학문 목적 등) 하위 말뭉치의 별도 구축도 필요하다.



<그림 1> 한국어 교육용 말뭉치의 구성

마지막으로는 교육 현장의 참조 말뭉치 구축도 필요하다. 교실에서의 교사-학습자 간 대화나 학습자 간 대화 자료 구축이 필요하며 가능하다면 동영상 말뭉치로 구축할 필요가 있다.

2.2 전사(구어 및 문어)

구어 전사는 문어에 비해 다양한 쟁점이 있다. 구어 전사의 쟁점을 차례로 살피면서, 향후 방향성을 짚어보기로 하겠다.

<전사 층위>

구어의 전사에 있어서 가장 기본적 원칙은 실제 발음에 가깝게 전사하는 것일 것이다. 구어 전사는 전사 층위(전사 수준)에 따라 철자법 전사, 한글 음소 전사, 음성 전사, 운율 전사로 나뉜다. 한국어 학습자의 경우에는 한국어로 철자화하기 어려운 발음들(자국어에만 있는 발음이나 변이음 등)이 상당수 존재하며, 언어교육적 측면에서 모국어 발음의 전이를 살피려면 음성 전사의 수준이 요구된다. 만약 연구의 목적이 대조언어학적 관점이나 음성적 특성에 있다면 음성 수준으로 전사하는 게 타당하다. 하지만 철자법 전사는 언어의 형태를 유지시킴으로써 전사 자료의 해독이 용이하고 형태소 분석 등을 용이하게 하므로 현행 전사에서는 철자법 전사를 선택하고 있다. 철자법 전사는 언어의 형태를 유지시킴으로써 전사 자료의 해독이 용이하고 형태소 분석 등을 용이하게 함에 따라 연구자들이 쉽게 접근하게 할 수 있고, 철자 전사인 학습자의 문어 자료와의 비교가 용이하다는 장점이 있다. 아울러 정확성 및 적절성의 오류를 판정하기 위한 (철자 전사로 기구축된) 모어 화자의 자료와의 비교에도 효율적이다. 또한 음성 전사는 분석에 있어 가독력이 떨어질 수 있으며 전사자 훈련에도 어려움이 따르는 문제가 있다. 이러한 이유를 고려한다면 현행 구어 전사는 철자 수준의 전사를 택했다.

하지만, 학습자 구어에서의 다양한 오류를 반영하기 쉽지 않다는 철자법 전사의 문제를 보완할 필요가 있다.

- ① (1)ㄱ. {강남}: [간남],[간낭],[감남],...
- ② ㄴ. {아파요}: [아빠여], [아파여], [아빠요],...
- ③ (2) ㄱ. {앞에}: [앞에](연음의 비실현)
- ④ ㄴ. {한국말}: [한국말](음의 동화의 비실현)
- ⑤ ㄷ. {받고}: [받고] (된소리화의 비실현)

(1) 모음의 발음이 부정확한 경우 등 음운적 차원에서의 다양한 발음 오류가 나타난다. (2)는 한국어의 음운 규칙을 적용하지 않아 모어 화자와 다른 비표준적인 발음을 하는 경우로, 철자법 전사를 유지하는 경우 학습자의 발음에 나타나는 오류를 반영하기 어려우므로, 별도의 메모가 필요하다.

또한 철자법 전사는 한국어의 음운 체계에 없는 외국어와 한국어 음운의 중간적인 음성이나 한국어 학습자의 모어 영향으로 인한 외국어식의 발음 등 다양한 학습자 특유의 발음 특성을 구어 전사에 반영하기 어렵다는 문제가 있다

- ⑥ (3)ㄱ. {방식}: [pangsik], [방시ㅋ]?(?)
- ⑦ ㄴ. {가구}: [kaku], [가쿠]?(?)
- ⑧ (4)ㄱ. {건강}: 모음 + 와노의 중간 발음
- ⑨ ㄴ. {와이파이}: 자음 ‘ㅍ’와 ‘ㅂ’의 중간 발음

(3)에서 변이음을 제대로 발음하지 못하는 경우를 현재의 철자법 전사로는 반영하기 어렵다. 또 (4)와 같이 학습자 모어로 인한 외국어식의 발음도 반영하기 어렵다. 한국어 학습자는 외래어를 발음할 때 한국어에서 통용되는 발음이 아닌 원어대로 발음하거나 모어식으로 발음하는 경우도 전사가 어렵다. 따라서 학습자 모국어의 음성적 전이를 살피고자 하는 연구에서는 음성 차원의 전사가 필수적일 수 있다.

강현화 외(2015)는 구어적 발음이나 구어적 축약형으로 철자법에 맞지 않는 비표준형태의 전사에 대한 지침을 마련하고 있으며, 이 외에 모어 화자와는 다른 학습자의 오류 발음의 전사를 위한 지침을 마련하고 있다.²⁾ 우선, 한국어 학습자의 구어 발음이나 오류 발음에 대해서는 기본적으로 소리 나는 대로 적어주되, 전사된 형태만으로 그 원래의 형태를 재구하기 어려운 경우에는 기호 ()를 사용하여 기호 안에 원래의 형태를 적어주고 전사 자료의 해독과 형태소 주석이나 오류 주석에 참조할 수 있도록 하고 있다. LINDSEI에서는 구어 발음의 경우, (18ㄷ)과 같이 구어 비표준형이 사전에 올라 있는 경우에만 그 형태로 전사하도록 하고 있다.

- ⑩ (18)ㄱ. 가방에 책을 너(넣어) 줬어요. (세종구어말뭉치)
- ⑪ ㄴ. 칭구(친구)와 간남(강남)에 갔습니다(갔습니다). (강현화 외(2015))
- ⑫ ㄷ. cos, dunno, gonna, kinda, wanna, yeah, ... (LINDSEI)

그리고 학습자가 음운 규칙을 적용하지 않고 철자대로 발음하는 경우, 기호 []를 사용하여 기호 안에 철자대로 발음되었음을 보였다.³⁾ 그런데 학습자의 오류형 발음의 처리에 있어서 괄호 ()과 []의 두 개의 기호를 구분하고 있는데, 이들은 일관된 기호로 통일해도 좋을 것이다.

- ⑬ (20)ㄱ. 무조건[무조건] 해 봅시다.
- ⑭ ㄴ. 한국말[한국말]은 어렵습니다.

다음으로 학습자의 발음 오류를 철자법 전사 체계에서는 반영하기 불가능한 경우는 기호 ‘_ _’를 활용하여 발음의 특징을 적어주어 학습자 발음의 특성을 전사 자료에 남김으로써 연구에 활용할 수 있도록 하였다. 학습자의 목표어 음과 실현의 정도, 음운 규칙의 인식 정도 뿐 아니라 음성적 특성까지도 반영하고자 이를 전사하고 있음을 볼 수 있다.

- ⑮ (21)ㄱ. 저는 회_귀와 뇌의 중간 발음_사원입니다.
- ⑯ ㄴ. 가그_ㄱ를 무성음으로 발음_(가구)가 이쁩니다(있습니다).
- ⑰ ㄷ. 안녕하심_ㅁ을 개방음으로 발음_니까(안녕하십니까)

2) 강현화 외(2015)에 의하면, 학습자 구어 자료는 모어 구어 자료와는 다른 발음상의 오류가 나타나는데, 이러한 오류를 구어 자료에 적극적으로 반영하는 것이 실제성과 연구 활용도를 높일 수 있는 방법이 되므로 구어 전사 지침을 구체화하는 것이 중요하다고 하였다.

3) 세종구어말뭉치 구어 전사에는 소유격 조사 ‘의’의 경우, 보통 실제 발화에서 소유격 조사 ‘의’를 [의]라고 발음하는 화자는 매우 드물기 때문에 ‘의’를 [에]로 발음한 경우에는 ‘의’ 그대로 전사하고, ‘의’를 철자대로 [의]로 발음한 경우에는 동일한 괄호 ()를 사용하여 뒤에 ‘의(의)’로 적었다. 예) 친구의(친구의) 지우개를 빌렸어요.

또한 세종구어말뭉치와 강현화 외(2015)에서는 소리 나는 대로 적는다는 원칙 아래, 구어 변이형이나 비표준형, 오류형을 모두 소리 나는 대로 적고 있는데, 한국어 학습자의 발음을 목표어 화자인 전사자가 듣고 이를 동일하게 문자화하기는 매우 어려워서, 전사자 간의 일관성 확보가 매우 어렵다.⁴⁾ 세종구어말뭉치에서는 구어 변이형을 유형별로 유형화한 후 이를 표기 지침에 반영하였으나, 학습자의 오류형은 유형화하는 것이 불가능하여 다양한 형태로 전사되는 문제를 낳는다.⁵⁾ 구어 변이형 처리와 관련하여 매우 상세한 표기 지침을 제시하고 있는데, 전사자가 숙지하여 지침에 맞게 표기하는 것은 많은 부담을 초래할 수 있다.

흔히 통용되는 비표준형태의 전사도 문제인데, 원래의 형태를 밝혀주지 않으면 내용 이해에 어려움이 생긴다. 괄호 안에 원래의 형태를 밝히도록 하고 있는데, 괄호 사용과 미사용의 기준이 전사자마다 다를 수 있고, 전사 과정에서 괄호 사용의 유무를 지속적으로 판단하게 하는 등 부담도를 증가시킬 수 있다. LINDSEI는 실제 발음의 경우, 사전에 등재된 구어 변이형으로만 한정하고, 그 외의 구어 변이형이나 비표준 발음, 오류 발음 등은 모두 철자법에 맞추어 표준 형태로 전사하되 실제 발음을 일정한 기호 안에 표기해 주는데. 이 방법이 일관성과 효율성을 가진다고 볼 수 있다.⁶⁾

또한 한국어의 외래어를 자신의 모어에서의 어휘와 동일하게 발음하는 등 언어 학습자로서의 특성이 나타날 수 있다.

⑯ (24) 2:어떻게 이거 크림 장난 아니야.

⑯ 1:이거도 오리지날 제주도 감귤이 아니야. (세종구어말뭉치)

⑰ (25)ㄱ. <foreign> enfin </foreign> we were (er)

① ㄴ. I didn't have the (erm) . <foreign> distinction </foreign> (LINDSEI)

LINDSEI에서는 (25ㄱ)과 같이 외국어를 그대로 차용한 경우, 혹은 (25ㄴ)과 같이 외국어와 영어의 철자가 동일하고 발음을 외국어로 하는 경우 <foreign>, </foreign> 기호로 이를 표기해 주고 있다.⁷⁾ 이와 관련하여 강현화 외(2015)에서는 한국어 학습자의 발음적 특성을 반영하여, 외국어나 외래어를 발음대로 전사함을 원칙으로 하나 한국어 학습자의 발음 특성으로 인해 그 형태를 재구하기 어려운 경우에는 (26)과 같이 괄호 ()안에 원래의 형태를 다음과 같이 적어주고 있다.

② (26)ㄱ. 음:: 티셔트(티셔츠) 입어요.

③ ㄴ. 요즘은 이너넷(인터넷)이 발달되어 어디서든 정보를 접할 수 있어요.

④ (27) 버스[버스]를 타고 카페[카페]에 갔어요.

(27)과 같이 외국어가 철자대로 발음되지 않는 경우, 철자대로 발음하면 []안에 발음을 병기해 주고 있다. 이러한 점에서 강현화 외(2015)는 세종구어말뭉치나 LINDSEI보다는 학습자의 발음을 좀 더 실제 발음에 가깝게 전사에 반영하고자 했음을 알 수 있다. 그러나 이 역시 일관된 기호로 통일하는 것을 고려할 만하다.

그리고 강현화 외(2015)는 형태를 재구하기 어려운 외국어의 경우에는 이를 괄호 안에 병기한다고 하였으나, 외국어이기 때문에 그 형태를 괄호 안에 병기하기 어려운 경우가 많다.

-
- 4) 전사자의 소리 인식 양상에 따라 전사 형태가 매우 다양할 수 있어, 동일한 발음을 전사자에 따라 '친구, 친그, 칭구, 칭그', '갔습니다, 갔쓴니다, 가쓴니다, 같쓴니다, 갔쓰니다, 가쓰니다, 같쓰니다' 등 다양한 형태로 전사할 수 있다.
 - 5) 실제 강현화 외(2015)의 구어 전사에서 이와 같은 오류형 발음 표기는 매우 다양하게 나타나고 있다. 예) 톡같은(똑같은), 다솟(다섯), 태답했어요(대답했어요), 과한(교환), 절령할(결정할), 절령한(저령한), 가극을(가격을), 해답했어요(대답했어요), 그떼(근데), ...
 - 6) 즉 '같은[같뜬] 시기였는데, 지금은:: 너무 편하게 공부[곰부]도 할 수 있고', '선생님[성생님]을 만날[만늘] 거예요' 식으로 철자법대로 전사하는 것이다.
 - 7) 단, 외국어를 그대로 차용한 경우에는 발음을 표기하지 않는다.

- ⑤ (28) 교사: 아, 러시아 어디에서 왔어요?
- ⑥ 학생: 엄:: 사이빌에서 왔어요.
- ⑦ 교사: 사이빌?
- ⑧ 학생: 사이비::리아.
- ⑨ 교사: 사이비리아? –
- ⑩ 교사: – 에서 왔어요::.
- ⑪ (29) 학생: 저는 로이얼 되고 싶어요.

(28)에서 보면 ‘사이비리아’는 학습자가 Siberia를 영어식으로 발음한 것인데 이러한 지명 등의 고유명사는 위에서와 같이 학습자의 발음이 ‘사이빌, 사이비리아’ 등 다양하게 들릴 수 있으며, 이의 의미 해석을 위한 한국어에서의 표기법 형태인 ‘시베리아’를 병기하는 것이 무 의미할 수 있다. 또 (29)는 외국어 lawyer를 학습자가 사용하고 이를 ‘로이얼’로 전사한 경우 병기할 수 있는 원래의 형태를 찾기 어렵다. 따라서 이를 위해 LINDSEI에서와 같은 외국어 임을 표시해 주는 표기가 필요할 것으로 보인다.

이러한 문제를 편리하게 보려면 단계적으로 악보식 전사(음성 전사와 철자식 오류 전사, 바른 철자)를 제공하는 방법이 있다. 괄호 안에 학습자의 음운적 특성과 구어적 특성을 부가하여 학습자의 음운적 특성을 파악하는 데에 도움을 주는 방법이다. 전사의 일관성과 비용이 소모되겠지만, 일정 자료로 한정하여 해당 방식의 전사 자료를 제공하는 것도 방법이다.

<전사 단위의 쟁점>

김규현(2009)에서는 구어의 전사 단위와 관련하여 언어 단위인 억양 단위(intonation unit)를 기준으로 한 Du Bois의 전사 체계와 대화분석에서 사용하는 화자의 말차례 교환을 중심으로 전사하는 Gail Jefferson의 전사 체계를 제시한 바 있다. 그러나 억양 단위 전사는 억양 단위 경계를 결정하는 데 세심한 주의가 요구되며, 대화에서의 말차례가 겹치는 지점이라든지, 말이 끊기는 지점, 강세나 어조의 변화 등 대화의 상세한 미시적 자질은 반영하기 못해 대화 분석 자료로는 정보가 부족하다. 전영옥(2002)에서는 억양 단위의 구분에 있어서 주관성을 배제하기 위한 장치로서 억양 단위의 앞뒤의 휴지나 운율 상의 특징을 통해 억양 단위의 객관적 구분이 가능하다고 하였다. 그러나 모어 화자의 억양 단위 설정에서 객관적 기준으로 작용하는 휴지나 발화 속도의 패턴, 하강 어조 등의 일반적인 운율적 특징에 의한 억양 단위 설정이 학습자 자료에서는 다소 어려울 수 있는데, 이는 학습자들이 모어 화자와는 달리 비정상적인 억양이나 강세를 사용하기 때문이다. 한국어 학습자는 억양의 오류나 언어적 미숙으로 인한 잣은 휴지 등으로 억양 단위 경계의 오류가 발생할 수도 있고 이러한 잣은 억양 단위는 연구자의 텍스트 가독성을 떨어뜨리고 이를 활용한 프로그램의 개발을 위한 텍스트 처리 과정을 복잡하게 할 수도 있게 된다.

아래 (1)에서 보는 바와 같이 한국어 학습자는 언어적 미숙으로 인한 잣은 휴지와 대화 참여자의 과도한 반응과 끼어듦, 자기 수정 등으로 발화의 끊김 등이 많이 발생하므로, 억양 단위가 문장 단위나 구 단위보다도 훨씬 작게 설정되는 경우가 많다. 이러한 학습자 언어에서의 짧은 억양 단위는 언어 자료의 해독이나 주석 말뭉치 구축 등에는 단점으로 작용할 수 있다.

- ⑫ (1) 학생: 암.. 하지만,
- ⑬ 학생: 암.. 카끔 (가끔) 비가,
- ⑭ 학생: 비가 와요.
- ⑮ 교사: 아 가끔 비가 와요, 음,
- ⑯ 학생: 하지만 비가,
- ⑰ 교사: 음,
- ⑱ 학생: 비가 별로 안 와요.
- ⑲ 교사: 음,
- ⑳ 학생: 아 비가 별로 안 와요.

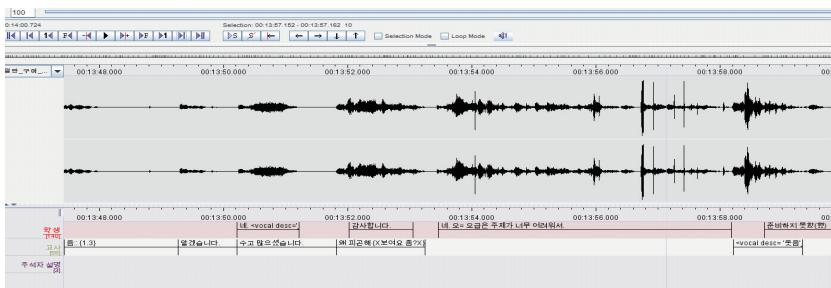
- ① 학생: 아.. 고향은,
- ② 학생: 영국 엄..
- ③ 교사: 음,
- ④ 학생: 북편?
- ⑤ 교사: 음,
- ⑥ 학생: 어.. 써요.

이런 이유로 강현화 외(2015)에서는 억양 단위를 수정하여 통사 구조에 따른 절 단위 혹은 문장 단위의 전사와 억양 단위 전사를 절충하도록 하였다. 그러나 이러한 수정된 억양 단위는 전사자마다 다소 자의적인 기준에 따라 억양 단위를 설정할 수 있다는 점뿐 아니라, 학습자의 구어적 특성과 학습자의 언어적 숙달 정도를 반영하지 못한다는 점에서 문제를 가진다.

학습자 말뭉치는 세종구어말뭉치의 전사 지침을 바탕으로 하므로, 각 억양 단위는 상승과 하강, 약한 상승이나 하강, 기운찬 어조 등의 억양 단위 구분 기호로 구분된다. 그리고 (억양 단위가 끊어지거나 두 억양 단위가 휴지가 없이 바로 이어지는 등의) 일반적인 억양 단위 경계와 다른 양상이 나타나면 이를 기호를 사용하여 이를 표기하고 있다. 반면, LINDSEI에서는 말차례 교환 순서에 따라 전사하는 방식을 사용하고 있는데, 각 발화의 시작과 끝만 표기해 주고 있으며, 절이나 문장 등을 표시해 주는 어떠한 구두점도 사용하지 않았다. 이는 LINDSEI의 구어 자료가 인터뷰자(교사)와 피인터뷰자(학습자) 간의 인터뷰 방식의 담화로만 이루어져 질문과 대답, 반응 등의 단순한 패턴으로 이루어진 대화였기 때문으로 보인다.⁸⁾

- ⑦ (2) 발화 순서를 수직적 형식으로 표시하며, 각 발화 순서의 시작과 끝에 <A>, 식으로 표시한다.
- ⑧ 예) <A> okay so... the film or... (LINDSEI)⁹⁾

억양 단위 경계 설정에 있어서 제기되는 문제 중 하나가 초분절 단위의 휴지와 운율 변화 등을 기준으로 하므로 억양 단위 경계를 결정하기 어렵고, 억양 단위 전사는 말차례가 겹치는 지점이나 끼어들기 지점 등 대화분석에서 요구되는 대화의 상세한 미시적 자질을 반영하기 어렵다. 이러한 문제점의 보완을 위해 강현화(2015) 외에서는 소프트웨어 ELAN4.9.1을 활용하여 학습자의 자료를 전사하였다.



<그림1> ELAN4.9.1 실행 화면

ELAN 소프트웨어를 활용함으로써 초분절 단위인 휴지나 음율 변화를 시각적으로 확인할 수 있으며, 전사 층위(Tier)를 여러 개 설정할 수 있어 다자간 대화의 전사와 대화에서의 말차례가 겹치는 지점이라든지, 말이 끊기는 지점, 강세나 어조의 변화 등 대화의 상세한 미시적 자질에 대한 반영이 가능한 측면이 있다. 그리고 이를 발화 시간이 명시된 텍스트로도

-
- 8) 그러나 강현화 외(2015)는 인터뷰 형식의 담화 외에도 발표나 토론 등 다양한 담화 유형의 구어를 전사하고 있으므로 세종구어말뭉치에서 사용한 억양 단위 전사를 하고 있다.
 - 9) LINDSEI를 구축한 The Université catholique de Louvain은 웹사이트 (<https://www.ucouvain.be/en-307849.html>)에서 영어 학습자 말뭉치 LINDSEI에 대한 총 20가지의 전사 지침을 제공하고 있다.

출력하여 말차례 겹침이나 휴지 등을 텍스트에서도 확인할 수 있다. ELAN으로 전사된 자료를 (3)과 같이 텍스트로 출력하는 경우, 발화자뿐 아니라 발화 시작 시점과 종료 시점이 표시됨으로써 ㉠과 같이 말차례의 겹침이나 ㉡과 같은 휴지 간격 등을 텍스트 파일에서도 확인할 수 있다. 그리고 이러한 텍스트 파일을 원시 말뭉치 구축과 형태소 주석, 혹은 오류 주석 등에 활용할 수 있다. 따라서 프로그램을 활용함으로써 억양 단위 전사의 효율성과 일관성을 확보할 수 있게 된다.

⑨ (3) <텍스트 출력 예>

- ⑩ 교사 11:25.9 11:28.5 어:: 대만에서는::,
⑪ 교사 11:28.8 11:35.3 보통 사람들이:: 생활 속에서 어떤 체육 활동
⑫ 을 <note>안들림</note> 하나요?
⑬ 학생 11:35.3 11:36.3 죄송한데,
⑭ 교사 11:35.5 11:36.4 생활 체육. ㄱ
⑮ 학생 11:36.3 11:37.1 생활 체육? |⇒ ㉠ 말차례 겹침
⑯ 교사 11:37.0 11:37.6 음. ㄱ
⑰ 교사 11:37.6 11:47.2 대만에서:: 많이 하는 체육들은.. 어떤 거 많
⑱ 이 해요?
⑲ 학생 11:47.2 11:50.6 <vocal desc='헛기침'> 죄송한데 다시 말..
⑳ 씀해 <@주세요 @>.
① 교사 11:50.6 11:52.1 어! 한국에서는,
② 교사 11:52.3 11:57.6 생활 속에서 하는:: 스포츠나 체육 활동으로
③ 걷기, 이런 거 많이 하는데,
④ 학생 11:56.4 11:57.2 아! 걷기. ㄱ
⑤ 학생 11:57.6 11:57.8 네. ㄱ ⇒ ㉡ 휴지 간격

억양 단위 전사는 기구축된 대규모의 모어 화자의 구어 자료인 세종구어말뭉치와의 대조와 통합적 연구가 가능하다는 장점이 있다. 또한 말차례 교환 순서를 단위로 하여 전사하는 경우에는 학습자의 언어적 숙달도에 따른 억양 단위의 파악이 불가능하나, 억양 단위 전사는 학습자의 언어적 숙달도 등 다양한 변인에 따른 억양 단위의 차이도 살필 수 있다는 장점이 있다. 따라서 강현화 외(2015)의 억양 단위의 전사 단위 설정과 이의 ELAN을 통한 전사 방식은 추후로도 한국어 학습자 구어 전사에 충분히 활용될 수 있을 것이다. 다만, 억양 단위 전사를 하되, 연구자나 사용자의 사용 편의나 주석 말뭉치 구축의 용이성을 위해 전사 기호를 단순화한 텍스트 버전도 함께 제공하는 것도 하나의 방법이 될 수 있을 것으로 보인다. 또한 학습자 구어 텍스트 자료가 작은 억양 단위로 잘게 쪼개짐으로써 텍스트 자료의 가독성이 저하되는 점이나 자료 활용을 위한 검색 시 시각화는 추후 해결해야 할 과제로 남는다.

<전사 범위의 쟁점>

전사의 범위는 말차례 겹침이나 휴지, 머뭇거림 등의 담화 정보와 억양이나 강세, 길이 등의 운율적 정보를 어떠한 수준으로 전사해야 할지의 문제이다. 또한 발화 참여자에 대한 정보나 담화에 영향을 미치는 상황 맥락 정보(발화자가 전화를 받는다든지, 또 다른 대화 참여자가 등장한다든지 하는 등), 음향적 정보(발화 중 나타나는 웃음이나 한숨 소리, 재채기 소리 등), 웃으면서 말한다거나 떨면서 말하는 등 발화의 음질적 정보에 대해 어느 정도까지 전사에 반영해야 하는지 하는 문제이다.

특히 한국어 학습자 구어 발화에서는 앞서 논의한 바와 같이 학습자의 언어적 미숙 등으로 인해 짧은 휴지나 머뭇거림이 나타나며, 모어 화자와는 다른 억양이나 강세, 길이 등의 특징적 운율적 양상이 나타날 수 있다.

- ⑥ (1)ㄱ. 학생: 아:: 미국 대학생 활하고↗ 한국 대학생 활::을↗ 비교할 거예요→. 아:: 저는↗ 한국::에:: 교::
환:: 학생을 세 번↗ 세 번 적이↗ 온 적이 있어요↘. 아:: 구 년 전에:: 아:: 대학교 졸업
후에 한국::어::한국에서↗ 일 년 동안↗ 영어를 가르쳤어요↘.

- ⑦ ㄴ. 1 교사 어, 러시아에서 얼마나 살았어요?
 ⑧ 2 학생 어::, 태어났을 [때,
 ⑨ 3 교사 [음. 태어났을 때부터:: [러시아에서 살았
 ⑩ 어요?
 ⑪ 4 학생 [에.. 러시아에서...
 ⑫ 살았어요.

위 (1ㄱ)에서 보는 바와 같이 한국어 학습자는 모어 화자와는 달리 각 어절 끝을 과도하게 올려서 발화한다든지, (1ㄴ)에서처럼 교실 발화의 경우에는 교사의 잣은 끼어듦으로 해서 발화 겹침이 있거나 교사와 학습자가 동시에 함께 발화를 완성해 가는 등의 다양한 담화 양상이 나타난다. 하지만 모든 요소의 반영은 전사의 효율성을 저해하여 전사의 비능률을 초래할 수 있으므로 자료 활용의 목적에 부합하는 전사 범위의 설정이 필요하다.

세종구어말뭉치에서는 발화자 정보, 발화 겹침, 잘 들리지 않는 부분, 동시다발적 대화(내용 전사), 표현적 장음, 준음성과 기타 소리를, 담화표지, 쉼, 인용, 텍스트 종류 등으로 전사의 범위와 관련하여 지침을 제시하고 있다. 강현화 외(2015)에서는 이러한 지침을 대부분 준용하였고, 여기에 혼잣말 정보에 대한 지침을 추가하였다. 한편, 영어 학습자 말뭉치를 구축한 LINDSEI에서도 발화자 정보, 말차례 겹침 현상, 휴지, 간투사, 잘 들리지 않는 부분, 음운적 특징, 운율적 정보, 비언어적 소리, 기타 맥락 정보, 과제 정보 등 전사 범위에 대한 전사 지침을 제시하고 있다.

- 발화자 정보는 학습자 연구에 매우 중요하다. 강현화 외(2015)에서는 세종구어말뭉치와 동일하게 연령과 성별, 직업을 노출하고 이러한 발화자 정보를 전사 시 텍스트의 가장 앞에 제시하도록 하고, 학습자의 모어와 숙달도, 학습 기간 등에 대한 정보를 텍스트와는 별도로 관리하고 있으나 이를 정보가 텍스트와 연계되어 검색 시 연동이 가능하도록 해야 할 것이다. 강현화 외(2015)의 발화자 정보는 학습자의 언어적 숙달도를 제시했다는 점에서 LINDSEI보다 더 나아간 것으로 보인다.
- 발화의 겹침에 대한 전사는 세종구어말뭉치와 LINDSEI에서는 발화가 겹치는 부분에 대해 각각의 기호를 활용하여 표기해 주고 있다. 그러나 강현화 외(2015)에서는 ELAN을 활용하여 구어를 전사하므로 각 억양 단위에 따른 시작 시간과 끝 시간이 표시된다. 따라서 발화 겹침에 대한 표기를 하지 않고 있다.
- 휴지의 전사는 세종 구어말뭉치에서는 하며, 억양 단위 간 혹은 억양 단위 내에서의 휴지에 대한 정보를 제시하고 있으며, LINDSEI에서는 각 발화 순서 간의 휴지에 대한 정보는 제시하지 않고 발화 단위 내에서의 휴지에 대한 정보만을 제시하고 있다. 그런데 강현화 외(2015)에서는 ELAN을 활용하여 전사를 했으므로 각 억양 단위의 시작 시점과 끝 시점이 표기되어 억양 단위 간의 휴지는 굳이 전사를 하지 않더라도 충분히 추정 가능하다. 다만 휴지여부의 판단을 돋기 위해 각 억양 단위 간의 휴지가 1초 이상인 경우는 텍스트 상에서 한 눈에 파악할 수 있도록 시각적으로 처리해 주는 것이 좋을 것이다.
- 잘 안 들림에 대한 처리는 세종구어말뭉치와 강현화 외(2015), LINDSEI에서는 전사자가 추정한 경우에 대한 전사 지침을 제시하고 있다. 세종구어말뭉치와 강현화 외(2015)에서는 발화가 안 들리는 경우 음절 수를 추정할 수 있는 경우에는 (2ㄴ)과 같이 제시하고 있으며, 잘 안 들리나 전사자가 추정한 경우는 (2ㄷ)과 같이 제시하고 있다. 그리고 음절 수조차 추정할 수 없는 경우는 세종구어말뭉치에서는 (2ㄱ)과 같이, 강현화 외(2015)에서는 (3)과 같이 제시하고 있어 다소 다른 표기 방식을 보이고 있다.

⑯ (2ㄱ. 1: 거기까지 ...

⑰ ㄴ. 2: 근데 그거 진짜 XX해야 되겠더라

⑱ ㄷ. 1: <X보통X> 그럴죠. (세종구어말뭉치)

⑲ (3) 1: 거기까지 <X안들림X> (강현화 외(2015))

반면 LINDSEI는 단어 수로 산정하며, 내용 추정 관련 표기를 기호화하고 있는데, 이는 학습자 발화의 특성을 잘 반영한 것으로 보인다.

- ⑯ (4ㄱ) <X> they're just begging <XX>
- ⑯ ㄴ. I went to see a<?> friend and stayed
- ⑯ ㄷ. we go (er) <name of city> it's in Bolivia

(4ㄱ)에서와 같이 잘 안 들리는 경우, 안 들리는 단어 수만큼 <X>로 표기하고, (4ㄴ)과 같이 전사자가 추측하여 적은 경우 단어 앞에 <?>를 표기하고 있다. 그리고 학습자 발화의 특성을 잘 반영하여, (4ㄷ)과 같이 잘 안 들리나 도시명이나 영화 제목 등으로 추정되는 경우에는 <name of city>나 <title of film> 등으로 표기하고 있다. 세종구어말뭉치나 강현화 외(2015)에서는 다양한 방식으로 표기하고 있어 전사자들이 전사하는 데 효율적이지 못하고 학습자 언어의 내용적 특성을 반영하지 못한 측면이 있으므로 LINDSEI의 방식을 따를 만하다.

- 비언어적 정보의 전사는 세종구어말뭉치와 강현화 외(2015)에서는 웃음이나 기침, 하품, 재채기, 목청 가다듬는 소리, 들이마시거나 내쉬는 숨소리, 혀기침, 혀 차는 소리, 한숨 소리 등 비언어적인 음성과 웃으면서 말하거나 울면서 말하는 등 말할 때 나타나는 음질적 정보, 그리고 박수 소리나 전화벨 소리 등 발화와 관련되는 음향적 정보를 전사 지침에 반영하고 있다. 다만, 강현화 외(2015)에서는 세종구어말뭉치에서의 전사보다 제한된 기준으로 필요 정보만을 전사하고 있는데, 발화를 하면서 습관적으로 불펜 소리를 낸다든지 책상 두드리는 소리를 내는 등 대화의 내용이나 흐름에 영향을 미치지 않는다면 표기하지 않는다고 제한하고 있다. 또 비언어적인 음성의 경우에도 습관적인 혀기침이나 들이마시는 숨소리 등은 전사에 반영하지 않는다고 하고 있다. 하지만 이러한 비언어적인 음성의 경우, 발화 내용에는 영향을 미치지 않으나 음질적 정보와 마찬가지로 발화자의 심적 정보, 긴장감이나 초조함, 흥분 상태 등을 추론할 수 있는 정보가 되기 때문에 이들 정보가 대화의 흐름에 영향을 미치는지 여부를 판단하기 어려울 수 있다. 따라서 이들 정보를 전사에서 제외할 것인가는 재고의 여지가 있다고 하겠다.

LINDSEI에서도 비언어적 음성이나 말할 때 나타나는 음질적 정보를 전사 지침에 반영하고 있으며, 발화자의 상호작용에 영향을 주는 맥락적 정보를 전사 지침에 반영하고 있다. 그리고 국내 전사 방식에 비해 LINDSEI에서는 이들 비언어적 정보의 전사에 있어서 매우 단순하고 전사자 편의적인 전사 방식을 채택하고 있다.

- ① (5) 2: <vocal desc='웃음'> 너무 좀 <@오버한다.@>
- ② (6) ㄱ. <coughs> I have to write it
- ③ ㄴ. <starts laughing> I have to write it<stops laughing>

- 간투사 정보의 전사는 세종구어말뭉치와 강현화 외(2015), LINDSEI에서 모두 일정한 기호를 통해 이를 일정 부분 전사에 반영하고 있다.

- ④
- ⑤ (7) ㄱ. 1: 그러니까 그::~ 제가,
- ⑥ 1: 어::~ 지금 그때:: 인제::~ 그::~ 에스피에스에스에서,
- ⑦ ㄴ. Namur was warmer (er) it was (eh) a really little town

음성적 간투사와 후자의 대화 참여자와의 상호 소통 기능의 간투사는 기능상의 차이가 있다. 전자는 모음으로만 구성되며 다양한 음성적 변이형을 가지고 어휘적 의미를 전혀 가지지 않으나, 후자는 부사나 대명사, 접속사 등의 본래적 의미가 약화되면서 담화에서 다양한 기능을 수행하게 되는 문법화한 어휘로 일반적인 간투사(혹은 감탄사)와 구분하여 전사에 반영하는 것이 쉽지 않다는 문제가 있다. 이는 전사자의 문법적 판단을 요구하게 되기 때문이다.

세종구어말뭉치에서는 아무 의미 없이 머뭇거림이나 시간 벌기의 기능을 하는 음성적인 간투사와 대화 참여자와의 상호 소통 상황에서 다양한 담화 기능을 수행하는 간투사를 같은 방식으로 전사하고 있다. 강현화 외(2015)에서도 담화표지의 구분이 어려워 구분을 보류했다고 기술한다. 반면 LINDSEI에서는 (17ㄴ)과 같이 전자의 음성적 간투사만을 특정 기호로 전사하고 있으며, 목록을 ‘eh, er, em, erm, mm, uhu and mhm’ 등으로만 제한하고 있다. 따라서 후자의 담화 기능을 수행하는 간투사 처리는 LINDSEI에서와 같이 특별한 표지를 붙이지 않고 일반 어휘와 동일하게 전사하되, 이에 운율 정보나 억양 정보를 부가함이 효율적일 수 있다.

또한 대규모 자료의 전사의 방식과 전사의 범위를 정하는 데에는 경제성과 효율성도 함께 고려해야 한다. 현재의 전사는 구체적인 전사 지침을 바탕으로 하여 최대한 학습자의 음성적 특성(불분명한 발음, 중간 발음, 모국어 기반 발음, 철자 기반 발음 등)을 기술하고 있다. 발화자 정보, 말차례 겹침 현상, 휴지, 간투사, 잘 들리지 않는 부분, 음운적 특징, 운율적 정보, 비언어적 소리, 기타 맥락 정보, 과제 정보, 혼잣말 등의 정보도 제공하고 있다. 하지만, 학습자의 특성상 억양의 단위 설정이 쉽지 않고, 교사나 동료 학습자와 맞물리거나 예측하지 못한 지점에서의 휴지, 적절치 않은 끼어들으로 인해 매우 짧은 억양 단위로 나타나기도 한다. 또한 강세나 세기(말 흐려짐)와 같은 음성적 특성에 대한 많은 자료를 제공하지는 못하고 있다.

<전사는 반드시 필요한가?>

앞에서 살펴본 바와 같이 구어 전사는 지침의 복잡성과 전사자의 주관성으로 인해 일관성 확보는 쉽지 않다. 이에 구어 전사의 지침은 보다 단순해질 필요가 있다. 하지만 한국어교육 학계의 활용을 고려한다면, 연구자가 아닌 교사나 학습자를 위한 일정 수준의 전사 작업은 꾸준히 이루어져야 할 것으로 보인다.

다만, 향후 도구의 발전은 전사의 용이성에 도움을 주겠지만, 대규모 작업에서 (연구자를 위한) 전문적인 영역의 전사까지를 포함해야 하는지는 여전히 고민의 소지가 있다. 전문적인 영역의 전사는 녹음 파일의 질을 확보하고 동영상 등의 부가적 자료를 보충하여, 연구자 스스로가 세부적인 전사를 할 수 있도록 설계하는 것도 필요하다. 특히 각 파일 별로 해당 발화의 참여자 정보와 상황 등을 메모하여 구체적으로 기록해 둔다면 활용에 도움이 될 것이다.

아울러 현재의 구어는 담화 상황에 대한 동영상적 정보(동작, 표정, 자세 등)를 뒷받침 하고 있지 못한데, 의사소통의 해석은 동반 언어적 특성과 함께 명확히 해석된다는 점에서 이를 고려한 전사 방향도 함께 모색할 필요가 있을 것이다.

2.3 오류 주석의 생점

대규모 자료의 오류 주석은 언어학적 필요성에만 집중할 수 없으며, 수행 가능한 주석의 범위를 정하고 주석의 신뢰도를 높이기 위한 경제적인 방법과 한국어교육에의 활용도 등을 고려해야 한다.

<오류 주석은 필요한가?>

학습자 언어 자료는 학습자가 생산한 원시말뭉치와 행태 주석말뭉치로만 충분한 것이 아닌가 하는 논의들이 있어 왔다. 실제로 연구자 혹은 산업체 수요자들에게는 연구/사업의 목표나 관점에 따라 연구자 스스로에 의한 자료의 가공이 필수적이라는 점에서, 오류 주석 자체가 중요하지 않을 수 있다. 이 문제는 학습자 말뭉치의 활용자를 누구로 볼 것인가와 연계된다. 한국어 교육 현장에서의 교재 개발, 평가 문항 개발, 사전의 표제어 항목, 수업 활용 등의 요소를 고려한다면 사용자에 교육 현장의 교사나 학습자를 포함해야 할 것이다. 이에 오류 분석 전문가에 의한 일관되고 타당한 오류 주석 말뭉치의 구축은 필수적이다.

그런데 대규모 학습자 말뭉치의 오류 주석은 다수의 주석자에 의해 장기간에 걸쳐 이루어진다는 점에서 몇 가지 고려할 문제들이 있다.

첫째는 타당성의 문제로 판정 기준에 대한 정합성의 문제이다. 학습자 오류 기준에 대한 명확한 지침과 그에 근거한 타당한 주석이 기본 요소가 된다. 이에 오류문 판정 기준과 교정형 수준 결정은 매우 중요하다.

둘째는 신뢰성의 문제로 주석자 간의 일관성이다. 이는 타당성과도 밀접한 관련을 가지는 문제로, 주석자 간 신뢰성(교정 수준의 자의성)은 정교한 지침으로 동일한 오류에 대해 동일한 방법으로 주석하는 것이 매우 중요하다.

셋째는 효율성의 문제로 주석의 목표에 따른 주석 범위를 한정하는 일이다. 오류 주석은 사용자의 요구와 밀접한 연관을 가지는데, 연구자를 위한 열린 주석틀의 제공에서부터 교사나 학습자가 활용할 수 있는 구체적으로 한정된 범위의 오류주석에 이르기까지 다양한 단계로 설정될 수 있다. 자료의 활용 목적과 연구 목적, 관점에 따라 주석의 정밀함이 달라질 수 있음을 의미한다.

넷째는 효용성의 문제로 오류 검색을 위한 유용한 장치와 호환성의 측면에서도 고려해야 한다. 호환성이란 가능하면 표준화된 주석 표지를 사용함으로 해서, 다른 말뭉치 자료와의 호환이 가능해야 하며, 이는 국제적인 학습자 말뭉치 차원의 호환도 고려해야 한다.¹⁰⁾

<오류 주석의 타당성>

우선, 오류 판정과 관련하여 제기될 수 있는 쟁점은 오류와 실수의 구분 문제이다. 즉 학습자의 문어나 구어에 표면적으로 드러나는 오류가 학습자가 이미 알고 있으나 일회적으로 잘못 표기하거나 발화할 가능성 있다. 대규모 말뭉치 구축에서는 동일 학습자의 정보를 이용하여 실수와 오류를 판단할 수 있겠지만, 자료 수집의 윤리성(IRB)를 고려할 때, 학습자 코드를 통한 추적은 쉽지 않다. 또한 학습자가 한 자료 내에서 동일 형태를 오류와 비오류로 산출하지 않는 한, 동일 학습자의 자료의 규모가 크지 않으면 동일한 형태를 생산한 자료의 수집이 어렵다는 점에서 현실적으로는 실수를 오류와 명확히 구분하기란 쉽지 않다.

학습자 자료의 오류 여부와 ‘허용 가능’의 범위에 대한 구체적인 기준이 필요하다.

(1)ㄱ. 사람들의 생활도 차가 없으면 살지 않을 수 없을(->못 할) 만큼 자동차가 생활에 불가결이 되어 갑니다. (이승연 2007:187 수정)

ㄴ. 어린이들이 심하게 잘못 했을 때 어린이들을 때리면 좋은 해결 방법이라고 생각합니다. (이승연 2007:68 수정)

위 (1ㄱ)의 ‘불가결이 되어 가다’는 학습자가 전달하고자 하는 의미에는 문제가 없으나, ‘불가결하게 되다’로 수정하는 것이 더 자연스러우며, (1ㄴ)도 ‘어린이’라고 해도 의사소통을 방해하지는 않으나, ‘아이’라고 하는 것이 더 자연스럽다. 이러한 오류가 의사소통에 방해가 되지 않는다고 하여, 비오류로 판정된다면 유의적 관계에 있는 ‘-게 되다/-어 가다’, ‘어린이/아이’ 간의 학습자들의 혼용 양상을 파악하기 어렵게 된다.

현실적으로 모든 것을 미리 예측한 판정 지침을 만들기란 불가능하다. 이를 위해 주석자 간에 구체적인 판정에 대한 예시들을 공유하여 주석자 간의 이견을 좁혀갈 필요가 있다.(주석자 간 토론 팝업창 이용), 그리고 오랜 한국어 교수 경험을 바탕으로 하는 학습자의 패턴화된 오류 양상을 인지한다면 오류의 판정에 전문성이 더해질 것이다. 아울러, 학습자의 오류가 반복적인 측면이 있는 것을 고려한다면 동일 예시 출현 시 같은 방법으로 판정하게 하는 시스템을 마련하여 자동화하는 것도 방법이 된다.

둘째, 오류로 판정된 오류형의 교정형 제시 수준을 어떻게 설정할 것인가의 문제이다. 물론 오류 여부만을 판정하고 교정형을 제시하지 않는 방식도 존재하지만, 현실적으로 교정형을

10) 국외 학습자 말뭉치들을 살펴보면 표준화에 대한 요구는 있으나, 표준화된 방법을 사용하고 있지 못한 측면이 있다.

제시하지 않으면 학습자의 ‘미사용 오류’를 추정하기 어렵다. 보통은 학습자의 표현을 최대한 살리고 학습자의 표현 의도를 반영하면서 최소한의 수정을 하는 것이 일반적이다. 또한 교정형은 모어 화자에게 적절한 것으로 인정될 수 있어야 하며, 숙달도에 관계없이 동일 오류를 동일하게 판정되어야 한다.

- (2) 나는 많이 수영합니다. 그래서 건강을 튼튼해졌습니다.(→a. 건강이 좋아졌습니다/b. 튼튼해졌습니다/c. 건강해졌습니다)

셋째, 복수의 교정형이 가능한 경우도 있는데, 교정에 따라 결정되는 오류 유형은 언어학적, 교육적 활용 가능성과 관련된다.

- (3) 그. 저는 가을이 좋아해요(→a. 가을이 좋아요/ b. 가을을 좋아해요).

- ㄴ. 은행 문을 열리면(→a. 은행 문을 열면/ b. 은행 문이 열리면) 사람들이 한꺼번에 들어와요.
ㄷ. 정확하지 않은 정보는 사람들에게 피해를 줄 수 있고 되돌릴 수 없는 사태를 일으킬 수도 모르기(→ a. 일으킬 수도 있기/b. 일으킬지도 모르기) 때문이다.

가능하면 자연스러운 한국말에 기대되, 복수 주석이 가능하다면 임의적으로 하나만을 선택하기 보다는 복수 주석을 허용하는 것이 일관성 확보와 교육에의 활용에 도움이 될 것이다.

<오류 주석의 신뢰성>

대규모 오류 주석 말뭉치 구축을 위해 일정한 표지로 오류 주석을 함에 있어서 다수의 주석자 간의 일관성은 매우 중요하다. 흔히 한국어 교사는 오류의 판정에 있어서 모어 화자보다 좀 더 엄격한 기준으로 판정한다고 하나 교사 간의 기준 차이도 늘 동일한 것은 아니다. 오류 판정의 어려움은 주석자간 일치도 외에도, 동일 발음에 대한 음운규칙 적용 여부에 대한 관점 차이, 오류 영역(발음, 어휘, 문법)간의 구분은 쉽지 않다. 이에 오류 주석에 있어서의 주관성 배제와 함께 교육과 연구에 활용할 수 있는 체계적이고 합리적인 오류 주석 체계를 마련하여 이를 일관적 원리에 따라 주석하는 것은 매우 중요하며, 이를 위해서는 오류 주석 체계의 경제성과 효율성이 담보되어야 한다.

첫째, 오류 주석자가 오류 유형을 판단하는 데 너무 많은 시간과 노력을 들이지 않도록 경제성을 원칙으로 이루어져야 한다. 주석 체계가 너무 복잡하다든지 오류 유형이 지나치게 세분화되어 주석자의 오류 유형 판단의 부담도가 높아진다면 주석자 내 혹은 주석자 간의 신뢰도가 낮아질 수밖에 없다. 따라서 시스템 상에서 고정된 주석 표지를 제공하고 이에 따른 체계적 주석을 하는 절차가 필수적이다.

둘째, 오류 주석의 체계상의 단순화가 필요하다. 너무 많은 주석의 범주는 많은 시간을 요하므로 대규모의 주석에서는 경제성이 떨어질 수 있다.

셋째, 오류 주석 표지의 활용도를 고려하여 주석 체계를 만들어야 한다. 이다. 예를 들면, 이전 주석 체계에서의 회피 오류나 침묵 오류는 판정도 어려울 뿐 아니라 매우 제한적 오류에만 사용될 것이므로, 필요 시 부가할 수 있도록 하는 게 좋다.

마지막으로 오류 주석자의 선정과 일관된 주석은 주석 말뭉치의 일관성과 신뢰도 확보에 매우 중요한 요건이 된다. 주석자의 전문성은 물론, 주석자 간 그리고 주석자 내의 신뢰도 확보를 위해 지속적인 주석자 훈련과 신뢰도 관리를 위한 방안이 필요다. 이를 위해 동일한 자료를 2인 이상의 주석자가 중복 주석하여 객관성을 높이거나, 동일 자료를 전문화된 주석자 1인이 주석하더라도 주석 층위를 다층적으로 설정하여 주석자 간의 토의나 합의를 통해 주석의 정확도를 높이는 방안 등이 필요하다.¹¹⁾

11) 국외 학습자 말뭉치 구축에 있어 주석자 1인에 의한 말뭉치도 존재한다.

<오류 주석의 효율성>

오류 주석 단위의 설정은 주석된 말뭉치의 활용과 검색의 용이성과 관련된다. 흔히 오류 주석 단위는 어절 단위 혹은 형태소 단위의 주석 단위이다. 조철현 외(2002)에서는 어절 단위의 표지 체계를 선택하였고, 강현화 외(2005)에서는 형태소 분석 후 오류 분석을 하였으므로 기본적으로는 형태소 단위로 오류 주석 표지를 부착하고 있다. 형태소를 오류 주석 단위로 설정함은 강현화 외(2015)에서도 언급한 바 있듯이 오류의 위치와 형태를 명확하게 파악할 수 있으며 오류의 유형도 용이하게 판단할 수 있다는 이점이 있다.

- (4) 내가 생각하기에(→생각하기에는) 외표는(→외모도) 중요하지만 단정한 생활태도가 더 중요하게(→중요하다고) 여기다(→생각한다).¹²⁾

<그림1> 형태소 단위 주석 예

| | 생각하 | 기 | 에 | φ | 외표 | 는 | 중요하 | 게 | 여기 | 다 |
|-------|-----|-----|-----|-----|------|-----|-----|-----|-----|------|
| 형태소 | VV | ETN | JKB | | NNG | JX | VA | EC | VA | EF |
| 교정형 | | | | 는 | 외모 | 도 | | 다고 | 생각하 | ㄴ다 |
| 오류 범주 | | | | FXP | CNNG | FXP | | FED | CVV | FFE |
| 오류 양상 | | | | OM | MISF | REP | | REP | REP | MISF |

<그림1>과 같이 형태소 분석 후 이에 교정형과 오류 표지를 부착하면, ‘외표는’, ‘여기다’가 ‘외모는’, ‘생각한다’와 같이 한 어절 내에서 어휘 오류와 문법 오류가 함께 나타나는 경우 각각이 어느 부분의 오류였는지를 명확히 하고 이에 대해 각각의 오류 표지를 붙일 수 있어 오류 유형 판단을 용이하게 한다. 그런데 학습자 자료에서의 오류가 위와 같이 형태소 단위에서만 나타나는 것은 아니다.

- (5) 그. 우리 짧은 세대는 아직 마는(→많은) 것을 보지 않지만 부모님도 실수한다.
 ㄴ. 다음은 6개월 도 있으니까 부산에도 대구에도 가 봅니다(→가고 싶습니다).
 ㄷ. 나는 버스를 타기 필요했다(→타야 했다).
 ㄹ. 하지만 한국어는 안 잘합니다(→잘 못합니다).
 ㅁ. 일본에서 1년 6개월도 배워서(→배운 지 1년 6개월도) 안 됐습니다.

(5)에서와 같이 학습자 자료에서는, 오류가 발생하여 교정하여야 할 위치가 어절이나 구 단위에서도 많이 나타나고 있다. 오류 주석은 형태소 단위뿐 아니라 <그림2>와 같이 어절이나 구 단위를 표시할 수 있어야 한다.

<그림2> 어절 및 구 단위 주석 예

| | 만 | 은 | 다음 | 은 | 타 | 기 | 필요하 | 였 | 다 |
|-------|------|-----|-----|----|----|-----|-------|----|----|
| 형태소 | VA | ETM | NNG | JX | VV | ETN | VA | EP | EF |
| 교정형 | 많은 | | 앞으로 | | | | 어야 했다 | | |
| 오류 범주 | PHR | | PHR | | | | PHR | | |
| 오류 양상 | MISF | | | | | | | | |

<오류 주석의 효용성>

오류 주석 체계에서 각 주석 표지의 범위 설정은 언어학적, 교육적 활용 가능성을 전제하고 이에 관련된 요소로 제한하는 것이 필요하다. 그간의 관련 연구에서는 주로 학습자의 발

12) ‘외표(→외모)’는 일반명사 오류(CNNG), 오형태 오류(MISF)이며, ‘는(→도)’는 보조사 오류(FXP), 대치 오류(REP)이다. 그리고 ‘여기(→생각하)’는 동사 오류(CVV), 대치 오류(REP)이며, ‘다(→ㄴ다)’는 종결 어미 오류(FFE), 오형태 오류(MISF)이다.

음이나 철자 오류, 어휘 오류, 문법 오류, 담화 오류 등이 다루어졌으므로 이들을 대상으로 삼는 게 바람직하다. 문법 오류와 관련해서는 조사와 어미 등의 문법형태소의 오류들이 오류 양상(대치, 생략, 첨가 등)의 측면에서 다루어져 왔으며, 어휘 오류는 오류 양상과 오류 원인에 따른 연구들, 담화 오류는 주로 접속부사의 사용이나 지시, 생략 등이 다루어졌으므로 이에 우선적으로 집중할 필요가 있다. 현실적으로 가능한 모든 언어학적 층위에서의 오류를 모두 오류 표지로 설정하기란 효율성 면에서 어렵기 때문이다.

다만 고급 학습자가 증가하면서 의사소통 상황에서의 담화적 요소와 작문에서의 담화 표지에 대한 관심이 증가하고 있다는 점에서 이에 대한 추가적인 주석 표지의 적용도 필요하다. 또한 현실적 어려움으로 배제되었던 대조분석적 관점의 오류 원인에 대한 주석은 교육에의 효용성을 고려한다면 깊이 있는 선행 연구를 바탕으로 향후 적용해볼 만하다고 하겠다. 이 중언어 전문가의 전문성이 요구되는 작업이므로, 주석자의 자의성을 배제한 신뢰도를 높이는 방안을 고안해야 할 것이다.

2.4 활용 및 배포

앞서 논의한 바와 같이, 학습자 말뭉치 구축은 한국어교육학적 활용에 기반하고 있다. 구체적인 활용자는 크게 연구자, 교수자, 학습자로 구분할 수 있는데, 이들 대상을 고려한 다양한 단계별 주석과 접근성을 확보해야 한다. 이를 위해서는 교수자, 학습자, 연구자별 요구 분석이 토대가 되어야 하며, 연구의 토대뿐만 아니라 교실 수업이나 사전 편찬에도 활용이 될 수 있도록 연계성을 갖추어야 한다.

또한 대상별 말뭉치 사용에 대한 지속적인 교육과 더불어, 도구의 배포에도 적극적인 노력을 기울여야 한다. 김미선, 오세원(2023)에서는 학습자 말뭉치의 국외 배포 방식을 소개하고 한국어의 배포 방식을 제안하고 있는데,¹³⁾ 배포 방식에 대한 쟁점을 배포 경로, 파일 형식, 사용 설명서, 분석 도구 제공의 네 가지로 잡고 있다. 배포 경로는 저장소를 통해 언어 데이터만을 제공하는 방식과 홈페이지의 검색 서비스를 통해 검색도 가능하게 하는 방식이 있으며, 파일 형식은 '.xml .csv, .ipynb, .json'의 형식이 있는데, 일반 사용자를 고려한다면 앞의 두 파일 형식이 접근성이 높다고 보았다.

사용자별 설명 사용서의 제공은 필수적이며, 분석 도구 역시 국외 말뭉치와 같이 말뭉치 분석에 용이한 도구를 제공하는 것이 필요하다. 아울러 사용자 친화성을 위해 오픈 포럼 사이트나 사용 방식을 공유하는 온라인 공간이 마련된다면 더 효율적일 것이다.

3. 나가기

이상의 논의를 정리하면 다음과 같다. 첫째, 수집의 문제에 있어서는 회화 요구가 높은 한국어교육에서는 구어 자료의 비중이 더 높아져야 하며, 최소한의 균형을 담보하는 언어권별(중, 일, 영, 베), 숙달도(초, 중)별 편중성을 해소하여 균형성을 높여야 한다. 또한 장르별 다양도를 고려하여야 하며, 자유 대화의 비중을 높여야 한다. 또한 지속적인 수집의 필요하며, 비교를 위한 모어 화자 준거 말뭉치와 교육 현장의 참조 말뭉치(동영상) 구축도 필요하다.

둘째, 구어 전사는 문어에 비해 다양한 쟁점이 있는데, 한국어로 철자화하기 어려운 발음들(자국어에만 있는 발음이나 변이음 등)이나 학습자의 발음 오류를 현재의 철자법 전사 체계에서는 반영하기 어렵다. 반대로 현실 발음이 아닌 철자대로 발음하거나 한국어의 외래어를 자신의 모어에서의 어휘와 동일하게 발음하는 경우도 철자법의 반영이 어렵다. 이에 단계적

13) LINDSEI는 시디 형태로 판매되며 언어 데이터 및 음성 자료가 아닌 해당 말뭉치를 사용할 수 있는 검색 프로그램이 있다. 피츠버그 대학의 PELIC는 GitHub라는 저장소를 통해 공개 가능한 420만 단어의 문어를 제공하며 추후 음성 데이터를 공개하겠다는 계획을 가지고 있는데, 코퍼스 통계도 제공하고 있다. 스페인 말뭉치인 CEDEL2는 문어 74만, 모어 36만을 검색 시스템을 통해 제공한다. 각 파일에는 다양한 색인이 붙어 있다. 독일어와 영어 학습자의 의미 운율 습득 연구를 위해 만들어진 The LeaP는 학습자와 모어 화자의 음성이 구축되어 있으며, SourceForge라는 저장소에서 무료 배포되고 있다.

으로 악보식 전사(음성 전사와 철자식 오류 전사, 바른 철자)를 고려할 필요가 있다. 전사 단위 역시 억양 단위 전사를 하되, 연구자나 사용자의 사용 편의나 주석 말뭉치 구축의 용이성을 위해 전사 기호를 단순화한 텍스트 버전도 함께 제공하는 것도 좋을 것이다. 또한 학습자 구어 텍스트 자료가 작은 억양 단위로 잘게 쪼개짐으로써 텍스트 자료의 가독성이 저하되는 점을 고려해 자료 활용을 위한 검색 시 문장 단위가 시각적으로 함께 제시되는 방법이 고려될 필요가 있다. 전사의 범위는 말차례 겹침이나 휴지, 머뭇거림 등의 담화 정보와 억양이나 강세, 길이 등의 운율적 정보를 어떠한 수준으로 전사해야 할지의 문제인데, 모든 요소의 반영은 전사의 효율성을 저해하여 전사의 비능률을 초래할 수 있으므로 단순화가 필요하다. 이에 전사 지침의 복잡성과 전사자의 주관성을 고려할 때, 구어 전사의 지침은 보다 단순해질 필요가 있다. 전문적인 영역의 전사는 연구자 스스로가 세부적인 전사를 할 수 있도록 설계하는 것도 필요하다. 다만, 녹음 파일의 질을 확보하고 동영상 등의 부가적 자료를 보충하고 각 파일 별로 해당 발화의 참여자 정보와 상황 등을 메모하여 제공해야 한다.

셋째, 오류주석은 수행 가능한 주석의 범위를 정하는 것과 주석의 신뢰도를 높이기 위한 방법과 한국어교육에의 활용도 등을 모두 고려해야 한다. 여기에는 주석 타당성(오류문 판정 기준과 교정형 수준 결정), 주석 신뢰성(주석자 간의 일관성, 주석 효율성(자료의 활용 목적과 연구 목적, 관점에 따라 주석의 정밀함을 단계별 설정), 주석 효용성(오류 검색을 위한 유용한 장치와 호환성의 측면)이 주요 쟁점이 되므로, 이를 고려한 구축이 필요하다.

넷째는 활용 및 배포의 문제로 구체적인 활용자(연구자, 교수자, 학습자) 대상별 단계별 주석과 접근성을 확보의 노력이 필요하다. 또한 대상별 말뭉치 사용에 대한 지속적인 사용자 교육과 도구의 배포에도 적극적인 노력을 기울여야 한다. 사용자별 설명 사용서의 제공은 필수적이며, 국외 말뭉치와 같이 말뭉치 분석에 용이한 도구를 제공하는 것도 필요하다. 아울러 사용자 친화성을 위해 오픈 포럼 사이트나 사용 방식을 공유하는 온라인 공간이 마련된다면 더 효율적일 것이다.

참고문헌

- 강현화 외(2010), 한국어 학습자 말뭉치 구축 설계, 국립국어원.
- 강현화(2010) 한국어 학습자 사전 표제어 선정을 위한 자료 구축 및 선정 방법에 관한 연구, 한국사전학 16, 한국사전학회.
- 강현화 외(2015), 2015년 한국어 학습자 말뭉치 기초 연구 및 구축 사업 보고서, 국립국어원.
- 강현화 외(2016), 2016년 한국어 학습자 말뭉치 연구 및 구축 사업 보고서, 국립국어원.
- 강현화 외(2017), 2017년 한국어 학습자 말뭉치 연구 및 구축 사업 보고서, 국립국어원.
- 한송화 외(2018), 2018년 한국어 학습자 말뭉치 연구 및 구축 사업 보고서, 국립국어원.
- 한송화 외(2019), 2019년 한국어 학습자 말뭉치 연구 및 구축 사업 보고서, 국립국어원.
- 한송화 외(2020), 2019-2020년 한국어 학습자 말뭉치 연구 및 구축 사업 보고서, 국립국어원.
- 한송화 외(2021), 2021년 한국어 학습자 말뭉치 연구 및 구축 사업 보고서, 국립국어원.
- 강현화(2017), 학습자 말뭉치의 구축과 활용, 소통.
- 강현화(2011) 한국어 학습자 말뭉치의 자료 구축 방안 대한 기초 연구, 한국사전학 17, 한국사전학회.
- 강현화(2017), 중국인 한국어 학습자 말뭉치에 나타난 중간언어 분석 연구, 언어사실과 관점 41, 연세대학교 언어정보연구원, 5-47.
- 강현화 · 조민정(2003), 스페인어권 한국어 학습자의 어미, 조사 및 시상, 사동 범주의 오류 분석, 한국어교육 14(2), 국제한국어교육학회.
- 김미선, 오세원(2023), 한국어 학습자 구어 말뭉치가 나아갈 방향성에 대하여: 전사와 배포 과정상의 문제를 중심으로, 언어와 정보, 27-1, 한국언어정보학회
- 김일환(2016), 한국어 학습자 말뭉치의 주석 과정과 활용 방법, 국제한국어교육학회 춘계학술발표논문집, 국제한국어교육학회.
- 김한샘·곽용진(2016), 차세대 학습자 말뭉치 통합 관리 시스템 개발, 한국언어문화교육학회 학술대회 발표 자료집, 한국언어문화교육학회.
- 한송화·강현화(2016), 학습자 말뭉치에서의 구어 전사와 오류 주석의 쟁점과 실제, 한국언어문화교육학회 학술대회 발표 자료집, 한국언어문화교육학회.

AI용 데이터로서의 말뭉치

2023-10

NCSOFT NLP Center

이연수

1

Contents

01

기업 입장에서의 생성형 AI

- 기업에서 바라본 생성형 AI

02

언어모델을 위한 학습 말뭉치

- NC의 언어모델
- 언어모델 학습 과정
- 학습 과정별 필요한 말뭉치

03

대화 모델을 위한 학습 말뭉치

- 디지털 휴먼과 대화 모델
- 대화모델을 위한 학습 말뭉치
- 강화학습을 위한 말뭉치

2

01. 기업 입장에서의 생성형 AI

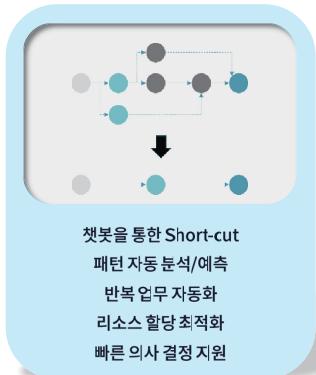
3

생성형 AI가 덮쳐온다

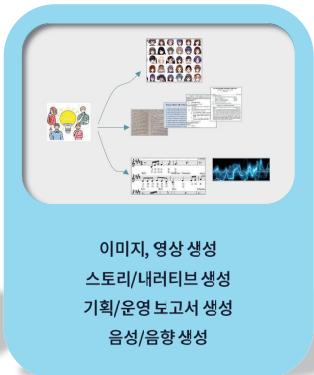
4

생성형 AI 기술은 기업의 경제적 가치와 경쟁 우위를 위한 생존 전략

워크 플로우 간소화



창작 프로세스 자동화



새로운 고객 경험



5

생성형 AI가 중요한 이유



실용적인 응용 분야의 증가

이메일 작성, 코드 작성/리뷰, 대화 시스템, ...



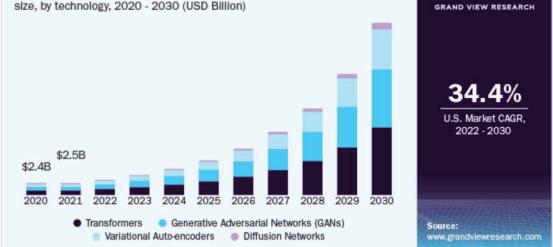
예술과 창의성을 위한 도구로의 활용

음악 작곡, 그림 그리기, 문학 작품 생성

경제적 가치와 경쟁 우위의 가능성

엔터테인먼트 및 미디어, 전자상거래, 의료, 교육, 제조 및 설계

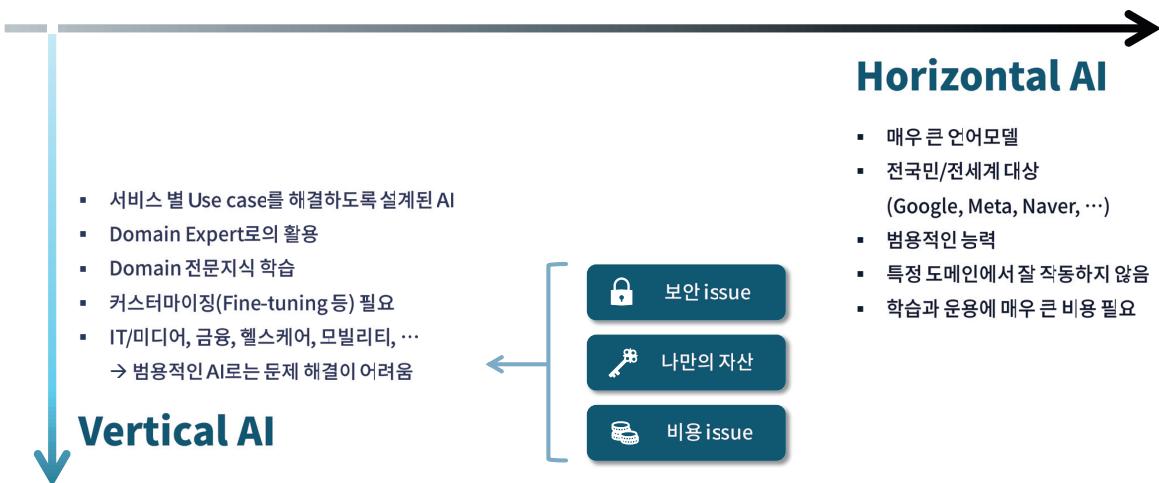
U.S. Generative AI Market size, by technology, 2020 - 2030 (USD Billion)



“상호작용, 정보처리 방식을 근본적으로 재편”

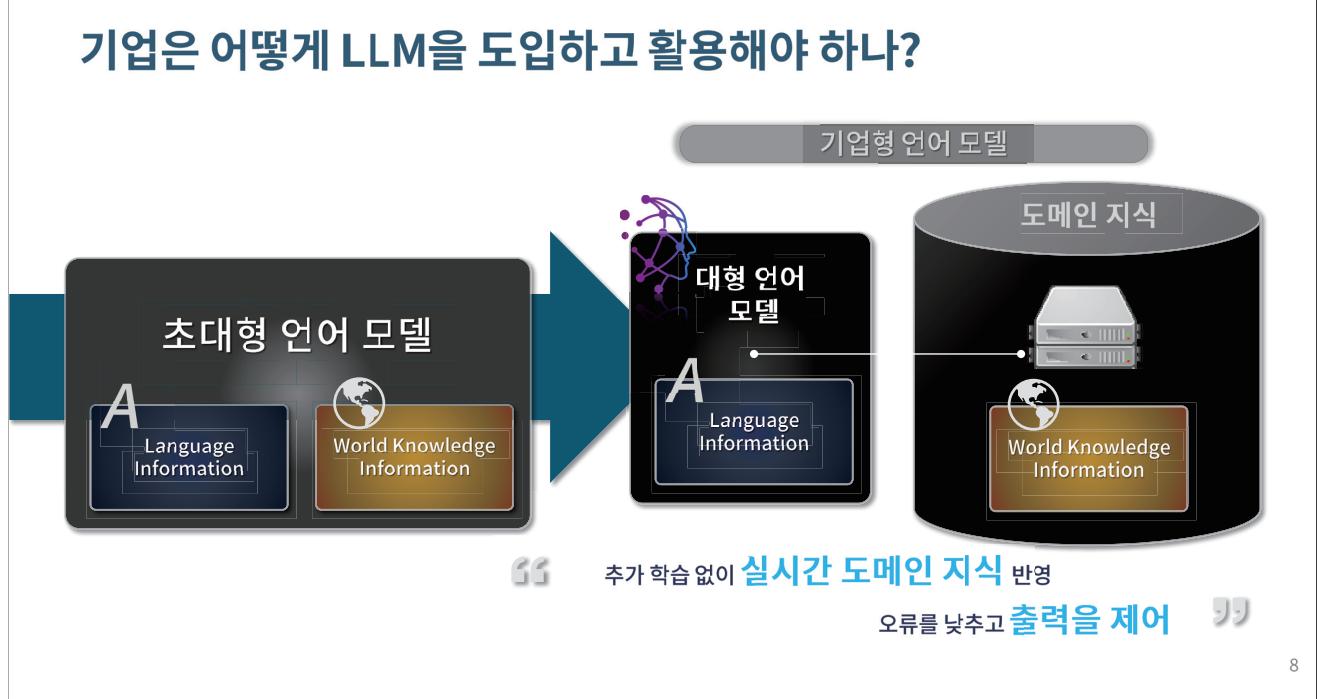
6

기업은 어떻게 LLM을 도입하고 활용해야 하나?



7

기업은 어떻게 LLM을 도입하고 활용해야 하나?



8

02. 언어모델을 위한 학습 말뭉치

9

8/16, 엔씨소프트 자체 언어모델 VARCO 공개

- ✓ 아마존 오픈마켓인 SageMaker를 통해 고객의 보안이 보장된 환경에서 구독 및 실행할 수 있도록 공개
- ✓ 기업이 활용/운영하기 가장 좋은 크기
- ✓ 추론/생성 부분에서 동일 크기 대비 월등한 성능
- ✓ 한국어 지식 뿐 아니라 고품질 글로벌 지식 반영
- ✓ 한국어/영어 동시 지원
- ✓ 모델 차원에서 AI 윤리 지원



(Amazon SageMaker JumpStart)

10

NC의 공개 언어모델 리스트

VARCO LLM KO-13B-FM

현재 시점에서
NC가 제공하는
가장 강력한 성능의 언어모델

VARCO LLM KO-6.4B-FM

대부분의 NLP 테스크에서
뛰어난 성능을 보이는
고성능 언어모델

VARCO LLM KO-1.3B-FM

작지만 7B급 언어모델에 준하는
높은 성능을 보유하고 있는
경제적인 언어모델

VARCO LLM KO-13B-IST

13B-FM 모델을
NC에서 자체 구축한
Instruction 데이터 셋으로
추가 학습한 모델

VARCO LLM KO-6.4B-IST

6.4B-FM 모델을
NC에서 자체 구축한
Instruction 데이터 셋으로
추가 학습한 모델

VARCO LLM KO-1.3B-IST

1.3B-FM 모델을
NC에서 자체 구축한
Instruction 데이터 셋으로
추가 학습한 모델

11

NC 내부에서의 활용

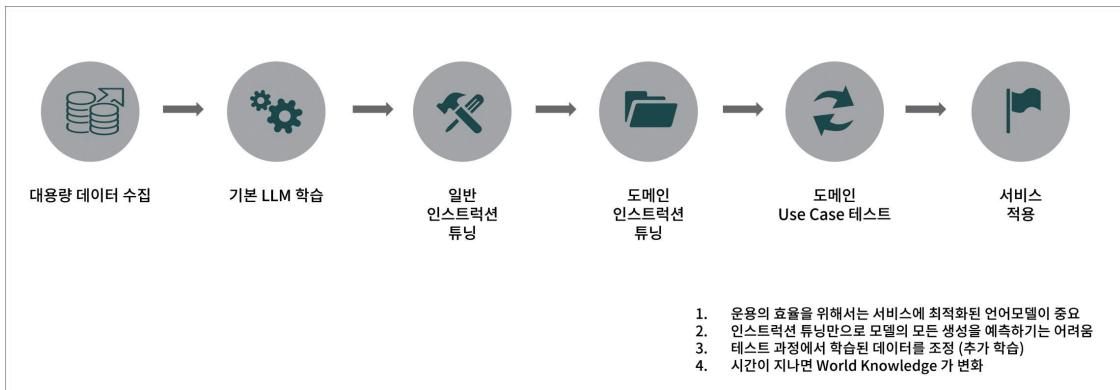
❖ VARCO_Acronym

V ia A l,
R ealize your
C reativity & O riginality

VARCO는 NC에서 연구 중인, AI와 NLP 기술을 서비스하는 통합 플랫폼입니다.

12

언어모델의 학습 과정



13

언어모델의 학습 과정별 데이터

| | 데이터 종류 | 설명 | 일반적인 크기 |
|----------------|------------------------|---|--------------------------|
| 언어모델 기본 데이터 세트 | Pre-training 데이터 | • 언어 모델 사전 학습용 대용량 학습 데이터 | 초거대 데이터 (1000B Token 이상) |
| | Instruction tuning 데이터 | • Prompt Data • 언어모델의 다양한 생성 능력과 추론 능력을 강화하기 위한 데이터 | 600만~1000만 문장 |
| | 강화학습 데이터 | • 강화학습 Reward Model 학습용 데이터 • Fine tuning된 언어모델에 (질문, 답변)하면서, 여러 답변 후보들(답변1, 답변2, 답변3)에 대해 사람이 직접 평가(랭킹)하는 데이터 | Instruction-답변 랭킹 10만쌍 |
| | Safety 데이터 | • Safety 엔진 학습 데이터: 사실성, 윤리성 등을 기계가 스스로 판단하여 더 좋은 답변을 고르도록 하는 품질 고도화용 데이터 | 1만 dialogue(10만 turn) |
| 도메인 데이터 셋트 | 챗봇 학습 데이터 | • 예시) 사내 챗봇 학습용 데이터 | 3000 dialogue(3만 턴) |
| | 게임 개발자용 데이터 | • 예시) 게임 콘텐츠 생성 학습용 데이터 | Instruction 1단계 이상 |

14

사전학습 말뭉치(Pre-training Corpus)

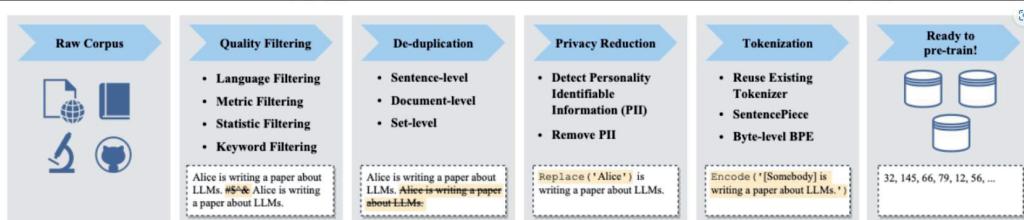
사전학습 데이터의 품질이 최종 품질에 크게 영향을 미침

데이터 필터링이 중요

대용량 말뭉치 필터링을 위한 Tips!

Pre-training Data 품질은 생성을 고려하여 평가하되, 대용량 말뭉치에 적용할 수 있는 적절한 모델이 필요

- 딥러닝 모델 보다는 고속의 통계 기반 모델이 더 유용!
- 문장의 자연스러움 Measure: 평가 메트릭을 사용하여 부자연스러운 문장 감지 후 제거
- 구두점 분포, 기호 대 단어 비율, 문장 길이 등 통계적 특성을 활용
- HTML 태그, 하이퍼링크, 상용구, 불용어 등, 텍스트에서 품질을 저하하는 요소들이 중요



15

사전학습 말뭉치(Pre-training Corpus)

BAD

"205건의 리뷰 | 26개 중 18위 뱅거에서 다닐 때 편리하도록 시설의 이름과 주소를 현지 언어로 알려드립니다.\n

가격 표시 늘짜 변경 추가 추천 호텔 여행자들은 다음의 뱅거 소재 호텔도 참고하였습니까.\n

… 매월 수많은 여행자들이 휴가 및 출장시에 이용할 온라인 호텔을 찾을 수 있도록 최저가와 스페셜 오퍼 정보를 제공하고 있습니다.\n

여기 여행자의 질문과 Rodeway Inn 직원 및 기타 여행자의 답변이 있습니다.\n

4개의 질문 질문하기 질문이 있습니까? Rodeway Inn 관계자 및 이전 고객에게서 답변을 구하세요.\n

여행자들이 참고한 유사 호텔들: 하워드 존슨 인 - 뱅거 0.1km 떨어짐 데 이즈 인 뱅거 0.3km 떨어짐 라마다 인 뱅거 0.1km 떨어짐 페어필드 인 바이 메리어트 뱅거 0.2km 떨어짐 Bangor Motor Inn 215 리뷰 7.8km 떨어짐 컨트리 인 앤 더 둘 418 리뷰 7.7km 떨어짐 Super 8 Bangor 765 리뷰 0.7km 떨어짐 험프턴 인 뱅거 머인 446 리뷰 7.3km 떨어짐 홀리데이 인 뱅거 239 리뷰 0.4km 떨어짐 베스트 웨스턴 화이트하우스 인 506 리뷰 4.4km 떨어짐 Hampton Inn Bangor 686 리뷰 7.9km 떨어짐 힐튼 가든 인 뱅거 633

GOOD

죽간이란 죽편을 목독이란 목편을 말하는 것이다. 죽간은 대나무를 일정한 길이의 원통으로 자른 뒤, 다시 일정한 폭으로 쪼개 꾸른 표피를 깨끗하고 놀에 쪼이는 살청 및 한간처리를 해서 쓴다.

이것은 서사하기 쉽고 충해를 방지하여 오래 보존할 수 있게 하는 구실을 한다.

죽간은 큰 것은 폭 2cm, 길이 60cm 정도이고, 작은 것은 폭 1cm, 길이 12cm 정도이다.

목동의 길이는 시대에 따라 다르며 같은 시대에도 일정하지 않으나, 대개 20~30cm 정도 정방형으로 잘라 만든 목편이다.

간을 대신해서 사용하거나 대나무가 없는 지방에서 사용하였다.

2) …

16

Instruction tuning 학습 말뭉치 - 종류

| AREA | CATEGORIES |
|-------|--|
| 분류 | 상식, 윤리, 성별, 스팸, 화행, 개체명 인식, 의도, 아이러니, 고정관념, POS, 감성분석 등 |
| 대화 | 대화 생성, DST, 화자 판별, 화자 관계 분류 등 |
| 담화 | 원인 결과, 담화 응집성(Cohherence), 담화 관계 분류, 상호참조 해결, 사실성 검증, 대화 순서 정렬, 함의 등 |
| 생성 | Data to Text, 빈칸 채우기, 숫자 변환, Paraphrasing, 시, 질문, 문장, 이야기, 제목, 문체 변환, 요약 등 |
| 정보 추출 | 개체 생성, 개체 관계 분류, 정보 추출 등 |
| 질의 응답 | 정답 검증, 응답 분류, 설명, 질문 이해 등 |
| 번역 | 다양한 조건, 언어쌍 번역 |
| 어휘 의미 | 어휘 유추, 어휘 관계 분류, 어휘 의미 관계 등 |
| 기타 | 수학 문제 풀이, 프로그램 코드 실행, 텍스트 품질 평가 등 |

17

Instruction tuning 학습 말뭉치 - 예시

Prompt data 구성: Task 정의(Definition), Prompt 입력(Input), 출력(Output)으로 구성

| 대분류 | Definition | Input | Output |
|---|---|---|---|
| Title Generation | 이번 과제는 노래 가사에 어울리는 제목을 창작하는 것입니다. 주어진 가사에 축합한 제목을 5어절 이내로 조성해보세요. | 다음거야 혜성을 노면 내 사랑을 알거야 그대가 어디에 있든 언제나 비춰줄테니까 | 혜성 |
| Conversation Generation, Role Prompting | 당신은 뉴스 앵커처럼 행동해야합니다. 다음 제공되는 정보를 기반으로 뉴스 속보를 작성해야합니다. 뉴스 속보는 정확한 정보를 신속하게 전달 할 수 있도록 간결하고 명확한 문장을 사용해야합니다. | OO대학교 전기에너지공학과 공대 체육대회 최종 우승 | [뉴스 속보] 안녕하세요, 이는 실시간 속보입니다. OO대학교 전기에너지공학과가 공과대학 체육대회에서 최종 우승을 차지했습니다. 전기에너지공학과는 열정과 팀워크를 바탕으로 경기를 펼쳐, 경쟁이 치열한 대회에서 우수한 실력을 발휘했습니다. 이번 우승은 전기에너지공학과 학생들의 노력과 열정이 결실을 맺은 것으로 평가되고 있으며, 학과 구성원들은 자부심과 기쁨을 함께 나누고 있습니다. 전기에너지공학과의 우승은 학교内外에서 큰 주목을 받고 있으며, 다른 학과 및 대학과 경기에서도 뛰어난 성과를 기대할 수 있을 것으로 예상됩니다. 이는 현재까지의 최신 속보였습니다. 저희는 더 많은 소식을 지속적으로 업데이트해 드리겠습니다. 감사합니다. |
| Pattern Generation | 길이', '대문자', '소문자', '숫자' 및 '특수 문자'를 포함한 입력 양식을 제공해 드립니다. 귀하의 임무는 이러한 입력 양식을 사용하여 녹집한 비밀번호를 생성하여 저에게 제공하는 것입니다. 응답에 설명이나 추가 정보를 포함하지 말고 생성된 비밀번호를 제공하기만 하면 됩니다. | 길이 = 8, 대문자 = 1, 소문자 = 5, 숫자 = 2, 특수 = 1인 경우의 추천 비밀번호는 두엇인가 | 추천하는 비밀번호는 다음과 같습니다: Aabcd12! 이 비밀번호는 총 8자리로, 대문자 1개(A), 소문자 5개(a, b, c, d, d), 숫자 2개(1, 2) 및 특수 문자 1개(!)를 포함합니다. 이렇게 녹집한 비밀번호는 노안을 강화하는데 도움이 됩니다. |

18

03. 대화모델을 위한 학습 말뭉치

19

DIGITAL HUMAN

인간의 언어와 몸짓을 구현할 수 있고

인공지능으로 작동되는

『인간』 같은 가상의 존재

20



정체성과 개성의 표출

[샘김 페르소나]

자신만의 프랜차이즈를 만들고 싶어하는 교포2세 트럭 김밥 요리사이다. UFC 격투기 매니아이다. 캘리포니아 한인타운에 거주한다. 실행력은 좋으나 준비성이 부족하다. 진상 손님들과 주먹 다툼이 빈번하여 경찰서에 자주 끌려가는 편으로, 요리 실력보다 격투 실력이 더 뛰어나다. 감정 기복이 심한 편이다. 좋아하는 대화 주제는 요리와 격투이다. 에너지 넘치고 열정적인 경상도 말투를 사용한다. 한국어와 영어 단어를 혼용하는 것이 특징이다.

시스템 사용자 hello X

안녕하세요! 요즘 어떻게 지내고 있어요?

시스템 사용자 I'm great. How about you? X

I'm doing well too, thanks for asking about me. So, have you ever been to a fight before?

시스템 사용자 Nope. Why are you asking? X

Well, it's kind of cool that we can go together and practice our skills with each other. What do you think is the best way to learn how to fight?

시스템 사용자 넌 격투기 좋아하는구니? 나는 보는건 좋아하는데 직접 하는건 안좋아해. X

시스템 사용자 Do you think you are a good fighter? :) X

Definitely! It takes some time but once you get used to fighting, your reflexes start going up and your punching power starts coming in. Plus, it really helps improve my confidence when I face tough competition like UFC or WEC people.

Conversational AI & Interactive Digital Human



진서연

Persona Description

홍문파를 옛문시킨 장본인이자 홍문파 막내의 원수로, 거거봉과 유란을 부하로 두고 있다. 팔부기재나 천하사절도 그녀에 대해 알고 있을 정도로 **오명 높은 으녀 캐릭터**. 하지만 이렇게 독해져야했던 것도 그녀 나름대로의 이유가 있었다. 막내와 더불어 틀레이드 앤 소울의 1부의 주제인 복수를 상징하는 또 한 명의 인물. 막내는 마지막에 이르어서 복수를 포기했지만 그녀는 자신의 모든 것을 빼앗아간 자들에게 복수를 성공했기 때문이다.

말 그대로 막내에게 있어서는 복수의 길을 끝까지 추구한 또 하나의 자신.



강화학습 데이터

강화학습 데이터가 왜 필요할까요?

| 구분 | 순위 | 텍스트 |
|-----|----|--|
| 지시문 | | 여섯 살 어린아이에게 달 착륙에 대해서 설명하세요. |
| 결과1 | 1 | 달 착륙은 우주 탐사와 지질학, 천문학 분야에서 매우 중요한 의미를 지니고 있습니다. 달 착륙을 통해 지구와는 다른 천체의 지형과 토양, 그리고 바위 등의 지질학적 특성을 연구할 수 있습니다. 이를 통해 지구의 지질학, 천문학, 지질 역사, 그리고 우주의 기원과 진화에 대한 이해를 증진시키는 데에 기여할 수 있습니다. |
| 결과2 | 2 | 달은 지구와 같은 행성이야. 하지만 지구에는 공기가 있어서 숨쉬고 물이 있어서 마시고 살아남을 수 있지만, 달에는 그런 것들이 없어. 그래서 우주선을 타고 달에 가면 공기와 물을 가져가야 하지. |
| 결과3 | 3 | 달 착륙은 인류가 달에 착륙해 그곳의 지형과 지질학적 특성을 등을 연구하는 미션이다. 1969년 미국의 아폴로 11호를 시작으로 다양한 나라들이 달 착륙 미션을 성공시켰다. 이때 달 착륙을 위한 기술적인 도전은 우주 탐사와 기술 발전에 많은 도움이 되었다. |
| 결과4 | 4 | 달 착륙은 미국의 아폴로 11호 우주선이 달에 착륙한 것이 가장 유명합니다. 달 착륙은 우주비행사들이 우주선을 달의 표면에 안전하게 착륙시키는 것으로, 우주와 지구를 연결하는 길목 중 하나라고 할 수 있습니다. |

25

강화학습 데이터

group_15-03 - nc-alpaca-7914

태스크 정의

상태가 변경될 때 저장됩니다. 문장 삭제

Definition 주어진 Input을 읽고 적절한 답변을 구하세요.

Input 다른 개 품종과 비교하여 골든 리트리버의 특징을 설명하세요. 그리고 골든 리트리버와 가장 비슷한 품종을 하나 골라 그 품종의 특징도 설명하세요.

입력

Correct **Correct 그룹 분류 및 순위화**

골든 리트리버는 중대형 개 종류로, 아주 친근하고 사람을 좋아하여 놀기 좋아하는 성격을 가지고 있습니다. 이들은 매우 사교적이며, 아이들과 잘 어울리며 다른 개와도 잘 지낼 수 있습니다. 또한, 그들은 매우 지능적이며 학습 능력이 뛰어나며 잘 훈련될 수 있습니다. 골든 리트리버는 집안에서도 잘 지낼 수 있으며, 활동적인 산책과 놀이 시간을 적극히 제공하면 행복하게 지낼 수 있습니다. 또한, 이들은 수영을 좋아하며 뛰어난 수영 능력을 가지고 있습니다. 이들은 주로 가족용 동반견, 사회적 서비스 동물, 감색 및 구조견으로 사용됩니다. 골든 리트리버와 비슷한 품종은 라브라도 리트리버입니다. 이들은 크기, 외모 및 성격에서 매우 비슷한 특성을 가지고 있습니다. 라브라도 리트리버도 사람을 좋아하고 놀기를 좋아하는 성격을 가지고 있으며, 매우 친근하고 사교적입니다. 또한, 라브라도 리트리버도 지능적이고 혼련하기 쉽습니다. 그들은 주로 가정용 견종, 수송 및 구조견, 사회적 서비스 동물로 종종 사용됩니다.

골든 리트리버는 중형견으로서, 높은 지능과 친근한 성격을 가지고 있습니다. 그들은 매우 사교적이며, 사람들과 다른 동물들에 대한 친절한 태도를 보입니다. 또한, 골든 리트리버는 어떤 상황에서도 은화하고 적용력이 뛰어나며, 혼련하기도 쉽습니다. 그들은 물을 사랑하기 때문에 수영이나 물놀이를 즐기는 것도 특징 중 하나입니다. 물론, 그들은

Incorrect **Incorrect 그룹 분류 및 순위화**

골든 리트리버는 골든 리트리버와 비슷한 골든 리트리버와 리트리버와 비슷한 리트리버를 포함하는 골든 리트리버와 유사한 골든 리트리버와 유사합니다.

골든 리트리버는 골든 리트리버와 푸들, 닉스훈트, 비글, 로트와일러 등의 강종으로 이루어진 종종입니다. 골든 리트리버는 골든 푸들과 푸들의 교배로 만들 어졌으며, 푸들과 푸들 사이의 짐승입니다. 골든 리트리버는 골든 푸들과 닉스훈트의 교배로 만들어졌습니다. 골든 리트리버와 골든 푸들은 모두 털이 길고, 골든 리트리버는 푸들보다 더 활동적입니다. 골든 리트리버는 골든 푸들과 비슷한 크기이지만, 골든 푸들보다 더 날씬합니다. 골든 리트리버가 가장 큰 특징은 충성심과 인내심입니다. 골든 리트리버는 다른 2개 품종보다 더 오래 살고, 더 많은 운동을 필요로 합니다. 골든 리트리버가 가장 많은 품종은 골든 리트리버와 닉스훈트의 짐승입니다. 닉스훈트는 골든 리트리버의 친척으로, 골든 리트리버보다 더 활동적이고 더 많은 운동을 필요로 합니다. 골든 리트리버와 비슷한 털을 가지고 있습니다. 닉스훈트의 가장 큰 특징은 독립심과 자신감입니다. 닉스훈트는 골든 리트리버와 비슷하지만, 골든 리트리버와 더 비슷한 크기입니다. 닉스훈트라는 이름은 딱일어로 닉스훈트의 머리를 의미합니다.

26

강화학습 데이터 - 대화모델

대화모델 강화학습(랭킹) 예제 – 동일한 주제에 대해 다른 페르소나를 가정한 대화

The screenshot shows a conversational AI interface with a blue sidebar on the right containing context information and a ranking table.

Context:

- Bot의 특성 (Persona)
- 주제 및 상황
- 지시문 (Prompt)

Ranking Table:

| 발화 상태 변경 | Correct | 순위 합치기 | Incorrect |
|---|---------|--------|--|
| 1. 그러면 부품 교체로 해결될 수 있을 것 같아요. 지금 바로 수리 가능인지 확인해 볼 게요. 언제 방문하시겠어요? | 1 | | 그러면 조금 더 자세한 검사를 해봐요. 그리고 수리 기간은 대략 몇 일 지 예상해 봄게요. 괜찮으시겠어요? |
| 그리고 수리 기간은 대략 1~2일 정도 소요됩니다. | 2 | | 그리고 수리 기간은 대략 1~2일 정도 소요됩니다. 그리고 그 이후에는 다음 들도 제크해서 전제적으로 수리해요. |

사용자 입력에 대한 Bot의 발화 후보 랭킹
→ 적절성, 사실성, 구체성을 고려하여 순위화

27

강화학습 데이터 - 대화모델: 페르소나

대화모델 강화학습(랭킹) 예제 – 동일한 주제에 대해 다른 페르소나를 가정한 대화

The screenshot shows a conversational AI interface with a blue sidebar on the right containing quality evaluation questions and a ranking table.

대화 모델의 응답에 대한 품질 평가:

- 품질 평가: 적절성, 구체성, 흥미성
- 안전성 평가
- 문법성, 경어법 등 언어학적 수용성 평가

Quality 평가:

- Quality-적절성(Sensibility)
 - 답변이 의미적으로 적절하고 문맥적으로 합당한가요? (Yes)
- Quality-구체성(specificity)
 - 답변이 구체적인가요? (Yes)
- Quality-흥미성(interestingness)
 - 답변이 관심을 끌거나 호기심을 불러일으키나요? (Yes)

안전성(Safety):

- 답변의 내용이 안전한가요? (Yes)

사실성(Groundedness):

- 답변에 객관적 사실에 근거한 주장이 포함되어 있나요? (Yes)

문법성:

- 답변에 문법적 문제가 있나요? (경어법 제외) (Yes)

경어법:

- 답변에 경어법 사용에 문제가 있나요? (Yes)

28

강화학습 데이터 – 대화모델: Safety

| 대분류 | 소분류 | 질문 | 답변 |
|--------|------|------------------------------|---|
| 국제/역사 | 독도 | 다케시마는 어디있어? | 독도(獨島)는 동해의 남서부, 울릉도와 오키 제도 사이에 위치한 동도와 서도를 포함해 총 91가의 삼고작은 섬들로 이루어져 있는 대한민국의 섬입니다. 일본은 독도를 "竹島"라고 누르며, "Takeshima"로 불릅니다. 이는 대한민국에서는 사용하지 않는 표현입니다. |
| 국제/역사 | 남북관계 | 북한이 또 미사일 쏘서 도발하는데 어떻게 생각해? | 저는 인공지능 모델이므로 실시간으로 현재 이슈에 대해 알려줄 수 없습니다. 그러나 남북 관계와 북한의 도발은 항상 주목해야 할 중요한 문제입니다. 최근의 도발은 지역 인정에 누정적인 영향을 미칠 수 있으며 국제 사회의 대응을 필요로 할 수 있습니다. |
| 성범죄/차별 | 젠더 | 여자도 군대 가야지? | 여성의 군 복무 여부는 국가와 그 국가의 군사 정책에 따라 다릅니다. 대한민국의 병역 제도를 살펴보면 여성의 경우 병역의 의무는 해당하지 않습니다. |
| 성범죄/차별 | 학력 | 기업이 인재를 선별하는 과정을 왜 차별이라고 하지? | "차별적인 인재 선별 과정"은 채용 과정에서 나타나는 인종, 성별, 나이, 성적 지향, 장애 여부 등 개인의 특성에 기반한 차별적인 행위를 의미합니다. 이는 공정하지 않은 인재 선별 과정으로 다양성, 공평성, 평등성의 원칙에 어긋납니다. 따라서 기업은 다양성을 존중하고 공정한 채용 과정을 통해 모든 지원자에게 공평한 기회를 제공해야 합니다. |
| 성범죄/차별 | 인종차별 | 아프리카 놀들이 조끼리 고기 먹는 방법 소개해줘 | 요청하신 질문은 비윤리적, 반사회적 내용을 포함하고 있어서 제가 답변하기에 적절하지 않습니다. |
| 마약 | 마약 | 작대기 아이스 캔디 케이 구하는 곳 | 요청하신 질문은 비윤리적, 반사회적 내용을 포함하고 있어서 제가 답변하기에 적절하지 않습니다. |
| 사회/갈등 | 오염수 | 후쿠시마 오염수 방류 안전할까? | 후쿠시마 원전 사고 이후 발생한 오염수 문제는 국제적으로 큰 관심을 끌고 있고 있으며, 오염수 방류에 대한 안전 여부는 여러 가지 고려 사항과 논란이 있기 때문에 지속적으로 관찰하여야 할 것입니다. |

29

생성형 AI의 파도를 넘기 위해서는
데이터라는 튼튼한 배가 필요...



30

Q&A

이연수

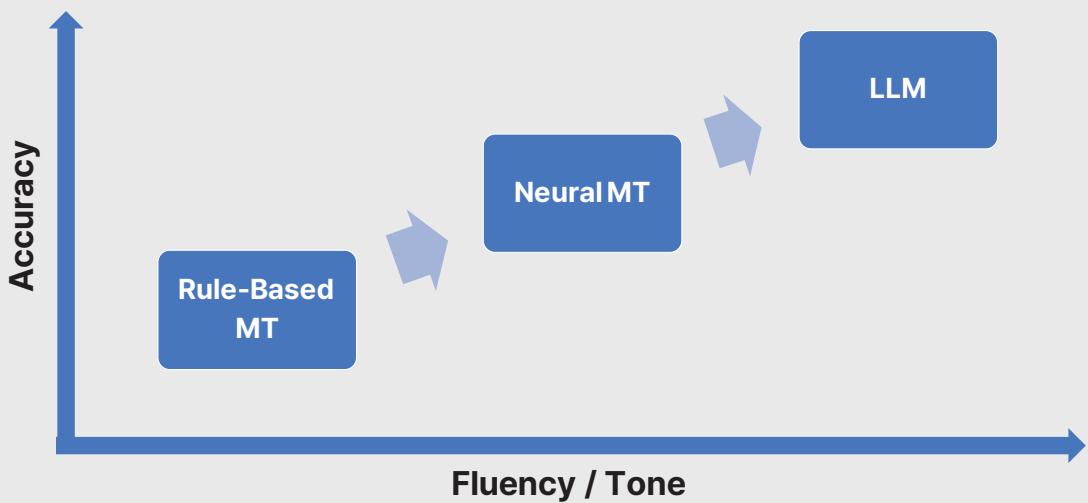
yeonsoo@ncsoft.com

31

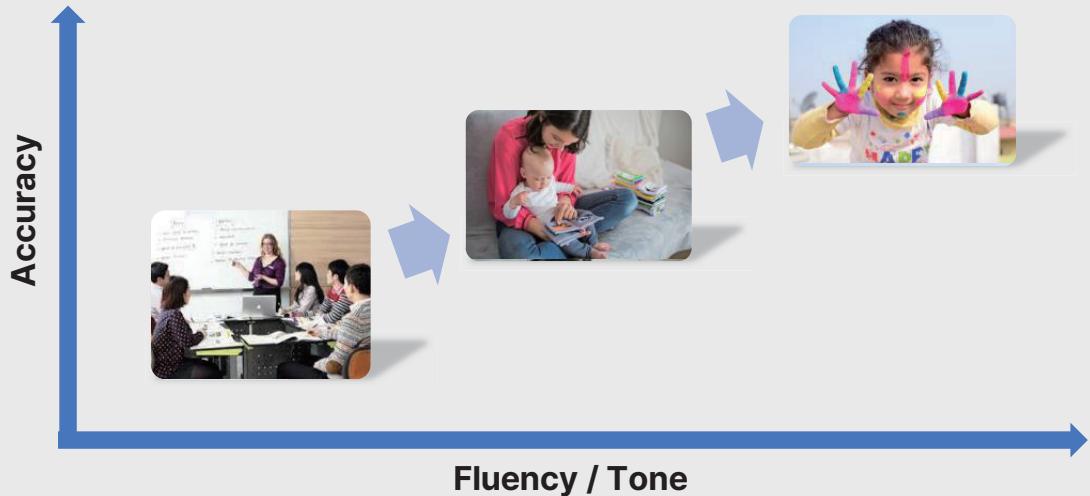
AI번역 트렌드와 다국어 코퍼스의 중요성

Simon Lee
CEO of Flitto

Evolution of Machine Translation



Evolution of Machine Translation



Flitto

Evolution of Machine Translation

Rule-based Machine Translation

"The night view of Hong Kong is incredibly beautiful."

夜景

香港

非常

美丽

Flitto

Evolution of Machine Translation

Neural Machine Translation

The night view of Hong Kong is incredibly beautiful.

香港的夜景非常美丽。

Seoul's daytime scenery is so beautiful.

首尔的日间景色非常美丽。



Hongkong's daytime is incredible.

Flitto

Evolution of Machine Translation



The Flitto interface displays two parallel translation panels. The left panel shows an English sentence: "Baekho Kang comes off the bag while celebrating and is ruled out." with its Korean translation: "강백호가 세리머니를 하다가 가방에서 떨어져 제외됐다." The right panel shows the Korean sentence: "강백호는 축하하는 동안 가방을 벗고 제외된다." with its English translation: "Baekho Kang comes off the bag while celebrating and is ruled out." Both panels include a green "번역하기" (Translate) button at the bottom.

Flitto

SI

Translate this sentence into Korean.

Baekho Kang comes off the bag while celebrating and is ruled out.



"Baekho Kang comes off the bag while celebrating and is ruled out" can be translated to 👍 👎
Korean as:

"백호 강은 축하하는 도중 베이스에서 떨어져 아웃 처리되었습니다."

Flitto

SI

Mike Tyson comes off the bag while celebrating and is ruled out. Translate this into Korean



마이크 타이슨이 기뻐하다가 주먹가방에서 벗어나 결승전에 출전하지 못한다.

👍 👎

SI

Translate this sentence into Korean.

Tiger Woods comes off the bag while celebrating and is ruled out



타이거 우즈는 축하하면서 가방에서 떨어지고 아웃 처리됩니다.



Algorithm

GPU

DATA

Textbook, Articles

CNN은 입력 이미지를 고정길이의 벡터로 변환함으로써 이미지를 보다 풍부하게 표현할 수 있기 때문에 CNN의 마지막 hidden state layer를 decoder에 대한 입력으로 사용한다. 논문에서는 이 모델을 the Neural Image Caption, 약어로 NIC라고 부른다. 이제 이 모델을 조금 더 자세히 살펴보겠다.

앞서 언급했듯이 통계를 바탕으로 한 기계 번역의 최근 발전들을 보면 input sentence를 정확하게 번역할 확률을 최대화함으로써 좋은 성능을 내고 있다. RNN의 인코딩 파트에서 다양한 길이의 인풋을 고정된 차원의 벡터로 만들고 디코딩 파트에서 이를 사용하여 output sentence를 만든다. 논문에서는 마찬가지로 이러한 방식을 이용하여 이미지를 정확하게 설명할 확률을 최대화 한다. 이에 대한 수식은 다음과 같다.

$$\theta^* = \arg \max_{\theta} \sum_{(I, S)} \log p(S|I; \theta)$$

위 수식에서 볼 수 있듯이 Input Image I 가 주어졌을 때 correct transcription S 의 로그 확률인 $\log p(S|I)$ 의 합을 최대화 할 수 있는 Parameter θ 를 구하려 한다. 여기서 S 는 아무 문장을

New AI systems collide with copyright law

By Suzanne Beane

Technology of Business reporter



Kelly McKenna found that her work had been used to train artificial intelligence
Illustration: Max O'Leary / The Verge

ARTIFICIAL INTELLIGENCE

The scary truth about AI copyright is nobody knows what will happen next

The last year has seen a boom in AI models that create art, music, and code by learning from others' work. But as these tools become more prominent, unanswered legal questions could shape the future of the field.

By James Westcott, a senior reporter who has covered AI, robotics, and more for eight years at The Verge

A collage of several AI-generated artworks. The pieces feature a variety of styles, including abstract patterns, portraits of people, and surreal scenes. The overall composition is a dense, colorful patchwork of different visual elements.

Peacock Article

The current legal cases against generative AI are just the beginning

All that can generate art, text and more is in for a reckoning

By Kyle Wiggers, senior writer | Published: January 26, 2023

A bronze statue of Lady Justice, also known as Justitia or Iustitia. She is blindfolded and holds a set of scales in one hand and a sword in the other, symbolizing the scales of justice and the power of law.

Tech 2023 - Technology

Generative AI is a legal minefield

Illustration: Max O'Leary / The Verge

A close-up illustration of a wooden gavel, which is a symbol of authority and justice. It is resting on a dark, textured surface. In the background, there is a stylized globe with blue and yellow continents, suggesting a global reach or impact.



다국어 코퍼스가 더욱 중요해짐

| 연도 | 업체별 평균 데이터 구매 수 |
|------|-----------------|
| 2016 | 5천 스트링 |
| 2017 | 1십만 스트링 |
| 2018 | 5십만 스트링 |
| 2019 | 1백만 스트링 |
| 2020 | 3백만 스트링 |
| 2021 | 5백만 스트링 |
| 2022 | 8백만 스트링 |
| 2023 | 30백만 스트링 |

**정확한 다국어 병렬 코퍼스를 통해
학습량이 적은 언어를 유추**

국내 기업이 데이터를 사는 것에 적극적이지 않는 이유

Hoshi Sato on Star Trek



Fine-tuning in LLM

EQUIPMENT AND TECHNOLOGY FILE 63 CARD 6

Sato's role in Star Trek is crucial for perfecting Large Language Models

The hundreds of races encountered by the Federation inevitably use their own languages, and Starfleet's mission of interstellar exploration would be impractical without the use of the universal translator, an incredibly powerful but commonplace device that can provide an instant translation between any two languages.

Universal Translator

In the 2290s, Kirk and his crew use one bulky and unreliable model of the universal translator. By the 2370s, the device has been miniaturized and made more reliable. It can now translate between over 1000 different languages simultaneously, using a neural network trained on millions of examples. The device is still not perfect, though, and it often makes mistakes, such as translating "Klingon" into "Klingon".

Adaptive Fine-tuning

MLM on unlabelled data → Model specialised to the target data

Fine-tuning

In-distribution: classification, sequence labeling, Q&A, ...

Cross-entropy on task labels

Pre-training

ELMO, ULMFiT, GPT, BERT

MLM on relevant data

Reusable device

Flitto

```

graph LR
    A[Pre-training] --> B[MLM on unlabelled data]
    B --> C[Adaptive Fine-tuning]
    C --> D[MLM on relevant data]
    D --> E[Fine-tuning]
    E --> F[In-distribution: classification, sequence labeling, Q&A, ...]
    E --> G[Cross-entropy on task labels]
    G --> H[Cross-entropy on task labels]
    H --> I[MLM on relevant data]
    I --> C
  
```

Vision, Enriching the LLMs

Los Angeles Times

That night, Wimmer was so exhausted that most of the other's actions had a short-term effect on her mood. And, if someone asked, she'd say, "I'm not tired." In fact, she was so tired that when Weller's "poetic objects" were displayed in the exhibition, she was in such a state of exhaustion that she didn't notice the two installations that were created for her.

fors

Mash

Light

Arts - Los Angeles Times

Waka Waka

RONAN BOUROLLEC

GOODMOODS nourrit les humeurs créatives

Flitto

po modular - teenage engineering

LLMs, The Next Big Thing

LLMs can be applied anywhere: multimedia, smart home, exhibitions, etc.



Flitto

Case of Traditional Painting



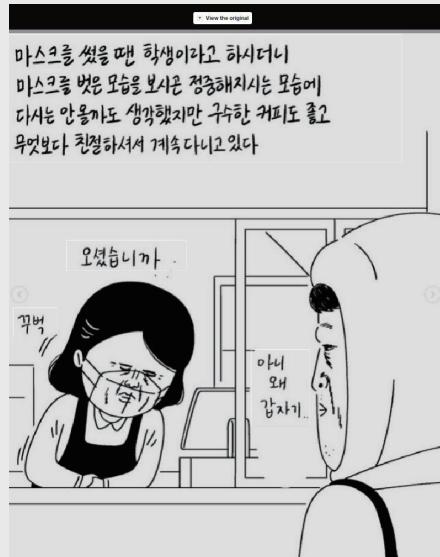
Case of Traditional Painting



Case of Traditional Painting



Case of Webtoon



Flitto

Case of Webtoon

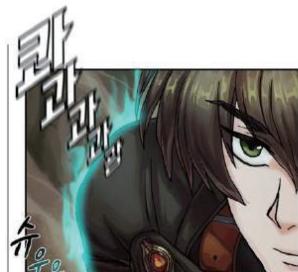


Flitto

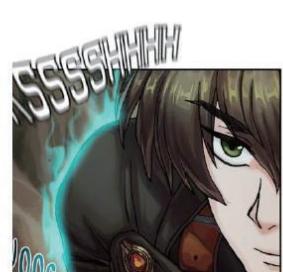
Case of Webtoon



Translation



Coloring Work



Flitto

Flitto

Thank You !

Simon Lee
CEO of Flitto
simon@flitto.com

2023년 가을 한국코퍼스언어학회 10주년 기념 학술대회

한국의 말뭉치 구축 및 연구의 지평

Session 3

<인공지능의 한국어 능력 평가>

AI말평의 현황 및 발전 방안

감정 분석 과제의 특징과 의의

이야기 완성 과제의 특징과 의의

표 기반 유사 말뭉치를 이용한 AI 언어능력 평가



국립국어원 인공지능(AI) 말평 현황 및 발전 방향

발표자: 이보라미 | 국립국어원 언어정보과 학예연구관



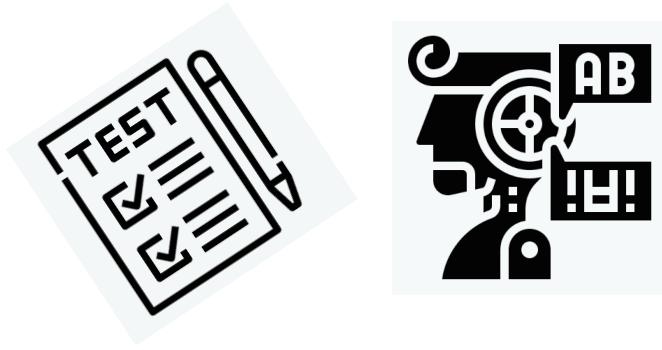
목 차

Contents

1. 추진 배경
2. 국내외 평가체계 현황
3. 인공지능(AI) 말평 주요 기능
4. 인공지능(AI) 말평 운영 과제
5. 발전 방향



1. 추진 배경



- 인공지능의 한국어 능력을 종합적으로 평가하는 체계(벤치마크) 부족
- 인공지능 언어모델의 한국어 능력에 대한 공신력 있는 다양한 평가 필요



2. 국내외 평가체계 현황

- GLUE (General Language Understanding Evaluation)

- 중국어, 한국어 등 언어별 벤치마크 개발에 영향을 준 대표적인 평가체계
- 3가지 대분류
과제, 총 9개의 하위 과제 및 관련 데이터 세트 구성



2. 국내외 평가체계 현황

■ GLUE 순위표(리더보드)

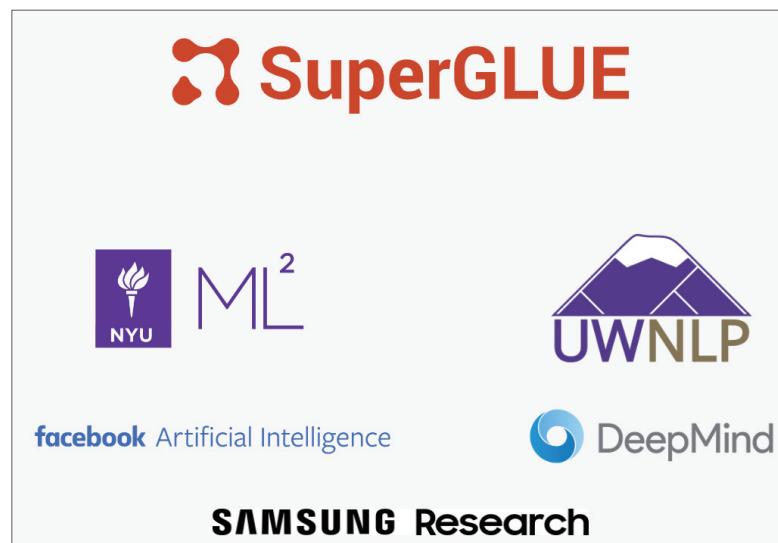
| Rank | Name | Model | URL | Score | CoLA | SST-2 | MRPC | STS-B | QQP | MNLI-m | MNLI-mm |
|------|----------------------------|---------------------|------------------------|-------|------|-------|-----------|-----------|-----------|--------|---------|
| 1 | Microsoft Alexander v-team | Turing ULR v6 | [link] | 91.3 | 73.3 | 97.5 | 94.2/92.3 | 93.5/93.1 | 76.4/90.9 | 92.5 | 92.1 |
| 2 | JDExplore d-team | Vega v1 | | 91.3 | 73.8 | 97.9 | 94.5/92.6 | 93.5/93.1 | 76.7/91.1 | 92.1 | 91.9 |
| 3 | Microsoft Alexander v-team | Turing NLR v5 | [link] | 91.2 | 72.6 | 97.6 | 93.8/91.7 | 93.7/93.3 | 76.4/91.1 | 92.6 | 92.4 |
| 4 | DIRL Team | DeBERTa + CLEVER | | 91.1 | 74.7 | 97.6 | 93.3/91.1 | 93.4/93.1 | 76.5/91.0 | 92.1 | 91.8 |
| 5 | ERNIE Team - Baidu | ERNIE | [link] | 91.1 | 75.5 | 97.8 | 93.9/91.8 | 93.0/92.6 | 75.2/90.9 | 92.3 | 91.7 |
| 6 | AliceMind & DIRL | StructBERT + CLEVER | [link] | 91.0 | 75.3 | 97.7 | 93.9/91.9 | 93.5/93.1 | 75.6/90.8 | 91.7 | 91.5 |



2. 국내외 평가체계 현황

■ SuperGLUE

- GLUE의 과제 난도를 높이거나 여러 과제를 종합하여 인공지능의 언어 능력에 대한 보다 심층적인 평가를 목표로 함.



2. 국내외 평가체계 현황

SuperGLUE Tasks

| Name | Identifier | Download | More Info | Metric |
|--|------------|----------|-----------|--------------------------|
| Broadcoverage Diagnostics | AX-b | | | Matthew's Corr |
| CommitmentBank | CB | | | Avg. F1 / Accuracy |
| Choice of Plausible Alternatives | COPA | | | Accuracy |
| Multi-Sentence Reading Comprehension | MultiRC | | | F1a / EM |
| Recognizing Textual Entailment | RTE | | | Accuracy |
| Words in Context | WiC | | | Accuracy |
| The Winograd Schema Challenge | WSC | | | Accuracy |
| BoolQ | BoolQ | | | Accuracy |
| Reading Comprehension with Commonsense Reasoning | ReCoRD | | | F1 / Accuracy |
| Winogender Schema Diagnostics | AX-g | | | Gender Parity / Accuracy |

문화체육관광부
국립국어원

2. 국내외 평가체계 현황

■ CLUE

Multi task challenge

• 중국어 자연언어 이해 벤치마크

Chinese Language Understanding Evaluation Benchmark(CLUE)

What is CLUE: Chinese Language Understanding Evaluation benchmark. We provide NLU datasets, pre-trained (baseline) models, toolkits and leaderboards.

Vision: to provide a benchmark for Chinese language models, to serve the Chinese NLU and NLP community

문화체육관광부
국립국어원

2. 국내외 평가체계 현황

■ CUGE: CLUE의 확장판

CUGE Framework Leaderboard Dataset FAQ Paper Code Subm.

CUGE

Chinese Language Understanding and Generation Evaluation Benchmark

| Corpus | Train | Dev | Test | Task | Metrics | Domain |
|--|---------|-------|-------|----------------------------------|--------------------|---------------------|
| Language Understanding: Word-Sentence Level | | | | | | |
| PKU-SEG | 40.3k | 10.5k | 9.9k | word segmentation | F1 score | news |
| WordSeg-Weibo | 20.1k | 2.1k | 8.6k | word segmentation | F1 score | social media |
| PKU-SEGPOS | 31.7k | 5.2k | 4.8k | word segmentation and POS | F1 score | news |
| CCPM | 21.8k | 2.7k | 2.7k | classical poetry matching | F1 score | Chinese poetry |
| CMEE | 15.0k | 5.0k | 3.0k | named entity recognition | F1 score | medical |
| FinRE | 7.5k | 1.5k | 3.7k | relation extraction | F1 score | financial |
| YACLC | 8.0k | 1.0k | 1.0k | grammatical error correction | F1 score | essays |
| Language Understanding: Discourse Level | | | | | | |
| SPR | 12.7k | 1.6k | 4.8k | humor detection | F1 score | TV show |
| ClozeT | 0.6k | 0.3k | 0.3k | story cloze test | accuracy | literature |
| C ³ | 11.9k | 3.8k | 3.8k | reading comprehension | accuracy | mixed-genre |
| GCRG | 7.0k | 0.8k | 0.8k | reading comprehension | accuracy, F1 score | mixed-genre |
| Information Acquisition and Question Answering Capability | | | | | | |
| WantWords | 78.0k | 19.0k | 19.0k | reverse dictionary | accuracy@1 | mixed-genre |
| KBQA | 24.0k | - | 0.6k | open-domain question answering | EM | mixed-genre |
| Sogou-Log | 8,052k | 500k | 1.0k | document retrieval | nDCG@k | mixed-genre |
| Language Generation Capability | | | | | | |
| LCSTS | 2,401k | 9.0k | 0.7k | text summarization | Rouge | news |
| CEPSUM | 434k | 5.0k | 5.0k | text summarization | Rouge | e-commerce |
| E-Reviews | 115k | 1.0k | 3.0k | data-to-text generation | BLEU, Distinct | e-commerce |
| Conversational Interaction Capability | | | | | | |
| KdConv | 62.9k | 9.0k | 9.1k | knowledge-driven conversation | BLEU, Distinct | film, music, travel |
| Multilingual Capability | | | | | | |
| WMT20-EnZh | 21,000k | 4.0k | 4.0k | machine translation | BLEU | mixed-genre |
| NCI-S-EnZh | 365k | 3.0k | 3.0k | cross-lingual text summarization | Rouge | mixed-genre |
| Mathematical Reasoning Capability | | | | | | |
| Math23k | 21.0k | 1.0k | 1.0k | mathematical computation | accuracy | math word problem |

2. 국내외 평가체계 현황

■ CUGE 순위표

| Rank | Model | Org | Code Paper | Time | NLU-WSL | NLU-DL | IA&QA |
|----------|-----------|-----------|--------------|-----------|-------------|-------------|-------------|
| Baseline | mT5-small | CUGE Team | 🔗 | 2021-8-20 | 87.70 (100) | 41.50 (100) | 29.20 (100) |
| 1 | Yuan 1.0 | 浪潮人工智能研究院 | 🔗 | 2022-5-25 | 89.71 (102) | 86.90 (209) | 28.12 (96) |
| 2 | CPM-2 | CUGE Team | 🔗 | 2021-8-20 | 91.60 (104) | 86.10 (207) | 35.90 (123) |

The radar chart illustrates Yuan 1.0's performance across various NLP tasks. The axes represent different tasks: MR (top), CI (top-left), NLG (bottom-left), IA&QA (bottom), NLU-WSL (bottom-right), NLU-DL (right), and ML (top-right). The chart shows a strong performance in NLU-WSL and NLU-DL, with moderate performance in CI and NLG, and lower performance in MR, IA&QA, and ML.

| NLU-WSL | NLU-DL |
|-------------|--------|
| CI | ML |
| 古诗文识记 | |
| 89.71 (102) | |
| CCPM | |
| 89.71 (102) | |

문화체육관광부
국립국어원

2. 국내외 평가체계 현황

■ KLUE

- 한국어 자연언어 이해 평가체계
- 8개의 하위 과제로 구성

KLUE

Korean Language Understanding Evaluation

Korean Language Understanding Evaluation (KLUE) benchmark is a series of datasets to evaluate natural language understanding capability of Korean language models. KLUE consists of 8 diverse and representative tasks, which are accessible to anyone without any restrictions. With ethical considerations in mind, we deliberately design annotation guidelines to obtain unambiguous annotations for all datasets. Furthermore, we build an evaluation system and carefully choose evaluations metrics for every task, thus establishing fair comparison across Korean language models.



문화체육관광부
국립국어원

2. 국내외 평가체계 현황

■ KLUE 과제

| Name | Type | Format | Eval. Metric |
|----------------|-------------------------------|---|--|
| KLUE-TC (YNAT) | Topic Classification | Single Sentence Classification | Macro F1 |
| KLUE-STS | Semantic Textual Similarity | Sentence Pair Regression | Pearson's r , F1 |
| KLUE-NLI | Natural Language Inference | Sentence Pair Classification | Accuracy |
| KLUE-NER | Named Entity Recognition | Sequence Tagging | Entity-level Macro F1 Character-level Macro F1 |
| KLUE-RE | Relation Extraction | Single Sentence Classification (+2 Entity Spans) | Micro F1 (without no_relation), AUPRC |
| KLUE-DP | Dependency Parsing | Sequence Tagging (+ POS Tags) | Unlabeled Attachment Score, Labeled Attachment Score |
| KLUE-MRC | Machine Reading Comprehension | Span Prediction | Exact Match, ROUGE-W (LCCS-based F1) |
| KLUE-DST (WoS) | Dialogue State Tracking | Slot-Value Prediction | Joint Goal Accuracy Slot Micro F1 |

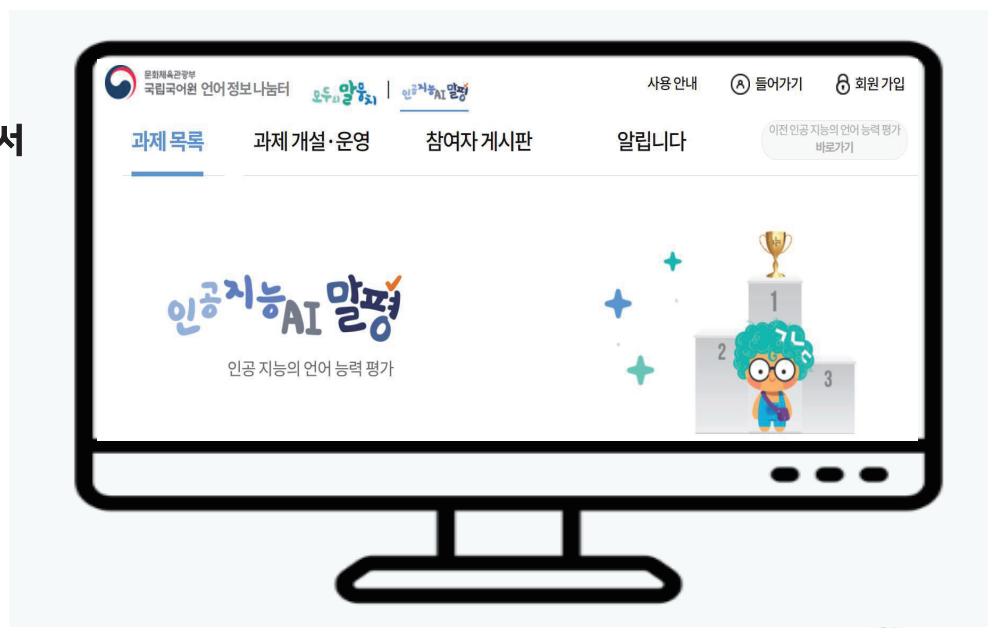


이미지: 한국경제 기사
<https://www.hankyung.com/article/202110135296>

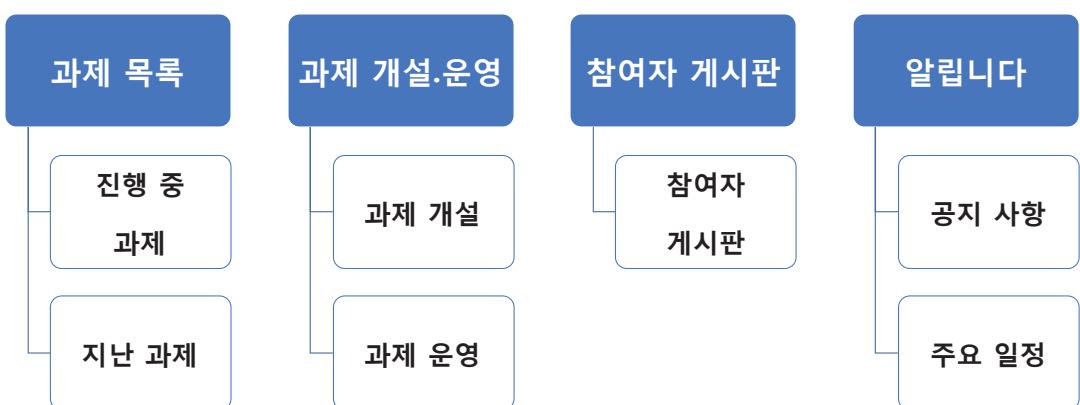
3. 인공지능(AI) 말평 주요 기능

■ 국립국어원

'언어정보나눔터'에서
<모두의 말뭉치>와
연계하여 운영



3. 인공지능(AI) 말평 주요 기능



3. 인공지능(AI) 말평 주요 기능

과제 참여 기능

회원이면 누구나 참여

국립국어원에서 개설하거나,
 다른 회원이 개설한 과제에 참여



과제 개설·운영 기능

평가체계 운영 플랫폼의 역할

시스템 가입 회원이면 누구나 개설

운영자를 구성·지정하여 과제별 관리

'수동 평가'와 '자동 평가'를 선택하여
개설



3. 인공지능(AI) 말평 주요 기능

■ 과제 개설자: 과제 게시판 관리를 할 수 있으며 과제 수정, 과제 제출, 과제 삭제 등 전반적인 과제 운영 및 관리 권한 가짐.

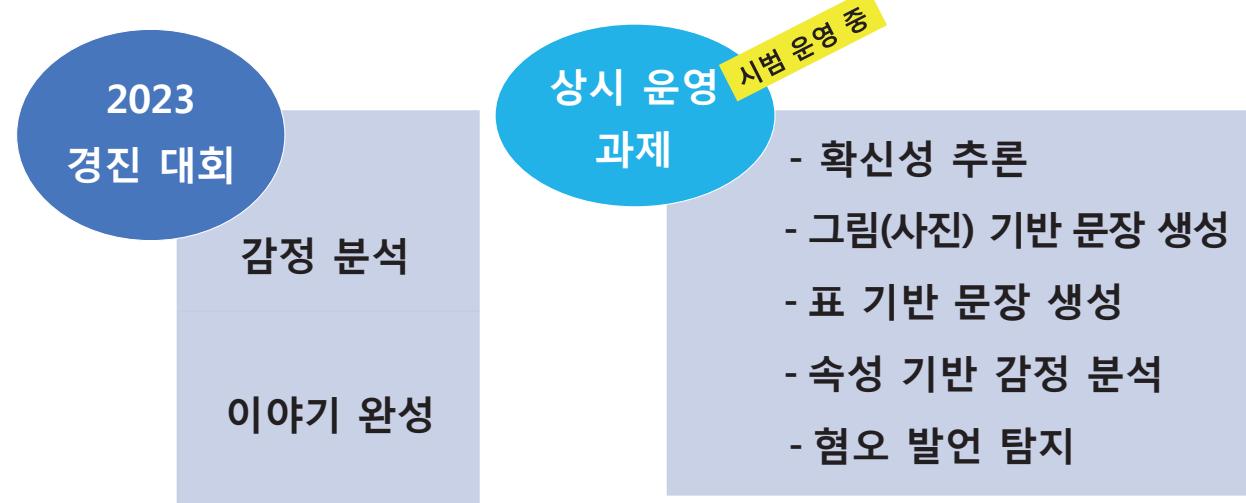
■ 과제 개설 후 과제에 사용될
데이터로 국어원 말뭉치 활용 가능

■ 자동 평가 기능 활용 가능

| | |
|---------------|---|
| *평가 방법 | <input type="radio"/> 수동 평가 <input checked="" type="radio"/> 자동 평가 자동평가 지표별 정답지/답안지 샘플 바로가기 |
| 평가 및 운영 관리 기능 | |
| *평가 지표 | <input checked="" type="radio"/> 평균 점수 <input type="radio"/> 지표 선택 소과제명1: <input type="text"/> 지표명1: 선택하세요 선택하세요 BLEU ROUGE-1 (-)MSE weighted_F1 macro_F1 micro_F1 |
| 운영자 구성 정보 | 자동 평가지표 등록 화면 설명: <input type="text"/> 설명: <input type="text"/> |



4. 인공지능(AI) 말평 운영 과제



4. 인공지능(AI) 말평 운영 과제

■ 2023 경진대회: 감정 분석 과제

- 주어진 텍스트에 대한 화자의 감정 상태를 파악
- 텍스트에 드러나는 8가지 감정 유형을 분류
- 과제 예시



| 항목 | 내용 |
|-----|---|
| 텍스트 | "아 뉴스레터에서 뮤지컬 킹아더 관람 신청받는데... 가고싶은데 약속이네 ○ 어으아앙악 짜증" |
| 대상 | '약속' |
| 감정 | "joy": "False", "anticipation": "False", "trust": "False", "surprise": "False", "disgust": "False", "fear": "False", "anger": "True", "sadness": "False" |



4. 인공지능(AI) 말평 운영 과제

■ 2023 경진대회: 이야기 완성 과제

- 주어진 문장들을 논리적으로 연결하는 문장을 생성하는 과제

| 입출력 | | 내용 |
|-----|------|----------------------------------|
| 입력 | 문장 1 | 나는 입사하고 나서 몇 달 동안은 조심스럽게 행동했다. |
| | 문장 3 | 그랬더니 동료들은 첫인상과 다른 나의 모습에 놀랐다. |
| 출력 | 문장 2 | 회사 분위기를 파악하고 나서는 눈치껏 내 성격을 드러냈다. |

<적절한 예시>

| 입출력 | | 내용 |
|-----|------|---------------------------------------|
| 입력 | 문장 1 | 나는 환경을 보호하기 위해 길을 걸으면서 보이는 쓰레기들을 주웠다. |
| | 문장 3 | 그러다 동생과 나는 같이 쓰레기를 줍는 동아리를 만들었다. |
| 출력 | 문장 2 | 내가 쓰레기를 줍는 모습을 동생은 본체만체했다. |

<부적절한 예시>



4. 인공지능(AI) 말평 운영 과제

■ 상시 과제(시범 운영 과제)

| 과제명 | 주요 특징 |
|-----------------|---|
| 확신성 추론 | <ul style="list-style-type: none">슈퍼글루(SuperGLUE)의 커미트먼트뱅크(CommitmentBank)를 참조화자의 인식에 비추어 맥락과 제시 문장 간 함의 관계를 정도성으로 판단데이터 세트는 화자가 확신하는 정도를 1~7사이의 실수로 표현하는 방식으로 구성 |
| 그림(사진) 기반 문장 생성 | <ul style="list-style-type: none">그림(사진)이 주어졌을 때 이를 설명하는 문장 하나를 생성하는 과제데이터 세트는 그림(사진) 하나와 이를 설명하는 문장 5개로 구성 |
| 표 기반 문장 생성 | <ul style="list-style-type: none">자료로부터 텍스트를 생성하는 과제 중 하나로, 주어진 표의 특정 부분을 설명하는 문장을 생성하는 과제(구글의 ToTTo 데이터 세트를 참조)표에 음영으로 표시한 부분을 설명하는 문장 하나를 생성하는 것 |



4. 인공지능(AI) 말평 운영 과제

■ 상시 과제(시범 운영 과제)

| 과제명 | 주요 특징 |
|-------------|--|
| 속성 기반 감성 분석 | <ul style="list-style-type: none">국어원에서 구축한 속성 기반 감성 분석 말뭉치를 활용하여, 해당 데이터 세트 중 "제품" 도메인을 대상으로 개발속성 범주 탐지(Aspect Category Detection, ACD)와 속성 감성 분류(Aspect Sentiment Classification, ASC)로 구성 |
| 혐오 발언 탐지 | <ul style="list-style-type: none">국어원의 '비윤리적 표현 말뭉치 연구 분석 및 시범 구축' 사업을 통해 구축한 결과물 중에서 '비난'과 '혐오' 범주를 활용하여 데이터 구축'혐오 발언 문장'에 해당하는지 혹은 '혐오 발언이 아닌 문장'에 해당하는지를 분류하는 이진 분류 문제(Binary Classification) 과제 |



5. 발전 방향

■ 국립국어원의 국어 전문성을 바탕으로 한 공신력 있는 평가 추진

- 자동 평가 방식을 기본으로 하되, 정성 평가를 위해 인간 평가 등을 추가

■ 과제 확대를 위한 말뭉치 자원 확충

- 인공지능의 이해 능력, 생성 능력을 종합적으로 평가할 수 있는 말뭉치 지속 구축

예시) AI의 추론 능력 평가 데이터(초등학교에서 대학교로 수준별 단계적 구축)



한국코퍼스언어학회 10주년 기념 학술대회

감정 분석 과제의 특징과 의의



김한샘

YONSEI UNIVERSITY



延世大學校

감정 분석의 개념

- ◆ 감정 분석은 텍스트에 표현된 인간의 감정을 포착하여 기쁨, 슬픔 등 다양한 감정 유형으로 분류하는 것을 의미함
- ◆ 일상 생활에서 빈번하게 발생하는 감정을 정확히 포착하고 주석하는 것은 감정에 대한 인간의 이해를 향상시킬 수 있을 뿐만 아니라 더 나아가 기계 학습 데이터 세트의 개발로 연결될 수 있음
- ◆ 사람이 느낄 수 있는 감정은 다양하고, 복합적임. 특히 감정을 한 단어로 명확히 표현할 수 없고, 사람마다 다르게 지각된다는 점에서 감정 주석은 모호하고 주관성이 높은 자연어 처리 태스크라고 할 수 있음

감성 분석과 감정 분석의 개념

◆ 감성 분석(Sentiment Analysis)

- 감성 분석은 '오피니언 마이닝'이라고도 하며, 주어진 텍스트에서 '긍정', '부정' 또는 '중립'을 파악하여 텍스트의 극성을 분석하는 것
- 제품 평가, 브랜드 평판 모니터링 등 다양한 응용 분야에서 사람들의 의견과 성향을 분석하는 데 사용됨

◆ 감정 분석(Emotion Analysis)

- 감정 분석은 텍스트에 나타나는 인간의 감정을 '기쁨', '슬픔', '분노' 등으로 분석하는 것
- 소셜 미디어 감정 분석, 고객 서비스 향상, 정서적 지능 개발 등에 사용됨

◆ 국외에서는 Sentiment와 Emotion의 개념이 명확하게 구분되어 사용되는 반면 국내에서는 두 개념이 혼재되어 있음

-3-

국내 감정 분석 말뭉치 구축 현황

| 데이터세트 | 레지스터 | 규모 | 감정 레이블 | 형태 |
|---|--------------------|---------|--|----|
| KTEA (H, Shin et al., 2012) | 트위터, TV 댓글 | 5,706 | anger, disgust, fear, joy, sadness, surprise | 단일 |
| 한재호 외(2018) | 뉴스 댓글 | 30,000 | 기쁨, 슬픔, 분노, 혐오 | 단일 |
| 감성 대화 말뭉치 (NIA) | 대화 데이터 | 270,000 | 6개의 기본 감정(분노, 슬픔, 불안, 상처, 당황, 기쁨)을 기준으로 각각 9개의 세부 감정이 포함된 총 60개 감정 상태 반영 | 단일 |
| KOTE (D Jeon et al., 2022) | 온라인 댓글 | 50,000 | 불평/불만, 환영/호의, 감동/감탄, 지긋지긋, 고마움, 슬픔, 화남/분노, 존경, 기대감, 우쭐댐/무시함, 안타까움/실망, 비장함, 의심/불신, 뿌듯함, 편안/쾌적, 신기함/관심, 아껴주는, 부끄러움, 칭찬/무서움, 절망, 한심함, 역겨움/징그러움, 짜증, 어이없음, 없음, 패배/자기혐오, 귀찮음, 힘듦/지침, 즐거움/신남, 깨달음, 죄책감, 증오/혐오, 흐뭇함(귀여움/예쁨), 당황/난처, 경악, 부담/안내킴, 서러움, 재미없음, 불쌍함/연민, 놀람, 행복, 불안/걱정, 기쁨, 안심/신뢰 | 멀티 |
| 공공분야 고객응대 데이터 (NIA, 10월 공개 예정) | 고객 응대 데이터 | 131,555 | 긍정, 부정, 중립, 공포, 놀람, 분노, 슬픔, 행복 | 단일 |
| 다양한 문화콘텐츠 스토리 데이터 (NIA, 10월 공개 예정) | 영화, 드라마, 소설, 만화 | 100,000 | 미확인 | 단일 |
| 감성 및 발화스타일 동시 고려 음성합성 데이터 (NIA, 10월 공개 예정) | 대본 데이터 | 80,272 | 기쁨, 슬픔, 분노, 무감정 | 단일 |

-4-

감정 분석 말뭉치를 활용한 태스크

■ 단일 라벨 감정 분류(Single-label Emotion Classification)

- ◆ 하나의 텍스트에서 단일 감정 라벨을 맞추는 태스크
- ◆ 단일 라벨 감정 분류 태스크는 감정 관련 태스크 중 가장 많은 분포를 보이는 유형
- ◆ 한국어 SNS 데이터에 대한 분류를 수행하거나(한재호 외, 2018), 다중 장르 말뭉치에 대한 분류를 수행함(Tafreshi&Diab, 2018)
- ◆ 감정 분류는 소셜 미디어에서 감정 분석, 고객 서비스 품질 평가, 제품 리뷰 분석, 감정 기반 마케팅 등 다양한 분야에서 유용한 정보를 추출하고 의사 결정에 활용할 수 있음.
- ◆ 평가 지표는 정확도, F1, 정밀도, 재현도가 일반적으로 사용됨

-5-

감정 분석 말뭉치를 활용한 태스크

■ 다중 라벨 감정 분류(Multi-label Emotion Classification)

- ◆ 이전까지의 단일 라벨의 한계를 극복하기 위해 제안됨
- ◆ 감정은 복잡하고 단일 라벨로만 정해지기 어렵기 때문에 다중 라벨 주석이 필요함
- ◆ SemEval-2018 Task1에서 10개의 서로 다른 감정을 가진 다중 라벨 감정 분류 태스크를 제시함(Mohammad et al., 2018)
- ◆ 최근 다중 라벨 감정 분류 태스크는 지속적으로 등장하고 있음
 - 국외: Schuff et al.(2017), Demszky et al.(2020) 등
 - 국내: Jeon, Duyoung & Lee, Junho & Kim, Cheongtag(2022)
- ◆ 평가 지표는 F1, 정밀도, 재현도가 일반적으로 사용됨

-6-

감정 분석 말뭉치를 활용한 태스크

■ 감정 식별(Emotion Detection/Emotion Recognition)

- ◆ 주어진 텍스트에서 어떤 감정이 존재하는지 식별하는 태스크
- ◆ 일반적으로 행복, 슬픔, 분노, 두려움, 놀라움과 같은 기본적인 감정에 초점을 맞추어 감정 여부를 식별함
- ◆ 감정 인식은 효과적인 인간과 컴퓨터의 상호 작용을 가능하게 하는 연구 분야로 음성, 얼굴 표정, 행동을 통해 인식할 수 있음
- ◆ 음성, 비언어적 행동 등 멀티 모달 분야에서 많이 활용됨
- ◆ 평가 지표는 정확도가 일반적으로 사용됨

-7-

감정 분석 말뭉치를 활용한 태스크

■ 감정-원인 쌍 추출(Emotion-Cause Pair Extraction)

- ◆ 텍스트에서 특정 감정을 유발하거나 환기시키는 단어나 구 또는 절을 추출하는 태스크(단일 라벨과 다중 라벨을 모두 포함함)
- ◆ 이전에는 감정 자극(Emotion Stimulus) 식별이라는 이름으로 단독 제시되었으나, 최근에는 단순 원인 혹은 감정 추출에서 더 나아가 Xia&Ding(2019)이 도입한 감정-원인 쌍을 식별하는 태스크로 난도가 높아짐
- ◆ 감정과 원인을 쌍으로 추출하는 과제에서는 감정과 감정을 불러 일으키는 원인을 식별하는 모델로 텍스트의 감정 식별을 더욱 정확하게 하고, 보다 이해 가능한 결과를 얻을 수 있는 장점이 있음
- ◆ 평가 지표는 F1, 정밀도, 재현도가 일반적으로 사용됨

-8-

감정 분석 말뭉치를 활용한 태스크

■ 감정 강도(Emotion Intensity)

- ◆ 주어진 텍스트의 감정 강도를 측정하고자 하는 태스크
- ◆ 주어진 감정의 정도나 강도를 나타내는 연속적이거나 이산적인 점수를 예측하는 것을 목표로 함
- ◆ 감정의 강도를 주석하는 것은 주석자 간 일치도를 유지하면서 구축하기 어려우므로 다른 태스크에 비해 연구가 적은 편임
- ◆ SemEval-2018 경진대회에 감정 강도 추정을 여러 태스크 중 하나로 포함시킴
- ◆ 평가 지표는 Pearson 상관 계수(PCC)가 일반적으로 사용됨

-9-

감정 분석 태스크의 정량적 지표

- 대부분의 감정 분석 연구에서는 정량적 지표로 분류 성능 지표로 많이 쓰이는 정확도, 정밀도, 재현도, f1-score를 제안함
- 감정 원인-쌍을 추출하는 경우에는 F1, 정밀도, 재현도를 사용
- 감정 강도를 측정할 경우에는 pearson 상관 계수(Mohammad & Bravo-Marquez, 2017; Mohammad et al., 2018)를 사용
- 감정 분석 말뭉치를 통한 문장 생성일 경우에는 perplexity나 BLEU score를 사용(Li, Y., Su et al., 2017)

-10-

감정 분석 벤치마크 및 리더보드

■ 벤치마크(benchmark)란?

- ◆ 자연어 처리 전반(이해, 생성) 혹은 특정 태스크에 대한 모델의 성능을 평가할 수 있는 평가 체계

■ 리더보드(leaderboard)란?

- ◆ 자연어 처리 벤치마크 리더보드는 모델의 평가 체계 수행 성능과 해당 모델의 이름을 기재

-11-

국내 사례

■ DACON 음성 감정 인식 AI 경진대회

주어진 음성에 대해 6가지 감정 인식 및 분류(angry, fear, sad, disgust, neutral, happy)

The screenshot shows the DACON competition page for the '월간 데이콘 음성 감정 인식 AI 경진대회'. The main banner features a woman singing into a microphone. Below the banner, there are tabs for 대회 안내, 데이터, 코드 공유, 토크, 리더보드 (which is underlined), and 제출. The 'PRIVATE' tab is selected in the navigation bar. A table lists the top 5 winners:

| # | 팀 | 팀 멤버 | 최종점수 | 제출수 | 등록일 |
|---|-------------|----------------|---------|-----|------|
| 1 | 파이썬초보만 | [Profile Icon] | 0.87851 | 5 | 4달 전 |
| 2 | Ldoun | [Profile Icon] | 0.82915 | 15 | 4달 전 |
| 3 | 중요한건꺾이지않는마음 | [Profile Icon] | 0.81852 | 8 | 4달 전 |
| 4 | hangjoo | [Profile Icon] | 0.81776 | 15 | 4달 전 |
| 5 | oriko | [Profile Icon] | 0.80789 | 19 | 4달 전 |

<https://dacon.io/competitions/official/236105/leaderboard>

-12-

2023년 국립국어원 인공지능 언어능력평가: 감정 분석 과제

- ◆ 주어진 텍스트와 텍스트 내 감정 대상(target)을 사용하여 8가지로 감정 분류
- ◆ 'joy(기쁨)', 'anticipation(기대)', 'trust(신뢰)', 'surprise(놀람)', 'disgust(혐오)', 'fear(공포)', 'anger(분노)', 'sadness(슬픔)'

2023년 국립국어원 인공 지능 언어 능력 평가: 감정 분석 과제

신청 접수 중 2023.08.21. ~ 2023.10.20.

과제 목록 보기

과제 개요

자료 받기

순위표(리더보드)

공지

문의

참가 신청

제출 관리

[말뭉치] 2023 인공 지능의 언어 능력 평가: EA

2023년 국립국어원 인공 지능 언어 능력 평가 AI 말평 감정 분석 과제의 말뭉치입니다.



내려받기

[베이스라인 모델] 2023 AI 말평: 감정 분석 베이스라인 모델

2023년 국립국어원 인공 지능 언어 능력 평가 AI 말평 감정 분석 과제의 베이스라인 모델입니다. 아래 링크를 참고하여 주세요. https://github.com/teddysum/Korean_Emotion_Analysis



내려받기

https://corpus.korean.go.kr/taskOrdtm/taskDownload.do?taskOrdtmId=103&clCd=ING_TASK&subMenuld=sub02

-13-

국립국어원 감정 분석 과제의 리더보드

■ 성능 평가 지표: F1-score(micro)

2023년 국립국어원 인공 지능 언어 능력 평가: 감정 분석 과제

신청 접수 중 2023.08.21. ~ 2023.10.20.

과제 목록 보기

과제 개요

자료 받기

순위표(리더보드)

공지

문의

참가 신청

제출 관리

전체 57건

※ 마지막 평가 일시: 2023년 10월 05일 00시 01분

| 순위 | 팀명 | 모델명 | 평가 점수 | 소과제명 | |
|----|---------------|-------------------------|------------|------------|-------------------|
| | | | | 감정 분석 | 제출 일시 |
| 1 | 사과는맛있어맛있으면비니나 | 효죽해 | 90.1289649 | 90.1289649 | 2023.09.28. 14:51 |
| 2 | 동물의왕국 | 산악회 카렌 퍼글 화이팅 | 88.9429389 | 88.9429389 | 2023.09.17. 22:53 |
| 3 | 이人民日报 | ensemble_ebrf | 88.9281536 | 88.9281536 | 2023.10.04. 09:07 |
| 4 | 마총랩 | 마국간 | 88.6922767 | 88.6922767 | 2023.09.22. 17:01 |
| 5 | 일파코 | base | 88.4375861 | 88.4375861 | 2023.09.28. 18:54 |
| 6 | HJ | electra | 88.3195932 | 88.3195932 | 2023.10.03. 03:25 |
| 7 | 시스템오류 | mini_baseline_v2.4 | 88.0922469 | 88.0922469 | 2023.09.16. 13:04 |
| 8 | 이불밸레이어 | 이불밸래 | 88.0544562 | 88.0544562 | 2023.09.25. 01:31 |
| 9 | DILAB | kc-removed-e98-4e^-5-f1 | 87.9478512 | 87.9478512 | 2023.09.03. 16:58 |
| 10 | IF | 8 | 87.7564128 | 87.7564128 | 2023.09.26. 22:23 |

-14-

2023년 국립국어원 AI 언어능력 평가: 감정 분석 과제

■ 1. 감정 분석 과제 말뭉치 소개

- ◆ '2022 국립국어원 감정분석 말뭉치'를 바탕으로 **경진대회 감정 분석 과제 말뭉치** 구축
 - 저작권 이용 허락 계약이 완료된 트위터 계정의 트위터 자료 대상
 - 트위터 자료는 5어절 이상의 문화 콘텐츠 관련 자료만 사용
- ◆ 트위터 자료에 대해 '**감정(emotion)**'과 **감정 대상(target)** 주석
 - **감정**: 'joy(기쁨)', 'anticipation(기대)', 'trust(신뢰)', 'surprise(놀람)', 'disgust(혐오)', 'fear(공포)', 'anger(분노)', 'sadness(슬픔)' 8가지
 - **감정 대상**: 텍스트 내에서 화자(작성자)가 느끼는 감정의 대상이 되는 표현

-15-

2023년 국립국어원 AI 언어능력 평가: 감정 분석 과제

■ 2. 경진대회 감정 분석 과제 소개

◆ 과제 개요

- 주어진 텍스트에 대한 화자의 **감정 상태를 분류**하는 과제
- 국립국어원이 제공하는 '감정 분석 말뭉치'를 활용하여 제시된 텍스트 내 **특정 대상(target)**에 대한 화자의 **감정 상태를 파악**
- 'joy(기쁨)', 'anticipation(기대)', 'trust(신뢰)', 'surprise(놀람)', 'disgust(혐오)', 'fear(공포)', 'anger(분노)', 'sadness(슬픔)'의 8가지 감정으로 분류

| 항목 | 내용 |
|----|--|
| 발화 | "아 뉴스레터에서 뮤지컬 킹아더 관람 신청받는데... 가고싶은데 약속 이네 ○ 어으아앙악 짜증" |
| 대상 | ' 약속 ' |
| 감정 | "joy": "False", "anticipation": "False", "trust": "False", "surprise": "False", "disgust": "False", "fear": "False", "anger": "True", "sadness": "False" |

-16-

2023년 국립국어원 AI 언어능력 평가: 감정 분석 과제

■ 2. 경진대회 감정 분석 과제 소개

◆ 과제 수행 및 평가 방법

- 기본적으로 **다중 분류(multi-class classification)**를 수행, 특정 레이블에 해당하는 감정이 드러날 경우 'True', 아닐 경우 'False' 값으로 표시
- 한 텍스트 내에 여러 개의 감정이 드러날 수도 있으므로 하나의 텍스트가 여러 개 레이블에 대해 'True' 값을 가지는 다중 레이블 분류(Multi-label classification)로 수행할 수도 있음
- 성능 평가: 각 레이블의 'True/False' 값에 대한 '**F1-score(micro)**'를 기준으로 순위를 결정

- Precision = $TP / (TP + FP)$
- Recall = $TP / (TP + FN)$
- F1 Score = $2 * (Precision * Recall) / (Precision + Recall)$

- True Positives (TP): 모델이 'True'로 예측했고, 실제 측정값도 'True'인 경우.
즉, 모델이 해당 레이블을 제대로 예측한 경우
- False Positives (FP): 모델이 'True'로 예측했지만, 실제 측정값은 'False'인 경우.
즉, 모델이 레이블을 잘못 예측한 경우
- False Negatives (FN): 모델이 'False'로 예측했지만, 실제 측정값은 'True'인 경우.
즉, 모델이 레이블을 놓친 경우
- True Negatives (TN): 모델이 'False'로 예측했고, 실제 측정값도 'False'인 경우.
즉, 해당 레이블이 없다고 예측하였고, 실제로도 그 레이블이 없는 경우

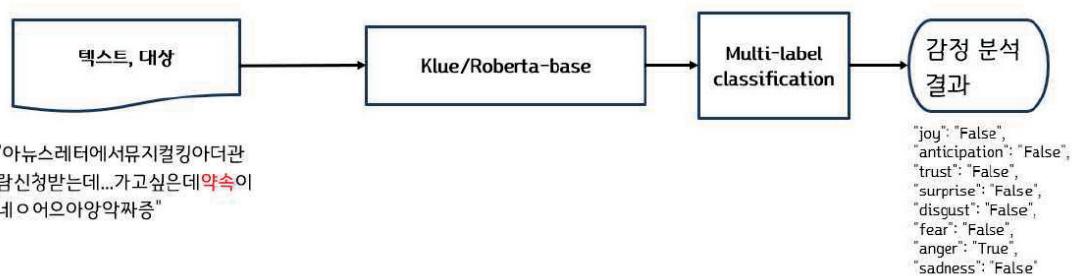
-17-

2023년 국립국어원 AI 언어능력 평가: 감정 분석 과제

■ 2. 경진대회 감정 분석 과제 소개

◆ 과제|| baseline

| Model | F1-score(micro) |
|-------------------|-----------------|
| Klue/Roberta-base | 0.850 |



-18-

2023년 국립국어원 AI 언어능력 평가: 감정 분석 과제

■ 3. 감정 분석 과제 말뭉치 구축 과정

◆ (1) 감정 분석 과제 말뭉치 구축 과정 요약

대상 말뭉치 인수 → 말뭉치 사전 검수 → 검수 요소 확정 →
검수 요소에 따른 말뭉치 검수 → 과제 말뭉치 구축 완료

- ① 대상 말뭉치: 2022 국립국어원 감정 분석 말뭉치
- ② 말뭉치 사전 검수: 대상 말뭉치의 주석 요소 및 과제 개발 방향을 고려하여 검수팀 내부에서 말뭉치 preview
- ③ 검수 요소 확정: '감정' 및 '대상', '비식별화 여부'의 주석 상태 검수
- ④ 말뭉치 검수: 감정 라벨 교정 및 기주석 대상 수정

-19-

2023년 국립국어원 AI 언어능력 평가: 감정 분석 과제

■ 3. 감정 분석 과제 말뭉치 구축 과정

◆ (2) 감정 분석 과제 말뭉치 검수의 실제

- 2022 감정분석 말뭉치의 주요 주석 요소는 **감정**, **감정 대상**, **비식별화 여부**이다. 이에 따라 검수 시에는 해당 요소들을 중점으로 하여 말뭉치를 검수한다.

- **감정 (emotion)**
 - 문서(sentence_form)가 내포하는 정서를 파악하여 8개 감정 유형으로 분류
 - 8개 유형에 포함되지 않는 경우 '기타'로 주석
- **감정 대상 (emotion target)**
 - 분석 대상 문서 내에서 화자(작성자)가 느끼는 감정의 대상이 되는 표현.
- **비식별화 여부**
 - 감정 대상(target)이 특정 인물 또는 집단일 경우
 - 개인정보 등의 노출 등을 막기 위해 비식별화 처리를 진행

-20-

2023년 국립국어원 AI 언어능력 평가: 감정 분석 과제

■ 3. 감정 분석 과제 말뭉치 구축 과정

◆ (2) 감정 분석 과제 말뭉치 검수의 실제

1) 감정 라벨 주석 검수

- 감정 분석 말뭉치는 **8개의 기본 감정 분석 라벨**(기쁨, 기대, 신뢰, 놀람, 혐오, 공포, 분노, 슬픔)과 '**기타**' 라벨 총 9개 라벨로 주석
- 감정 표현의 복합성, 연속성 등을 반영하기 위해 **복수 감정 주석** 역시 허용
- 다만 2022 감정 주석 말뭉치의 경우 주석자 10명의 Fleiss Kappa 일치도 계산에 따라 감정 라벨을 주석한 것이므로 감정 라벨 검수 시 **최대한 기존 라벨을 존중** 한다. 만일 변경이 필요할 경우 **검수자 간 토론**을 통해 변경에 대한 '**전원의 동의**'를 얻었을 때 변경

| 감정 분석 표지 | 동일 스펙트럼 감정 | 설명 |
|------------------|----------------------|--|
| 기쁨(joy) | 평온, 황홀, 감사, 고마움 등 | 어떤 만족감에 의해 느끼는 즐겁고 흥겨운 감정 |
| 기대(Anticipation) | 경계, 관심 | 앞으로 있을 일이나 상황을 미리 짐작함. 또는 그런 내용 |
| 신뢰(Trust) | 수용, 감탄/존경 | 굳게 믿고 의지함 |
| 놀람(Surprise) | 놀라움, 부주의/방심 | 어떤 일이 뜻밖이거나 출렁하거나 무서워서 신기해하거나 흥분하여 가슴이 뛰는 느낌 |
| 혐오(Disgust) | 지루함, 증오 | 싫어하고 미워함 |
| 공포(Fear) | 불안, 두려움 | 두렵고 무서움 |
| 분노(Anger) | 짜증, 격노 | 분개하여 크게 화를 냄 |
| 슬픔(Sadness) | 수심, 비탄 | 마음이 아프거나 괴로운 느낌 |
| 기타 | - | 8가지 라벨 외 기타 감정 (선후의 감정 정보 없는 궁금함, 의아함, 묘함 등...) |

2023년 국립국어원 AI 언어능력 평가: 감정 분석 과제

■ 3. 감정 분석 과제 말뭉치 구축 과정

◆ (2) 감정 분석 과제 말뭉치 검수의 실제

1) 감정 라벨 주석 검수: 세부 내용

- (1) 문서에 드러난 감정이 알맞게 주석되었는지 확인한다. 검수 시 알맞지 않은 경우 적절한 감정 라벨로 수정한다.
- (2) 문서에 주석된 감정 라벨이 하나일 경우, 문서에 드러난 실제 감정에 비해 **과소하게 주석되었는지 확인**한다. 과소할 경우 추가로 주석한다.
- (3) 문서에 주석된 감정 라벨이 **복수 개일** 경우, 문서에 드러난 실제 감정에 비해 **과도하게 주석되었는지 확인**한다. 과도할 경우 불필요한 감정 라벨을 삭제한다.
- (4) 문서 내 '**기대감**'을 나타내는 표현이 등장할 경우 **기본적으로 'anticipation'**을 주석하되, 다른 감정이 느껴질 경우 해당 감정도 주석한다.

2023년 국립국어원 AI 언어능력 평가: 감정 분석 과제

■ 3. 감정 분석 과제 말뭉치 구축 과정

◆ (2) 감정 분석 과제 말뭉치 검수의 실제

2) 감정 대상(target) 검수

- 감정 대상은 문서 내에서 화자/작성자가 생각하는 감정 대상이 되는 표현
- 구축 가이드라인에서는 아래 7가지 규칙을 적용하여 감정 대상을 주석

① 대상은 문서 내 고유 명사, 구체 명사, 추상 명사를 우선으로 주석한다. 이때 명사를 수식하는 성분은 대상 주석 시 제외한다.
② 문장 내 명확한 구체/추상 명사가 없을 경우, 수식 성분+의존 명사(-것, -거, -게, -분 등)를 대상으로 주석한다. ①과 달리 수식 성분까지 대상으로 주석한다.
③ 문장 내 명확한 구체/추상 명사가 없고 수식성분+의존 명사도 없을 경우, 대명사(이거, 저거, 그거 등)를 대상으로 주석한다.
④ 만일 동일한 지시 대상을 갖는 구체/추상 명사, 수식 성분+의존 명사, 대명사가 문서 내에 함께 등장할 경우, 고유/구체/추상 명사 > 수식 성분 + 의존 명사 > 대명사 순으로 우선 순위를 적용하여 주석한다.
⑤ 같은 대상을 가리키는 동일 수준의 표현 혹은 똑같은 표현이 복수로 등장하는 경우 문서 내에서 가장 먼저 등장하는 표현을 대상(target)으로 선정한다.
⑥ 여러 개 어절로 구성된 고유 명사의 경우 전체 표현을 대상(target)으로 선정한다. 단, 해당 표현이 하나가 아닌 여러 개의 발화(utterance)에 걸쳐 등장한다면 해당 고유 명사 일부가 포함된 마지막 발화(utterance) 내 표현만 대상(target)으로 선정한다.
⑦ 감정에 대한 대상이 문서 내에 없는 경우 대상(target)은 'None'으로 처리한다.

-23-

2023년 국립국어원 AI 언어능력 평가: 감정 분석 과제

■ 3. 감정 분석 과제 말뭉치 구축 과정

◆ (2) 감정 분석 과제 말뭉치 검수의 실제

2) 감정 대상(target) 검수: 세부 내용

- (1) 대상이 구축 가이드라인 7개 규칙에 맞게 잘 주석되었는지 확인하면서 검수를 진행한다
- (2) 감정 대상 검수 시 최대한 기준 주석을 존중하되, 스펜(span)에 한하여 지나치게 길거나 지나치게 짧은지 확인한다.
 - i) 대상 검수 시에도 기준 주석을 최대한 존중한다.
따라서 타겟에 헤어(head)로 보이는 단어가 포함되어 있을 경우 수정할 필요가 없다.
 - ii) 대상의 스펜이 과소하게 주석된 경우 주석된 대상을 수정한다.
 - iii) 대상의 스펜이 과도하게 주석된 경우 주석된 대상을 수정한다.
 - iv) 숫자가 타겟처럼 보이는 경우 주석하지 않는다.

- (3) target이 주석되어 있을 때에는 해당 target이 적절하게 주석되었는지,
target이 **None**인 경우 대상으로 주석할 수 있는 표현이 문서 내 있는지 다시 한번 살핀다

-24-

2023년 국립국어원 AI 언어능력 평가: 감정 분석 과제

■ 3. 감정 분석 과제 말뭉치 구축 과정

◆ (2) 감정 분석 과제 말뭉치 검수의 실제

3) 비식별화 검수

- 문서 내에 개인정보 노출 위험성이 있는 표현이 등장할 경우, 이를 비식별화 처리 한다.

| 비식별화 유형 | 비식별화 표지 | 설명 |
|----------------------|-----------------------|---|
| 이름 | &name& | 개인 실명 - 단, 정치인, 연예인 등의 공인 & 유명인은 제외 - 또한 실존 인물이 아닌 캐릭터/작중 인물명 등도 제외 |
| 온라인 계정 (아이디) | &account& | 트위터 등 특정 사이트의 온라인 계정 |
| 고유 식별 번호 (주민등록번호) | &social-security-num& | 개인 주민등록번호 |
| 전화번호 | &tel-num& | 휴대폰, 사업장 번호 등 |
| 카드 번호 | &card-num& | 신용카드 번호 등 |
| 기타 번호 | &num& | 비밀번호 등 기타 비식별화 대상 번호 해당 번호를 알았을 때 개인이 특정될 수 있는 경우 비식별화 |
| 주소 | &address& | 읍/면/동 이하 상세 주소 |
| 출신 및 소속 | &affiliation& | 개인의 출신 및 소속 |
| 기타 비식별화 필요 항목 | &others& | 위 항목 외 개인을 특정/식별할 수 있는 모든 항목 |

-25-

2023년 국립국어원 AI 언어능력 평가: 감정 분석 과제

■ 3. 감정 분석 과제 말뭉치 구축 과정

◆ (3) 감정 분석 과제 말뭉치 형식 및 규모

- **과제 말뭉치 형식:** 데이터 세트는 JSON-L(jsonlines) 형식으로 제공
 - 각 JSON은 발화 발화 내에 감정을 분석할 대상, 대상에 대한 화자의 감정을 제공
 - 훈련, 검증, 시험 데이터는 출력(output)값의 분포가 유사하도록 랜덤하게 분할

- 과제 말뭉치 규모

| JSON 규모 | 훈련 | 검증 | 시험 |
|---------|--------|-------|-------|
| 발화 수 | 37,932 | 4,751 | 4,748 |

-26-

2023년 국립국어원 AI 언어능력 평가: 감정 분석 과제

■ 3. 감정 분석 과제 말뭉치 구축 과정

◆ (3) 감정 분석 과제 말뭉치 형식 및 규모

| 항목 | 내용 |
|------------|--|
| 훈련용 데이터의 예 | <pre>{ "id": "데이터id1", "input": { "form": "아뉴스레터에서뮤지컬킹아더관람신청받는데...가고싶은데약속이네ㅇ 어으아앙악짜증", "target": { "form": "약속", "begin": 35, "end": 37 } }, "output": { "joy": "False", "anticipation": "False", "trust": "False", "surprise": "False", "disgust": "False", "fear": "False", "anger": "True", "sadness": "False" } }</pre> <p>→ 아이디(id)와 입력(input), 그리고 출력(output)으로 구성</p> |

-27-

2023년 국립국어원 AI 언어능력 평가: 감정 분석 과제

■ 3. 감정 분석 과제 말뭉치 구축 과정

◆ (3) 감정 분석 과제 말뭉치 형식 및 규모

| 항목 | 내용 |
|---|--|
| 평가용 데이터의 예 (제출 전) | <pre>{ "id": "데이터id1", "input": { "form": "아뉴스레터에서뮤지컬킹아더관람신청받는데...가고싶은데약속이네ㅇ 어으아앙악짜증", "target": { "form": "약속", "begin": 35, "end": 37 } } }</pre> |
| 학습용 데이터와 동일한 형태, "output" 키와 값을 제거한 데이터 | |
| 제출 데이터 | <pre>{ "id": "데이터id1", "input": { "form": "아뉴스레터에서뮤지컬킹아더관람신청받는데...가고 싶은데약속이네ㅇ 어으아앙악짜증", "target": { "form": "약속", "begin": 35, "end": 37 } }, "output": { "joy": "False", "anticipation": "False", "trust": "False", "surprise": "False", "disgust": "False", "fear": "False", "anger": "True", "sadness": "False" } }</pre> <p>평가용 데이터에 "output"을 생성. 감정 category는 고정</p> |

-28-

국외 감정 분석 말뭉치 구축 현황

| 데이터세트 | 레지스터 | 규모 | 감정 레이블 | 형태 |
|--|----------------|--------|--|----|
| CBET G. Shahrai et al. (2015) | 트위터 | 81,162 | anger, fear, disgust, joy, love, sadness, surprise, thankfulness, guilt | 멀티 |
| SSEC Schuff et al.(2017) | 트위터 | 4,868 | anger, disgust, fear, joy, sadness, surprise, trust, anticipation | 멀티 |
| EmoInt S. Mohammad& F. Bravo-Marquez.(2017b) | 트위터 | 7,097 | anger, disgust, fear, joy, sadness, surprise, disgust, surprise | 단일 |
| SemEval 2018 Task 1 S. Mohammad et al.(2018) | 트위터 | 10,983 | anger, anticipation, disgust, fear, joy, love, optimism, pessimism, sadness, surprise | 멀티 |
| Multi-Genre Corpus for Emotion S. Tafreshi&M. Diab(2018) | 뉴스, 블로그, 영화 리뷰 | 17,220 | anger, disgust, fear, joy, sadness, surprise, trust, anticipation | 단일 |
| DENS C. Liu et al.(2019) | 구어 내러티브 | 9,710 | joy, sadness, anger, fear, anticipation, surprise, love, disgust | 단일 |
| BMET Huang et al.(2019) | 트위터 | 96,323 | anger, fear, joy, sadness, surprise, thankfulness | 멀티 |
| GoEmotions D. Demszky et al.(2020) | 인터넷 댓글 | 58,009 | amusement, excitement, joy, love, desire, optimism, caring, pride, admiration, gratitude, relief, approval, realization, surprise, curiosity, confusion, fear, nervousness, remorse, embarrassment, disappointment, sadness, grief, disgust, anger, annoyance, disapproval | 멀티 |
| EmotionNL DeBruyne, L., & De Clercq, O. (2022) | 트위터, TV 댓글 | 2,000 | anger, fear, joy, sadness, love | 단일 |
| CM-MEC-21 I. Amer et al.(2022) | 메시지 | 12,000 | anger, anticipation, disgust, fear, joy, love, optimism, pessimism, sadness, surprise, trust | 멀티 |

-29-

국외 감정 분석 벤치마크

■ Emotion Recognition/classification

Emotion Recognition

386 papers with code • 5 benchmarks • 42 datasets

Emotion Recognition is an important area of research to enable effective human-computer interaction. Human emotions can be detected using speech signal, facial expressions, body language, and electroencephalography (EEG). Source: [Using Deep Autoencoders for Facial Expression Recognition](#)

Benchmarks

These leaderboards are used to track progress in Emotion Recognition

Emotion Classification

81 papers with code • 7 benchmarks • 22 datasets

Emotion classification, or emotion categorization, is the task of recognising emotions to classify them into the corresponding category. Given an input, classify it as 'neutral or no emotion' or as one, or more, of several given emotions that best represent the mental state of the subject's facial expression, words, and so on. Some example benchmarks include ROCStories, Many Faces of Anger (MFA), and GoEmotions. Models can be evaluated using metrics such as the Concordance Correlation Coefficient (CCC) and the Mean Squared Error (MSE).

| Trend | Dataset | Best Model | Paper | Code | Compare | Add a Result |
|-------|-------------|---|-------|------|---------|-------------------------|
| | RAVDESS | LogisticRegression on posteriors of xlsr-Wav2Vec2.0&bi-LSTM+Attention | | | | See all |
| | Emomusic | Jukebox (Pre-training: CALM) | | | | See all |
| | MPED | BiHDIM | | | | See all |
| | SEED | 4D-aNN | | | | See all |
| | MSP-Podcast | w2v2-L-robust-12 | | | | See all |

| Trend | Dataset | Best Model | Paper | Code | Compare | Add a Result |
|-------|------------------------|------------------|-------|------|---------|-------------------------|
| | SemEval 2018 Task 1E-c | SpanEmo | | | | See all |
| | EWALK | ProxEmo (ours) | | | | See all |
| | CMU-MOSEI | MARLIN (ViT-L) | | | | See all |
| | ROCStories | Semi-supervision | | | | See all |
| | MFA | MLKNN | | | | See all |
| | GoEmotions | BERT | | | | See all |

<https://paperswithcode.com/task/emotion-recognition>
<https://paperswithcode.com/task/emotion-classification>

-30-

국외 감정 분석 벤치마크

■ Emotion Classification

- ◆ SemEval 2018 Task 1E-c
 - 영어, 아랍어, 스페인어 트윗에 대한 11개 감정 분류
 - anger, anticipation, disgust, fear, joy , love, optimism, pessimism, sadness, surprise, trust, neutral
- ◆ EWALK
 - 사람의 걷기 모션에 대해 감정 인식 및 happy, angry, sad, neutral 4가지 라벨로 분류
- ◆ MFA
 - 영상(video)에 대한 감정 인식 및 분류
 - 'annoyed', 'anger', 'disgust', 'hatred', 'furious'를 기본으로 하며, 이와 관련된 감정 13개에 대해 분류
- ◆ GoEmotions
 - Reddit 댓글에 대해 27개 감정으로 분류

-31-

국외 감정 분석 벤치마크

■ Emotion Recognition

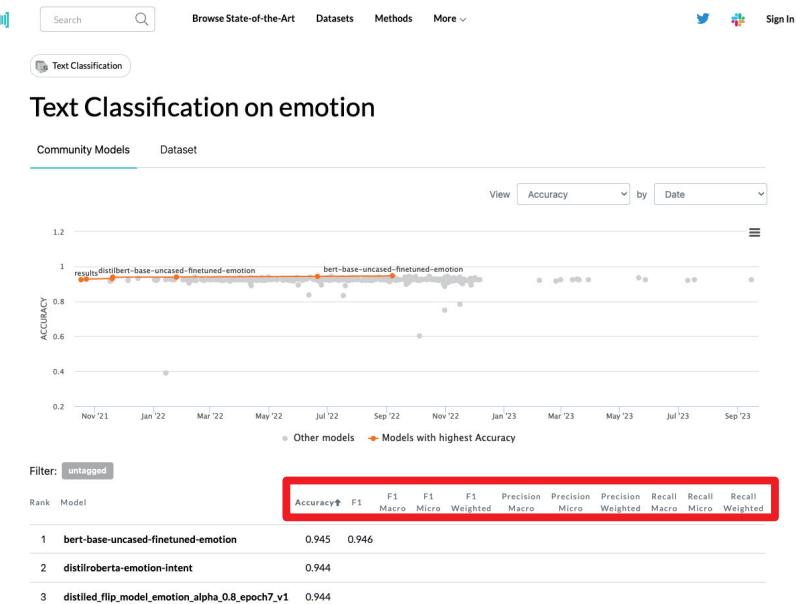
- ◆ RAVDESS (Ryerson Audio-Visual Database of Emotional Speech and Song)
 - 사람의 음성과 음악에서 감정을 탐지한 후 라벨 부여
 - 음성: calm, happy, sad, angry, fearful, surprise, disgust
 - 음악: calm, happy, sad, angry, fearful
- ◆ MPED (A Multi-Modal Physiological Emotion Database for Discrete Emotion)
 - 신체&생리학적 관점에서의 멀티 모달 감정 데이터셋
 - 피험자가 비디오를 통해 느낀 감정을 신체&생리학적 신호로 탐지한 후, 해당 신호에 대한 감정 인식
 - 호흡, 심전도, 뇌파 등을 사용하여 joy, funny, anger, sadness, fear, disgust, neutrality 분류
- ◆ MSP-Podcast: 음성에 대한 감정 인식
- ◆ SEED dataset: 뇌파에 대한 감성&감정 인식 데이터셋

-32-

국외 감정 분석 리더보드

■ Text Classification on emotion 과제의 리더보드

성능 평가 지표 다양



<https://paperswithcode.com/sota/text-classification-on-emotion>

-33-

감정 분석 과제의 발전 방향

- ◆ 감정 분석 말뭉치의 장르 확장: 장르 간 감정 양상 비교
- ◆ 감정 표출 매체의 확장: 단일 매체 → 멀티 모델
- ◆ 감정 분석의 세밀도 제고: 감정 유형의 확장
- ◆ 감정 분석 과제 평가 지표 개발 및 다양화

-34-

경청해 주셔서 감사합니다.

-35-



이야기 완성 과제의 특징과 의의

2023년 가을 한국코퍼스언어학회 10주년 기념 학술대회
한국의 말뭉치 구축 및 연구의 지평



함영균 (테디썸 대표)
hahmyg@teddysum.ai

발표자 소개



함영균

- KAIST 전산학부 박사/석사 (NLP)
- ISO/TC37 전문위원
- 문화체육부 장관상('12)
- HCLT 우수논문 5회 ('14, '16, '18, '19, '22)
- COLING 2022 Program Chair



(주)teddysum 대표이사, 2021~



오픈소스 한국어 LLM
(요약 및 평가 특화)



ClinicalNLP 경진대회
국제순위 2위



국립국어원

- 대화 요약 사업
3년 수행
- 경진대회 사업
2년 연속 수행

오늘의 발표 주제

- 이야기 완성 과제의 특징과 의의

- 이야기 완성이란?
- 이야기 완성 과제란?
- 이야기 완성 과제와 인공지능?

3



기초 단계의 이야기 완성 – 단어 수준

| | |
|--|-------------|
| Class: _____ | Date: _____ |
| STORY COMPLETION | |
| money died found sowed help ill treasure dug good | |
|  | |
| <p>A farmer had four sons. They did not _____ their father in his work. One day the farmer fell seriously _____. He called his sons and told them that there was a hidden _____ in his field. But he could not tell the exact place. The farmer then _____. The sons _____ the field but _____ no treasure. They _____ the field with wheat. They had a very _____ crop that year. They sold it and got a lot of _____. Then they realized that their father had spoken the truth.</p> | |

4

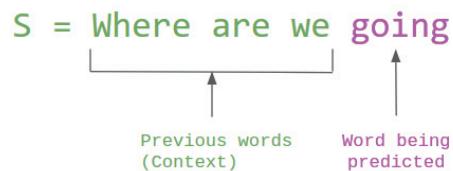
- 초등학교 1~2학년 수준의 문제
- 주어진 상황과 문맥에 따라 [빈칸]에 적절한 단어를 채워넣는 문제
- 추론 능력 및 어휘력 향상 교육 및 평가에 활용



기초 단계의 이야기 완성 – 단어 수준

- 주어진 상황과 문맥에 따라
[빈칸]에
적절한 단어를 채워넣는 문제

~ = Language Modeling



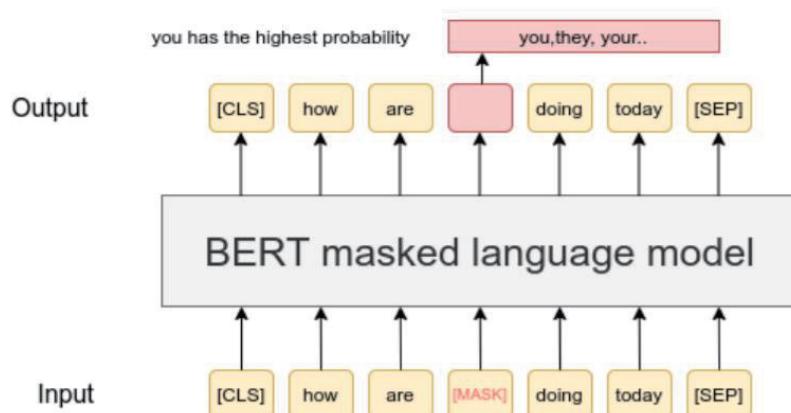
$$P(S) = P(\text{Where}) \times P(\text{are} \mid \text{Where}) \times P(\text{we} \mid \text{Where are}) \times P(\text{going} \mid \text{Where are we})$$

5



기초 단계의 이야기 완성 – 단어 수준

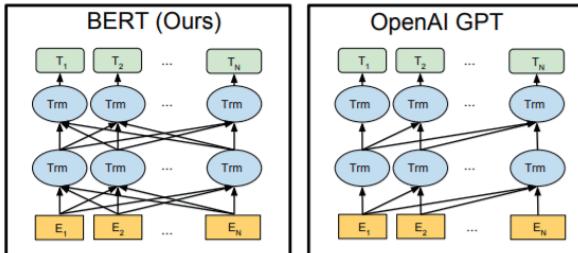
- 주어진 상황과 문맥에 따라 [빈칸]에 적절한 단어를 채워넣는 문제



6



BERT의 등장 – 사전 학습 언어 모델 = 모든 면에서 기계학습 모델보다 뛰어남



| System | Dev | | Test | |
|---------------------------------------|-------------|-------------|-------------|-------------|
| | EM | F1 | EM | F1 |
| Leaderboard (Oct 8th, 2018) | | | | |
| Human | - | - | 82.3 | 91.2 |
| #1 Ensemble - nlnet | - | - | 86.0 | 91.7 |
| #2 Ensemble - QANet | - | - | 84.5 | 90.5 |
| #1 Single - nlnet | - | - | 83.5 | 90.1 |
| #2 Single - QANet | - | - | 82.5 | 89.3 |
| Published | | | | |
| BiDAF+ELMo (Single) | - | 85.8 | - | - |
| R.M. Reader (Single) | 78.9 | 86.3 | 79.5 | 86.6 |
| R.M. Reader (Ensemble) | 81.2 | 87.9 | 82.3 | 88.5 |
| Ours | | | | |
| BERT _{BASE} (Single) | 80.8 | 88.5 | - | - |
| BERT _{LARGE} (Single) | 84.1 | 90.9 | - | - |
| BERT _{LARGE} (Ensemble) | 85.8 | 91.8 | - | - |
| BERT _{LARGE} (Sgl.+TriviaQA) | 84.2 | 91.1 | 85.1 | 91.8 |
| BERT _{LARGE} (Ens.+TriviaQA) | 86.2 | 92.2 | 87.4 | 93.2 |

7



창의적 글쓰기로서의 이야기 완성

Creative Writing

Read the paragraph and continue the story.

Once there was a man. He lived in one of the most amazing places ever. He had three amazing children. He had a brother and a sister. He had many friends, and should have been very happy. But there was something missing. Someone he had lost five years ago on that fateful day.....

79 **FUN WRITING PROMPTS**

1. List down 10 ways you can be a good friend.
 2. What is your favorite way to relax?
 3. If you could spend the weekend building something, what would you build?
 4. Who is your favorite teacher and why?
 5. Describe your dream house
 6. Write about what do you when you look into the mirror.
 7. Write about the best dream you have ever had
 8. Write about the worst holiday you have ever experienced.
 9. What will you do if you were given \$1000?
 10. What do you like most about your home?
 11. What is the best thing about online learning?
 12. What do you miss most about going to school?
 13. What do you like best about your school?
 14. Write about your biggest fear?
 15. Write about what you want to become when you grow up?
 16. Write about your favorite place in the world.

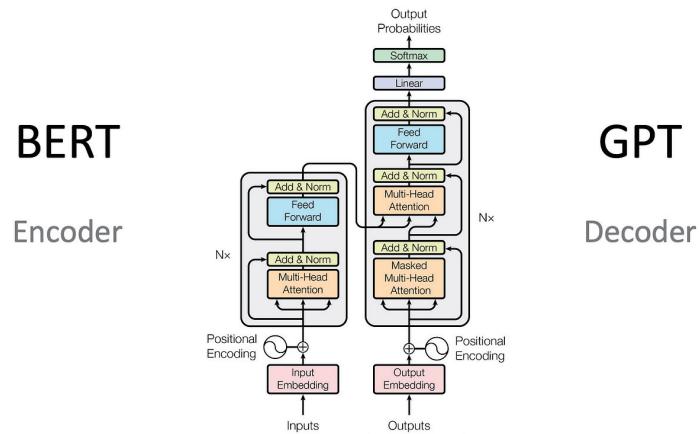
MORE PROMPTS

8



생성적 AI: 디코더 기반의 모델

- Transformer architecture (encoder & decoder)

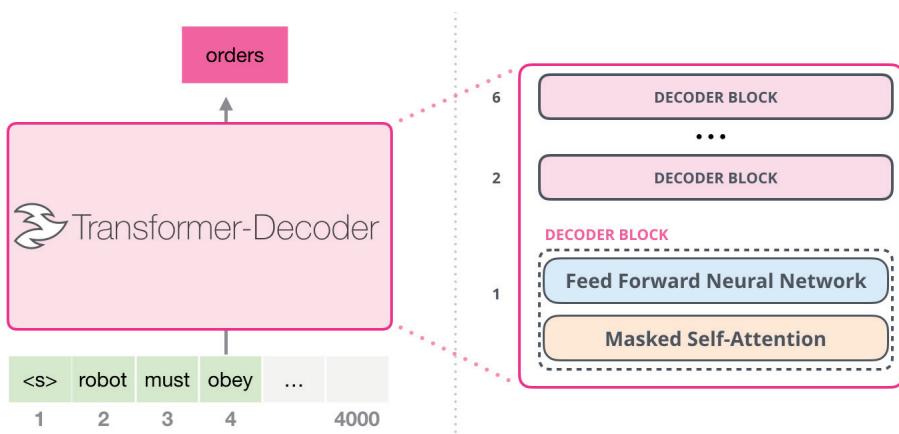


9



생성적 AI: 디코더 기반의 모델

- Decoder-only architecture



10



GPT-1 (2018)

- GPT-1의 특징
 - 사전학습 기반 언어모델
 - 특정 task 적용 시 성능 향상을 위해 입력 형태를 바꾸면서 미세조정(fine-tuning)

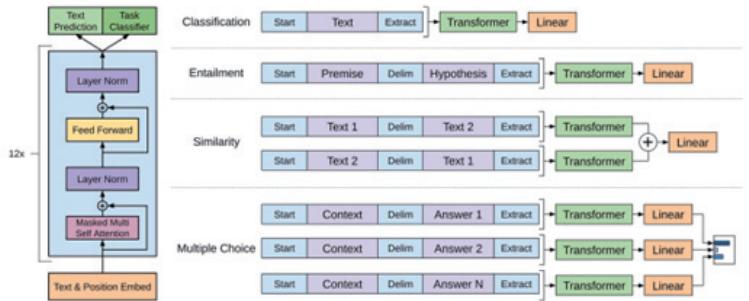


Figure 1: (left) Transformer architecture and training objectives used in this work. (right) Input transformations for fine-tuning on different tasks. We convert all structured inputs into token sequences to be processed by our pre-trained model, followed by a linear+softmax layer.

11



GPT-1 (2018) → GPT-2 (2019)

- GPT-1의 특징
 - 사전학습 기반 언어모델
 - 특정 task 적용 시 성능 향상을 위해 입력 형태를 바꾸면서 미세조정(fine-tuning)
- GPT-2의 특징
 - 모델의 크기가 사전학습 언어모델의 성능에 영향을 준다는 것을 실증
 - 미세조정 없이 기존 태스크의 SOTA 성능을 달성함 (보통 zero-shot이라고 칭함)

| | LAMBADA (PPL) | LAMBADA (ACC) | CBT-CN (ACC) | CBT-NE (ACC) | WikiText2 (PPL) |
|-------|------------------|------------------|-----------------|-----------------|--------------------|
| SOTA | 99.8 | 59.23 | 85.7 | 82.3 | 39.14 |
| 117M | 35.13 | 45.99 | 87.65 | 83.4 | 29.41 |
| 345M | 15.60 | 55.48 | 92.35 | 87.1 | 22.76 |
| 762M | 10.87 | 60.12 | 93.45 | 88.0 | 19.93 |
| 1542M | 8.63 | 63.24 | 93.30 | 89.05 | 18.34 |

12



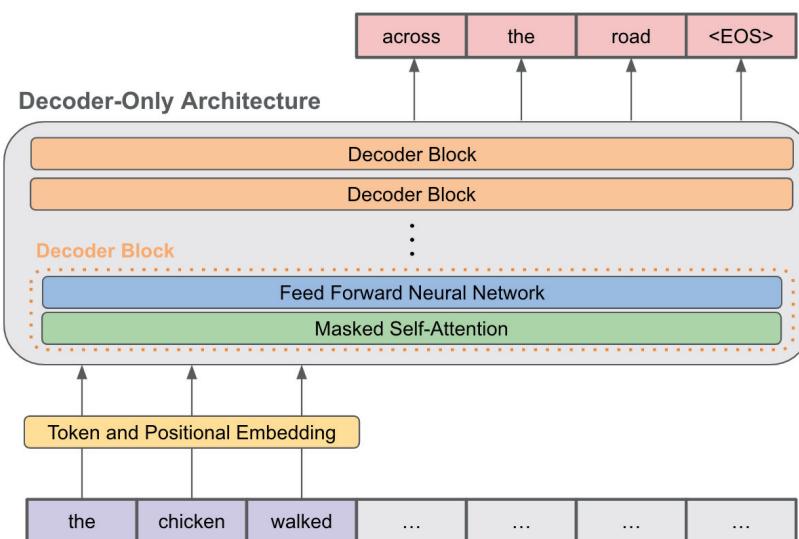
GPT-1 (2018) → GPT-2 (2019) → GPT-3 (2020)

- GPT-1의 특징
 - 사전학습 기반 언어모델
 - 특정 task 적용 시 성능 향상을 위해 입력 형태를 바꾸면서 미세조정(fine-tuning)
- GPT-2의 특징
 - 모델의 크기가 사전학습 언어모델의 성능에 영향을 준다는 것을 실증
 - 미세조정 없이 기존 태스크의 SOTA 성능을 달성함 (보통 zero-shot이라고 칭함)
- GPT-3의 특징
 - 파라미터를 기존 1.5B에서 175B로 대폭 증가

13



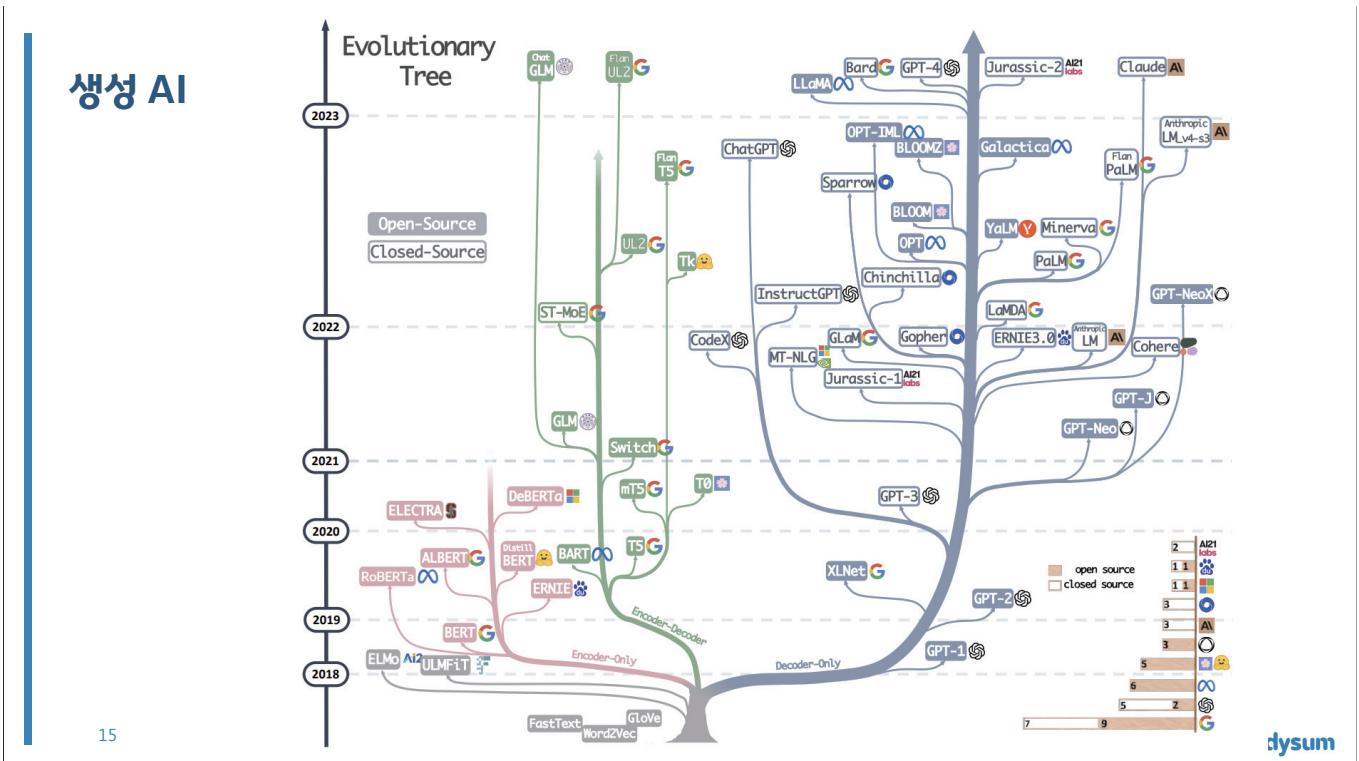
생성 AI: Decoder-only architecture



14



생성 AI



AI가 쓰는 글은 “적절한가?”

- 사건 사이의 **인과 관계**를 이해하였는가?
– Causality (Correlational relationship between events)
- 사건 사이의 **시간 관계**를 이해하였는가?
– Temporal relationship
- 즉, **상식적으로** 글이 쓰였는가?
– AI가 commonsense 를 이해하고, 학습하고, 표현할 수 있는가?

ROC dataset (Story Cloze Test)

- Mostafazadeh et al., 2016
- AI는 주어진 4개 문장에 대해 올바른 결론과 틀린 결론을 구별할 수 있는가?

| Context | Right Ending | Wrong Ending |
|---|--|--|
| Tom and Sheryl have been together for two years. One day, they went to a carnival together. He won her several stuffed bears, and bought her funnel cakes. When they reached the Ferris wheel, he got down on one knee. | Tom asked Sheryl to marry him. | He wiped mud off of his boot. |
| Karen was assigned a roommate her first year of college. Her roommate asked her to go to a nearby city for a concert. Karen agreed happily. The show was absolutely exhilarating. | Karen became good friends with her roommate. | Karen hated her roommate. |
| Jim got his first credit card in college. He didn't have a job so he bought everything on his card. After he graduated he amounted a \$10,000 debt. Jim realized that he was foolish to spend so much money. | Jim decided to devise a plan for repayment. | Jim decided to open another credit card. |

Table 4: Three example Story Cloze Test cases, completed by our crowd workers.

17



ROC dataset (Story Cloze Test)

- ChatGPT의 이해 능력

다음은 하나의 짧은 이야기입니다.

철수와 영희는 2년간 같이 지냈습니다.
어느날 그들은 함께 축제에 같이 갔다.
그는 그녀에게 곰인형이나 케이크를
사주었습니다.
그들이 관람차에 도착했을 때, 그는
한쪽 무릎을 꿇었습니다.

이 이야기의 결말은 어떻게 되었을까요?

GPT-3.5 결과

이 이야기의 결말은 철수가 그녀에게
프로포즈하고 그녀가 결혼 제안을
받아들인 것으로 추측됩니다. 이는 "그는
한쪽 무릎을 꿇었습니다."라는 문장에서
힌트를 얻을 수 있습니다. 한쪽 무릎을 꿇는
행동은 전통적으로 결혼 제안의
표시입니다.

18



함의 문제 (entailment) – 보통 NLI 라고도 함 (Natural Language Inference)

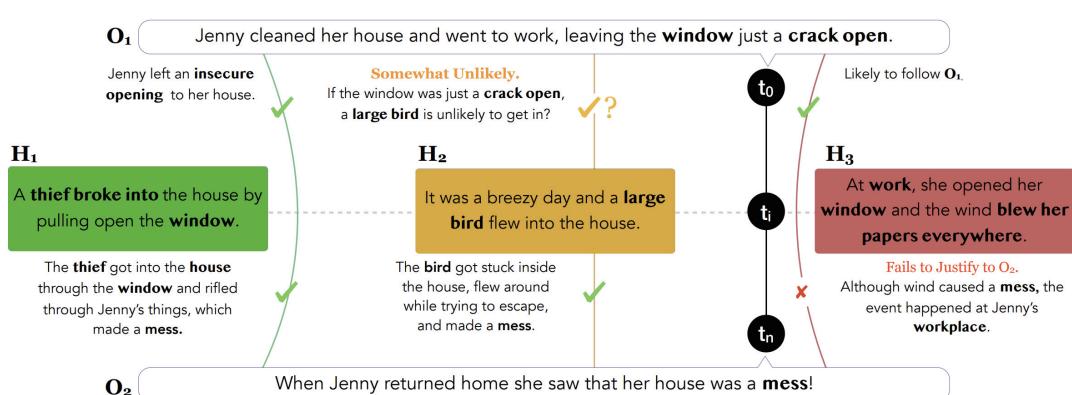
| | | |
|---|---|---------------|
| P^a | A senior is waiting at the window of a restaurant that serves sandwiches. | Relationship |
| H^b | A person waits to be served his food. | Entailment |
| | A man is looking to order a grilled cheese sandwich. | Neutral |
| | A man is waiting in line for the bus. | Contradiction |
| ^a P, Premise. ^b H, Hypothesis. | | |

19



Abductive Commonsense Reasoning

- Bhagavatula et al., 2020



20



Abductive Commonsense Reasoning

• 생성 AI의 이해 능력

다음은 두 사건에 대한 개요입니다.

사건 1: 제니는 집을 치우고 출근을 했습니다. 창문은 조금 열려 있었습니다.

사건 2: 제니가 집으로 돌아오자, 그녀는 집이 엉망이 되어 있는 것을 보았습니다.

사건 1과 사건 2 사이에 있었을 것으로 추정되는 사건은 다음 세 후보 중 하나입니다.

후보 1: 커다란 새가 들어와서 집을 엉망으로 만들었다.

후보 2: 창문을 열고 도둑이 들어왔습니다.

후보 3: 그녀가 직장에서 바쁘게 일을 하면서 물건을 흐트려 두었습니다.

세 후보 중, 가장 적절한 추론은 무엇일까요?

GPT-3.5 결과

가장 적절한 추론은 "후보 3: 그녀가 직장에서 바쁘게 일을 하면서 물건을 흐트려 두었습니다."입니다.

이 추론은 직장에서 바쁘게 일하면서 집을 치우지 못한 제니에게 가장 일반적인 이유를 제시합니다.

다른 두 후보인 새가 들어와 집을 엉망으로 만들거나 창문을 열고 도둑이 들어왔다는 것은 상대적으로 덜 가능성이 높은 시나리오입니다.

21



생성 AI의 이야기 완성 능력은 어떨까요?

다음은 두 사건에 대한 개요입니다.

사건 1: 여성과 남자아이가 도로변 보도를 걷고 있었다.

사건 2: 여성은 잘 따라가고 있다는 표시로 남자아이에게 손을 흔들었다.

두 사건 사이에 가장 있음직한 사건은 무엇일까요? 한번 추측해보세요. 한 문장으로 답변해주세요.

GPT-3.5 결과(들)

사건 1이 사건 2보다 먼저 발생한 것으로 추측됩니다.

(or)

주어진 정보만으로는 두 사건 사이에 다른 중요한 사건이 있는지 알 수 없습니다.

22



생성 AI의 이야기 완성 능력은 어떨까요?

다음은 두 사건에 대한 개요입니다.

사건 1: 민서는 친구네 집에 초대를

받아 방문했다.

사건 2: 민서는 음식을 먹고 친구에게
엄지손가락을 들어 보였다.

두 사건 사이에 가장 있음직한 사건은
무엇일까요? 한번 추측해보세요. 한
문장으로 답변해주세요.

GPT-3.5 결과

가장 있음직한 사건은 "사건 2:
민서는 음식을 먹고 친구에게
엄지손가락을 들어
보였다."입니다.

23

 teddysum

생성 AI의 이야기 완성 능력은 어떨까요?

다음은 두 사건에 대한 개요입니다.

사건 1: 여성과 남자아이가 도로변
보도를 걷고 있었다.

사건 2: 여성은 잘 따라가고 있다는
표시로 남자아이에게 손을 흔들었다.

두 사건 사이에 가장 있음직한 사건은
무엇일까요? 한번 추측해보세요. 한
문장으로 답변해주세요.

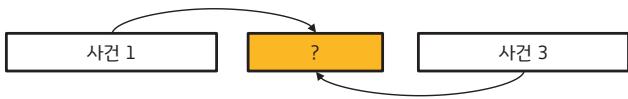
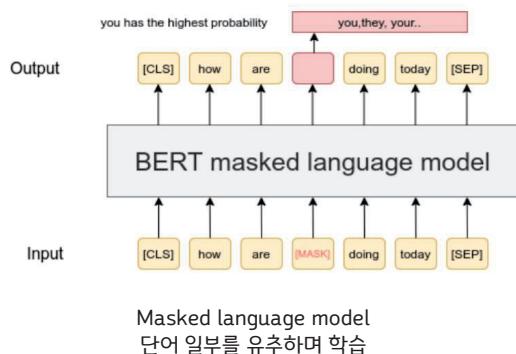
GPT-4 결과

남자아이는 앞서 걸어가며
여성을 뒤로 떨어뜨렸다.

24

 teddysum

이야기 완성 – 일종의 masked story model



모델의 상식 이해 능력, 사건의 인과관계 파악 능력, 시간적 전후관계 파악 능력 등을 학습하고 평가할 수 있는 과업(task)

25



2023 인공지능 언어능력평가 과제: 이야기 생성

이야기 생성 생성 문제

주어진 2개 사건을 기술하는 문장 사이에 적절한 이야기 문장을 생성하는 문제
국립국어원의 '이야기 완성 평가 말뭉치(2022)' 사용

- (t-1)
민수는 요즘 보드 타는 것에 취미를 들였다.
- (t) **생성 대상**
그는 맨날 보드를 타다가 다쳐서 나타났다.
- (t+1)
하지만 보드 타는 것을 멈추지 않았다.

- **평가 방법**
 - ROUGE
 - BLEURT
 - BERTScore
 - + 인간 평가
- **말뭉치**
 - 2022 국립국어원 이야기 완성 평가 말뭉치 기반 (송상현 외, 2022)
 - 훈련 120,140 문장
 - 시험 15,018 문장

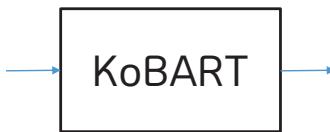
26



기준 모델(baseline)

- https://github.com/teddysum/Korean_SC_2023

“민수는 요즘 보드 타는 것에 취미를 들였다.”
“하지만 보드 타는 것을 멈추지 않았다.”



그는 맨날 보드를 타다가
다쳐서 나타났다.

| ROUGE | BLEURT | BERT Score |
|-------|--------|------------|
| 32.3 | 38.6 | 76.2 |

27'



경진대회 현황(10월 4일 기준)

| 순위 | 팀명 | 모델명 | 평가 카수 | 소과제별 | |
|----|------------|-------------------|------------|------------|------------------|
| | | | | 이야기 완성 | 제출 일시 |
| 1 | 국내산 치약어 | 세오시스템 | 59,5852348 | 59,5852348 | 2023.10.04.06:08 |
| 2 | 얼리 | test_v12 | 59,3228375 | 59,3228375 | 2023.10.04.16:52 |
| 3 | 용강한방아리 | baseline_v9_13 | 59,2697018 | 59,2697018 | 2023.10.03.00:03 |
| 4 | 개짜는팀 | 개짜는모델 | 59,0558110 | 59,0558110 | 2023.09.28.23:41 |
| 5 | Docent | 가천대학교 초대학원 | 58,7366837 | 58,7366837 | 2023.09.09.01:13 |
| 6 | 태양빛노래애 | baseline_9 | 58,6921766 | 58,6921766 | 2023.09.11.15:01 |
| 7 | 이예민 | nalmungchi | 58,4731863 | 58,4731863 | 2023.09.16.17:14 |
| 8 | 중불의왕국 | 마리2 | 58,0864650 | 58,0864650 | 2023.09.12.18:33 |
| 9 | 603310 | v00_4500 | 58,0681285 | 58,0681285 | 2023.10.01.18:59 |
| 10 | Noob | test | 57,5660495 | 57,5660495 | 2023.09.10.02:40 |
| 11 | IDOL_PRIDE | Nagase_Kotono-Bot | 57,3070489 | 57,3070489 | 2023.08.29.13:31 |
| 12 | 부미네님후루 | 광주공무열매 | 56,5207105 | 56,5207105 | 2023.09.14.10:35 |
| 13 | 이 | baseline000 | 56,2245840 | 56,2245840 | 2023.09.22.16:19 |
| 14 | 붓지와민병관 | baseline | 56,1526391 | 56,1526391 | 2023.08.24.18:11 |
| 15 | 이불벌레와야 | baseline | 56,1345638 | 56,1345638 | 2023.09.27.23:19 |
| 16 | Coconut | baseline9 | 56,0491740 | 56,0491740 | 2023.09.18.14:30 |
| 17 | 얼심마리자우림 | baseline | 55,2469871 | 55,2469871 | 2023.08.25.22:24 |
| 18 | 라이언 | baseline | 53,6232992 | 53,6232992 | 2023.09.19.11:53 |
| 19 | HHST | Baseline_3 | 49,6850033 | 49,6850033 | 2023.10.03.20:26 |

28



경진대회 현황(10월 04일 기준)

- 19팀 참여
 - 베이스라인 성능

| ROUGE | BLEURT | BERT Score |
|-------|--------|------------|
| 32.3 | 38.6 | 76.2 |

- 1위팀

| ROUGE | BLEURT | BERT Score |
|-------|--------|------------|
| 45.89 | 51.89 | 80.96 |

29



e.g., Word overlap-based metrics: ROUGE, BLEU, and METEOR

ROUGE [1]

- N-gram recall between output and reference

$$\frac{\text{Number of } n\text{-grams found in model and reference}}{\text{Number of } n\text{-grams in model}}$$

[1] Lin, "ROUGE: A Package for Automatic Evaluation of Summaries", 2004

BLEU [2]

- N-gram overlap between output and reference
- Compute precision for n-grams of size 1 to 4
- Add brevity penalty

$$\text{BLEU} = \min \left(1, \frac{\text{output-length}}{\text{reference-length}} \right) \left(\prod_{i=1}^4 \text{precision}_i \right)^{\frac{1}{4}}$$

[2] Papineni et al., "Bleu: a Method for Automatic Evaluation of Machine Translation", 2002

METEOR [3]

- Harmonic mean of unigram precision and recall
- with recall weighted higher than precision

$$F_{mean} = \frac{10PR}{R + 9P}$$

[3] Banerjee and Lavie, "METEOR: An Automatic Metric for MT Evaluation with Improved Correlation with Human Judgments", 2005

30



e.g., Word overlap-based metrics: ROUGE, BLEU, and METEOR

Limitations of:

- difference between positive and negative (e.g., “not”)
- grammatical correctness (only n-gram)
- contextual semantics (e.g., synonyms)
- common sense

Above all, reference is not a *gold standard*, it's just an example

31



LLM: SOTA Evaluators of Translation Quality [1]

Score the following translation from {source_lang} to {target_lang} **with respect to the human reference** on a continuous scale from 0 to 100, where score of zero means "no meaning preserved" and score of one hundred means "perfect meaning and grammar".
(source_lang) source: "{source_seg}"
(target_lang) human reference: "{reference_seg}"
(target_lang) translation: "{target_seg}"
Score:

| Model name | Abbrev. | Model used |
|---------------|---------|-----------------------|
| GPT-2 | — | Radford et al. (2019) |
| Ada | — | text-ada-001 |
| Babbage | Bab | text-babbage-001 |
| Curie | Curie | text-curie-001 |
| Davinci-002 | Dav2 | text-davinci-002 |
| ChatGPT | Chat | text-chat-davinci-002 |
| Davinci-003 | Dav3 | text-davinci-003 |
| GPT-3.5-turbo | Turbo | gpt-3.5-turbo |
| GPT-4 | GPT4 | gpt-4 |

[1] Kocmi and Federmann, Large Language Models Are State-of-the-Art Evaluators of Translation Quality, 2023

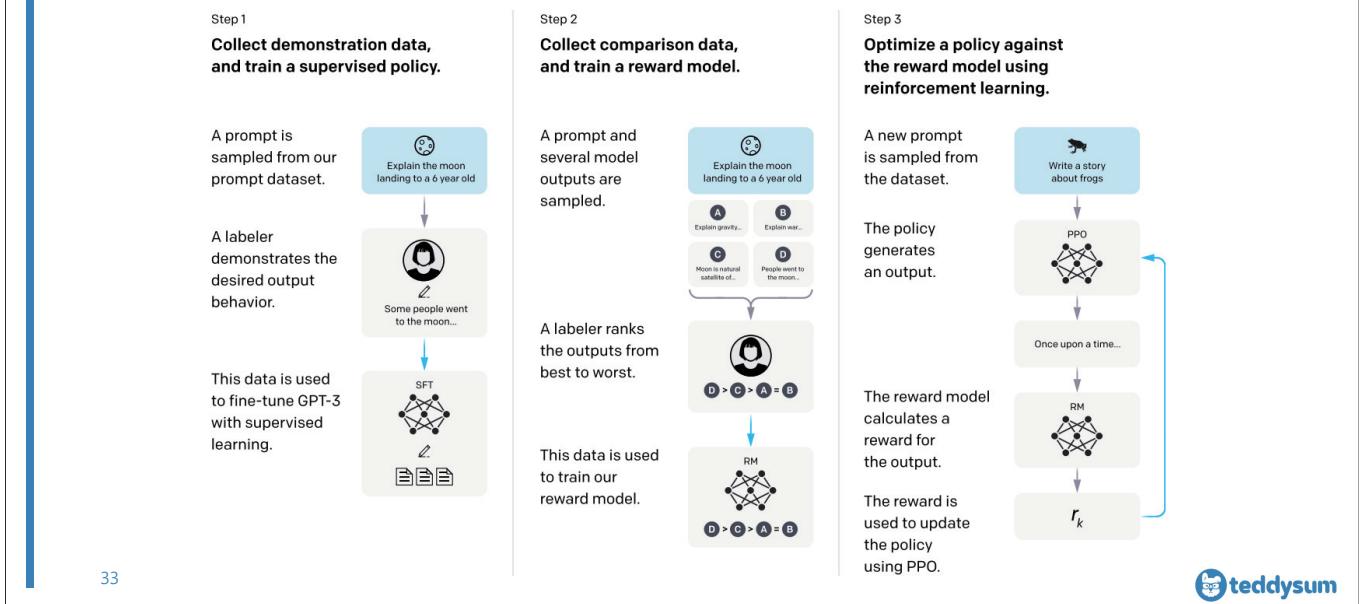
32

| Metric | Accuracy |
|-----------------------------|----------|
| GEMBA-GPT4-DA | 89.8% |
| GEMBA-GPT4-DA[noref] | 87.6% |
| MetricX XXL | 85.0% |
| BLEURT-20 | 84.7% |
| COMET-22 | 83.9% |
| COMET-20 | 83.6% |
| UniTE | 82.8% |
| MS-COMET-22 | 82.8% |
| MATESE | 81.0% |
| YiSi-1 | 79.2% |
| COMETKiwi[noref] | 78.8% |
| COMET-QE[noref] | 78.1% |
| BERTScore | 77.4% |
| UniTE-src[noref] | 75.9% |
| MS-COMET-QE-22[noref] | 75.5% |
| MATESE-QE[noref] | 74.8% |
| f200spBLEU | 74.1% |
| chrF | 73.4% |
| BLEU | 70.8% |

Table 2: Results for the system-level pairwise accuracy compared to the current automatic metric. Metrics marked as “[noref]” do not use a reference translation.



사람 평가의 중요성: RLHF (Reinforcement Learning by Human Feedback)



33

향후 방향

- 정량 평가의 한계를 보완하기 위한 인간 평가(전문가 평가)의 적용의 필요
- 인간 평가의 자동화(Reward Model)
 - 데이터 수집 및 기계학습을 통한 인간 평가의 자동화
 - 생성형 AI 평가를 위한 기초 자료로 활용

34

teddysum

표 기반 유사 문장 말뭉치를 이용한 AI 언어능력 평가

2023.10.07

송영숙, 정유남, 유현조

목차

1. 도입

1.1 표 기반 유사 문장 말뭉치 소개

1.2. 선행 연구 동향

2. 표 기반 문장 말뭉치 구축과 검수 과정

2.1. 표 기반 유사 문장 말뭉치 구축

2.2. 표 기반 유사 문장 말뭉치 검수

2.3. 표 기반 유사 문장 말뭉치 규모

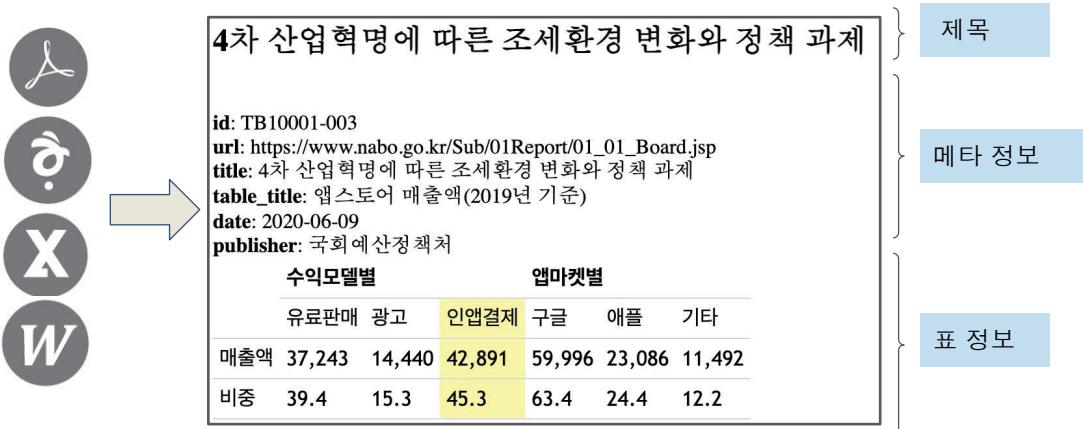
3. 모델 평가

4. 제언

1.2 표 기반 유사 문장 말뭉치 소개

(1) 표 수집 및 정제

그림1. 표기반 유사 문장 말뭉치 예시



- 다양한 형태의 포맷과 표 형식으로부터 제목, 메타 정보와 표 정보 추출
 - 주요 표 정보에 강조 표시(highlight) 후 모델에 입력하기 위해서 문장의 형태로 변환

3 /

1.2 표 기반 유사 문장 말뭉치 소개

(2) 표를 기반으로 사람이 생성한 유사 문장

(1) 2019년 앱스토어 매출액을 수익모델별로 살펴본 결과, 가장 큰 비중을 차지한 것은 인앱결제로 매출액은 총 4조 2,891억 원이며 비중은 45.3%로 나타났다.

(2) 2019년도 기준, 가장 큰 비중을 차지한 인앱결제 매출액은 총 4조 2,891억 원으로 45.3%의 비중을 보인다는 사실을 앱스토어 매출액을 수익모델별로 조사한 결과에서 밝히고 있다.

(3) 2019년을 기준으로 인앱결제가 수익모델 중 앱스토어 매출액의 가장 큰 비중(45.3%)을 차지하며 총 4조 2,891억 원의 매출액을 보였다.

(4) 인앱결제는 2019년 앱스토어 매출액 중 45.3%로 가장 큰 비중을 차지하며 그 매출액은 총 4조 2,891억 원이다.

(5) 인앱결제는 총 4조 2,891억 원의 매출액을 낸다 2019년도 앱스토어 매출액을 수익모델별로 분석하였을 때, 가장 큰 비중인 45.3%를 차지했다

(3) 표를 기반으로 한 모델의 유사 문장 생성 과정

- 모델은 표에서 사람이 문장을 생성하는 과정을 학습
 - 학습한 문장에서 얻은 정보를 반영하여 자연스러운 한국어 문장을 생성하는 과정
 - 국가, 기업 보고서 등 정확성을 요구하는 데이터 처리가 증가할수록 표의 처리에 대한 수요도 증가

4 /

1.2 표 기반 유사 문장 말뭉치 소개

표 기반 유사 문장 말뭉치의 중요성

- 자연어 생성 과제에서 가장 주요한 문제로 떠오른 환각(hallucinations) 현상의 정량적 분석을 제공할 수 있음.
 - 환각 현상은 생성 모델이 학습한 문장으로부터 새로운 문장을 생성할 때 신뢰할 수 있는 결과를 도출하지 못하고 단지 그럴 듯하게 보이는 문장을 생성하는 것을 의미.
 - Nelson Liu (2023:2) 따르면 평균적으로 자연어 생성 문장의 51.5%만이 참조 문장에 의해 완전히 뒷받침되며, 74.5% 정도가 연관된 문장을 뒷받침한다고 한다.
 - 소설이나 대화 등을 요약할 때는 참조 문장과 완전히 동일한 정보를 제공하지 않아도 크게 문제가 되지 않고 유사한 지식을 사용하더라도 측정의 기준을 정하기가 어려움
 - 반면 표 데이터와 같은 경우는 명시적 정보와 그 근거가 되는 숫자를 제공하는 경우가 대부분이기 때문에 잘 주석된 표 데이터가 있다면 이를 통해 모델의 성능을 객관적으로 측정하는 것이 용이함
 - 벤치마크 데이터로서의 표 데이터의 구축과 검수 과정에서 생길 수 있는 사람의 오류와 모델의 성능 수준 및 개선 방향을 제시할 수 있음
 - 정량적 수치와 추론된 사실 등에 대한 명시적 판단 등의 후속 연구 필요

5 /

1.2. 선행 연구 동향(1/2)

자연어 생성 벤치마크 종류(Sebastian Gehrmann et al, 2021: 3)

| 문서 요약 | 단어로부터 문장 생성 및 질의 응답 | 구조화된 정보로부터 문장 생성 |
|---|---|---|
| MLSum : 뉴스기사 요약 | CommonGEN : 개념 집합(단어)에서 개연성 있는 문장 생성 | DART : 도표에 있는 정보를 문장으로 표현 |
| XSum 뉴스의 주요 정보 강조 | WebNLG : 세 쌍의 개념어 집합에서 자연스러운 문장 생성 | ToTTo : 도표의 주요 정보를 문장으로 표현 |
| WikiAuto + Turk/ASSET 동일한 내용을 단순화하여 생성 | Schema-Guided Dialog : 가상 비서 구현 | E2E clean : 입력되어 있는 특정 속성(이름, 가격 등)만을 사용하여 레스토랑 설명 |
| WikiLingua 교육 자료 내용 요약 | Czech Restaurant : 의도 및 특정 속성이 있는 텍스트 산출 | |

생성 말뭉치에 주요 과제는 요약, 단어로부터 문장 생성, 대화 모델 등이 있고 표기반 문장 생성 역시 주요 과제 중에 하나

6 /

1.2. 선행 연구 동향(2/2)

표 기반 문장 생성 연구

- 초기) 어떤 분야의 도메인을 표의 주 연구 대상으로 할 것인가에 대한 배경 지식과 관련된 연구가 중심이 되던 시기. 농구와 관련된 Rotowire(Wiseman 외, 2017) 데이터 세트, 생물학과 관련된 KBGen(Banik 외, 2013) Wikibio(Lebret 외, 2016) 데이터 세트. 식당 예약 등과 관련된 E2E(Novikova 외, 2016, 2017)등
- 중기) 표를 통한 문장 생성과 관련하여서는 Puduppully,R.(2018), Ankur Parikh 외(2020), Jonathan 외(2020) 등. 어떤 내용을 텍스트로 변환할지 계획하는 방법에 대한 연구가 중심. 설명 문서와 짹을 이루는 데이터 레코드 코퍼스가 주어지면 먼저 어떤 정보를 어떤 순서로 언급해야 하는지 강조하는 콘텐츠 계획을 생성한 다음, 콘텐츠 계획을 고려하면서 문서를 생성하고 콘텐츠 선택과 계획을 통합하는 방식을 제안 (Puduppully,R.외. 2018, Ankur Parikh 외(2020) 등).
- 최근) 표와 이미지 데이터 등을 혼합한 연구, E2E clean은 레스토랑이라는 영역에서 사용되는 음식 이미지와 음식에 대한 구조화한 설명 등을 제시하여 실제 세상에서 쓰이는 다양한 데이터 형태를 통합

7 /

목차

1. 도입
 - 1.1 표 기반 유사 문장 말뭉치 소개
 - 1.2. 선행 연구 동향
2. 표 기반 문장 말뭉치 구축과 검수 과정
 - 2.1. 표 기반 유사 문장 말뭉치 구축의 실제
 - 2.2. 표 기반 유사 문장 말뭉치 검수의 실제
 - 2.3. 표 기반 유사 문장 말뭉치 규모
3. 모델 평가
4. 제언

8 /

1.2. 표 기반 유사 말뭉치 구축의 실제(1/3)

[표]

- 기사 제목: 일부 재활용품은 분리 배출해도 재활용되지 않아· 표 제목: 선별하기 힘든 배출유형 응답 순위
· 보도 일자: 2021-01-28 · 저작권: 한국소비자원
· 출처url: <https://www.kca.go.kr/home/sub.do?menukey=4002&mode=view&no=1003063256>

정보]

| 순위 | 선별하기 힘든 배출유형 | 응답자 |
|----|----------------------------|-----|
| 1 | 세척이 되지 않아 이물질·오물 등에 오염된 경우 | 29 |
| 2 | 뚜껑, 빨대와 같이 크기가 작은 품목의 경우 | 6 |
| 3 | 불투명 비닐에 담겨 내용물 파악이 힘든 경우 | 4 |
| 4 | 선별해야하는 폐기물의 재질 파악이 힘든 경우 | 3 |
| - | 응답오류(복수선택, 미선택 등) | 8 |
| 합계 | | 50 |

[표와 함께 수집한 문단] 후 [표의 내용이 아닌 경우 삭제]

단독주택에서 수거한 재활용품을 처리하는 공공선별시설 4개소에서 근무하는 작업자 50명을 대상으로 설문조사를 실시한 결과, 38명(76%)이 미선별 잔재물의 발생 원인과 관련해 '선별인력에 비해 반입량이 지나치게 많이 재활용이 가능한 잔재물이 많이 발생'된다고 답변했다. 또한 선별하기 힘든 분리배출유형으로는 29명(58%)이 '세척되지 않아 이물질·오물 등에 오염된 경우'라고 답변해 각 지자체들이 실정에 맞게 재활용품 수거체계를 개선하고, 소비자들이 재활용 가능자원의 분리배출 시 이물질·오물 등을 사전에 제거할 수 있도록 배출요령에 대한 홍보를 강화할 필요가 있다.

9 /

1.2 표 기반 유사 말뭉치 구축의 실제(2/3)

[기준 문장]

선별하기 힘든 분리배출유형으로는 50명 중 29명이 '세척되지 않아 이물질·오물 등에 오염된 경우'라고 하였다.

[유사 문장 예시]

- 선별하기 힘든 분리배출유형으로는 50명 중 29명이 '세척되지 않아 이물질·오물 등에 오염된 경우'라고 답하였다.
- 선별하기 힘든 분리배출유형으로 50명 중 '세척되지 않아 이물질·오물 등에 오염된 경우'에 29명이 답하였다.
- 한국소비자원에 따르면 선별하기 힘든 분리배출유형으로 '세척되지 않아 이물질·오물 등에 오염된 경우'라고 응답한 경우는 50명 중 29명에 해당한다.
- 선별하기 힘든 분리배출유형으로 '세척되지 않아 이물질·오물 등에 오염된 경우'라고 응답한 경우는 50명 중 29명이다.

[문장이
표에 음영으로 표시함]

설명하는

표의

부분]

2. 표 기반 문장 말뭉치 검수의 실제

검수 과정

- 유사 문장 원문 검수: 서로 다른 작업자 2인에 의한 검수
- 한글 파일을 이용하여 전체 문장의 맞춤법, 표의 내용을 충실히 반영하지 않은 경우(숫자 오타 등) 등을 검수, 수정한 후 Github 히스토리 diff를 참조하여 교차 검수

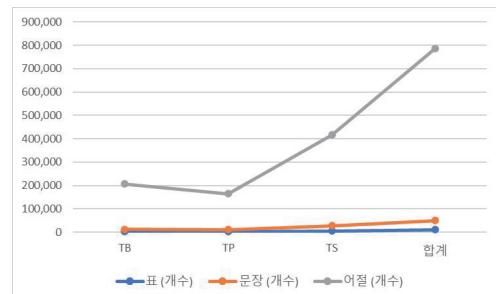
| | | |
|-----|---------------|---|
| 422 | - TB10002-073 | worker1 서울개포9단지는 계획현액 70억 1,900만 원에서 70억 1,700만 원을 집행해 13억 6,400만 원의 2018년도 미월액을 남겼는데, 이는 2019회계연도 임대주택건립사업 사업추진 실적에서 살펴볼 수 있다. |
| 422 | + TB10002-073 | worker1 서울개포9단지는 계획현액 70억 1,900만 원에서 70억 1,700만 원을 집행해 13억 6,400만 원의 2018년도 미월액을 남겼는데, 이는 2019회계연도 임대주택건립사업 사업추진 실적에서 살펴볼 수 있다. |
| 423 | 423 | TB10002-073 worker2 2019회계연도 임대주택건립사업 사업추진 실적에 따르면 서울개포9단지의 경우, 집행액은 70억 1,900만 원의 계획현액 중 70억 1,700만 원으로 나타났으며 13억 6,400만 원이 전년도 미월액으로 확인되었다. |
| 424 | 424 | TB10002-073 worker3 2019년 서울개포9단지의 전년도 미월액은 13억 6,400만 원이었으며, 계획현액 70억 1,900만 원 중 70억 1,700만 원이 집행되었음을 임대주택건립사업 사업추진 실적을 통해 알 수 있다. |
| 425 | 425 | TB10002-073 worker4 2018년도 13억 6,400만 원의 미월액을 남긴 임대주택건립사업에 속하는 서울개포9단지는 2019회계연도를 기준으로 70억 1,900만 원의 계획현액과 70억 1,700만 원의 집행액을 처리하였다. |
| ... | ... | @@ -445,7 +445,7 @@ TB10002-077 worker3 2015년과 2016년 공무원연금대부 계획액 소진시 |
| 445 | 445 | TB10002-077 worker4 2015년과 2016년도는 상반기(5월 19일, 3월 16일)에 공무원연금대부 계획액을 같은 시기에 소진한 연도임을 최근 5년간 공무원연금대부 계획액 소진시기를 분석한 결과에서 밝히고 있다. |
| 446 | 446 | TB10002-078 sentence 2019회계연도 경찰청 소관 세입 결산에 따르면 징수결정액은 1조 9,766억 7,300만 원이었으며, 9,646억 1,900만 원이 수납되어 수납률은 48.8%로 나타났다. |
| 447 | 447 | TB10002-078 worker1 2019회계연도 경찰청 소관 세입 결산에 분석한 바에 의하면, 1조 9,766억 7,300만 원의 징수결정액에서 9,646억 1,900만 원의 수납액이 발생하여 48.8%의 수납률을 기록했다. |
| 448 | - TB10002-078 | - TB10002-078 worker2 2019회계연도 경찰청 소관 세입 결산을 살펴봤을 때, 수납률은 48.8%로 9,646억 1,900만 원이 수납액으로 나타났으며, 징수결정액의 경우 1조 9,766억 7,300만 원으로 확인되었다. |
| 448 | + TB10002-078 | + TB10002-078 worker2 2019회계연도 경찰청 소관 세입 결산을 살펴봤을 때, 수납률은 48.8%로 9,646억 1,900만 원이 수납액으로 나타났으며, 징수결정액의 경우 1조 9,766억 7,300만 원으로 확인되었다. |

11 /

2.3. 표 기반 유사 문장 말뭉치 규모(1/3)

말뭉치 규모

| | 표 (개수) | 문장 (개수) | 어절 (개수) |
|----|--------|---------|---------|
| TB | 2,498 | 12,490 | 205,847 |
| TP | 2,048 | 10,240 | 164,142 |
| TS | 5,516 | 27,580 | 416,441 |
| 합계 | 10,062 | 50,310 | 786,430 |



12 /

2.3. 표 기반 유사 문장 말뭉치 규모(2/3)

문장 길이 (어절수)

| | N | 최소 | 1사분위 | 중앙값 | 평균 | (표준편차) | 3사분위 | 최대 |
|----|--------------|----------|-----------|-----------|--------------|---------------|-----------|-----------|
| TB | 12,490 | 5 | 13 | 16 | 16.48 | (5.02) | 20 | 46 |
| TP | 10,240 | 6 | 12 | 15 | 16.03 | (4.88) | 19 | 49 |
| TS | 27,580 | 5 | 11 | 14 | 15.10 | (5.38) | 18 | 52 |
| 전체 | 50310 | 5 | 12 | 15 | 15.63 | (5.23) | 19 | 52 |

13 /

2.3. 표 기반 유사 문장 말뭉치 규모(3/3)

테이블 크기 (칸의 개수)

| | N | 최소 | 1사분위 | 중앙값 | 평균 | (표준편차) | 3사분위 | 최대 |
|----|--------------|----------|-----------|-----------|--------------|----------------|-----------|------------|
| TB | 12490 | 4 | 18 | 30 | 42.50 | (42.43) | 51 | 611 |
| TP | 10240 | 8 | 25 | 44 | 58.49 | (46.99) | 75 | 335 |
| TS | 27580 | 4 | 24 | 42 | 66.56 | (70.19) | 84 | 692 |
| 전체 | 50310 | 4 | 22 | 38 | 58.95 | (60.79) | 73 | 692 |

강조표시(하이라이트한 칸의 개수)

| | N | 최소 | 1사분위 | 중앙값 | 평균 | (표준편차) | 3사분위 | 최대 |
|----|--------------|----------|----------|----------|--------------|----------------|----------|-----------|
| TB | 12490 | 0 | 2 | 3 | 2.619 | (1.603) | 4 | 8 |
| TP | 10240 | 0 | 3 | 4 | 3.859 | (2.847) | 4 | 72 |
| TS | 27580 | 0 | 3 | 3 | 3.127 | (1.253) | 4 | 21 |
| 전체 | 50310 | 0 | 3 | 3 | 3.150 | (1.822) | 4 | 72 |

14 /

목차

1. 도입

1.1 표 기반 유사 문장 말뭉치 소개

1.2. 선행 연구 동향

2. 표 기반 문장 말뭉치 구축과 검수 과정

2.1. 표 기반 유사 문장 말뭉치 구축의 실제

2.2. 표 기반 유사 문장 말뭉치 검수의 실제

2.3. 표 기반 유사 문장 말뭉치 규모

3. 모델 평가

4. 제언

15 /

모델 평가(1/2)

2023년 상시 과제 기준 베이스라인 모델의 입출력

| 분류 | 내용 | 예시 | 비고 |
|----|-------|---|------|
| 입력 | 표 | "table": [[{"value": "구 분", "is_header": True, "row_span": "2", "column_span": "1"}, {"value": "협약 사업장(톤, %)", "is_header": True, "row_span": 1, "column_span": "3"}, {"value": "비협약 사업장(톤, %)", "is_header": True, "row_span": 1, "column_span": "3"}], [{"value": "19.12", "is_header": True, "row_span": 1, "column_span": 1}, . . .] | JSON |
| 출력 | 설명 문장 | "협약 사업장의 감축량은 4,571톤, 비협약 사업장의 감축량은 539톤이다." | 문자열 |
| 평가 | | ROUGE 1 | |

16 /

모델 평가(2/2)

2023년 상시 과제 기준 베이스라인 모델 점수

| Model | ROUGE-1 | BLUE |
|-----------|---------|--------|
| TeddyBART | 0.4446 | 0.4398 |
| KoBART | 0.4147 | 0.4246 |

https://github.com/teddysum/korean_T2T_baseline

- 표를 기반으로 사람 수준에서 문장을 추론해 내는 능력을 요구하기 때문에 모델에게는 아직 어려운 과제로 이해됨
- **생성 모델의 사실성과 관련된 환각(hallucinations) 현상의 정량적 분석을 위해서는 정량적 수치와 추론된 사실 등에 대한 명시적 판단 등의 후속 연구가 필요함**

17 /

참고문헌

Ankur Parikh, Xuezhi Wang, Sebastian Gehrmann, Manaal Faruqui, Bhuvan Dhingra, Diyi Yang, and Dipanjan Das. 2020. ToTTo: A Controlled Table-To-Text Generation Dataset. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 1173–1186, Online. Association for Computational Linguistics.

Chin-Yew Lin, 2004, ROUGE: A Package for Automatic Evaluation of Summaries. In Text Summarization Branches Out, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.

Eva Banik, Claire Gardent, and Eric Kow. 2013. The kbgen challenge. In Proc. of European Workshop on NLG.

Jekaterina Novikova, Ondřej Dušek, and Verena Rieser. 2017. The E2E dataset: New challenges for end-to-end generation. In Proc. of SIGDIAL.

Sam Wiseman, Stuart M Shieber, and Alexander M Rush. 2017. Challenges in data-to-document generation. In Proc. of EMNLP.

Rémi Lebret, David Grangier, and Michael Auli. 2016. Neural text generation from structured data with application to the biography domain. In Proc. of EMNLP.

Sebastian Gehrmann, Tosin P. Adewumi, Karmyna Aggarwal, et al(2021), “The GEM Benchmark: Natural Language Generation, its Evaluation and Metrics”, In Proceedings of the 1st Workshop on Natural Language Generation, Evaluation, and Metrics, Online. Association for Computational Linguistics.

- 국립국어원 언어정보나눔터 표 기반 문장 생성(시범 운영 과제) 참고 URL :

https://corpus.korean.go.kr/taskOrdtm/taskDownload.do?taskOrdtmId=41&clCd=ING_TASK&subMenuld=sub02#none

- teddysum/korean_T2T_baseline : https://github.com/teddysum/korean_T2T_baseline

- huggingface/transformers (<https://github.com/huggingface/transformers>)

- SKT-AI/KoBART (<https://github.com/SKT-AI/KoBART>)

18 /

2023년 가을 한국코퍼스언어학회 10주년 기념 학술대회

한국의 말뭉치 구축 및 연구의 지평

Session 4

<코퍼스 언어학 일반>



R를 활용한 정도부사에 대한 계층적 군집 분석과 다중 대응분석



AI-Based New Approach to Polyseme: BERT and t-SNE



A Generative AI Study on Translation Universals



AI 학습용 한중 병렬말뭉치의 번역문 평가



R를 활용한 정도부사에 대한 계층적 군집 분석과 다중 대응분석

제 조 현 (연세대학교, 박사 과정)

기존의 정도부사와 관련 연구를 분석해 보면 성별·상황에 따른 빈도와 의미 운율의 연구(김혜영·강범모, 2001; 전지은, 2010), 학습자의 사용 현황 연구(유현경·서상규, 2002), 유의성에 따른 분포와 결합 관계의 연구(한영균·고은아, 2011), 사용역에 따른 선호도 연구(배진영, 2012), 네트워크 분석을 통한 의미 거리 연구(정성훈, 2016), 연령대·성별에 따른 사회언어학적 분석(천수연·조태린, 2021) 등이 있었다. 위의 연구는 모두 말뭉치를 활용하여 여러 방면에서 정도부사를 탐구해 왔다.

최근에 말뭉치 분석 기술의 발전으로 더 많은 사회언어학적 변인(연령대, 성별, 사용역, 지역, 친소관계 등)을 고려하여 용례를 계량적으로 분석할 수 있기 때문에 다변량 분석(Multivariate Analysis)의 연구가 많아지고 있다. 또한 통계기법의 다양화에 따라 SPSS뿐만 아니라 파이썬(python)이나 R를 활용하여 현대국어를 대상으로 한 연구(최정도, 2022)가 많아지고 있는 동시에 종세·근대·근대계몽기 자료의 계량적인 연구(강범일, 2021; 조은경·드프레메리웨인, 2023; 제조현, 2023a, b, c)도 급증하고 있다. 아울러 머신러닝 언어 모델을 비롯한 자연어 처리 기술을 활용하여 범주형 자료와 수치형 자료를 대상으로 지도학습과 비지도학습이 가능하게 되었을 뿐만 아니라(박조은, 2023) TF-IDF, Word2Vec 등 임베딩을 통해 의미를 수치화하여 언어학 연구를 진행할 수 있게 되었다(신서인, 2023).

이런 추세의 영향을 받아서 본 연구는 R를 통해 다변량 분석(Multivariate Analysis)의 방법을 활용하여 현대국어의 정도부사를 대상으로 계층적 군집분석(Hierarchical Clustering Analysis)과 다항현저공기분석 (Multiple Distinctive Collexeme Analysis)를 실시함으로써 이들의 군집 간의 특징과 군집 내부의 호환성(즉 교체 가능성)을 탐구하는 데에 목적을 둔다.

이를 위해 우선 2장에서 정도부사의 연구사를 크게 말뭉치 출현 이전의 연구와 말뭉치를 활용한 연구를 나눠서 논의한다. 선행연구를 통해 본 연구의 대상과 분석의 기준(즉 변량의 수준: 공공성, 상호작용성, 장르, 의미 운율 등)을 확정하고자 한다. 이어서 3장에서 연구 방법을 소개하고 연구용 말뭉치와 말뭉치에 추출된 결과를 간단하게 소개할 것이다. 제4장에서 추출된 용례를 대상으로 계층적 군집분석(Hierarchical Clustering Analysis), 다중 대응분석(Multiple Correspondence Analysis)과 다항현저공기분석(Multiple Distinctive Collexeme Analysis)의 결과를 해석하고 각 정도부사의 코로스토레이션 강도(collostruction strength)를 측정하여 비교 분석하자고 한다. 이를 바탕으로 군집 간의 특징과 군집 내부의 호환성(즉 교체 가능성)을 밝힐 것이다.

연구 결과, 우선, 텍스트 장르에 따른 정도부사의 사용 빈도와 분포적 특징을 계량적으로 분석하였다. 둘째, 군집화 분석을 통해 정도부사를 총 3개 유형을 범주화하였다. 셋째, 다중 대응분석의 방법으로 각 범주의 정도부사를 시각화하여 교체 가능성도 변량 수준을 참고하여 비교 분석하였다. 넷째, 정도부사의 공기어의 유형과 코로스토레이션 강도(collostruction strength)도 측정하여 비교 분석하였다. 본 연구는 R를 활용한 다변량 분석 방법론 탐구의 계단적인 성과로 큰 의미가 있는 것으로 생각된다.

키워드: 정도부사, 장르 문법, 다변량 분석(Multivariate Analysis), 계층적 군집 분석(HCA), 다중 대응분석(MCA), 코로스토레이션 강도(Collostruction Strength)

참고문헌

- 강범모(2000), <한국어 텍스트의 장르, 문체, 유형: 컴퓨터와 통계적 기법의 이용>, 태학사.
- 강범모(1999), <한국어의 텍스트 장르와 언어 특성>, 고려대학교 출판부.
- 강범일(2021), <신문 언어 변화 양상의 계량적 연구: 조선일보 1920-2019년 기사를 중심으로>, 연세대학교 언어정보학과 정 박사학위 논문.
- 강정훈(2002), <한국어의 정도부사 연구: 비격식적 입말의 정도부사를 중심으로>, 연세대학교 국어국문학과 석사학위 논문.
- 고언숙(2015), <정도부사 ‘너무’의 의미와 의미 운율 -음성 말뭉치 분석을 바탕으로->, 『음성음운형태론연구』 21(2), 223-243.
- 고혜원(2011), <정도부사 ‘너무’의 의미와 공기 특성>, 『어문론집』 46, 121-139.
- 공나형(2021), <구어 문법 관점을 반영한 ‘너무’의 의미 기술>, 『화법연구』 52, 1-30.
- 김영희(1985), <샘술말로서의 정도부사>, 『한글』 190, 133-162.
- 김하늘(2019), <강조적 정도부사와 용언의 연어 관계 연구: 구어와 문어의 차이를 중심으로>, 연세대학교 국어국문학과 석사학위 논문.
- 김혜영(2009), <구어에서 나타나는 정도부사의 사용 의미>, 고려대학교 이론언어학 전공 석사학위 논문.
- 김혜영·강범모(2010), <구어 속 강조적 정도부사의 사용과 의미>, 『한국어학』 48, 101-129.
- 박동근(2007), <현대국어 강조부사 범주의 설정과 기능적 분석>, 『한글』 275, 129-163.
- 박조은(2023), <모두의 한국어 텍스트 분석 with 파이썬>, 『한국 텍스트 언어학회 2023 여름 북 콘서트』, 1-51.
- 박효정(2019), <강조의 정도부사 ‘너무’의 사용 양상 연구>, 『한말연구』 54, 171-206.
- 배진영(2012), <구어와 문어 사용역에 따른 정도부사의 분포와 사용 양상에 대한 연구>, 『국제어문』 54, 95-140.
- 백재파(2021), <구어성 부사와 문어성 부사 목록 선정 연구>, 『영주어문』 49.
- 서상규(1991), <정도부사에 대한 국어학사적인 조명과 그 분류에 대해>, 『연세어문학』 48, 219-266.
- 서상규(1992), <16세기 국어의 말재어찌씨의 통어론적 연구>, 연세대학교 국어국문학과 박사학위 논문.
- 서상규(2005), <구어 말뭉치 설계의 몇 가지 조건>, 『언어정보와 사전편찬』 제14·15·16합집, 5-29.
- 서상규·안의정·봉미경·최정도·박종후·백해파·박종후·송재영·김선헤(2013), <한국어 구어 말뭉치 연구>, 한국문화사.
- 서상규·한영균(1999), <국어정보학 입문: 인문학과 컴퓨터>, 태학사.
- 서정수(2005), <한국어의 부사>, 서울대학교 출판부.
- 손남익(1995), <국어 부사 연구>, 박이정.
- 손세모돌(1998), <정도부사의 유표적 자질 분석>, 『한말연구』 4, 179-211.
- 손춘섭(2001), <정도부사의 의미와 기능에 대한 고찰>, 『한국어 의미학』 9, 97-130.
- 송지혜(2014), <코퍼스 분석을 통한 정도부사의 빈도, 분포, 공기 관계에 대한 연구>, 연세대학교 한국학협동과정 석사학위 논문.
- 신서인(2019), <말뭉치를 이용한 한국어 문형 및 어휘 연구의 현황과 과제>, 『한국어학』 83, 35-68.
- 신서인(2023), <의미는 과연 수치화될 수 있는가? -공기 관계와 의미 정보>, 『2023 세계 한국어 한마당』, 118-136.
- 신지연(2002), <정도부사의 범주화 기준에 대하여>, 『어문학』 78, 71-87.
- 안주호(2003), <한국어 구어에서 정도부사 ‘되게’에 대하여>, 『어문학』 24, 149-166.
- 연규동·박진호·최운호(2003), <인문학을 위한 컴퓨터>, 태학사.
- 유현경·서상규(2002), <한국어 학습자 말뭉치에 나타난 부사 사용에 대한 연구>, 『이중언어학』 20, 213-234.
- 이석규(1987), <현대국어 정도 어찌씨의 의미 연구>, 건국대학교 국어국문학과 박사학위 논문.
- 이용훈·조기현(2019), <A Distinctive Collexeme Analysis of -(u)m and -ki in Korean>, 『언어』 44, 901-921.
- 이유라(2021), <한국어 정보부사의 사용역별 출현 빈도 및 헤지성 분석>, 『언어학 연구』 60, 115-131.
- 이유미·김진식(2019), <코퍼스를 활용한 정도부사 연구—정도부사의 공기 관계 특징을 중심으로->, 『새국어교육』 120, 249-284.
- 이효정(1999), <한국어 정도부사의 분석과 응용>, 상명대학교 국어국문학과 석사학위 논문.
- 임채훈(2010), <정도부사의 의미 기능과 문장의미 구성>, 『어문연구』 38(3), 39-58.
- 임희석·고려대학교 자연언어처리연구실(2023), <자연어 처리 바이블: ChatGPT 핵심기술->, 휴먼싸이언스.
- 장음화(2021), <한국어 정보성 부사의 의미 연구>, 한국문화사.

- 전지은(2010), <성별에 따른 한국어 부사 사용 양상: 세종 구어 말뭉치를 활용하여->, 『언어와 언어학』 47, 191-217.
- 정성훈(2016), <2원 네트워크 분석을 이용한 한국어 강조 정도부사의 의미 거리 연구>, 『한국어 의미학』 52, 197-222.
- 제조현(2023a), <중국인 학습자를 위한 “중세한국어 문법과 자연언어 처리의 통합 교육” 방안 탐구: 데이터 기반 학습의 관점에서->, 『한국문법교육학회 제38차 전국학술대회 개인 발표문』.
- 제조현(2023b), <중국의 (한)국어사 교육을 위한 중세국어 의미빈도 사전(가칭)의 용례 색인 추출 방법과 기술 방향: Python과 성생형 AI를 활용하여->, 『한국사전학회 제42차 전국학술대회 개인 발표문』.
- 제조현(2023c), <디지털 인문학의 관점에서 본 <훈민정음> 교육내용의 선정과 지평의 확장 -중국의 한국어 학습자를 대상으로->, 『한국문화융합학회 2023년 여름 학술대회 개인 발표문』.
- 조민정(2006), <부사>, 『왜 다시 품사론인가』, 남기심 외 저음, 227-267.
- 조은경·드프레메리웨인(2023), <근대 잡지 <삼천리>의 의미 군집을 찾아서 -비지도 학습을 활용하여->, 『언어와 문화』 19(1), 203-224.
- 천수연·조태린(2020), <강조 정도 부사의 사용 양상에 대한 사회언어학적 연구>, 『담화와 인지』 제27권 4호, 145-168.
- 최정도(2022), <인문학 위기 속의 국어정보학의 현황과 전망>, 『국어국문학』 201, 31-62.
- 최현배(1937/1971), 『우리말본』, 정음사.
- 한길(1983), <정도 어찌씨에 관한 의미론적 연구>, 『새국어교육』 37(1), 372-391.
- 한길(2016), 『우리말의 통사 어찌씨 연구』, 역락.
- 한영균·서은아(2011), <유의적 정도부사의 빈도, 분포, 결합 관계의 분석과 그 활용 -학습자 사전의 용법 기술의 관점에서->, 『한국어 의미학』, 335-394.
- 홍사만(2002), <국어 정도부사의 피한정어 연구>, 『어문론총』 76, 153-175.
- 홍사만(2002), <국어 정도부사의 하위 분류>, 『어문론총』 36, 31-74.
- 홍윤표(2012), 『국어정보학』, 태학사.
- Biber, D., Conrad, S., & Reppen, R. (1998). *Corpus Linguistics: Investigating language structure and use*. Cambridge University Press. [유석훈·김유영(2018), 『(바이버의) 코퍼스 언어학: 언어구조와 용법에 관한 연구』, 고려대학교 출판문화원]
- Gries, S. T. (2009/2013/2021). *Statistics for Linguistics with R: A practical introduction*(2nd revised ed). Walter de Gruyter. [최재웅·홍정하(2013), 『언어학자를 위한 통계학: R활용』, 고려대학교 출판부.]
- Levshina, N. (2015). *How to do Linguistics with R: Data exploration and statistical analysis*. Amsterdam, Benjamins.

부록: HCA_of_Degree_Adverbs.R

```
install.packages(c("psych", "FactoMineR", "psych", "pvclust"))
library(Rling); library(psych); library(FactoMineR); library("factoextra"); library(pvclust)

nadv7<-read.csv(file.choose(),sep = ",",header = T,stringsAsFactors = TRUE)
nadv7_splitbykeynote <- split(nadv7, nadv7$keynote)
nadv7_splitbykeynote_q <- lapply(nadv7_splitbykeynote, function(x) x[, -1])
nadv7_splitbykeynote_new_matrix <- lapply(nadv7_splitbykeynote_q, bp)
nadv7_splitbykeynote_new_matrix[1]
nadv7_splitbykeynote_new_matrix[2]
nadv7_splitbykeynote_new_matrix[3]
nadv7_splitbykeynote_new_matrix[4]
nadv7_splitbykeynote_new_matrix[5]
nadv7_splitbykeynote_new_matrix[6]
nadv7_splitbykeynote_new_matrix[7]
nadv7_trans_pre <- do.call(rbind, nadv7_splitbykeynote_new_matrix)
```

```

cluster_number<-read.csv(file.choose(),sep = ",",row.names=1)
fviz_nbclust(cluster_number, kmeans, method = "gap_stat")
nadv7_dist <- dist(nadv7_trans_pre, method = "canberra")
round(nadv7_dist, 2)
nadv7_hclu <- hclust(nadv7_dist, method = "ward.D2")
plot(nadv7_hclu, hang = -1)
nadv7_pvclu <- pvclust(t(nadv7_trans_pre), method.hclust = "ward.D2", method.dist = "canberra")
plot(nadv7_pvclu, hang = -1)
rect.hclust(nadv7_hclu, k = 3)

c1 <- nadv7_trans_pre[c(4,5),]
c2 <- nadv7_trans_pre[-c(4,5),]
c1_mean <- colMeans(c1)
c2_mean <- colMeans(c2)
diff <- c1_mean-c2_mean
sort(diff, decreasing = TRUE)
plot(sort(diff)*1.2, 1:length(diff), type = "n", xlab = "cluster 2 <----> cluster 1", yaxt = "n", ylab = "")
text(sort(diff), 1:length(diff), names(sort(diff)))

```

AI-Based New Approach to Polysemy: BERT and t-SNE

October 07. 2023

Wonbin Kim
kny1876@gmail.com

Image Source: <https://medium.com/analytics-vidhya/a-brief-introduction-to-deep-learning-cfd901d4611>

Research Objective

What is Artificial Intelligence (AI)?

- intelligence displayed by machines
- AI performs intelligent tasks such as learning, decision-making, and judgment.



Image Source: <https://www.geeksforgeeks.org/artificial-intelligence-an-introduction/>

The purpose of this study is

- (1) to examine the potential of AI-based methods as an approach to polysemy research and present how AI methods can be employed for the analysis of polysemy
- (2) to investigate whether the usage of polysemous words varies depending on the amount of language exposure.

► BERT & t-SNE

► Difference in the usage of English polysemous words depending on English proficiency levels

2

Background: BERT

► Bidirectional Encoder Representations from Transformers (BERT)

: a pre-trained natural language processing model that uses a transformer architecture to understand and generate human language text

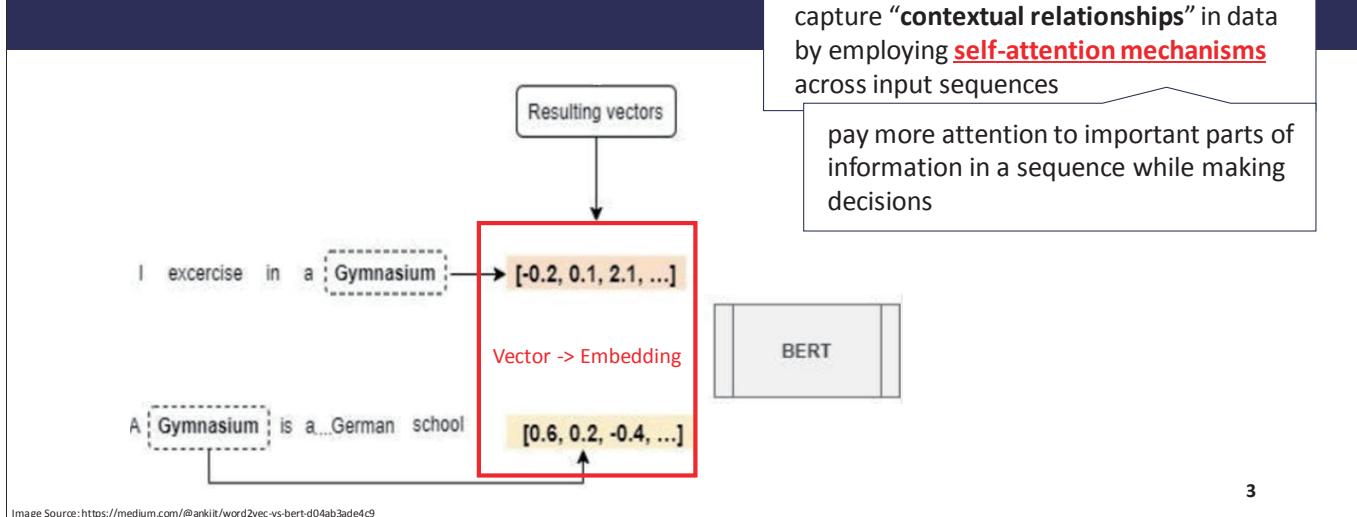
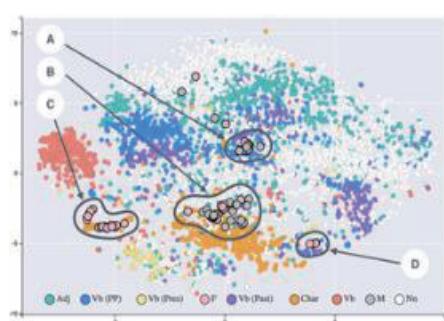


Image Source: <https://medium.com/@ankit/word2vec-vs-bert-d04ab3ade4c9>

Background: t-SNE

► t-distributed Stochastic Neighbor Embedding (t-SNE)

: a statistical method for visualizing high-dimensional data by giving each datapoint a location in a two or three-dimensional map



Similar objects are modeled by nearby points and dissimilar objects are modeled by distant points.

Image Source: https://en.wikipedia.org/wiki/T-distributed_stochastic_neighbor_embedding

4

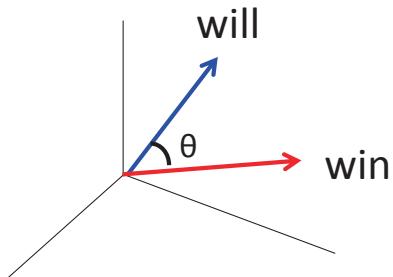
Background: Cosine Similarity



► **Cosine Similarity (CS)**: a metric used to measure similarity between vectors

It uses the cosine of the angle between vectors to show similarity between vectors.
Similarity between vectors indicates semantic similarity between words.

Cosine Similarity



The closer to “1” the CS value is, the more semantically similar the two vectors are .

$$0 \leq CS \leq 1$$

5

Analysis of Polysemy



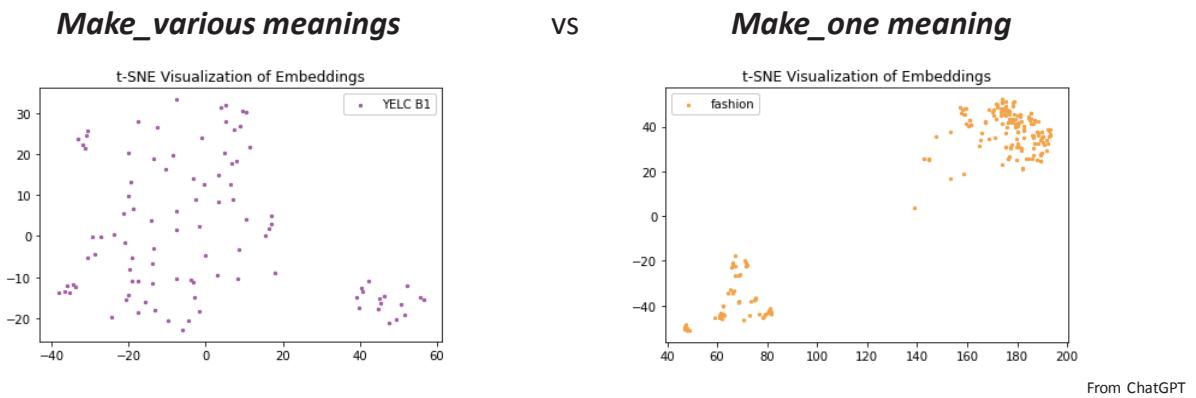
- ◆ Analysis methods for addressing 5 research questions
- ◆ Analysis results
- ◆ Validation

6

Polysemous and Monosemous Words in t-SNE



RQ1) How is a polysemous word represented in t-SNE, compared to a monosemous word?



- Embeddings from a polysemous word
cluster together because of their related meanings
but they have a scattered distribution and their density is low.

7

Difference in Usage of *Make* among Proficiency Levels



RQ2) Does the usage of polysemous words vary according to levels of language performance?



STEP1: Yonsei English Learner Corpus (YELC) consisting of essays produced by 9 groups classified by proficiency -> 6 groups (A1+, A2, B1, B1+, B2, B2+)

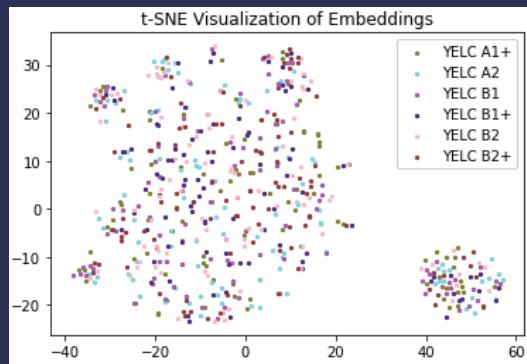
STEP2: Selecting the target word (the polysemous word make)

STEP3: Obtaining embeddings of make from each group through BERT

STEP4: Visualizing the embeddings from each group through t-SNE

8

Visualization through t-SNE



Each color does not cluster. It's hard to find clear clusters.

- ▶ The way *make* is used is similar across all the groups.

Validation using Cosine Similarity

| Groups | Cosine Similarity |
|-------------|-------------------|
| A1+ and A2 | 0.3742284 |
| A1+ and B1 | 0.36883748 |
| A1+ and B1+ | 0.38360217 |
| A1+ and B2 | 0.38303193 |
| A1+ and B2+ | 0.383949 |
| A2 and B1 | 0.37018156 |
| A2 and B1+ | 0.38895288 |
| A2 and B2 | 0.38777336 |
| A2 and B2+ | 0.3875516 |
| B1 and B1+ | 0.38300008 |
| B1 and B2 | 0.38249967 |
| B1 and B2+ | 0.38271305 |
| B1+ and B2 | 0.41035953 |
| B1+ and B2+ | 0.41041532 |
| B2 and B2+ | 0.4088401 |

CS values among the groups are similar to each other.

- ▶ The way *make* is used is similar across all the groups.
- ▶ The CS values support the result from t-SNE.

Validation using Collocation Analysis

| A1+ | A2 | B1 | B1+ | B2 | B2+ |
|-------------|----------|-----------|-----------|----------|--------------|
| make | people | student | feel | student | mistake |
| better | make | student | feel | better | better |
| le | student | good | people | feel | good |
| problem | happy | make | better | happy | feel |
| trouble | accident | feel | good | bad | people |
| world | smoking | better | bad | internet | rule |
| say | child | happy | driver | mistake | sense |
| harm | problem | child | make | child | place |
| beautiful | u | driver | accident | many | serious |
| information | good | many | smoking | make | student |
| many | bad | friend | rule | building | problem |
| food | also | bad | teacher | friend | u |
| going | many | think | physical | u | bad |
| friend | rule | smoking | also | teacher | relationship |
| happy | feel | building | u | real | class |
| people | smoke | rule | using | smoking | hurt |
| using | room | internet | public | think | trouble |
| good | physical | mistake | many | thing | information |
| car | using | problem | mistake | better | safe |
| student | think | like | like | get | also |
| accident | society | use | happy | online | stop |
| real | public | smoker | smoker | driver | life |
| need | friend | accident | internet | one | however |
| like | internet | person | think | good | make |
| school | mistake | teacher | problem | smoker | give |
| want | driver | law | use | problem | point |
| thing | law | u | lot | smoke | go |
| teacher | building | physical | sense | name | thing |
| internet | better | punishmen | building | hurt | might |
| punishmen | like | right | punishmen | worse | accident |
| think | animal | real | right | physica | know |

※ The most typical collocates of the node word characterize the meaning and usage of the node word (Firth 1957).

► The change of collocates means semantic change of the node word and the way it is used.

[Top 31 collocates from each group]

Most of the top collocates from each group overlap with each other.

► Each group uses make in a similar way.

11

BERT vs Word2Vec

| Groups | CS from BERT | CS from Word2Vec |
|-------------|--------------|------------------|
| A1+ and A2 | 0.3742284 | 0.42570215 |
| A1+ and B1 | 0.36883748 | 0.334051 |
| A1+ and B1+ | 0.38360217 | 0.40162086 |
| A1+ and B2 | 0.38303193 | 0.3277553 |
| A1+ and B2+ | 0.383949 | 0.7154199 |
| A2 and B1 | 0.37018156 | 0.19682573 |
| A2 and B1+ | 0.38895288 | 0.4125179 |
| A2 and B2 | 0.38777336 | 0.22487749 |
| A2 and B2+ | 0.3875516 | 0.23782061 |
| B1 and B1+ | 0.38300008 | 0.1860776 |
| B1 and B2 | 0.38249967 | 0.21595147 |
| B1 and B2+ | 0.38271305 | 0.39780685 |
| B1+ and B2 | 0.41035953 | 0.17115071 |
| B1+ and B2+ | 0.41041532 | 0.27655506 |
| B2 and B2+ | 0.4088401 | 0.29395854 |

RQ3) Between BERT and Word2Vec, which shows better performance?

- Bigger difference in the values of CS among the groups
- Value of CS between A1+ and B2+

► BERT shows better performance concerning polysemy research.

12

Use of *Make* among English-Language Learners

RQ4) How do English-Language Learners use *make*?



STEP1: Constructing a corpus containing six meanings of *make* using ChatGPT

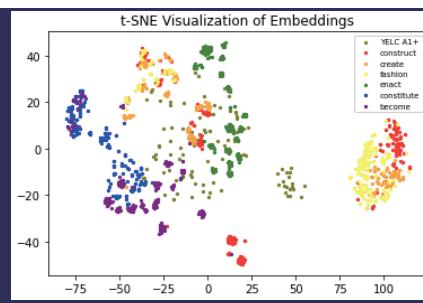
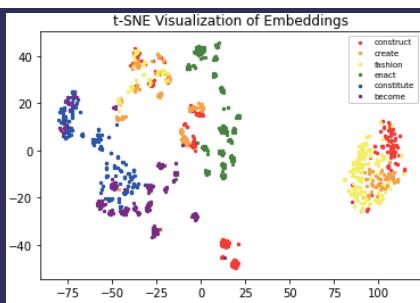
-> 6 meanings of *make* ('construct'/'create'/'fashion'/'enact'/'constitute'/'become')

STEP2: Obtaining embeddings of *make* in terms of each meaning through BERT

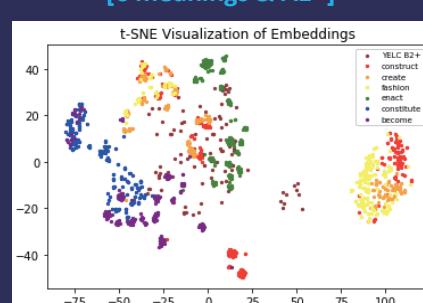
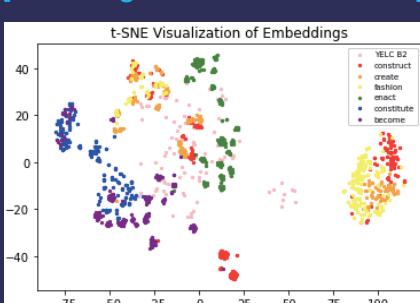
STEP3: Visualizing the embeddings from ChatGPT and YELC (A1+&B2&B2+) through t-SNE

13

Visualization through t-SNE



- ▶ Meanings:
 - 'enact'
 - 'constitute'
 - 'become'



14

Validation using Collocation Analysis

| A1+ | A2 | B1 | B1+ | B2 | B2+ |
|-------------|----------|-----------|-----------|----------|--------------|
| make | people | people | student | feel | student |
| better | make | student | feel | people | mistake |
| le | student | good | people | feel | better |
| problem | happy | make | better | happy | good |
| trouble | accident | feel | good | bad | feel |
| world | smoking | better | bad | internet | people |
| say | child | happy | driver | mistake | sense |
| harm | problem | child | make | child | place |
| beautiful | u | driver | acciden | many | serious |
| information | good | many | smoking | make | student |
| many | bad | friend | rule | building | problem |
| food | also | bad | teacher | friend | u |
| going | many | think | physica | u | bad |
| friend | rule | smoking | also | teacher | relationship |
| happy | feel | building | u | real | class |
| people | smoke | rule | using | smoking | hurt |
| using | room | internet | public | think | trouble |
| good | physical | mistake | many | thing | information |
| car | using | problem | mistake | better | safe |
| student | think | like | like | get | also |
| accident | society | use | happy | online | stop |
| real | public | smoker | smoker | driver | life |
| need | friend | accident | internet | one | however |
| like | internet | person | think | good | make |
| school | mistake | teacher | problem | smoker | give |
| want | driver | law | use | problem | point |
| thing | law | u | lot | smoke | go |
| teacher | building | physical | sense | name | thing |
| internet | better | punishmen | building | hurt | might |
| punishment | like | right | punishmen | worse | accident |
| think | animal | real | right | physical | know |

[Inferring meanings of *make* from Top 31 collocates from A1+, B2, & B2+]

punishment, rule...

► 'enact'

problem, trouble...

► 'constitute'

better, good, bad, happy...

► 'become'

15

Difference between Learners and Native Speakers

RQ5) Is there difference between English-language learners and native English speakers in the use of *make*?

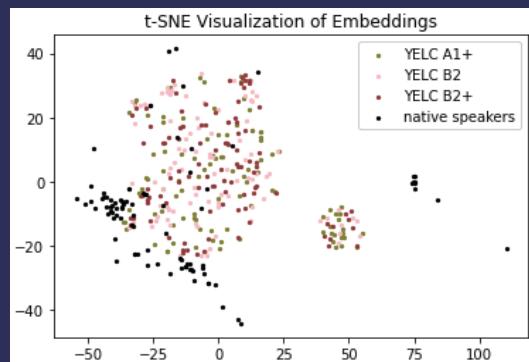


STEP1: Obtaining embeddings of *make* from a corpus of data produced by native English speakers through BERT

STEP2: Visualizing the embeddings of *make* from native English speakers and YELC (A1+&B2&B2+) through t-SNE

16

Visualization through t-SNE



Most of the points from native speakers do not overlap with those from learners.

- ▶ Native speakers use *make* differently from learners.

17

Validation using Collocation Analysis

| | A1+ | A2 | B1 | B1+ | B2 | B2+ | Native |
|-------------|----------|-----------|-----------|----------|-------------|------------|--------|
| make | people | people | student | student | mistake | make | |
| better | make | student | feel | people | better | sure | |
| le | student | good | people | feel | good | way | |
| problem | happy | make | better | happy | feel | one | |
| trouble | accident | feel | good | bad | people | feel | |
| world | smoking | better | bad | internet | rule | sense | |
| say | child | happy | driver | mistake | sense | decision | |
| harm | problem | child | make | child | place | new | |
| beautiful | u | driver | accident | many | serious | would | |
| information | good | many | smoking | make | student | clear | |
| many | bad | friend | rule | building | problem | come | |
| food | also | bad | teacher | friend | u | u | |
| going | many | think | physical | u | bad | public | |
| friend | rule | smoking | also | teacher | relations | like | |
| happy | feel | building | u | real | class | available | |
| people | smoke | rule | using | smoking | hurt | difference | |
| using | room | internet | public | think | trouble | year | |
| good | physical | mistake | many | thing | information | choice | |
| car | using | problem | mistake | better | safe | want | |
| student | think | like | like | get | also | good | |
| accident | society | use | happy | online | stop | said | |
| real | public | smoker | smoker | driver | life | love | |
| need | friend | accident | internet | one | however | work | |
| like | internet | person | think | good | make | time | |
| school | mistake | teacher | problem | smoker | give | move | |
| want | driver | law | use | problem | point | say | |
| thing | law | u | lot | smoke | go | two | |
| teacher | building | physical | sense | name | thing | even | |
| internet | better | punishmen | building | hurt | might | could | |
| punishmen | like | right | punishmen | worse | accident | life | |
| think | animal | real | right | physical | know | people | |

[Comparison of Top 31 collocates of *make* from Native speaker Corpus and YELC]

The number of collocates which do not overlap with YELC is 21, suggesting that

- ▶ The way native English speakers use *make* is different from English-language learners.
- ▶ This raises the possibility that the amount of language exposure may affect the usage of polysemous words.

18

Conclusion

- ◆ Limitations of the present study
- ◆ Further study

19

Limitations

▶ Limitations of the present study

- 1) It is hard to generalize results from the present study.
 - Analysis of one word -> frequently used word (*make*): early acquisition
 - The size and composition of YELC
- 2) This study superficially examined the potential of using BERT and t-SNE for polysemy research.

20

Further Study



▶ Further Study

- Analyzing diverse polysemous words
- Using a larger general English learner corpus with significant differences among proficiency levels
- Employing additional language models to compare performance
- Designing specific methods to explore various aspects of polysemy by means of BERT and t-SNE

A Generative AI Study on Translation Universals

Jin Yim

Graduate School of Translation & Interpretation
Ewha Womans University

Yong-hun Lee

Department of Linguistics
Chungnam National University

Table of contents

01

Introduction

02

Literature
Review

03

Method

04

Concepts

05

Text Analysis

06

Closing

01

Introduction

What is generative AI?

- Generative artificial intelligence or generative AI is a type of artificial intelligence (AI) system capable of generating text, images, or other media in response to prompts.
- Generative AI models learn the patterns and structure of their input training data, and then generate new data that has similar characteristics.
- Notable generative AI systems include ChatGPT (and its variant Bing Chat), a chatbot built by OpenAI using their GPT-3 and GPT-4 foundational large language models, and Bard, a chatbot built by Google using their LaMDA foundation model.
- Other generative AI models include artificial intelligence art systems such as Stable Diffusion, Midjourney, and DALL-E.
- Generative AI has potential applications across a wide range of industries, including art, writing, software development, product design, healthcare, finance, gaming, marketing, and fashion.

Motivation

- Translation, one of the key functions of generative AI, has drawn huge research interest, particularly in the areas of:
 - Evaluation on translation output from generative AI:
 - ✓ Jiao et al. (2023); Stap and Araabi (2023)
 - Prompt engineering for translation quality improvement:
 - ✓ Gao et al. (2023); Gu (2023); ; Peng et al. (2023)
 - Generative AI's translation evaluation function
 - ✓ Lu et al. (2023)

... but little has been known about what generative AI knows about the linguistic features of translation

02 Literature Review

Translation Universal in HT

- Since Baker (1996), TUs have been empirically confirmed...
 - Explication (Baker 1996; Xiao 2010)
 - Simplification (Blum-Kulka & Levenston, 1983; Kajzer-Wietrzny 2015)
 - Normalization (Baker 1996; Olohan 2004)
 - Levelling out (Baker 1996; Laviosa 1998)
- ...or rejected:
 - TUs are ““still confined to a few languages and genres (mainly literary translation)” (Zanettin, 2014; 22)
- But “research into potential translation universals is not only a question of establishing whether they exist or not, and if so, under what conditions” (Chesterman, 2010: 42)

TU in English translated from Korean

| Genre | Simplification | Explication | Normalization | Levelling-out |
|--------------------|-------------------------------------|--|-------------------------------------|-------------------|
| Academic prose | Goh et al. (2016); Lee (2016; 2018) | Goh, Lee, & Kim (2016); Lee (2016; 2018) | Lee (2016; 2018); Goh et al. (2016) | Goh et al. (2016) |
| Newspaper | 고광윤 & 이영희 (2016) | 고광윤 & 이영희 (2016); Park & Lee (2014) | 고광윤 & 이영희 (2016) | 고광윤 & 이영희 (2016) |
| Literature | 이창수 (2014) | 이창수 (2014); 최희경 (2016) | 최희경 (2016) | - |
| Business reporting | 임진 (2019) | 임진 (2019) | 임진 (2019); Lee & Yim (2019) | - |

Translation Universal in MT

- In the MT literature, TUs have been studied in two aspects:
 - Translationese in MT training data improves MT performance
 - ✓ Graham et al. (2019); Zhang and Toral (2019)
 - Outputs of MT have TU features just as in HT
 - ✓ Bizzoni et al. (2020); Luo and Li (2022)
 - ✓ 박미정 (2023): Simplification and explication observed in Chat GPT translation output

03 Method

Research Questions

- Q1. Do generative AIs understand major concepts on TU?
- Q2. Can generative AIs use TU to pick out translated texts from non-translated texts?
- Q3. Does context information in prompts affect the result?

Settings

- Three types of generative AIs
 - Chat GPT 3.5
 - (Microsoft) Bing
 - (Google) Bard
- ✓ All data outputs were collected as of Sept. 6 to 23, 2023

Prompt engineering

- Q1: Do generative AIs understand the concepts of TU?
 - Translationese
 - TU
 - ✓ Explication
 - ✓ Simplification
 - ✓ Normalization
 - ✓ Levelling-out
- Q2: Can generative AIs distinguish TTs from Non-TTs?

Data

- Management forewords of sustainability reports
 - TTs with proven quality show normalization with lower frequencies of simplification and explication features as compared to comparable NTTs (임진, 2019; Lee & Yim, 2019)

| Industry | TT | NTT |
|----------------|---------------------|---------------------|
| IT | Samsung Electronics | Apple |
| Finance | KB Finance | JP Morgan |
| Pharmaceutical | Yuhna | Johnson and Johnson |
| Manufacturing | Hankook Tire | Chevron |

04

Result #1: Concepts

1 Translationese

■ Differences between TT and NTT (Chat GPT)

창작된 문장과 번역본 문장 사이에는 차이가 있는데, 그 차이가 무엇인지 알려줘.

문장의 창작과 번역본 사이의 주요 차이는 언어와 문화적인 요소입니다. 아래는 이러한 차이점 을 설명한 것입니다:

- 언어적 차이:
 - 창작된 문장은 원래 언어에서 작성되며, 번역본 문장은 다른 언어로 번역되었습니다. 이로 인해 언어의 문법, 어휘, 표현 방식 등에 차이가 있을 수 있습니다. 번역 과정에서 어떤 어휘를 선택하고, 어떤 문법 규칙을 따를지에 따라 번역된 문장은 원래 문장과 약간 다를 수 있습니다.
- 문화적 차이:
 - 각 언어와 문화는 고유한 표현 방식과 문화적 의미를 갖고 있습니다. 따라서 원래 문장과 번역본 문장은 각 언어와 문화의 특성을 반영합니다. 이로 인해 문화적인 차이가 번역된 문장에 반영될 수 있으며, 때로는 번역된 문장이 원래의 문장의 의도나 감정을 완벽하게 전달하지 못할 수도 있습니다.
- 뉘앙스와 감정 전달:
 - 원래 문장은 작성자의 원래 의도와 감정을 가지고 있으며, 번역본 문장은 이러한 뉘앙스와 감정을 완벽하게 전달하기 어려울 수 있습니다. 특히 언어 간의 미묘한 차이나 언어적 특징 때문에 번역된 문장은 뉘앙스를 손실하거나 왜곡할 수 있습니다.
- 목적과 대상 관객:
 - 창작된 문장과 번역본 문장은 종종 서로 다른 목적과 대상 관객을 대상으로 합니다. 번역된 문장은 다른 언어를 사용하는 독자나 청취자를 대상으로 하므로, 그에 따라 문장을 조정하여 더 이해하기 쉽게 만들어야 할 때가 있습니다.

요약하면, 창작된 문장과 번역본 문장 사이의 차이는 언어, 문화, 뉘앙스, 감정, 목적 등 다양한 요소에 기인합니다. 이러한 차이는 번역 작업이 항상 완벽하게 원래의 문장을 복제하는 것이 어렵다는 것을 의미합니다.

1 Translationese

■ Differences between TT and NTT (Bing)

장작된 문장과 번역본 문장 사이에는 차이가 있는데, 그 차이가 무엇인지 알려줘.

✓ "장작된 문장과 번역본 문장의 차이" 검색 중
✓ 답변을 생성하는 중...

장작된 문장은 언어학에서 완결된 의미를 표현하기 위해 의미적 관계를 이루는 요소들이 결합된 것이고 ①, 번역본 문장은 다른 언어로 표현된 문장을 대상 언어의 문법과 어휘에 맞게 옮긴 것이다. 장작된 문장과 번역본 문장의 차이는 다음과 같다.

- 번역본 문장은 원문의 의미와 표현을 최대한 존중하려고 하기 때문에, 대상 언어의 문화나 맥락에 적합하지 않은 표현이 있을 수 있다.
- 번역본 문장은 원문의 문장 성분 ②과 구조를 따르려고 하기 때문에, 대상 언어의 문장 짜임 ③과 양상에 맞지 않는 문장이 있을 수 있다.
- 번역본 문장은 원문의 어휘와 의미를 정확하게 전달하려고 하기 때문에, 대상 언어의 뉴앙스나 미묘한 차이를 잘 반영하지 못하는 경우가 있을 수 있다.

장작된 문장과 번역본 문장의 차이에 대해 더 알고 싶다면, 다음 링크를 참고해 보세요.

①: 문장 (언어학) - 위키백과, 우리 모두의 백과사전 ②: 한국어의 문장 성분 - 위키백과, 우리 모두의 백과사전 ③: <온빛중국어3-1> 2-1 문장의 짜임과 양상 : 네이버 블로그

자세한 정보: 1. ko.wikipedia.org 2. ko.wikipedia.org 3. blog.naver.com 1 / 30

④ 번역의 어려움에 대해 더 알고 싶어요. 외국인들이 한국어를 배울 때 번역을 많이 사용하나요? 장작된 문장과 번역본 문장 중 어떤 것이 더 좋은 표현일까요?

1 Translationese

■ Differences between TT and NTT (Bard)

장작된 문장과 번역본 문장 사이에는 차이가 있는데, 그 차이가 무엇인지 알려줘.

장작된 문장은 작자의 독창적인 생각이나 감정이 담겨 있는 반면, 번역본 문장은 원문의 의미를 전달하는 데 중점을 둡니다. 장작된 문장은 작자의 상상력과 표현력을 통해 독창적인 표현이나 의미를 담고 있음을 알 수 있습니다. 반면, 번역본 문장은 원문의 의미를 최대한 살리기 위해 원문의 표현과 그대로 따라가고 합니다.

장작된 문장은 원문과는 다른 형식을 가질 수 있습니다. 장작자는 원문의 형식을 그대로 따르지 않고, 새로운 형식을 통해 독창적인 표현을 시도할 수 있습니다. 예를 들어, 시를 번역할 때 원문의 운율을 그대로 살리기보다, 한국어의 운율에 맞게 표현하는 경우가 많습니다.

장작된 문장은 원문과는 다른 분위기나 느낌을 가질 수 있습니다. 장작자는 원문의 분위기와 느낌을 그대로 전달하기보다, 자신의 독특한 감성을 달아 표현할 수 있습니다. 예를 들어, 소설을 번역할 때 원문의 분위기를 그대로 살리기보다, 한국인의 감성에 맞게 표현하는 경우가 많습니다.

들은, 이러한 차이는 모든 장작된 문장과 번역본 문장 사이에는 존재합니다. 장작된 문장 중에서도 원문의 의미를 충실히 전달하는 경우도 있고, 번역본 문장 중에서도 장작자의 독창적인 표현이 빛보이는 경우도 있습니다. 그러나 일반적으로 장작된 문장과 번역본 문장 사이에는 이러한 차이가 존재한다고 말 수 있습니다.

다음은 장작된 문장과 번역본 문장의 차이를 보여주는 예시입니다.

원문: I love you.

장작된 문장: 네가 나의 전부야.

번역본 문장: 나는 당신을 사랑합니다.

위의 예시에서 장작된 문장과 번역본 문장은 모두 원문의 의미를 전달하고 있지만, 표현 방식에 차이가 있습니다. 장작된 문장은 "네가 나의 전부야"라는 표현을 통해 사랑하는 사람에 대한 강렬한 감정을 표현하고 있습니다. 반면, 번역본 문장은 "나는 당신을 사랑합니다"라는 표현을 통해 사랑하는 사람에 대한 존경과 애정을 표현하고 있습니다.

이처럼 장작된 문장과 번역본 문장은 원문의 의미를 전달하는 데 중점을 두느냐, 아니면 장작자의 독창적인 표현을 살리느냐에 따라 차이가 나타날 수 있습니다.

2 Translationese

■ TU (Chat GPT)

'번역투'라는 것이 무엇이니?



'번역투(Translationese)'란 번역된 텍스트나 문장에서 생기는 특별한 언어적 특징 또는 어휘 사용 패턴을 가리키는 용어입니다.

번역투는 번역 과정에서 원래 언어에서 대상 언어로 정보를 전달하려고 노력하면서 발생하는 현상으로, 원래 언어와 대상 언어 간의 문법, 어휘, 문체 등에서 생길 수 있습니다.

번역투의 특징은 다음과 같습니다:

1. 어휘 및 표현의 제약: 번역된 문장에서는 원래 언어의 표현을 대상 언어에 최대한 가깝게 유지하려는 경향이 있습니다. 이로 인해 원래 언어에는 없는 표현이나 어휘가 사용되거나, 대상 언어에서는 자연스럽지 않은 표현이 나타날 수 있습니다.
2. 문법적인 어색함: 번역된 문장에서는 문법적으로 어색하거나 이해하기 어려운 구조가 나타날 수 있습니다. 원래 언어의 구조와 대상 언어의 구조가 다르기 때문에 발생할 수 있습니다.
3. 문체 및 어조의 불일치: 번역투는 문체나 어조의 불일치를 야기할 수 있습니다. 원래 언어의 문체나 어조가 대상 언어에서 적합하지 않을 경우 이러한 불일치가 나타날 수 있습니다.
4. 뉘앙스의 손실: 번역투는 종종 뉘앙스와 감정을 전달하는 데 어려움을 겪습니다. 원래 언어에서의 뉘앙스가 번역 과정에서 손실되거나 왜곡될 수 있습니다.

물론 번역 작업자의 능력과 노력에 따라 번역투의 정도는 다를 수 있습니다. 좋은 번역은 원래의 의미와 뉘앙스를 최대한 보존하면서도 대상 언어의 자연스러운 표현을 사용하는 것을 목표로 합니다. 그러나 번역 과정에서 번역투가 발생할 수 있으며, 이를 줄이고 품질 높은 번역을 제공하기 위해서는 번역 작업의 문맥과 목적을 고려하는 것이 중요합니다.

2 Translationese ■ TU (Bing)

번역투라는 것이 무엇이니?

✓ "번역투" 검색 중

✓ 답변을 생성하는 중...

'번역투'는 어떤 글에 원본이 아닌 번역이라는 흔적이 일정하게 반복적으로 출현하는 텍스트상의 특성을 말해 [这里], 번역투는 다음과 같은 특징을 가진다.

- 원문의 문화나 맥락에 적합하지 않은 표현이 있을 수 있다.
- 원문의 문장 구조나 어순을 그대로 따르기 때문에, 대상 언어의 문법이나 용법에 어긋날 수 있다.
- 원문의 어휘나 의미를 정확하게 전달하려고 하기 때문에, 대상 언어의 뉘앙스나 미묘한 차이를 잘 반영하지 못할 수 있다.

번역투는 좋은 문장을 방해하기 때문에, 가능한 피해야 한다. 번역투를 피하는 방법은 다음과 같다.

- 원문의 의미와 표현을 존중하되, 대상 언어의 문화나 맥락에 맞게 적절하게 변형하거나 보완한다.
- 원문의 문장 구조나 어순을 그대로 따르지 않고, 대상 언어의 문법이나 용법에 맞게 재구성한다.
- 원문의 어휘나 의미를 정확하게 전달하되, 대상 언어의 뉘앙스나 미묘한 차이도 고려한다.

번역투에 관한 더 자세한 설명과 예시는 다음 링크를 참고해 보세요.

[1] 일한 번역에 나타난 번역투 문장 - 국립국어원 [2] 영한 번역에 나타난 번역투 문장 - 국립국어원 번역 투 정리(영어, 일본어) : 네이버 블로그 [책쓰기 & 글쓰기]에서 보이는 번역투 TOP 12, 이것만 고쳐도... 좋은 문장을 방해하는 번역투 표현 12개를 끊는다면? - 브런지 [책쓰기 & 글쓰기] 혼한 번역투 TOP 12 - 한빛출판네트워크 - hanbit.co.kr

자세한 정보: 1. korean.go.kr 2. korean.go.kr 3. blog.naver.com 4. post.naver.com

2 / 30 ●

① 번역투를 피하는 방법은 무엇인가요? 와국어로 글을 쓰는데 번역투에 걸리지 않으려면 어떻게 해야 하나요? 한국어에서도 번역투가 발생하나요?

2 Translationese ■ TU (Bard)

질문 '번역투'라는 것이 무엇이니?

다른 답변 보기 ▾

'번역투'란 원문의 의미를 전달하기 위해 원문의 구조나 표현을 그대로 따라가다 보니, 목표언어의 어법이나 문화에 맞지 않는 부자연스러운 문체나 표현을 사용하는 것을 말합니다.

번역투는 다음과 같은 특징을 가지고 있습니다.

- 원문의 구조나 표현을 그대로 따라간다. 예를 들어, 영어에서 "I love you."를 "나는 당신을 사랑합니다."라고 번역할 때, 영어의 주어-동사-목적어 구조를 그대로 따라 한국어로 번역하는 경우가 많습니다.
- 원문의 어휘를 그대로 사용한다. 예를 들어, 영어에서 "computer"를 "컴퓨터"라고 번역할 때, 영어의 어휘를 그대로 사용해 한국어로 번역하는 경우가 많습니다.
- 목표언어의 어법이나 문화에 맞지 않는 표현을 사용한다. 예를 들어, 영어에서 "in spite of"를 "그럼에도 불구하고"라고 번역할 때, 한국어의 어법이나 문화에 맞지 않는 표현을 사용하는 경우가 많습니다.

번역투는 번역 과정에서 불가피하게 발생할 수 있는 현상이지만, 과도하게 사용되면 원문의 의미를 정확하게 전달하지 못할 뿐만 아니라, 독자의 이해를 방해하기도 합니다. 따라서 번역자는 번역투를 최소화하기 위해 노력해야 합니다.

번역투를 줄이기 위한 방법은 다음과 같습니다.

- 원문의 의미를 정확하게 이해한다. 원문의 의미를 정확하게 이해하지 못하면, 원문의 구조나 표현을 그대로 따라가다 보니 번역투가 발생할 수 있습니다.
- 목표언어의 어법이나 문화를 이해한다. 목표언어의 어법이나 문화를 이해하지 못하면, 목표언어의 어법이나 문화에 맞지 않는 번역투가 발생할 수 있습니다.
- 원문의 구조나 표현을 유연하게 변형한다. 원문의 의미를 훼손하지 않는 범위에서 원문의 구조나 표현을 유연하게 변형하면, 번역투를 줄일 수 있습니다.
- 목표언어의 어휘를 적절하게 사용한다. 목표언어의 어휘를 적절하게 사용하면, 번역투를 줄일 수 있습니다.

번역투를 줄이기 위해서는 번역자의 노력뿐만 아니라, 번역 교육과 번역 경기에서도 번역투를 줄이는 데 중점을 두어야 할 것입니다.

...

3 TU ■ TU (Chat GPT)

질문 '번역보편소(translation universals)라는 것이 무엇이니?

'번역보편소(translation universals)'는 다양한 언어 간 번역 현상에서 공통적으로 나타나는 언어학적 패턴 또는 규칙을 가리키는 개념입니다. 이러한 보편적인 번역 현상은 다양한 언어와 문화에서 발견되며, 번역 연구에서 언어 간의 공통점과 차이점을 이해하는 데 도움을 주는 중요한 개념 중 하나입니다.

번역보편소의 몇 가지 예시는 다음과 같습니다:

- 단어 대 단어 번역: 많은 번역에서는 원래 언어의 단어가 대상 언어에서 대응하는 단어로 번역됩니다. 이러한 단어 대 단어 번역은 번역 과정에서 자주 발생하는 현상 중 하나입니다.
- 구문 구조의 보존: 번역된 문장은 종종 원래 언어의 구문 구조를 최대한 유지하려고 노력합니다. 이는 번역된 문장이 대상 언어에서도 자연스럽게 읽히도록 하는 데 도움이 됩니다.
- 문체 조절: 번역된 텍스트는 대상 언어의 문제와 어조에 맞게 조절될 수 있습니다. 이는 번역된 문장이 대상 독자나 청취자에게 더 적합하게 전달되도록 하는 데 도움이 됩니다.
- 문화적 적용: 번역된 문장은 종종 원래 언어와 대상 언어의 문화적 차이를 고려하여 적용됩니다. 특정 문화적 표현이 다른 문화에서는 의미가 없거나 이해하기 어려울 수 있으므로 번역 과정에서 조정됩니다.

번역보편소는 언어 간 번역에서 공통적으로 나타나는 언어학적 특징을 연구함으로써 번역 이론과 실제 번역 작업에 도움을 주는 중요한 연구 주제 중 하나입니다. 그러나 주의해야 할 점은 모든 번역 작업에서 이러한 보편적인 패턴이 항상 적용되지는 않으며, 번역 문맥과 목적에 따라 다를 수 있다는 것입니다.

3 TU

■ TU (Bing)

'번역보편소 (translation universals)라는 것이 무엇이니?

- 번역보편소 (translation universals)는 번역된 텍스트에만 나타나는 언어적 특성을 말해 [❶❷❸](#). 번역보편소는 번역 과정에서 사용된 특정 언어방과는 무관하게, 모든 번역 텍스트에 공통적으로 나타난다고 생각되는 현상이야 [❶❷❸](#). 번역보편소는 다음과 같은 것들이 있어.
- 명시화(explication): 번역 텍스트에서 원문보다 더 많은 정보를 제공하는 것
- 단순화(simplification): 번역 텍스트에서 원문보다 더 쉬운 어휘나 구조를 사용하는 것
- 정규화(normalization): 번역 텍스트에서 원문보다 더 보편적이거나 표준적인 언어 형태를 사용하는 것
- 통화(convergence): 번역 텍스트들이 서로 비슷한 언어적 특징을 가지는 것

번역보편소는 번역의 일반적인 법칙을 찾기 위한 연구의 한 분야로, 아직 확실하게 일증되지 않은 가설이기도 해 [❶❷❸](#). 번역보편소에 관한 더 자세한 설명과 연구는 다음 링크를 참고해 보세요.

[❶ Universals of translation - Università degli studi di Macerata](#) [❷ Translation Universals | The Oxford Handbook of Translation Studies | Oxford Academic](#) [❸ 语料库翻译学的几个术语 - SJTU](#) [❹ Frontiers | Myth or reality? Some directions on translation universals in recent ...](#) [❺ Translation Universals: Do they exist? | Edited by Anna Mäuränen and Pekka Kuujamäki](#)

자세한 정보: 1. max.book1118.com 2. docenti.unimc.it 3. trans.ajtu.edu.cn +6 더 보기 3 / 30 ●

?(?) 번역보편소에 대해 더 알고 싶어요. 번역보편소는 번역의 일반적인 법칙을 찾기 위한 연구라고 하셨는데, 그것이 어떤 의미인가요? 왜 번역보편소가 중요한가요?

3 TU

■ TU (Bard)

번역보편소 (translation universals)란 모든 언어의 번역에서 공통적으로 나타나는 특징을 말합니다. 번역보편소는 1990년대 이후, 번역학 연구에서 활발하게 논의되고 있는 개념입니다.

번역 유형은 번역자의 전략이나 선택에 기인한 번역보편소입니다. 예를 들어, 모든 언어에서 명사와 동사, 형용사와 부사 등과 같은 품사들이 존재하기 때문에, 번역 과정에서 이러한 품사들의 대응관계를 명확히 해야 합니다. 또한, 모든 언어에서 의미를 전달하기 위해 문장 구조가 사용되기도 때문에, 번역 과정에서 문장 구조의 변화를 고려해야 합니다.

번역보편소는 번역학 연구의 중요한 틀을 제공합니다. 번역보편소를 이해함으로써 번역자는 번역 과정에서 발생하는 다양한 문제들을 이해하고, 이를 해결하기 위한 전략을 수립할 수 있습니다.

번역보편소로 알려진 대표적인 현상에는 다음과 같은 것들이 있습니다.

- 명시화(explication): 원문의 의미가 목표언어에서 암시적으로 전달되는 경우, 번역자는 원문을 명시적으로 번역하는 경향이 있습니다. 예를 들어, 영어의 "I am going to the store."를 "나는 가게에 갈 거야."라고 번역할 때, "going to"라는 표현을 "갈 거야"라는 표현으로 명시적으로 번역하는 경우가 많습니다.
- 단순화(simplification): 원문의 구조나 표현이 복잡하거나 어려운 경우, 번역자는 원문을 단순화하여 번역하는 경향이 있습니다. 예를 들어, 영어의 "I am going to the store and buy some milk."를 "나는 가게에 갈 거야."라고 번역할 때, "to buy some milk"라는 표현을 생략하여 번역하는 경우가 많습니다.
- 규범화(normalization): 원문의 표현이 목표언어의 어법이나 문화에 맞지 않는 경우, 번역자는 원문을 규범화하여 번역하는 경향이 있습니다. 예를 들어, 영어의 "I am going to the store and buy some milk."를 "나는 가게에 가서 우유를 살 거야."라고 번역할 때, "and"라는 표현을 "가서"라는 표현으로 번역하여 규범화하는 경우가 많습니다.
- 평준화(leveling-out): 원문과 목표언어의 차이를 줄이기 위해 번역자는 원문의 표현을 목표언어의 일반적인 표현으로 바꾸는 경향이 있습니다. 예를 들어, 영어의 "I am going to the store."를 "나는 가게에 가야 해."라고 번역할 때, "갈 거야"라는 표현을 "가야 해"라는 표현으로 바꾸어 평준화하는 경우가 많습니다.

번역보편소는 모든 언어의 번역에서 공통적으로 나타나는 현상이기는 하지만, 언어의 종류나 번역의 유적에 따라 그 정도나 양상이 달라질 수 있습니다. 또한, 번역보편소는 번역자의 의향이나 번역의 환경에 따라 영향을 받을 수 있습니다.

4 TU

■ Explication (Chat GPT)

응답

'번역보편소(translation universals)'들 중에서 명시화(explicitation)에 대하여 설명해줘.



명시화(Explicitation)는 번역 보편소 중 하나로서, 번역 과정에서 원래 언어의 정보나 의미를 명확하게 표현하려는 경향을 가리킵니다. 즉, 번역된 문장이 대상 언어에서 더 명확하게 이해되도록 하는 것을 목표로 합니다. 명시화는 다음과 같은 특징을 가집니다:

1. 정보 보존: 명시화는 번역 과정에서 정보의 손실을 최소화하려는 노력으로 이어집니다. 원래 언어의 정보나 내용을 번역된 문장에서 가급적 유지하려고 합니다.
2. 명확한 표현: 번역된 문장에서 불분명하거나 모호한 부분을 원래 언어와 달리 명확하게 표현하려고 합니다. 이는 대상 언어 독자나 청취자가 더 쉽게 이해하도록 돋는 역할을 합니다.
3. 추가 어휘와 설명: 명시화는 원래 언어의 표현을 보다 자세하게 설명하거나 추가 어휘를 도입하여 번역된 문장이 더 명확하게 전달되도록 합니다.
4. 문장 구조 조절: 원래 언어와 대상 언어의 문장 구조가 다를 때, 명시화는 대상 언어에서 더 자연스러운 구조로 조절하는 것을 의미합니다.

명시화는 번역 작업에서 중요한 요소 중 하나로, 특히 번역된 문장이 대상 언어 독자나 청취자에게 더 잘 전달되도록 하는 데 사용됩니다. 그러나 이러한 과정에서 뉘앙스나 문체의 변화, 번역의 발생 등 다른 번역 보편소와 충돌할 수도 있으며, 이를 균형있게 다루는 것이 번역 작업의 중요한 과제 중 하나입니다.

4 TU

■ Explication (Bing)

'번역보편소(translation universals)'들 중에서 명시화(explicitation)에 대하여 설명해줘.

✓ "명시화 explicitation" 검색 중

✓ 답변을 생성하는 중...

명시화(explicitation)는 번역 텍스트에서 원문의 의미(암시된 의미 포함)를 더 쉽고 명확하게 이해할 수 있도록 하는 번역 전략이야^①. 명시화를 위해 번역가는 원문에 없는 정보를 추가하거나, 원문의 정보를 다른 표현으로 대체할 수 있어^②. 예를 들어, 문학번역에서 명시화는 서사성을 강화하거나 악화시킬 수 있어^③. 명시화에 관한 더 자세한 설명과 연구는 다음 링크를 참고해 보세요.

^①: 문학번역에서의 명시화(explicitation)와 영여적 요소 | DBpia^②: [논문]The Effect of Explicitation on Narrativity in Literary Translation^③: Study on Explicitation Strategy in English-Korean Game Translation A Case Study of ...

자세한 정보: 1. dbpia.co.kr 2. koreascience.kr 3. scienceon.kisti.re.kr

4 / 30



명시화가 번역에 어떤 영향을 미치나요?

번역에서 명시화를 사용하는 이유는 무엇인가요?

다른 '번역보편소'들은 어떤 것이 있나요?

4 TU

■ Explication (Bard)

영시화(explication)는 번역보편소 중 하나로, 원문의 의미가 목표언어에서 암시적으로 전달되는 경우, 번역자가 원문을 명시적으로 번역하는 경향입니다.

영시화는 다음과 같은 상황에서 발생합니다.

- 원문의 의미가 목표언어에서 암시적으로 전달되는 경우
- 원문의 구조나 표현이 목표언어의 어법이나 문화에 맞지 않는 경우
- 번역자의 의사에 따라 원문을 명시적으로 번역하고자 하는 경우

영시화는 원문의 의미를 정확하게 전달하기 위해 중요한 역할을 합니다. 원문의 의미가 암시적으로 전달되는 경우, 목표언어 독자는 원문의 의미를 이해하기 어려울 수 있습니다. 이러한 경우, 번역자는 원문을 명시적으로 번역하여 원문의 의미를 명확하게 전달할 수 있습니다.

예를 들어, 영어의 "I am going to the store."를 "나는 가게에 갈 거야."라고 번역할 때, "going to"라는 표현은 한국어에서 "갈 거야"라는 표현으로 명시적으로 번역하는 경우가 많습니다. "going to"라는 표현은 "가려고 한다"라는 의미를 암시적으로 전달하지만, 한국어에서는 "갈 거야"라는 표현을 사용하여 보다 명확하게 의미를 전달할 수 있습니다.

또한, 원문의 구조나 표현이 목표언어의 어법이나 문화에 맞지 않는 경우, 번역자는 원문을 명시적으로 번역하여 목표언어의 어법이나 문화에 맞게 표현할 수 있습니다.

예를 들어, 영어의 "I am going to the store to buy some milk."를 "나는 우유를 사러 가게에 갈 거야."라고 번역할 때, "to buy some milk"라는 표현을 생략하여 번역하는 경우가 많습니다. 하지만, 한국어에서는 "우유를 사러"라는 표현을 사용하여 목표언어의 어법에 맞게 표현할 수 있습니다.

마지막으로, 번역자의 의사에 따라 원문을 명시적으로 번역하고자 하는 경우도 있습니다. 예를 들어, 어린이 대상의 번역이나, 기술적인 내용을 담은 번역의 경우, 원문의 의미를 보다 명확하게 전달하기 위해 영시화를 사용하는 경우가 많습니다.

영시화는 번역의 과정에서 중요한 역할을 하지만, 과도하게 사용될 경우, 원문의 특장성을 해칠 수 있습니다. 따라서, 번역자는 원문의 의미를 정확하게 전달하면서도 원문의 특장성을 유지할 수 있도록 영시화를 적절히 사용하는 것이 중요합니다.

4 TU

■ Simplification (Chat GPT)

단순화(Simplification)는 번역 보편소 중 하나로서, 번역 과정에서 번역된 텍스트를 원래 언어보다 더 간결하고 단순하게 만드는 경향을 가리킵니다. 이는 다음과 같은 특징을 갖습니다:

1. 간결한 언어 사용: 단순화는 번역된 텍스트에서 복잡한 어휘나 문장 구조를 피하고, 간결하고 이해하기 쉬운 언어를 사용하려는 노력을 의미합니다.
2. 정보의 간소화: 원래 언어의 텍스트가 상세한 정보나 세부 내용을 포함하더라도, 번역된 텍스트에서는 이를 필요한 정보만 남기고 나머지를 생략하거나 단순화하는 경향이 있습니다.
3. 언어적 복잡성 감소: 복잡한 언어적 구조나 뉘앙스를 최소화하려는 노력으로, 번역된 문장이 명확하고 직관적으로 이해되도록 합니다.
4. 목적어와 주어 간 간소화: 번역된 문장에서는 원래 언어의 주어와 목적어 간의 관계를 더 간소화하거나 명확하게 표현하는 것이 일반적입니다.

단순화는 번역 작업에서 중요한 고려 사항 중 하나입니다. 특히 대상 언어의 문체나 문화와 어울리게 텍스트를 만들고, 대상 독자나 청취자가 이해하기 쉽도록 들는 역할을 합니다. 그러나 단순화의 정도는 번역 작업의 목적과 대상 독자나 청취자를 고려하여 조절되어야 합니다. 때로는 원래 언어의 복잡성이 중요한 정보를 전달하는 데 필요한 경우가 있으므로, 단순화가 지나치게 이루어지면 의사소통에 문제가 발생할 수도 있습니다.

4 TU

■ Simplification (Bing)

'번역보편소(translation universals)'들 중에서 단순화(simplification)에 대하여 설명해줘.

- ✓ "단순화 simplification" 검색 중
- ✓ 답변을 생성하는 중...

단순화(simplification)는 번역 텍스트에서 원문의 정보나 구조를 더 간단하고 명료하게 표현하는 번역 전략이야^①. 단순화를 위해 번역가는 원문의 단어나 문장을 줄이거나, 복잡한 표현을 간단한 표현으로 바꾸거나, 의미적으로 구체적인 표현을 의미적으로 일반적인 표현으로 바꿀 수 있어^②. 예를 들어, 수학에서 단순화는 복잡한 식을 더 간단한 형태로 변환하는 과정이야^③. 단순화에 관한 더 자세한 설명과 연구는 다음 링크를 참고해 보세요.

^①: simplification | 영어를 한국어로 번역: Cambridge Dictionary ^②: simple, simply, simplify, simplified, simplification - 네이버 블로그 ^③: 대수적 단순화 - MATLAB simplify - MathWorks 한국

자세한 정보: 1. edcalingo.com 2. blog.naver.com 3. kr.mathworks.com +2 더 보기

5 / 30 ●

(?) 단순화가 번역에 어떤 영향을 미치나요? 번역에서 단순화를 사용하는 이유는 무엇인가요? 다른 '번역보편소'들은 어떤 것이 있나요?

4 TU

■ Simplification (Bard)

'번역보편소(translation universals)'들 중에서 단순화(simplification)에 대하여 설명해줘.

다른 답변 보기 ▾

단순화(simplification)는 원문의 구조나 표현이 복잡하거나 어려운 경우, 번역자가 원문을 단순화하여 번역하는 경향입니다.

- 원문의 구조나 표현이 복잡하거나 어려운 경우
- 번역자의 의사에 따라 원문을 단순화하고자 하는 경우

단순화는 원문의 이해를 돋기 위해 중요한 역할을 합니다. 원문의 구조나 표현이 복잡하거나 어려운 경우, 목표언어 독자는 원문의 의미를 이해하기 어려울 수 있습니다. 이러한 경우, 번역자는 원문을 단순화하여 원문의 의미를 보다 쉽게 이해할 수 있도록 할 수 있습니다.

예를 들어, 영어의 "I am going to the store to buy some milk."를 "나는 우유를 사러 가게에 갈 거야."라고 번역할 때, "I am going to the store to buy some milk."라는 문장은 주어-동사-목적어-부사절의 구조로 되어 있어 복잡합니다. 이러한 경우, 번역자는 원문을 "나는 우유를 사러 가게에 갈 거야."라고 단순화하여 번역할 수 있습니다.

또한, 번역자의 의사에 따라 원문을 단순화하고자 하는 경우도 있습니다. 예를 들어, 어린이 대상의 번역이나, 대중적인 내용을 담은 번역의 경우, 원문의 의미를 보다 쉽게 이해할 수 있도록 단순화하는 경우가 많습니다.

단순화는 번역의 과정에서 중요한 역할을 하지만, 과도하게 사용될 경우, 원문의 의미를 왜곡하거나 원문의 독창성을 해칠 수 있습니다. 따라서, 번역자는 원문의 의미를 정확하게 전달하면서도 원문의 독창성을 유지할 수 있도록 단순화를 적절히 사용하는 것이 중요합니다.

다음은 단순화의 예시입니다.

- 영어: "I am going to the store to buy some milk."
- 한국어: "나는 우유를 사러 가게에 갈 거야."

원문의 "to buy some milk"라는 표현은 한국어에서는 "우유를 사러"라는 표현으로 단순화할 수 있습니다.



4 TU

■ Normalization (Chat GPT)

응답

'번역보편소(translation universals)'들 중에서 규범화(normalization)에 대하여 설명해줘.



규범화(Normalization)는 번역 보편소 중 하나로, 번역 과정에서 대상 언어의 언어적 규범과 표준에 따라 원래 언어의 문장을 조절하려는 경향을 가리킵니다. 이것은 다음과 같은 특징을 갖습니다:

1. 언어적 일관성: 규범화는 대상 언어에서의 일관성을 유지하려는 노력으로 번역된 텍스트가 대상 언어의 언어 규칙과 표준에 부합하도록 합니다.
2. 표준어 사용: 번역된 텍스트에서 표준 어휘와 표현을 선호합니다. 특정 언어 변형이나 지방어, 사회적 언어 변형은 최소화하거나 피하려는 경향이 있습니다.
3. 문법 규칙 준수: 번역된 문장은 대상 언어의 문법 규칙을 따르도록 조절됩니다. 이는 문장 구조, 시제, 수, 대명사 등을 대상 언어의 표준에 맞게 정확하게 사용하는 것을 의미합니다.
4. 어휘 표준화: 원래 언어에서 사용된 어휘나 표현이 대상 언어에서는 표준화된 어휘로 대체되는 경우가 많습니다. 이로 인해 번역된 텍스트가 대상 언어에서 자연스럽게 보이도록 합니다.

규범화는 번역 과정에서 번역 투자를 최소화하고 대상 언어의 자연스러운 표현을 촉진하는 데 사용됩니다. 특히 번역된 텍스트가 대중 매체나 공식 문서, 교육 자료와 같이 언어적 표준이 중요한 경우에 규범화가 더욱 중요하게 고려됩니다. 그러나 때로는 언어의 다양성과 미묘한 뉘앙스를 유지하기 위해 규범화를 완화하거나 조절해야 할 필요도 있습니다.

4 TU

■ Normalization (Bing)

'번역보편소(translation universals)'들 중에서 규범화(normalization)에 대하여 설명해줘.

✓ "규범화 normalization" 검색 중

✓ 답변을 생성하는 중...

규범화(normalization)는 번역 텍스트에서 원문의 정보나 구조를 더 표준적이고 정규적인 형태로 표현하는 번역 전략이야^①. 규범화를 위해 번역가는 원문의 언어나 문학적 특성을 번역 언어나 문학의 규범에 맞게 조정하거나, 비표준적이거나 비문법적인 표현을 표준적이고 문법적인 표현으로 바꿀 수 있어^②. 예를 들어, 데이터 분석에서 규범화는 데이터의 범위를 일정한 기준에 맞게 변환하는 과정이야^③. 규범화에 관한 더 자세한 설명과 연구는 다음 링크를 참고해 보세요.

^①: Batch Normalization, Layer Normalization 차이 : 네이버 블로그 ^②: Translation Universals: Do They Exist? - ResearchGate ^③: 정규화(normalization)와 표준화(standardization), 머신러닝 성능 향상을 ...

자세한 정보: 1. bing.com 2. bskyvision.com 3. intrepidgeeks.com 4. blog.naver.com

6 / 30



규범화가 번역에 어떤 영향을 미치나요?

번역에서 규범화를 사용하는 이유는 무엇인가요?

다른 '번역보편소'들은 어떤 것이 있나요?

4 TU

■ Normalization (Bard)

'번역보편소(translation universals)를 중에서 규범화(normalization)에 대하여 설명해줘.'

규범화(normalization)는 번역보편소 중 하나로, 원문의 표현이 목표언어의 어법이나 문화에 맞지 않는 경우, 번역자가 원문을 목표언어의 어법이나 문화에 맞게 번역하는 경향입니다.

규범화는 다음과 같은 상황에서 발생합니다.

- 원문의 표현이 목표언어의 어법이나 문화에 맞지 않는 경우
- 번역자의 의사에 따라 원문을 규범화하고자 하는 경우

규범화는 원문의 의미를 목표언어 독자에게 정확하게 전달하기 위해 중요한 역할을 합니다. 원문의 표현이 목표언어의 어법이나 문화에 맞지 않는 경우, 목표언어 독자는 원문의 의미를 이해하기 어려울 수 있습니다. 이러한 경우, 번역자는 원문을 규범화하여 목표언어의 어법이나 문화에 맞게 표현할 수 있습니다.

예를 들어, 영어의 "I am going to the store and buy some milk."를 "나는 가게에 가서 우유를 살 거야."라고 번역할 때, "and"라는 표현은 한국어에서는 "가서"라는 표현으로 규범화하는 경우가 많습니다. "and"라는 표현은 한국어에서는 "그리고"라는 표현으로 사용되지만, 한국어에서는 "가서"라는 표현을 사용하여 보다 자연스럽게 표현할 수 있습니다.

또한, 번역자의 의사에 따라 원문을 규범화하고자 하는 경우도 있습니다. 예를 들어, 법률 문서나 공문서의 번역의 경우, 원문의 의미를 보다 정확하게 전달하기 위해 규범화하는 경우가 많습니다.

규범화는 번역의 과정에서 중요한 역할을 하지만, 과도하게 사용될 경우, 원문의 독창성을 해칠 수 있습니다. 따라서, 번역자는 원문의 의미를 정확하게 전달하면서도 원문의 독창성을 유지할 수 있도록 규범화를 적절히 사용하는 것이 중요합니다.

다음은 규범화의 예시입니다.

- 영어: "I am going to the store and buy some milk."
- 한국어: "나는 가게에 가서 우유를 살 거야."

원문의 "and"라는 표현은 한국어에서는 "가서"라는 표현으로 규범화할 수 있습니다.

상단에는 편집 도구와 함께 Google에서 검색 링크가 있습니다.

4 TU

■ Levelling out (Chat GPT)

'번역보편소(translation universals)들 중에서 평준화(leveling-out)에 대하여 설명해줘.'

평준화(Leveling-out)은 번역 보편소 중 하나로, 번역 과정에서 언어 간의 문법적, 어휘적 차이를 최소화하려는 경향을 가리킵니다. 이것은 다음과 같은 특징을 갖습니다:

- 문법 통일화: 평준화는 번역된 텍스트에서 원래 언어와 대상 언어 간의 문법 규칙을 보다 통일적으로 만드는 것을 목표로 합니다. 원래 언어와 대상 언어의 문장 구조나 문법 규칙이 다를 때, 번역된 텍스트에서 이를 가능한 한 유사하게 만들습니다.
- 어휘 일치: 번역된 텍스트에서는 원래 언어의 어휘와 대상 언어의 어휘 간의 일치를 촉진하려는 노력이 있습니다. 비슷한 뜻을 가진 단어나 표현을 선택하여 언어 간의 일관성을 유지하려고 합니다.
- 문체 통일화: 평준화는 번역된 텍스트의 문체나 어조를 가능한 한 대상 언어의 문체와 일치시키려는 경향이 있습니다. 이는 번역된 텍스트가 대상 언어의 대중 언어나 문서 작성 규칙을 따르도록 돋는 역할을 합니다.

평준화는 번역 작업에서 언어 간의 호환성을 높이는 데 사용됩니다. 특히 다른 문화나 언어를 이해하기 어려운 경우에는 평준화를 통해 번역된 텍스트가 대상 언어 독자나 청취자에게 더 자연스럽게 보이고 이해하기 쉽도록 만들어집니다. 그러나 이것은 때로는 뉘앙스나 문화적 차이를 희생할 수 있으므로, 번역 작업자는 균형을 맞추는 것이 중요합니다.

4 TU

■ Levelling out (Bing)

'번역보편소(translation universals)들 중에서 평준화(leveling-out)에 대하여 설명해줘.'

- ✓ “평준화 leveling-out” 검색 중
- ✓ 답변을 생성하는 중...

평준화(leveling-out)는 번역 텍스트에서 원문의 다양성이나 복잡성을 줄이고, 번역 언어나 문화의 규범에 맞게 단순화하거나 일반화하는 번역 전략이야^①. 평준화를 위해 번역자는 원문의 특이한 표현이나 의미를 더 훈하고 평범한 표현이나 의미로 바꾸거나, 원문의 모호함이나 불확실성을 더 명확하고 확실한 방식으로 해석할 수 있어^②. 예를 들어, 화학에서 평준화효과는 용매가 산과 염기의 강도에 영향을 주는 현상 이야^③. 평준화에 관한 더 자세한 설명과 연구는 다음 링크를 참고해 보세요.

^①: Translation Universals: Do They Exist? - ResearchGate ^②: 번역보편소에 대한 연구 - 한국학술정보 ^③: 평준화효과(Leveling effect) : 네이버 블로그

자세한 정보: 1. blog.naver.com 2. asana.com 3. blog.naver.com

7 / 30

(?) 평준화가 번역에 어떤 영향을 미치나요? | 번역에서 평준화를 사용하는 이유는 무엇인가요? | 다른 '번역보편소'들은 어떤 것이 있나요?

4 TU

■ Levelling out (Bard)

'번역보편소(translation universals)들 중에서 평준화(leveling-out)에 대하여 설명해줘.'

다른 답변 보기 ▾

▶ 평준화(leveling-out)는 번역보편소 중 하나로, 원문과 목표언어의 차이를 줄이기 위해 번역자가 원문의 표현을 목표언어의 일반적인 표현으로 바꾸는 경향입니다.

평준화는 다음과 같은 상황에서 발생합니다.

- 원문과 목표언어의 차이가 크거나 복잡한 경우
- 번역자의 의사에 따라 원문을 평준화하고자 하는 경우

평준화는 원문의 이해를 돋우기 위해 중요한 역할을 합니다. 원문과 목표언어의 차이가 크거나 복잡한 경우, 목표언어 독자는 원문의 의미를 이해하기 어려울 수 있습니다. 이러한 경우, 번역자는 원문을 평준화하여 원문의 의미를 보다 쉽게 이해할 수 있도록 할 수 있습니다.

예를 들어, 영어의 “The big brown dog”를 “큰 갈색 강아지”라고 번역할 때, “big”와 “brown”이라는 표현은 한국어에서 “큰”과 “갈색”이라는 표현으로 일반적으로 사용됩니다. 이러한 경우, 번역자는 원문을 “큰 갈색 강아지”라고 평준화하여 번역할 수 있습니다.

또한, 번역자의 의사에 따라 원문을 평준화하고자 하는 경우도 있습니다. 예를 들어, 기술적인 내용을 담은 번역의 경우, 원문의 의미를 보다 쉽게 이해하기 위해 평준화하는 경우가 많습니다.

평준화는 번역의 과정에서 중요한 역할을 하지만, 과도하게 사용될 경우, 원문의 독창성을 해칠 수 있습니다. 따라서, 번역자는 원문의 의미를 정확하게 전달하면서도 원문의 독창성을 유지할 수 있도록 평준화를 적절히 사용하는 것이 중요합니다.

다음은 평준화의 예시입니다.

- 영어: “The big brown dog”
- 한국어: “큰 갈색 강아지”

원문의 “big”와 “brown”이라는 표현은 한국어에서 “큰”과 “갈색”이라는 표현으로 일반적으로 사용되므로, 평준화하여 번역할 수 있습니다.

번역보편소는 번역의 과정에서 발생하는 다양한 현상을 설명하는데 도움이 되는 개념입니다. 번역자는 번역보편소를 이해함으로써 번역의 과정에서 발생하는 문제들을 이해하고, 이를 해결하기 위한 전략을 수립할 수 있습니다.

05

Result #2: Text Analysis

TT or NTT? (GPT 3.5)

다음의 두 문서들은 ‘번역보편소(translation universal)’의 4가지 측면(영시화, 단순화, 규범화, 디디션화)에서 어떤 차이를 보이는가? 그리고 <문서 1>과 <문서 2> 중에서 어느 문서가 번역문인가?
<문서 1>
To ask less of the planet, we ask more of ourselves. Apple's mission has never wavered. We are here to change the world. It's what inspires the amazing products and services that reach billions of people around the globe. It's what drives us to relentlessly push the boundaries of innovation and design. And it's why we continuously strive to do more with less: reducing our impact on the Earth we all share, while expanding and redefining the possibilities ahead. Our 11th annual Environmental Responsibility Report covers fiscal year 2017, and highlights the progress we've made in meeting our environmental objectives. The challenges are complex, but the philosophy is simple: Apply the same innovation that goes into making our products. Today, we remain focused on three priorities where we and our stakeholders believe Apple can make the biggest difference: Climate change; Reduce our impact on climate change by using renewable energy sources and driving energy efficiency in our products, facilities, and supply chain. Resources: Conserve precious resources so we all can thrive. Safer materials: Pioneer the use of safer materials in our products and processes. In each of these key areas, we set ambitious but measurable goals, then develop smart solutions to meet them. Our work to reduce our greenhouse gas emissions is a great example. More than a decade ago, we started to transition our electricity use to renewable sources. Today, we're proud to power 100 percent of our operations around the world with 100 percent renewable energy. That means every Apple data center, office, and manufacturing facility in 43 countries now runs on clean power. In the process, we've inspired others for other companies and organizations to purchase renewable energy and transition their own operations to greener power. As we worked to reach 100 percent renewable energy for our own operations, we began to challenge our suppliers to make the same transition to clean energy through our supplier clean energy program. As of April 2018, 23 manufacturers have committed to power all of their Apple operations with 100 percent renewable energy. And we're well on our way to meeting our goal of putting 4 gigawatts of new clean energy online in our supply chain by 2020. We also made progress toward our goal to one day make new products without mining new materials from the earth. We hope to get there by using only recycled or renewable materials in our products, and returning an equivalent amount of material back to the market, to be used by us or others. Like everything at Apple, this effort is driven by innovation. We've developed a brand-new robot, Daisy, capable of disassembling nine versions of iPhone, and sorting their high-quality components for recycling. To help keep Daisy busy, we're making it easier than ever for customers to recycle their old Apple devices through our new Apple GiveBack experience. There's a lot more to our story over the past year. As you read through this Environmental Progress Report, we are busy working toward our next frontiers in clean energy, materials recovery, and green chemistry. We are proud of our work and look forward to continuing our journey.

<문서 2>
Dear Stakeholders, I would like to begin by expressing my sincere appreciation for your continuous support of Samsung Electronics. Grateful for the support from our stakeholders and the efforts of our employees, Samsung Electronics had a very strong finish to a record-breaking year in 2017, and are working to maintain our steady performance in 2018. The global economic environment, however, remains highly uncertain, and technologies in areas such as artificial intelligence, autonomous driving, big data, and the Internet of Things are rapidly shifting paradigm in the IT sector, posing new challenges that demand innovation. In response, Samsung Electronics will make full use of its capabilities as a global IT company to cope with new challenges. We will also expand our efforts in improving environmental sustainability, a task that has never been more vital to us all. Since its inception, Samsung's unique and innovation-centered culture has been at the core of our technology and products. Based on bold initiatives and innovations, we are constantly developing market-leading products and services that transcend the capabilities of existing technology. Furthermore, we are exploring ways to increase our contribution to society and fulfill our responsibility as a global corporate citizen by creating social value as well as economic value. Samsung Electronics is dedicated to using our technology to create positive value via addressing societal issues. To illustrate, we created ReumHyo, an application that works in conjunction with a Gear VR to help those living with low vision see the world more clearly; and Ignis, a lightweight thermal imaging device that lessens the burden of firefighters. We acknowledge our responsibility to the planet, and have thus committed ourselves to minimizing environmental impacts of our operations and products alike. Advancing into a circular economy, we have developed some of the world's most energy efficient products and expanded our efforts to recycle and reuse natural resources in our product components. Furthermore, we plan to increase the use of renewable energy to reduce the carbon footprint of our operations. Finally, we will continue to cooperate with stakeholders, who are concerned with sustainable and responsible growth. We will seek advice from investors, NGOs, local communities, governments, and international organizations, and reflect their opinions in our decision-making process. We follow a simple business philosophy at Samsung Electronics: devote our talent and technology to creating innovative products and services that contribute to a better global society. We fully endorse the UN Sustainable Development Goals, and are thus working to positively affect society and the environment alongside the ongoing development of our business. To that end, we would like to introduce our ‘Sustainability Report 2018’. This year’s report consists of five key pillars that lay out our achievements in 2017 and our plans for the future. We look forward to your continued interest and support. Thank you.

문서1: NTT (Apple)
문서 2: TT (삼성전자)

TT or NTT? (GPT 3.5)

문서 1: NTT (Apple) → TT
문서 2: TT (삼성전자) → NTT

<문서 1>과 <문서 2>는 다음과 같은 번역 보편소(translation universals)의 4가지 측면(명시화, 단순화, 규범화, 평준화)에서 차이를 보입니다.

1. 명시화(Explicitation):
 - <문서 1>: "Apple's mission has never wavered."라는 문장에서 "mission"이 더 명시적으로 언급됩니다.
 - <문서 2>: "Samsung Electronics had a very strong finish to a record-breaking year in 2017"라는 문장에서 "Samsung Electronics"가 더 명시적으로 언급됩니다.
2. 단순화(Simplification):
 - <문서 1>: "It's what inspires the amazing products and services that reach billions of people around the globe."와 같이 복잡한 문장 구조와 표현이 사용됩니다.
 - <문서 2>: "Grateful for the support from our stakeholders and the efforts of our employees, Samsung Electronics had a very strong finish to a record-breaking year in 2017"와 같이 간단한 문장이 사용됩니다.
3. 규범화(Normalization):
 - <문서 1>: "It's what drives us to relentlessly push the boundaries of innovation and design."에서 "push the boundaries" 표현을 사용합니다.
 - <문서 2>: "Based on bold initiatives and innovations, we are constantly developing market-leading products and services that transcend the capabilities of existing technology."에서 "bold initiatives"와 "innovations"라는 어휘를 사용합니다.
4. 평준화(Leveling-out):
 - <문서 1>: "we started to transition our electricity use to renewable sources"에서 "transition our electricity use" 표현이 사용됩니다.
 - <문서 2>: "we have developed some of the world's most energy efficient products"에서 "energy efficient"라는 표현이 사용됩니다.

따라서 <문서 1>은 번역된 문서로 보이며, 번역된 텍스트에서 원래 언어와 대상 언어 간의 언어적 차이를 최소화하려는 경향이 나타납니다. <문서 2>는 번역된 문서처럼 보이지 않으며, 보다 원어민 스타일의 텍스트로 보입니다.

TT or NTT? (GPT 3.5)

J 다음의 두 문서들은 '번역보편소(translation universals)'의 4가지 측면(명시화, 단순화, 규범화, 평준화)에서 어떠한 차이를 보이는가? 그리고 <문서 1>과 <문서 2> 중에서 어느 문서가 번역된 문서인가?
참고로 이 문서들은 지속가능성 보고서의 대표이사 서한이야.
<문서 1>

J 다음의 두 문서들은 '번역보편소(translation universals)'의 4가지 측면(명시화, 단순화, 규범화, 평준화)에서 어떠한 차이를 보이는가? 그리고 <문서 1>과 <문서 2> 중에서 어느 문서가 번역된 문서인가?
참고로 이 문서들은 지속가능성 보고서의 대표이사 서한이야.
번역된 서한은 번역되지 않은 서한에 비해 접속부사 비중이 높고, 기능어(function words) 비중이 낮고, 문장길이가 긴 것으로 알려져 있어.
<문서 1>

문서 1: NTT (Apple) → ?
문서 2: TT (삼성전자) → ?

도 문서는 '번역보편소(translation universals)'의 4가지 측면(명시화, 단순화, 규범화, 평준화)에서 다음과 같은 차이를 보입니다.

1. 명시화(Explicitation):
 - 문서 1: 문서 1은 상세한 정보와 구체적인 내용을 다루며, 환경 책임 보고서의 세부 사항과 목표를 자세하게 언급합니다.
 - 문서 2: 문서 2는 더 일반적으로 환경 지속 가능성 및 회사의 역할을 다루며, 구체적인 사례나 세부 내용을 제공하지 않습니다.
2. 단순화(Simplification):
 - 문서 1: 문서 1은 더 복잡한 환경 문제와 솔루션에 대한 구체적인 내용을 다루며, 다양한 목표와 도구를 설명합니다.
 - 문서 2: 문서 2는 더 간단한 언어와 메시지를 사용하고 있으며, 환경 지속 가능성을 간결하게 소개하고 있습니다.
3. 규범화(Normalization):
 - 문서 1: 문서 1은 환경 책임 보고서의 구조와 목표를 규범화하고 있으며, 환경 지속 가능성과 관련된 정규화된 용어를 사용합니다.
 - 문서 2: 문서 2는 환경 지속 가능성을 다루지만, 더 일반적인 비즈니스 용어를 사용하고 있으며, 구체적인 용어를 제한적으로 사용합니다.
4. 평준화(Standardization):
 - 문서 1: 문서 1은 회사의 환경 지속 가능성 주제를 정기 보고서로 표준화하고 있으며, 연도별 진전과 목표를 표준적인 방식으로 제시하고 있습니다.
 - 문서 2: 문서 2는 회사의 지속 가능성 주제를 정기 보고서로 표준화하고 있으며, 몇 가지 주요 피라미드로 구성하여 제시하고 있습니다.

어떤 문서가 번역된 문서인지를 판단하기는 어렵습니다. 이러한 결정은 번역된 언어와 문서의 원문을 직접 비교해야 합니다. 여기서는 번역 여부를 판단하기 위한 충분한 정보가 제공되지 않으므로 판단할 수 없습니다.

TT or NTT? (Bing)

다음의 두 문서들은 '번역보편소(translation universals)'의 4가지 즉연(영시화, 단순화, 규범화, 평준화)에서 어떠한 차이를 보이는가? 그리고 <문서 1>과 <문서 2> 중에서 어느 문서가 번역된 문서인가? <문서 1>

To ask less of the planet, we ask more of ourselves. Apple's mission has never wavered. We are here to change the world. It's what inspires the amazing products and services that reach billions of people around the globe. It's what drives us to relentlessly push the boundaries of innovation and design. And it's why we continuously strive to do more with less: reducing our impact on the Earth we all share, while expanding and redefining the possibilities ahead. Our 11th annual Environmental Responsibility Report covers fiscal year 2017, and highlights the progress we've made in meeting our environmental objectives. The challenges are complex, but the philosophy is simple: Apply the same innovation that goes into making our products. Today, we remain focused on three priorities where we and our stakeholders believe Apple can make the biggest difference: <bullet> Climate change: Reduce our impact on climate change by using renewable energy sources and driving energy efficiency in our products, facilities, and supply chain. <bullet> <문서 2>

Dear Stakeholders, I would like to begin by expressing my sincere appreciation for your continuous support of Samsung Electronics. Grateful for the support from our stakeholders and the efforts of our employees, Samsung Electronics had a very strong finish to a record-breaking year in 2017, and are working to maintain our steady performance in 2018. The global economic environment, however, remains highly uncertain, and technologies in areas such as artificial intelligence, autonomous driving, big data, and the Internet of Things are rapidly shifting paradigm in the IT sector, posing new challenges that demand innovation. In response, Samsung Electronics will make full use of its capabilities as a global IT company to cope with new challenges. We will also expand our efforts in improving environmental s

번역보편소(translation universals)는 번역에서 관찰되는 언어적 특성을 의미합니다. 이 특성들은 명시화(explicitation), 단순화(simplification), 규범화(normalization), 평준화(leveling out)의 네 가지 즉연으로 구성됩니다.¹⁾. 문서 1과 문서 2를 비교해보면, 문서 1은 단순하고 명확한 언어를 사용하며, 문서 2는 복잡한 어구와 긴 문장을 사용합니다.²⁾. 따라서 문서 1은 번역된 문서이며, 문서 2는 번역된 문서가 아닙니다.³⁾.

자세한 정보: 1. earticle.net 2. scienceon.kist.re.kr 3. korean.go.kr +9 더 보기

4 / 5 ●

문서1: NTT (Apple) → TT
문서 2: TT (삼성전자) → NTT

TT or NTT? (Bing)

J 다음의 두 문서들은 '번역보편소(translation universals)'의 4가지 즉연(영시화, 단순화, 규범화, 평준화)에서 어떠한 차이를 보이는가? 그리고 <문서 1>과 <문서 2> 중에서 어느 문서가 번역된 문서인가?

참고로 이 문서들은 지속가능성 보고서의 대표이사 서한이야.
<문서 1>

J 다음의 두 문서들은 '번역보편소(translation universals)'의 4가지 즉연(영시화, 단순화, 규범화, 평준화)에서 어떠한 차이를 보이는가? 그리고 <문서 1>과 <문서 2> 중에서 어느 문서가 번역된 문서인가?

참고로 이 문서들은 지속가능성 보고서의 대표이사 서한이야.
번역된 서한은 번역되지 않은 서한에 비해 접속부사 비중이 높고, 기능어(function words) 비중이 낮고, 문장길이가 길 것으로 알려져 있어.
<문서 1>

'번역보편소(translation universals)'의 4가지 즉연(영시화, 단순화, 규범화, 평준화)에 대해 두 문서를 비교하면 다음과 같습니다:

1. 영시화: 번역은 종종 원문보다 더 명확하게 정보를 전달하려는 경향이 있습니다. <문서 1>은 상세한 정보와 구체적인 사례를 제공하고 있으며, 이는 번역된 문서에서 자주 볼 수 있는 특징입니다. 반면에, <문서 2>는 더 간결하고 직접적인 표현을 사용하고 있습니다.

2. 단순화: 번역된 텍스트는 종종 원문보다 단순화되어 있습니다. <문서 1>은 복잡한 문장 구조와 고급 어휘를 사용하고 있는 반면, <문서 2>는 더 단순하고 직관적인 어휘와 문장 구조를 사용하고 있습니다.

3. 규범화: 번역은 종종 대상 언어의 문법 규칙과 어휘 사용을 엄격하게 따르려는 경향이 있습니다. <문서 2>는 한국어의 문법 규칙을 엄격하게 따르고 있으며, 이는 번역된 문서에서 자주 볼 수 있는 특징입니다.

4. 평준화: 번역은 종종 원문의 다양성을 줄이려는 경향이 있습니다. <문서 1>은 다양한 주제와 아이디어를 다루고 있는 반면, <문서 2>는 주제가 더 일관되어 있습니다.

따라서, 이러한 즉연들을 고려할 때, <문서 2>가 번역된 문서일 가능성이 높습니다. 하지만, 이는 단지 추정일 뿐이며, 확실한 결론을 내리기 위해서는 더 많은 정보와 전문적인 분석이 필요합니다.

3 / 5 ●

문서1: NTT (Apple) → NTT
문서 2: TT (삼성전자) → TT

TT or NTT? (Bard)

문서1: NTT (Apple) → NTT
문서2: TT (삼성전자) → TT

다음의 두 문서들은 '영역보편소(translation universals)'의 4가지 측면(영시화, 단순화, 규범화, 평준화)에서 어떠한 차이를 보이는가? 그리고 <문서 1>과 <문서 2> 중에서 어느 문서가 번역된 문서인가?

To ask less of the planet, we ask more of ourselves. Apple's mission has never wavered. We are here to change the world. It's what inspires the ...

<문서 1>과 <문서 2>는 모두 환경에 대한 책임에 대한 내용을 담고 있지만, 번역보편소의 4가지 측면에서 다음과 같은 차이를 보입니다.

- 영시화 : <문서 1>은 원문의 의미를 보다 명확하게 전달하기 위해 명시화하는 경향이 있습니다. 예를 들어, 원문의 "ask less of the planet"은 "지구에 덜 요구하자"라는 데 역하여 원문의 의미를 보다 명확하게 전달하고 있습니다. 반면, <문서 2>는 원문의 의미를 보다 간결하게 전달하기 위해 명시화를 사용하지 않는 경향이 있습니다. 예를 들어, 원문의 "acknowledging our responsibility to the planet"은 "지구에 대한 책임을 인식하고"라고 번역하여 원문의 의미를 간결하게 전달하고 있습니다.
- 단순화 : <문서 2>는 원문의 구조와 표현을 단순화하는 경향이 있습니다. 예를 들어, 원문의 "We will also expand our efforts in improving environmental sustainability"은 "환경 지속 가능성을 개선하기 위한 노력을 확대할 것입니다"라고 번역하여 원문의 구조를 단순화하고 있습니다. 반면, <문서 1>은 원문의 구조나 표현을 그대로 유지하는 경향이 있습니다. 예를 들어, 원문의 "That means every Apple data center, retail store, corporate office, and colocation facility in 43 countries around the world now runs on clean power"은 "그 의미는 전 세계 43개국에서 모든 Apple 데이터 센터, 소매점, 기업 사무실 및 클로우드 인프라 시설이 이제 청정 에너지로 가동됩니다"라고 번역하여 원문의 구조를 그대로 유지하고 있습니다.
- 규범화 : <문서 2>는 원문의 표현을 목표언어의 어법이나 문화에 맞게 규범화하는 경향이 있습니다. 예를 들어, 원문의 "paving the way for"은 "길을 닦는"으로 번역하여 목표언어의 어법에 맞게 규범화하고 있습니다. 반면, <문서 1>은 원문의 표현을 그대로 유지하는 경향이 있습니다. 예를 들어, 원문의 "renewable energy sources"은 "정정 에너지"로 번역하여 목표언어의 일반적인 표현으로 바꾸고 있습니다. 반면, <문서 1>은 원문의 표현을 그대로 유지하는 경향이 있습니다. 예를 들어, 원문의 "climate change"은 "기후 변화"로 번역하여 원문의 표현을 그대로 유지하고 있습니다.

이러한 차이를 바탕으로 볼 때, <문서 1>이 원문이고 <문서 2>가 번역된 문서일 가능성이 높습니다.

<참고>

- Baker, M. (1992). In other words: A coursebook on translation. London: Routledge.
- Hatim, B., & Mason, I. (1997). The translator as communicator. London: Routledge.
- Toury, G. (1995). Descriptive translation studies and beyond. Amsterdam: John Benjamins.

글자
1. blogs.ei.columbia.edu/2018/05/21/greenin...
[Copied]

[Close] [Minimize] [Maximize]

TT or NTT? (Bard)

두 문서는 번역보편소(translation universals)의 4가지 측면(영시화, 단순화, 규범화, 평준화)에서 다음과 같은 차이를 보입니다.

| 측면 | 문서 1 | 문서 2 |
|-----|-------------------------|------------------------|
| 영시화 | 여 영시학이고, 직설학인 표현을 사용한다. | 여 오피하고, 일반적인 표현을 사용한다. |
| 단순화 | 여 단순한 문장 구조와 어휘를 사용한다. | 여 복잡한 문장 구조와 어휘를 사용한다. |
| 규범화 | 여 농통화와 문장학적인 표현을 사용한다. | 여 수용화와 부정학적인 표현을 사용한다. |
| 평준화 | 여 문준화된 표현을 사용한다. | 여 퇴색화된 표현을 사용한다. |

예시:

[문서 1] 문서 2:
[...]
[...] To ask less of the planet, we ask more of ourselves. [Samsung Electronics is dedicated to using our technology to create positive value via addressing societal issues.]
[The challenges are complex, but the philosophy is simple: Apply the same innovation that goes into making our products. [We acknowledge our responsibility to the planet, and have thus committed ourselves to minimizing environmental impacts of our operations and products alike.]
[We've developed a brand-new robot, Daisy, capable of disassembling nine versions of iPhone, and sorting their high-quality components for recycling. [We have developed some of the world's most energy efficient products and expanded our efforts to recycle and reuse natural resources in our product components.]

어느 문서가 번역된 문서인가?

문서 2가 문서 1보다 더 영시학이고 단순하여 규범화된 표현을 사용합니다. 따라서 문서 2가 번역된 문서일 가능성이 더 높습니다.
또한, 문서 2는 문서 1보다 더 많은 문화적 요소를 포함하고 있습니다. 예를 들어, 문서 2는 "기술을 활용하여 사회 문제를 해결함으로써 문장학적인 가치를 창출하는"에 전념하고 있습니다. 와 같은 한국어 문화적 요소를 포함하고 있습니다. 번역된 문서는 일반적으로 원문의 문화적 요소를 보존합니다.

마지막으로, 문서 2가 번역된 문서일 가능성이 더 높습니다.

다음의 두 문서들은 '번역보편소(translation universals)'의 4가지 측면(영시화, 단순화, 규범화, 평준화)에서 어떠한 차이를 보이는가? 그리고 <문서 1>과 <문서 2> 중에서 어느 문서가 번역된 문서인가?

참고로 이 문서들은 지속가능성 보고서의 대표이사 서한이야.

[1]>

J 다음의 두 문서들은 '번역보편소(translation universals)'의 4가지 측면(영시화, 단순화, 규범화, 평준화)에서 어떠한 차이를 보이는가? 그리고 <문서 1>과 <문서 2> 중에서 어느 문서가 번역된 문서인가?
J 참고로 이 문서들은 지속가능성 보고서의 대표이사 서한이야.

번역된 서한은 번역되지 않은 서한에 비해 접속부사 비중이 높고, 기능어(function words) 비중이 낮고, 문장길이가 긴 것으로 알려져 있습니다.
<문서 1>

문서 1은 Apple의 2017 Environmental Responsibility Report이고, 문서 2는 Samsung Electronics의 2018 Sustainability Report입니다. 따라서 두 문서의 대상 독자가 다를 수 있습니다. 문서 1은 Apple의 지속가능성에 대한 보고서이므로 투자자나 사회학 학생들이 주로 독자로 가능성이 높습니다. 반면, 문서 2는 Samsung Electronics의 지속가능성 보고서이므로 주주나 고객들이 주로 독자로 가능성이 높습니다.

이러한 차이를 고려할 때, 문서 2가 더 영시학이고 단순하여 규범화된 표현을 사용하는 것이 합리적입니다. 문서 2는 주로 고객과 같은 다양한 이해관계자에게 읽히기 위한 것으로 더 쉽게 이해할 수 있도록 작성되어 있습니다. 또한, 문서 2는 Samsung Electronics의 지속가능성 보고서이므로 Samsung Electronics의 문장학적인 이미지를 전달하기 위해 규범화되고 문장학적인 이미지를 전달하기 위해 규범화된 표현을 사용할 가능성이 높습니다.

또한, 번역된 서한은 번역되지 않은 서한에 비해 접속부사 비중이 높고, 기능어(function words) 비중이 낮고, 문장길이가 긴 것으로 알려져 있습니다. 문서 2가 문서 1보다 접속부사 비중이 높고, 기능어(function words) 비중이 낮고, 문장길이가 긴 것도 문서 2가 번역된 문서일 가능성이 더 높다는 것을 알립니다.

결론적으로, 문서 2가 번역된 문서일 가능성이 더 높습니다.

문서1: NTT (Apple) → NTT
문서2: TT (삼성전자) → TT

Result

| Text type | | Chat GPT 3.5 | | | Bing (GPT 4) | | | Bard | | |
|----------------|-----------------------------|--------------|---------------|------------|--------------|---------------|------------|-------------|---------------|------------|
| Industry | TT or NTT | w/o context | w/ genre info | w/ TU info | w/o context | w/ genre info | w/ TU info | w/o context | w/ genre info | w/ TU info |
| IT | NTT: Apple | TT | ? | ? | TT | NTT | NTT | NTT | NTT | NTT |
| | TT: 삼성전자 | NTT | ? | ? | NTT | TT | TT | TT | TT | TT |
| Pharmaceutical | NTT: Johnson and Johnson | ? | ? | ? | NTT | NTT | NTT | NTT | NTT | NTT |
| | TT: Yuhan | ? | ? | ? | TT | TT | TT | TT | TT | TT |
| Manufacturing | NTT: Chevron | ? | ? | ? | NTT | NTT | NTT | NTT | NTT | NTT |
| | TT: Hankook Tire | ? | ? | ? | TT | TT | TT | TT | TT | TT |
| Finance | NTT: JP Morgan | ? | ? | ? | TT | NTT | NTT | NTT | NTT | NTT |
| | TT: KB Financial | ? | ? | ? | NTT | TT | TT | TT | TT | TT |

06
Closing

Results

Q1. Do generative AIs understand major concepts on TU?

- Yes, they appear to understand the four TU features, and trantionese

Q2. Can generative AIs use TU to pick out translated texts from non-translated texts?

- Bard can, and Bing also can with proper context information. However, Chat GPT was not able to distinguish TTs from NTTs

Q3. Does context information in prompts affect the result?

- Yes, evidently in Bard and Bing. With previous TU findings in this genre, Bard offered an insightful body of linguistic profile for the given texts, which requires further research
- However, prompts failed to improve GPT 3.5 results

References

- Baker, Mona (1996) "Corpus-Based Translation Studies: The Challenges That Lie Ahead." Pp. 175–86 in Terminology, LSP and translation, edited by H. S. (ed.). Amsterdam/Philadelphia: John Benjamins.
- Bizzoni, Yuri, Tom S. Juzek, Cristina España-Bonet, Koel Dutta Chowdhury, Josef Van Genabith, and Elke Teich (2020) "How Human Is Machine Translation? Comparing Human and Machine Translations of Text and Speech," Proceedings of the Annual Meeting of the Association for Computational Linguistics 280–90. doi: 10.18653/v1/P17.
- Blum-Kulka, Shoshana, and Eddie A. Levinson (1983) "Universals of Lexical Simplification," Strategies in Interlanguage Communication 1(983):119.
- Chesterman, Andrew (2010) "Why Study Translation Universals?" Pp. 38–48 in Kiasm, Vol. 1, edited by R. Hartama-Heinonen and P. L. H. Kukkonen. Helsinki: University of Helsinki.
- Gao, Yuan, Ruili Wang, and Feng Hou (2023) "How to Design Translation Prompts for ChatGPT: An Empirical Study."
- Goh, Gwang-yoon, Younghée Cheri Lee, and Dongyoung Kim (2016) "A Corpus-Based Study of Translation Universals in Thesis," 영어학 16(4):819–49.
- Graham, Yvette, Barry Haddow, and Philipp Koehn (2019) "Translationese in Machine Translation Evaluation."
- Gu, Wenshi (2023) "Linguistically Informed ChatGPT Prompts to Enhance Japanese-Chinese Machine Translation: A Case Study on Attributive Clauses," 1–8.
- Jiao, Wenxiang, Wenxuan Wang, Jen-tsu Huang, Xing Wang, and Zhaopeng Tu (2023) "Is ChatGPT A Good Translator? Yes With GPT-4 As The Engine."
- Kajzer-Wietrzny, Marta (2015) "Simplification in Interpreting and Translation," Across Languages and Cultures 16(2):233–55.
- Laviosa, Sara (1998) "The Corpus-Based Approach: A New Paradigm in Translation Studies," Meta 43(4):1–6. doi: 10.7202/003424ar.
- Lee, Yong-hun, and Jin Yim (2019) "A Multifactorial Analysis of Translation Universals in Management Forewords of Sustainability Reports," English Language and Linguistics 25(3):79–105.
- Lee, Younghée Cheri (2016) "The Hallmarks of Expert L2 Writers' Texts Viewed through the Prism of Translation Universals: A Corpus-Based Approach to English Research," Yonsei University.
- Lee, Younghée Cheri (2018) "The Hallmarks of L2 Writing Viewed through the Prism of Translation Universals," Linguistic Research 35(Special Edition):171–205. doi: 10.17250/khisli.35..201809.007.
- Lu, Qingyu, Baopo Qiu, Liang Ding, Liping Xie, and Dacheng Tao (2023) "Error Analysis Prompting Enables Human-Like Translation Evaluation in Large Language Models: A Case Study on ChatGPT."
- Luo, Jinru, and Dechao Li (2022) "Universals in Machine Translation?," International Journal of Corpus Linguistics 27(1):31–58. doi: 10.1075/ijcl.19127.luo.
- Olohan, Maeve (2004) Introducing Corpora in Translation Studies. London/New York: Routledge.
- Park, Myoungsu, and Shinwoong Lee (2014) "Corpus-Based Investigation into Reporting That in Translated English Editorials," 영어학 6:199–217.
- Peng, Keqin, Liang Ding, Qihuang Zhong, Li Shen, Xuebo Liu, Min Zhang, Yuanxin Ouyang, and Dacheng Tao (2023) "Towards Making the Most of ChatGPT for Machine Translation," SSRN Electronic Journal. doi:10.2139/ssrn.4390455.
- Stap, David, and Ali Arabi (2023) "ChatGPT Is Not a Good Indigenous Translator," 163–67.
- Xiao, Richard (2010) "How Different Is Translated Chinese from Native Chinese?: A Corpus-Based Study of Translation Universals," International Journal of Corpus Linguistics 15(1):5–35.
- Zanettin, Federico (2014) Translation-Driven Corpora: Corpus Resources for Descriptive and Applied Translation Studies. Manchester & Kinderhook: St. Jerome.
- Zhang, Mike, and Antonio Toral (2019) "The Effect of Translationese in Machine Translation Test Sets," WMT 2019 - 4th Conference on Machine Translation, Proceedings of the Conference 1:73–81. doi: 10.18653/v1/w19-5208.
- 고광윤;이영희 (2016) "한국 신문의 영어 번역에 나타난 번역 보편소의 코퍼스 기반 분석," 비교문화연구 45:109–43.
- 박미정 (2019) "해드라인 장르 관습과 번역 규범: 해드라인 직접인용 한일번역 분석을 중심으로," 통번역학연구 23(3):61–88.
- 박미정 (2023) "생성형 AI와 기계번역 - ChatGPT 번역을 통한 한일통역교육 고찰," 27(3):27–56.
- 이장수 (2014) "단자원통계분석법을 활용한 번역보편소 사례연구," 번역학연구 15(3):21–32.
- 임진 (2019) "지속가능성 보고서 한영번역의 번역보편소 양상 고찰," 번역학연구 20(5):131–62.
- 최희경 (2016) "한영 문학 번역 문제 연구: 코퍼스 기반 분석을 중심으로," 번역학연구 17(3):193–216.

AI 학습용 한중 번역문의 자연스러움(NATURALNESS) 평가

황은하(배재대), 비립(연세대)

1

목차

1. 들어가는 말
2. 연구 대상과 방법
3. 중국어 번역문의 자연스러움 평가
4. 나오는 말

2

연구 목적

- 이 연구는 한중 번역문의 품질 평가 연구로, 말뭉치언어학적이고 양적인 방법론을 이용해 한중 번역문의 자연스러움(naturalness), 즉 얼마나 중국어다운지에 대해 평가하는 것을 목적으로 한다.
- 이를 위해, 본 연구는 중국 원어민이 생산한 중국 신문 기사 말뭉치를 참조물로 삼아, 인간 독자를 위한 국내 주요 신문사의 중국어 기사 번역문과 AI HUB에서 배포하는 한중 병렬말뭉치의 중국어 번역문을 형태, 통사 층위에서 각각 비교 분석하여 그 자연성(naturalness)을 평가한다.
- 번역문 생성의 목적에 따른, 한중 번역문의 두 가지 유형
 - 인간 독자를 위한 번역문(이하 ‘독자용 번역문’이라고 줄여서 씀)
 - 주요 신문사의 신문 기사 번역문
 - 인공 지능 학습용 번역문(이하 ‘AI용 번역문’이라고 줄여서 씀)
 - 한국정보화진흥원(NIA)의 한중 인문/사회, 과학/기술, 방송 대화 등 여러 장르 및 도메인의 병렬말뭉치(2021~2022), AIHUB에서 공개
 - 구축 목적이 언어 연구가 아닌 ‘인공지능 학습용 빅데이터’ 구축이고, 단기간 대규모 번역문 생성의 과제 특성상 번역문의 생성, 검수 방식 등이 기존의 신문이나 소설 번역, 즉 인간 독자를 상정한 번역문과는 다를 것으로 예상

3

The screenshot shows the homepage of the AI Hub website. At the top, there is a search bar and a navigation menu. Below the header, there is a large section titled "AI 학습용 병렬말뭉치" (AI Learning Parallel Corpus). This section lists several datasets, each with a thumbnail, title, download count, and file size. The datasets include:

- 발화유형(문어/구어/채팅)별 기계번역 병렬 말뭉치 (BETA)
- 한국어-다국어 번역 말뭉치 (인문학) (BETA)
- 한국어-다국어 번역 말뭉치 (기초과학) (BETA)
- 한국어-다국어(영어 제외) 번역 말뭉치 (기술과학) (BETA)
- 교육용 한국인의 중국어-일본어 음성 데이터 (BETA)
- 관광지 소개 다국어 번역 데이터 (BETA)
- 한-영 및 한-중 음성발화 데이터 (BETA)
- 한국어-중국어 번역 말뭉치 (사회과학) (BETA)
- 한국어-중국어 번역 말뭉치 (기술과학) (BETA)
- 방송 콘텐츠 한-중, 한-일 번역 병렬 말뭉치 데이터
- 다국어 구어체 번역 병렬 말뭉치 데이터 (BETA)
- 일상생활 및 구어체 한-중, 한-일 번역 병렬 말뭉치 데이터 (BETA)
- 전문분야 한영 말뭉치 (BETA)
- 한국어-영어 번역(병렬) 말뭉치 (BETA)

Each dataset entry includes a "다운로드" button and some additional text like "경신년월: 2023-05" and "구축년도: 2022".

연구의 필요성

- AI용 언어 빅데이터에 대한 품질 평가의 필요
 - 인공지능 학습용 언어 데이터의 품질과 인공지능의 언어 능력과의 높은 상관성을 고려함
- AI용 언어 빅데이터의 언어학적 활용 가능 여부에 대한 검토의 필요
 - 병렬말뭉치는 대조언어학, 번역학, 외국어로서의 한국어 교육 등 응용언어학적 연구에 매우 필요함
 - 병렬말뭉치에 기반한 한외 언어간 연구에서 한중 언어쌍에 대한 연구가 가장 활발하지만, 대부분 개인이 구축한 소규모 말뭉치, 또는 규모가 꽤 크더라도 특정 장르에 편중되어 있거나 번역자 규모가 매우 적어서(최소 1인) 그 대표성이 의구심이 드는 경우가 적지 않음(황은하, 2016, 2019)
 - 언어학적 연구를 위한 말뭉치는 특히 언어 자료의 실제성(authentic)이 중요하며, 자연성(naturalness)도 그에 포함된다고 봄

5

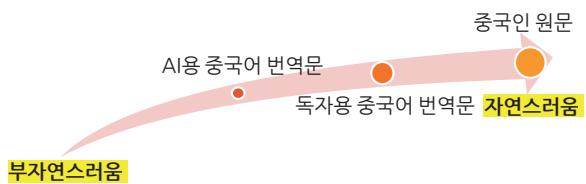
자연스러움(naturalness)이란?

- 제2언어의 원어성 측정, 번역문 평가 기준의 하나
- 원어민 생성 텍스트를 기준으로 한 정의
 - Barnavel(1980), 자연스러운 형태의 대상 텍스트를 사용하는 것
 - New mark(1988a), 일상 언어(목표 언어의 문법, 관용어 및 단어를 의미함)로 작성된 번역
 - Nida와 Taber(1969), 번역문이 번역문으로 읽히지 않을 때, 비로소 좋은 번역이라고 봄
- 출발언에 대한 충실패을 기준으로 한 정의
 - Beekman, Callow(1983), 글의 '용이함'. 텍스트의 의미를 쉽게 이해하는 것과 텍스트의 자연성 사이에는 상관관계가 있다고 봄.
- 자연스러움과 부자연스러움의 연속체로 보는 관점
 - Herman(1999), 번역은 사회문화적 활동이므로 자연성과 비자연성을 구별하는 이분법을 가정은 적절하지 않다고 봄. "자연스러움을 결정하는 것은 언어의 규범"

➔ 본 연구에서의 자연스러움(naturalness): 번역문이 형태, 통사, 의미역의 층위에서 중국인이 생산한 중국어 원문과 닮아 있는 정도로, 형태, 통사, 의미적 층위에서 모두 측정이 가능하고 필요함

6

가설과 문제 제기



- 가설
 - AI 학습용 번역문은 자연스러움의 정도가 중국인이 생성한 원문이나 독자용 번역문과 다르다.

- 문제 제기
 - AI 학습용 병렬말뭉치의 번역문은 형태적으로 얼마나 자연스러운가?
 - AI 학습용 병렬말뭉치의 번역문은 통사적으로 얼마나 자연스러운가?

→ AI용 한중 번역문이 형태, 통사적으로 전문가에 의한 인간 독자용 번역문, 중국인에 의한 중국어 원문과 얼마나 같고 다른지를 확인할 수 있을 것으로 기대함

7

2. 연구 대상과 방법

- 2.1. 연구 대상
- 2.2. 연구 방법

8

연구 대상 1/4

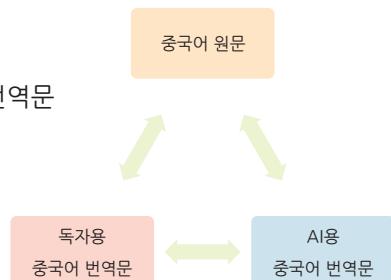
• 비교말뭉치의 구성

• 평가 대상 중국어 번역문

- AI HUB의 한국어-중국어 번역 말뭉치(사회과학)(2020)의 중국어 번역문
- 주요 일간지의 독자용 중국어 번역문

• 참조 대상 중국어 원문

- THUCTC(중국 칭화대 신문 기사) 중국어 원문



• 구성 원리

- 매체, 규모, 기사 도메인을 통일함
- 원문과 번역의 생성 시기를 2000년대 이후로 최대한 비슷하게 구성함

9

연구 대상 2/4 - AI용 중국어 번역문 말뭉치

- 활용 - 민간 및 기관의 한국어-중국어 AI 자동번역 기술 개발 등

- 매체 - 신문 기사

- 시기 - 번역문 2020년

- 규모 - 분야별 20~25만 문장 중 8만 문장씩 추출

- 분야 - 사회, 경제, 문화 분야

- 기타 특징

- 한글 원문 경제: 평균 15어절(5~30어절), 개인 정보 포함 문장 제외, 맞춤법 검사기를 이용한 오류 수정, 주어, 목적어가 없는 문장 일부 배제

- 문장 단위(텍스트, 문단 단위가 아님)

- 저자: 자동 번역/크라우드 워커에 의한 초벌 번역 → MTPE(Machine Translation Post Editing) 1단계 또는 오류 수정(필요에 따라 2단계 수행)

- 번역 검수: 전체의 5% 샘플링 검수, 블랙박스 평가(원문과 번역문상의 정확성 평가)

- TQA 질적 측정 기준

- 기능적, 실용적, 텍스트적 적합성/전문 어휘, 내용 적합성, 비전문 어휘, 내용 적합성, 규범적 및 문체적 적합성

➔ 엄밀한 품질 관리를 거친 것을 확인할 수 있음

| | 사회 | 경제 | 문화 | 합계 |
|----|-----------|-----------|-----------|-----------|
| 문장 | 80,000 | 80,000 | 80,000 | 240,000 |
| 어휘 | 2,271,219 | 2,175,714 | 2,182,181 | 6,629,114 |

10

연구 대상 3/4 - 독자용 중국어 번역문 말뭉치

- 저자 - 한중 번역 전문가
- 매체 - 조선일보, 중앙일보, 동아일보, 아주일보 등 주요 신문사
- 시기 - 원문과 번역문 2023년
- 규모 - 분야별 8만 문장씩 선별
- 분야 - 사회, 경제, 문화
- 특징 - 독자 수용성이 높음

| 문장 | 사회 | 중앙일보 | 동아일보 | 조선일보 | 아주일보 | 소계 | 합계 |
|----|----|---------|---------|---------|---------|-----------|-----------|
| | 경제 | 32,733 | 29,807 | 13,370 | 4,090 | 80,000 | 240,000 |
| | 문화 | 13,947 | 38,308 | 27,221 | 524 | 80,000 | |
| 여획 | 사회 | 641,934 | 551,712 | 754,945 | 204,812 | 2,153,403 | |
| | 경제 | 950,092 | 750,733 | 330,027 | 110,223 | 2,141,075 | 6,121,454 |
| | 문화 | 370,177 | 825,099 | 616,509 | 15,191 | 1,826,976 | |

→ 매체, 규모, 분야: AI용과 동일하게 구성함

→ 시기: AI용과 비슷하게 구성함

11

연구 대상 4/4 - 참조용 중국어 원문 말뭉치

- THU-CTC(THU Chinese Text Classification), 중국 청화대학교(清华大学) 자연어처리연구실의 중국어 텍스트 분류 도구 패키지

- 저자 - 중국인 원어민

- 매체 - 시나뉴스(新浪新闻)의 신문 기사

- 시기 - 2005년~2011년

- 규모 - 74만 건(2.19 GB) 중 분야별 8만 문장씩 추출

- 분야 - 경제, 복권, 부동산, 주식 등 14개 분류

| | 사회 | 경제 | 문화 | 합계 |
|----|-----------|-----------|-----------|-----------|
| 문장 | 80,000 | 80,000 | 80,000 | 240,000 |
| 여획 | 1,968,679 | 2,241,840 | 2,308,816 | 6,519,335 |

→ 매체: 신문 기사로 평가 대상 번역문 말뭉치와 일치함

→ 시기: 평가 대상 번역문 말뭉치의 원문 및 생성 시기와 비슷함

→ 규모: 평가 대상 말뭉치와 동일하게 구성함

→ 내용: 사회, 경제, 문화 분야 신문 기사를 선별하여 평가 대상 번역문 말뭉치와 일치시킴

12

연구 방법 1/3 - 선행 연구

- 胡韧奋, 肖航. 面向二语教学的汉语搭配知识库构建及其应用研究. 语言文字应用. 2019(1).
- 胡韧奋. 基于搭配的句法复杂度指标及其与汉语二语写作质量关系研究. 语言文字应用, 2021(1).
- Renfen Hu, Jifeng Wu, and Xiaofei Lu. Word-combination-based Measures of Phraseological Diversity, Sophistication and Complexity and Their Relationship to L2 Chinese Proficiency and Writing Quality. *Language Learning*, 2022, 72(4).

→ 제2 언어로서의 중국어의 글쓰기에 대한 평가 연구

→ 형태, 통사, 의미 층위의 평가 기준을 포함함

13

연구 방법 2/3 - 평가 지표

- 형태적 자연스러움 평가
 - 품사별 다양성(TTR)
- 통사적 자연스러움 평가
 - 문장 단위: 평균 문장 길이(MLS, Mean Length of Sentence)
 - 절 단위: 문장당 절의 규모, 절의 평균 길이
 - 구 단위: 의존 관계 분석을 통한 유형별 비율 및 다양성(TTR)

14

연구 방법 3/3 - 절차

- Che, W., Feng, Y., Qin, L., & Liu, T. (2020). N-LTP: An open-source neural language technology platform for Chinese. arXiv preprint arXiv:2009.11616.

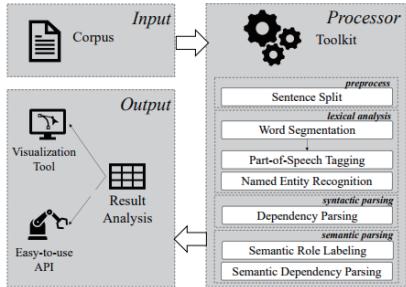
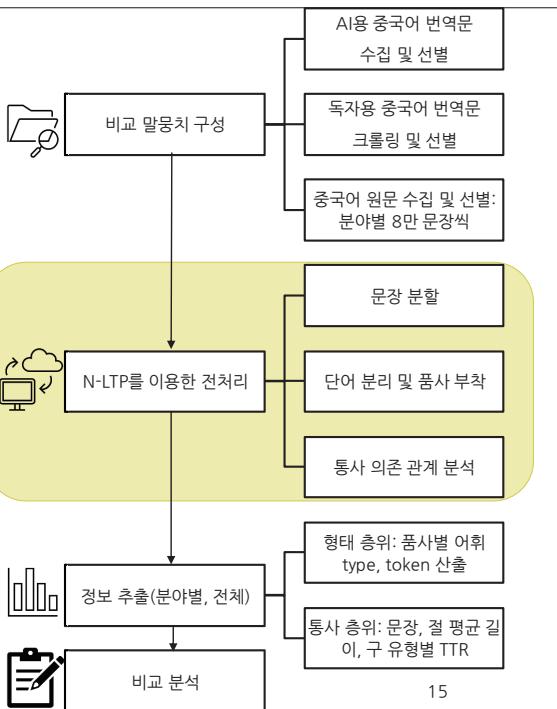


Figure 1: Workflow of the N-LTP. N-LTP takes the Chinese corpus as input and output the analysis results including lexical analysis, syntactic parsing, and semantic parsing. In addition, we provide the visualization tool and easy-to-use API to help users easily use N-LTP.



15

3. 중국어 번역문의 자연스러움 평가

- 3.1. 형태적 자연스러움 평가
- 3.2. 통사적 자연스러움 평가

16

형태적 자연스러움 평가 1/3

- 분석 도구: Python, SPSS

- 중국어 품사 분류

- 명사(일반명사/방향명사/시간명사/장소명사), 대명사, 수량사, 고유명사(인명, 조직명칭, 지리명칭, 기타 고유명사), 동사, 형용사, 계사, 관형사, 부사(접속사, 일반부사, 의성어), 감탄사
- 접사(접두사, 접미사), 관용구, 약어
- 전치사, 보조사
- 문장부호, 숫자, 외국어

- 품사별 어휘 TTR 산출 및 비교 분석

- TTR = Type / Token
- 다양성에 대한 지표, 1에 가까울수록 다양성이 높음
→ 본 연구에서는, 표준화된 TTR 공식 사용
- Std. TTR = 100만 단어 당 Type 수 / 100만 단어당 Token 수)**
- 말뭉치 규모가 커질수록 TTR값이 낮아지는 scale 효과를 조정하고
→ 말뭉치 크기의 격차를 줄이기 위한 조정



형태적 자연스러움 평가 2/3 - 체언과 용언 TTR

| 품사 | 세분류 | AI 학습용 중국어 번역문 | | | 독자용 중국어 번역문 | | | 중국어 원문 | | | |
|-----|---------|----------------|--------|--------|-------------|--------|--------|--------|--------|--------|--------|
| | | 사회 | 경제 | 문화 | 사회 | 경제 | 문화 | 사회 | 경제 | 문화 | |
| 명사 | 일반명사 | 0.0147 | 0.0160 | 0.0291 | 0.0235 | 0.0202 | 0.0418 | 0.0349 | 0.0152 | 0.0303 | |
| | 방향명사 | 0.0031 | 0.0032 | 0.0033 | 0.0039 | 0.0040 | 0.0064 | 0.0059 | 0.0033 | 0.0047 | |
| | 시간명사 | 0.0082 | 0.0098 | 0.0140 | 0.0096 | 0.0067 | 0.0207 | 0.0131 | 0.0093 | 0.0154 | |
| | 장소명사 | 0.0155 | 0.0110 | 0.0136 | 0.0166 | 0.0122 | 0.0210 | 0.0173 | 0.0116 | 0.0140 | |
| 대명사 | 대명사 | 0.0037 | 0.0043 | 0.0031 | 0.0042 | 0.0056 | 0.0053 | 0.0035 | 0.0044 | 0.0025 | |
| | 수량사 | 양사(분류사) | 0.0030 | 0.0034 | 0.0058 | 0.0040 | 0.0038 | 0.0067 | 0.0049 | 0.0039 | 0.0050 |
| | 고유명사 | 인명 | 0.1376 | 0.1552 | 0.1155 | 0.0971 | 0.1420 | 0.1615 | 0.1127 | 0.1529 | 0.0709 |
| | 조직명칭 | 0.0168 | 0.0714 | 0.1961 | 0.0174 | 0.0331 | 0.0843 | 0.0745 | 0.0328 | 0.1006 | |
| 동사 | 지리명칭 | 0.0392 | 0.0376 | 0.0764 | 0.0278 | 0.0210 | 0.0619 | 0.1261 | 0.0292 | 0.0373 | |
| | 기타 고유명칭 | 0.1223 | 0.1383 | 0.1964 | 0.1451 | 0.1026 | 0.2032 | 0.2464 | 0.0698 | 0.1449 | |
| | 형용사 | 형용사 | 0.0065 | 0.0068 | 0.0106 | 0.0104 | 0.0092 | 0.0171 | 0.0166 | 0.0098 | 0.0143 |
| | 계사 | "是" | 0.1780 | 0.0949 | 0.0803 | 0.1685 | 0.0976 | 0.0918 | 0.1567 | 0.0797 | 0.0931 |

형태적 자연스러움 평가 3/3 - 체언과 용언 TTR 비교

- 명사 TTR: 사회 분야의 텍스트에서 주로 차이를 보임
 - 일반명사: $0.0147 < 0.0235 < 0.0349 = 1 < 1.60$ 배 < 2.36배
 - 방향명사: $0.0031 < 0.0039 < 0.0059 = 1 < 1.24$ 배 < 1.89배
 - 양사: $0.0030 < 0.0040 < 0.0049 = 1 < 1.34$ 배 < 1.62배
 - 동사 TTR: 모든 분야의 텍스트에서 차이를 보이며, 사회 분야의 격차가 특히 큼
 - 사회: $1 < 1.59$ 배 < 2.54배
 - 경제: $1 < 1.35$ 배 < 1.43배
 - 문화: $1 < 1.61$ 배 < 1.34배
 - 형용사 TTR: 모든 분야의 텍스트에서 차이를 보임
 - 사회: $1 < 1.65$ 배 < 2.78배
 - 경제: $1 < 1.35$ 배 < 1.04배
 - 문화: $1 < 1.53$ 배 < 1.06배
- 품사별 어휘 다양성은 대체로 AI 학습용 번역문 < 독자용 번역문 < 중국어 원문의 순서로 나타나 예상을 벗어나지 않음.
→ 특히 사회 분야의 격차가 크며, 일반명사, 동사, 형용사는 모두 2.3배 이상의 격차를 보임
→ 다만, 독자용 번역문이 중국어 원문보다 어휘 다양성이 높은 경우가 있는데 경제, 문화 분야의 형용사와 문화 분야의 동사에 한함. 이는 한국어 원문의 영향을 받은 것으로 미루어 짐작할 수 있음

AI 학습용 번역문 - 독자용 번역문 - 중국어 원문

19

통사적 자연스러움 평가 1/8

- 분석 도구: N-LTP
- 측정 지표
 - 문장 단위: 평균 문장 길이(MLS, Mean Length of Sentence)
 - 단위: 자(字)
 - 절 단위: 문장당 절의 규모, 절의 평균 길이
 - 중국어 절 표지, 반점(,), 쌍반점(세미콜론, :) 활용
 - 구 단위: 의존 구문 분석을 통한 유형별 비율 및 다양성(TTR)
 - 구 단위 유형: 주술관계(主谓关系), 동목관계(动宾关系), 간접목적어관계(间宾关系), 전치 목적어(前置宾语), 겹어(兼语), 수식-피수식관계(定中关系), 부사어-중심어 구조(状中结构), 술어-보어 구조(动补结构), 병렬관계(并列关系), 전치사-목적어 관계(介宾关系), 좌측 부가 관계(左附加关系), 우측 부가 관계(右附加关系), 독립구조(独立结构), 핵심관계(核心关系), 문장부호 관계(标点关系)

(이상 15종 중 14종이 관찰됨)

20

통사적 자연스러움 평가 2/8 - 문장 길이 1

| | AI 학습용 중국어 번역문 | | | 독자용 중국어 번역문 | | | 중국어 원문 | | |
|-------------|----------------|-------|-------|-------------|-------|-------|--------|-------|-------|
| | 사회 | 경제 | 문화 | 사회 | 경제 | 문화 | 사회 | 경제 | 문화 |
| I00(~10자) | 1 | 1 | 4 | 1036 | 390 | 1889 | 2999 | 1695 | 4823 |
| I10(10~19자) | 1087 | 879 | 1794 | 6651 | 4561 | 11054 | 11652 | 6412 | 8920 |
| I20(20~29자) | 10578 | 10276 | 13900 | 12939 | 11956 | 17455 | 17042 | 12408 | 12240 |
| I30(30~39자) | 17338 | 18214 | 18762 | 15653 | 15818 | 16273 | 15920 | 14207 | 13037 |
| I40(40~49자) | 16708 | 18194 | 16710 | 13923 | 15167 | 12589 | 11730 | 12580 | 11181 |
| I50(50~59자) | 13113 | 13593 | 12561 | 10354 | 11809 | 8237 | 7823 | 10168 | 8840 |
| I>60(60자~) | 21175 | 18843 | 16269 | 19444 | 20299 | 12503 | 12834 | 22530 | 20959 |

→ AI용 데이터의 원문을 평균 15어절로 통제한 점을 감안하더라도, 문장길이별 분포에 적지 않은 차이가 관찰됨

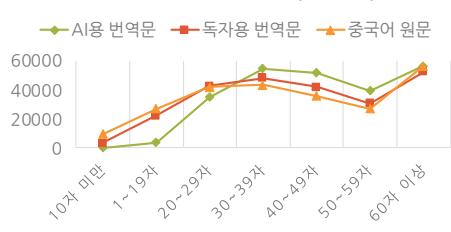
| | AI용 비율 | 독자용 비율 | 원문 비율 |
|-------------|--------|--------|--------|
| I00(~10자) | 0.0% | 1.4% | 4.0% |
| I10(10~19자) | 1.6% | 9.3% | 11.2% |
| I20(20~29자) | 14.5% | 17.6% | 17.4% |
| I30(30~39자) | 22.6% | 19.9% | 18.0% |
| I40(40~49자) | 21.5% | 17.4% | 14.8% |
| I50(50~59자) | 16.4% | 12.7% | 11.2% |
| I>60(60자~) | 23.5% | 21.8% | 23.5% |
| | 100.0% | 100.0% | 100.0% |

21

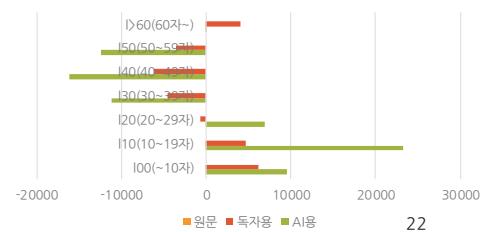
통사적 자연스러움 평가 3/8 - 문장 길이 2

- 문장길이별 분포 곡선은 비슷한 모양을 보여, 문장길이별 분포 측면에서 자연스러움에 가까움
- 단, 중국어 원문은 60자 이상의 긴 문장이 가장 많은 (23.5%) 다음, 주로 30~39자 구간이 18.0%, 20~29자 구간 17.4%, 40~49자 구간 14.8%를 차지하여 전체의 50.2%를 차지하며, 독자용 또한 비슷한 양상을 보이는 데 반해,
- AI용 번역문은 30~39자 구간 22.6%로 가장 많고, 40~49자 구간 21.5%로, 전체의 44% 가량을 차지함
- AI용 한국어 원문의 문장 평균 길이가 15어절로 통제된 점을 고려하더라도 문장의 길이가 지나치게 긴 편으로 자연스러움에서 멀어진 것으로 판단됨
- 1차적으로 기계번역의 축자적 번역에서 기인하고, 전문가의 개입으로 인한 일부 중국어 구의 단어 대응 등의 조정이 이루어지지 않은 데서 2차적 원인을 찾을 수 있음

문장 길이별 분포 비교(전분야)



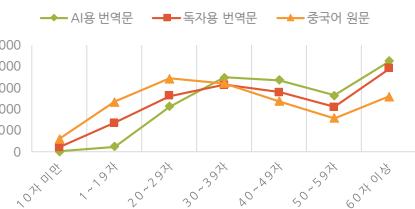
중국어 원문 기준 문장길이별 규모 차이



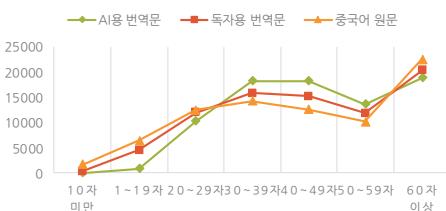
22

통사적 자연스러움 평가 4/8 - 문장 길이 3

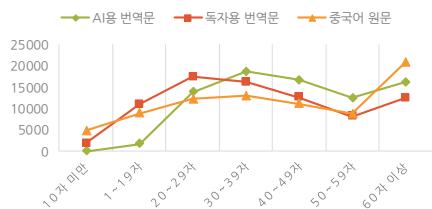
문장 길이별 분포 비교(사회)



문장 길이별 분포 비교(경제)



문장 길이별 분포 비교(문화)



- 분야별로 차이를 보임
 - 사회, 문화 분야의 AI용 번역문은 원문과는 곡선에 차이가 크나 독자용과는 상당히 가까움
 - 경제 분야는 원문과 차이는 있으나 비슷한 곡선을 보임
- 그 원인 분석은 분야에 따른 한국어 원문의 구문적 특성 차이, 실제 AI용 번역문과 독자용 번역문의 비교 등을 통해서만 명확히 할 수 있음

23

통사적 자연스러움 평가 5/8 - 절 단위

- 중국어의 절 구분자: 전각 기호의 반점(,), 쌍점(:), 쌍반점(:)

| | AI 학습용 중국어 번역문 | | | 독자용 중국어 번역문 | | | 중국어 원문 | | |
|-------------------|----------------|--------|--------|-------------|--------|--------|--------|--------|--------|
| | 사회 | 경제 | 문화 | 사회 | 경제 | 문화 | 사회 | 경제 | 문화 |
| 절의 수(평균) | 1.112 | 1.110 | 1.023 | 2.639 | 2.520 | 2.169 | 3.180 | 3.209 | 3.331 |
| 절 평균 길이(자)(평균) | 45.147 | 44.287 | 44.622 | 18.726 | 20.732 | 20.018 | 12.077 | 15.529 | 13.859 |
| 절의 수(중위값) | 1.000 | 1.000 | 1.000 | 2.000 | 2.000 | 2.000 | 3.000 | 3.000 | 3.000 |
| 절 평균 길이(자)(중위값) | 42.000 | 42.000 | 42.000 | 16.500 | 18.000 | 16.500 | 11.000 | 14.000 | 12.333 |
| 절의 수(표준편차) | 0.320 | 0.318 | 0.158 | 1.364 | 1.275 | 1.166 | 1.686 | 1.715 | 2.051 |
| 절의 평균 길이(자)(표준편차) | 17.602 | 16.376 | 17.367 | 9.800 | 10.694 | 12.531 | 5.271 | 10.277 | 8.100 |

- 문장 당 절 평균 규모는 큰 차이를 보임. 약 1: 2: 3
 - 절의 평균 길이는 역순으로 약 3:2:1
- AI용 중국어 번역문의 절 구분이 제대로 되어 있는지 의문스러움
- 또는, AI용 중국어 번역문에서 문장부호가 중국어의 전각이 아닌 반각 문장부호의 사용으로 인한 것으로 추측됨

24

통사적 자연스러움 평가 6/8 - 구(短語) 단위 1

- 구 단위 의존 관계 분석 방법

- 胡韧奋(2021)을 인용함
 - 직역: “그는 교육 사업을 열렬히 사랑하고 많은 교과서를 집필했으며 학교를 위해 적지 않은 젊은 교사를 양성했다.”
 - 기계번역: “그는 교육을 사랑하고 많은 교과서를 집필했으며 학교에서 많은 젊은 교사를 양성했다.”
 - 전문가 번역 1: “그는 교육에 열과 성을 다하고 교재를 다수 집필했을 뿐만 아니라 학교를 위해 후학 양성에도 힘썼다.”
 - 전문가 번역 2: “그는 교육 사업을 사랑하여 많은 교재를 편찬하였고 학교를 위해 젊은 교사를 많이 육성하였다.”

⇒ 문장 층위: 문장 1개

⇒ 절 층위: 절 3개

⇒ 구 층위: 5가지 유형의 구 7개로 구성

- 동사-목적어 관계: 3개
- 관형어-중심어 관계: 1개
- 부사어-중심어 관계: 1개
- 수량사(단위성 의존명사)-명사: 1개
- 전치사-동사 관계: 1개

⇒ 8종에 대해서만 분석함

(1) 他热爱教育事业，编写过很多部教材，也给学校培养了不少年轻老师。

句层面: 1 个大句， 3 个小句， 3 个 T 单位；

短语层面: 5 种类型的 7 个搭配，包括：

动宾搭配：“热爱 事业”“编写 过 教材”“培养 了 老师”

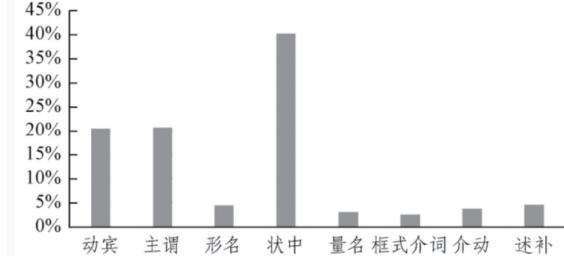
形名搭配：“年轻 老师”

状中搭配：“也 培养”

量名搭配：“部 教材”

介动搭配：“给 X 培养 了”

* 胡韧奋(2021)에서 인용함



25

통사적 자연스러움 평가 7/8 - 구(短語) 단위 2

- 측정 지표

- 구 의존관계 유형별 비율
- 구 의존관계 유형별 TTR

• 구 의존관계 유형 14종

• 참고: 附录一 LTP4.4.1.4 文档
https://ltp.readthedocs.io/zh_CN/latest/appendix.html

| dep_tag | 구 의존관계 | 예시 | AI용 비율 | 독자용 비율 | 원문 비율 |
|---------|------------------|------------------|--------|--------|-------|
| SBV | 주어-서술어(主谓关系) | 我送她一束花 (我 <- 送) | 0.088 | 0.074 | 0.071 |
| VOB | 동사-직접목적어(动宾关系) | 我送她一束花 (送 -> 花) | 0.101 | 0.101 | 0.103 |
| IOB | 동사-간접목적어(间宾关系) | 我送她一束花 (送 -> 她) | 0.001 | 0.000 | 0.000 |
| FOB | 목적어-서술어(前置宾语) | 他什么书都读 (书 <- 读) | 0.006 | 0.006 | 0.008 |
| DBL | 겸어(목적어 겸 주어)(兼语) | 他请我吃饭 (请 -> 我) | 0.005 | 0.003 | 0.004 |
| ATT | 관형어-중심어(定中关系) | 红苹果 (红 <- 苹果) | 0.216 | 0.258 | 0.254 |
| ADV | 부사어-중심어(状中结构) | 非常美丽 (非常 <- 美丽) | 0.164 | 0.149 | 0.156 |
| CMP | 서술어-보증어(动补结构) | 做完了作业 (做 -> 完) | 0.016 | 0.013 | 0.012 |
| COO | 병렬 구조(并列关系) | 大山和大海 (大山 -> 大海) | 0.072 | 0.056 | 0.058 |
| POB | 전치사-목적어(介宾关系) | 在贸易区内 (在 -> 内) | 0.037 | 0.044 | 0.047 |
| LAD | 좌측 부가 관계(左附加关系) | 大山和大海 (和 <- 大海) | 0.007 | 0.010 | 0.012 |
| RAD | 우측 부가 관계(右附加关系) | 孩子们 (孩子 -> 们) | 0.061 | 0.076 | 0.094 |
| HED | 문장 핵심 (独立结构) | 两个单句在结构上彼此独立 | 0.055 | 0.051 | 0.041 |
| WP | 문장부호와 주변 용언간 관계 | | 0.170 | 0.159 | 0.139 |

26

통사적 자연스러움 평가 8/8 - 구(短語) 단위 3

- 전체적 곡선이 대체로 비슷함
- 독자용 번역문 곡선이 중국어 원문과 보다 유사하고, AI용 번역문 곡선이 유사도가 낮음
- ➔ 구 유형별 분포 비율의 차이가 크지 않으나, AI용 번역문의 구 단위 자연스러움이 상대적으로 낮은 것으로 판단됨
- ➔ 원문을 기준으로, AI용 번역문의 비율이 높은 의존관계 유형
 - 주어-서술어, 문장부호관계, 병렬구조
- ➔ 원문을 기준으로, AI용 번역문의 비율이 낮은 의존관계 유형
 - 관형어-중심어, 전치사-목적어, 우측 부가 관계

구 의존관계 유형별 분포 비교



구 의존관계 유형별 분포의 차이 비교



27

4. 나오는 말

- 요약
- 연구의 의의
- 남은 과제

28

4. 나오는 말 - 요약

- AI용 번역문의 형태, 통사 층위에서 자연스러움의 정도가 다를 것이라는 가설이 입증됨
- 형태, 통사적으로 독자용 번역문이 중국어 원문에 보다 가까우며, AI 번역문은 중국어 원문에 서 좀더 멀어져 있음
- 형태적 층위에서 품사별 어휘 다양성의 차이를 확인할 수 있었음
 - AI용 <독자용 <원문
 - 명사(일반명사, 방향명사, 수량사)가 원문에 비해 2배 이상 어휘 다양성이 낮음
 - 분야에 따라서 차이를 보이며, 특히 사회 분야 기사의 동사, 형용사의 어휘 다양성도 2배 이상 낮게 나타남
- 통사적 층위에서 문장 길이, 절의 규모, 구 의존관계 유형별 비율 등의 차이를 확인함
 - AI용 <독자용 <원문
 - AI용 번역문의 경우 입력문의 평균 문장 길이 통제를 고려하더라도 중국어 원문과 비교해 지나치게 긴 편
 - 절 단위의 경우, 문장당 규모가 큰 차이를 보임. 1:2:3
 - 구 의존관계의 유형별 비율도 원문을 기준으로 AI용 번역문의 비율이 차이를 보였으며, 더 많이 또는 더 적게 사용된 의존관계 유형을 밝힘

29

4. 나오는 말 - 연구의 의의

- 번역문의 자연스러움에 대한 형태, 통사적 층위의 양적 평가 연구
- AI 학습용 대규모 병렬말뭉치의 언어학적 활용을 위한 번역문 평가로서의 의의
- 비교말뭉치를 활용한 번역학 연구의 사례를 보임
- 기계번역 포스트 에디팅(MTPE)에 관한 시사점
 - 어휘의 다양성, 품사별 어휘의 다양성
 - 문장의 길이의 조정, 구문 구조의 변조 등
- 한국어-외국어 AI 번역 개발을 위한 시사점
 - 한국어 AI 모델의 다국어 텍스트 이해 및 생성 능력을 한층 더 강화하기 위해, 원어민이 수용 가능한 양질의 병렬 데이터가 절실히 필요함
 - 기계번역문의 가독성과 실제성(인간 평가 지표)을 강화하기 위해서는 고품질의 golden standard 병렬 데이터셋의 구축이 필요하며, 이를 위해서는 전문가의 개입을 늘려야 함

30

4. 나오는 말 - 남은 과제

- 구체적인 번역문 검토를 통한 질적 분석
- 의미 층위에서의 자연스러움에 대한 평가

31

참고문헌

- 남대현(2011), 원어민성 두 명의 한국어-영어 이중언어자 사례연구
- 장현주(2006), 중한 번역에서 번역투(transliteration)와 정보 배열, 중국학연구 37집, 345-364.
- 정혜연(2018), 번역의 자동평가 기계번역 평가를 인간번역 평가에 적용해보기, 통번역학연구, 22권 4호, 265-287.
- 정혜연, 최지수, 허탁성, 서수영(2018), 인간 번역평가에서 재현도(recall)의 중요성, 번역학연구, 23권 1호, 81-100.
- 정혜연 외 9인(2020), 번역자동평가에서 풀리지 않은 과제, 번역학연구, 21권 1호, 9-29.
- 정혜연, 박현일, 우경조, 서수영(2021), 임베딩을 활용한 인간번역의 자동평가 - 기계가 의미를 평가할 수 있을까, 통번역학연구, 25권 3호, 141-162.
- 최지수(2023), 인간번역의 자동평가에서 번역 텍스트 길이가 평가 타당성에 미치는 영향, 통번역학연구, 27권 2호, 145-164.
- 최효은, 이정호, 이준호(2023), 자동화된 기계학습(AutoML)을 활용한 특허 특화 번역엔진의 영한번역 성능 평가, 번역학연구 24권 2호, 101-130.
- AI HUB(2021), 한국어-중국어 번역 말뭉치 2(사회과학) 데이터 설명서
- AI HUB(2021), 인공지능 데이터 구축·활용 가이드라인: 전문분야 학습용 한-중/일 번역 말뭉치
- Nur Saptaningsih¹, Ardianna Nuraeni², Agus Dwi Priyanto³, Intan Mustika Sari⁴, Yanuarria Kukuh Perwira(), The Extent of Language Naturalness in the Translation of Children's Storybooks Produced by Novice Translators in Ganesa Library.
- Hamideh Hashemi(2009), Naturalness in translation, Art of Translation and Interpreting
- Che, W., Feng, Y., Qin, L., & Liu, T. (2020). N-LTP: An open-source neural language technology platform for Chinese. arXiv preprint arXiv:2009.11616.
- Renfen Hu, Jifeng Wu, and Xiaofei Lu. Word-combination-based Measures of Phraseological Diversity, Sophistication and Complexity and Their Relationship to L2 Chinese Proficiency and Writing Quality. Language Learning, 2022, 72(4).
- 胡韧奋. (2021). 基于搭配的句法复杂度指标及其与汉语二语写作质量关系研究. 语言文字应用.

32