

국립국어원 2023-01-07

발간등록번호
11-1371028-000938-01

2022년 말뭉치 비윤리성 분석 및 연구

연구책임자
조 태 린



제 출 문

국립국어원장 귀하

국립국어원과 체결한 연구용역 계약에 따라 '2022년 말뭉치 비윤리성 분석 및 연구'에 관한 연구 보고서를 작성하여 제출합니다.

■ 사업 기간: 2023년 8월 ~ 2023년 2월

2022년 2월 21일

연구책임자: 조태린(연세대학교)

연구 기관: 연세대학교 산학협력단

연구책임자: 조태린

공동연구원: 공나형, 김미숙, 한용운

박승희, 변순용, 김봉제

이청호, 윤기현, 김성진

보조연구원: 박예슬, 심주희, 김승래

배승희, 최민경, 이재엽

김미영

<국문 요약>

2022년 말뭉치 비윤리성 분석 및 연구

이 사업의 목적은 말뭉치의 부적절성 판정 기준을 설정하고 식별된 부적절성의 분석을 위한 세부 지침을 수립하는 한편, 이를 바탕으로 기구축 말뭉치를 검토함으로써 궁극적으로 내용 분석 표지가 부착된 부적절성 말뭉치를 구축하는 데 있다.

이를 위해 본 사업은 부적절 표현의 탐색 및 선정을 위하여 기존 선행 연구를 검토하고 현대 사회의 윤리적·도덕적 기준을 고려하여 ‘부적절성’을 재개념화하였으며 그 결과 공격성, 비하성, 편향성 중 하나라도 해당하는 표현을 ‘부적절성’으로 판정하고자 하였다. 다음으로 선정된 부적절 문장의 분류 및 분석을 위한 주석 범주로서 부적절성 실현의 명시성 여부(명시/비명시), 맥락의 긍·부정성(긍정성/부정성), 부적절성의 강도(강/약), 부적절성의 의미 영역을 제시하였다. 의미 영역의 경우 선행 연구 및 국가인권위원회 법 제2조 제3항을 참고하여 성(gender), 연령/세대, 출신(인종, 국가, 지역), 신체(장애, 질병, 외모), 문화(종교 정치), 기타 등으로 세분화하였다.

또한 상기의 분류 기준에 덧붙여 분석 경계 및 범위 설정, 중첩 시 주석 방안 등과 관련한 세부 지침을 상세화함으로써 분석 및 주석 작업에 정밀성을 기하고자 하였다. 이 외에도 대상 자료의 개인정보에 대한 비식별화 방안을 수립하고 작업함으로써 자료 개방 혹은 공유 시 제기될 수 있는 보안 문제 또한 고려하였다. 해당 지침에 근거한 작업 수행에는 2020년 이후 7개 사업에서 작업 도구로 사용되어 그 성능과 기능이 확인된 ‘아이달고나(AI달고나)가 본 과제의 필요에 맞게 개선되어 활용되었다. 특히 해당 도구를 활용하여 작업자들 간 작업 후 검수를 진행함으로써 부적절성 판정 및 주석에 통일성과 신뢰성을 도모하였다. 이러한 방법론 및 세부 지침에 근거하여 본 사업은 2020년 기구축된 ‘개체명 분석 말뭉치(웹 자료)’ 500만 어절 대상으로 작업을 수행하였고, 부적절성 영역의 비중을 고려하여 16,240개의 문장으로 구성된 부적절성 말뭉치를 가공 및 정비하였다.

본 사업은 ‘부적절성’을 재개념화하고 맥락을 고려한 작업을 통하여 명시적으로 표상된 욕설이나 혐오 표현은 물론 우회적이고 간접적으로 이루어지는 차별 및 혐오, 편향적 표현 등을 주석할 수 있는 구체적 방안을 제안하였다. 이는 범용적으로 활용 가능한 대화체 텍스트 부적절성에 대한 검증 틀의 개발 및 고도화에 기여할 수 있다는 의의를 지니며, 부적절성 강도에 대한 타당한 검증 기준을 제공함으로써 추후 부적절성 문제의 관리 역량을 강화하고자 하는 산업계에도 기여할 수 있을 것이다.

주요어: 말뭉치 비윤리성, 말뭉치 부적절성, 부적절성 판정, 부적절성 분석, 기계 학습, 인공지능 언어 능력 평가

<Abstract>

2022 Analysis and Research of Corpus Unethicality

This project aims to establish criteria for determining the inappropriateness of the corpus and to establish detailed guidelines for analyzing identified inappropriateness. Moreover, based on that, by reviewing the existing corpus, we ultimately tried to build an inappropriate corpus labeled with indicators for analyzing inappropriate content.

To this end, this project reviewed prior studies and reconceptualized ‘inappropriateness’ regarding modern society's ethical and moral standards. As a result, in this project, the expression corresponding to aggression, degradation, and bias was determined as ‘inappropriateness’. Next, as annotation categories for classification and analysis of selected inappropriate sentences, whether or not the realization of inappropriateness is explicit (disambiguation/non-disambiguation), the context (positive/negative), the intensity of inappropriateness (strong/weak), and semantic areas of inappropriateness presented. In the case of semantic areas, they were subdivided into gender, age/generation, origin (race, country, region), body (disability, disease, appearance), culture (religious politics), and others, referring to previous studies and Article 2 (3) of the National Human Rights Commission Act.

In addition to the classification criteria mentioned above, the work process was carried out by detailing the annotation plan in the case of the boundary of the analysis, the setting of inappropriate ranges, and the overlapping of categories. In addition, the security issue of the data was also considered by establishing a de-identification plan for the personal information of the target data and performing tasks accordingly. In carrying out these tasks, AI Dalgona, which has guaranteed its performance and function in seven projects since 2020, has been improved and utilized to meet the needs of this task. In particular, by conducting peer inspection of the tool, workers tried to promote unity and reliability in the judgment of inappropriateness and annotation. Based on these methodologies and detailed guidelines, this project worked on 5 million words of the ‘Named Entity Tagged Corpus (web data)’ constructed in 2020 and processed and reorganized an inappropriate corpus consisting of 16,240 sentences considering the proportion of inappropriate areas.

This project proposed a specific plan to reconceptualize ‘inappropriateness’ and comment on explicitly revealed swear words, hate expressions, and indirect expressions by considering the context of its selection and analysis. This work is meaningful because it can contribute to developing and upgrading a framework for verifying universally available text inadequacy. In addition, this project will be able to contribute to the industry that wants to strengthen its management capabilities for inappropriate issues in the future by providing reasonable verification criteria for the intensity of inappropriateness.

Key-words: Corpus unethicality, Corpus inappropriateness, Inappropriateness judgment, Inappropriateness analysis, Machine learning, AI language proficiency evaluation

차례

제1장 사업 개요

1.1. 사업의 목적 및 필요성	3
1.2. 사업의 범위	3
1.3. 연구진의 구성	5
1.4. 사업 수행 전략	7
1.5. 기대 효과	8

제2장 사업 추진 경과

2.1. 사업 수행 절차와 일정	13
2.2. 분석 대상 말뭉치와 분석 요소 선정	14
2.3. 말뭉치 부적절성 분석 작업 지침 작성	20
2.4. 말뭉치 부적절성 분석 작업 도구 개발	22
2.5. 말뭉치 부적절성 분석 작업 수행	26
2.6. 부적절성 말뭉치 구축	32

제3장 사업 주요 내용

3.1. 부적절성의 개념 설정	35
3.2. 부적절성의 분석 및 주석(태깅) 단위 설정	36
3.3. 부적절성의 명시성 주석	37
3.4. 부적절성의 맥락 주석	39
3.5. 부적절성의 영역 주석	40
3.6. 부적절성의 강도 주석	42
3.7. 부적절성 말뭉치의 비식별화	43

제4장 사업 결과와 논의 및 제언

4.1. 사업 결과	47
4.2. 논의 및 제언	53

참고 문헌	57
-------------	----

부록

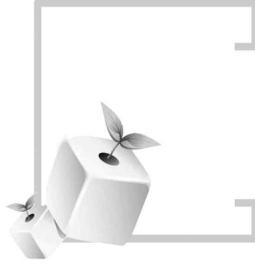
[붙임 1] 말뭉치 부적절성 분석 작업 지침_최종	63
[붙임 2] 부적절성 관련 어휘(표현) 목록	105

표 차례

<표 1> 분석/구축팀 구성 명단	6
<표 2> 사업 수행 일정	13
<표 3> ‘명시성’의 유형과 범위	37
<표 4> ‘맥락’의 유형과 범위	39
<표 5> ‘영역’의 유형과 범위	40
<표 6> ‘강도’의 유형과 범위	42
<표 7> 부적절성 문장의 ‘명시성’, ‘강도’, ‘맥락’ 관련 종합적 분포	47
<표 8> 부적절성 문장의 ‘명시성’, ‘영역’ 관련 종합적 분포 (단일 영역 기준)	49
<표 9> 부적절성 문장의 ‘명시성’, ‘영역’ 관련 종합적 분포 (복합 영역 기준)	51

그림 차례

[그림 1] 과제 수행 조직 및 인원 현황	5
[그림 2] 사업 수행 절차	13
[그림 3] 아이달고나 로그인 화면	22
[그림 4] 아이달고나 1단계 부적절성 문장 선별 화면 예시	23
[그림 5] 아이달고나 2단계 부적절성 분석 요소 주석 화면 예시	24
[그림 6] 아이달고나 2단계 검수 작업 화면 예시	25
[그림 7] 작업 도구 상 총괄 공동연구원 메모 예시	26
[그림 8] 작업자 분석 작업 현황(1월 25일 기준)	27
[그림 9] 작업자 분석 작업 현황(2월 18일 기준)	31
[그림 10] 부적절성 말뭉치 json 파일 예시	32
[그림 11] 부적절성 문장의 '명시성' 관련 분포	48
[그림 12] 부적절성 문장의 '강도' 관련 분포	48
[그림 13] 부적절성 문장의 '맥락' 관련 분포	49
[그림 14] 부적절성 문장의 '영역' 관련 분포(단일 영역 기준)	50
[그림 15] 부적절성 문장의 '영역' 관련 분포(단일 영역 기준)	52



제 1 장

사업 개요



1.1. 사업의 목적 및 필요성

- 본 사업은 말뚱치의 부적절성¹⁾ 분석을 위한 분석 방법론 및 세부 지침을 수립하고, 수립된 지침을 바탕으로 기구축 말뚱치를 검토한 후 부적절성 분석 말뚱치를 구축하는 것을 목적으로 한다. 인공지능 기술 개발 및 서비스 활용을 위한 대규모 말뚱치 구축이 활성화되고 있으나, 구축된 말뚱치 자체는 물론이고 말뚱치 활용에 있어 윤리성을 비롯한 부적절성 문제가 지속적으로 제기되고 있기 때문이다. 이에 부적절성 문제 해결을 위해 국내 표준화 및 참조 기반 자료가 될 수 있는 비윤리적 표현 관련 언어 정보를 부가한 말뚱치를 구축할 필요가 있다.

1.2. 사업의 범위

- 본 사업의 범위와 주요 내용은 다음과 같다.
 - 1) 부적절성 분석을 위한 분석 방법론 및 세부 지침 수립
 - 한국어의 특성을 반영한 분석 대상 선정 및 분석 방법론 수립
 - : 사회적 관습, 태도 등에서 비롯될 수 있는 비윤리성을 포함하는 부적절성이 드러난 대상 문장 탐색 및 선정 방안 수립
 - ※ 맥락에 따라 부적절성이 발생하는 경우의 효과적 제시 방안 포함
 - 예) 선·후행 문장 추가, 대화 또는 지시 대상 정보 추가 등
 - : 분석 경계 및 범위 설정 등에 대한 지침 상세화
 - : 대상 자료 개인정보(실명, 계정명, 소속 등) 비식별화 방안 수립
 - 부적절성 분석 세부 지침 수립

1) 본 사업은 당초 말뚱치의 '비윤리성'을 분석하고 비윤리성 분석 말뚱치를 구축하는 것으로 발주되었으나, 착수 보고 이후 분석 및 구축 대상 관련 개념이 '비윤리성'에서 '부적절성'으로 변경되었으므로 이하에서는 '부적절성'이라는 용어를 사용하되, 해당 개념과 용어가 변경되기 이전인 기술 협상 및 계약, 착수 보고 관련 사업 추진 경과를 기술하는 부분에서는 비윤리성이라는 용어를 그대로 사용한다. 부적절성 개념 관련 자세한 설명은 2장 2.3절과 붙임의 '작업 지침'을 참조할 수 있다.

: 부적절성 분석을 위한 비윤리성 영역(domain) 정의 및 범주 구체화

※ 성(gender), 연령/세대, 출신(인종, 국가, 지역), 신체(장애, 질병, 외모), 문화(종교, 정치) 등

※ 대상 문장의 부적절성 영역이 중첩되는 경우의 분석 방안 포함

: 대상 문장의 내용을 바탕으로 한 부적절성 내용 분류 표지 정의

2) 부적절성 분석 말뭉치 구축

- 분석 자료 검토 및 선별

: '20년 구축 '개체명 분석 말뭉치(웹 자료)' 500만 어절 대상

- 말뭉치 대상 부적절성 정보 분석

: 수립한 방법론에 따라 대상 문장의 부적절성 내용 분석 표지 부착

- 말뭉치로 가공 및 정비(15,000 문장 이상)

- 영역별 구축 규모 및 작업 체계 제안

: 부적절성 영역별 균등한 비윤리성 말뭉치 구축

※ 1개의 영역이 전체 구축량의 35%를 초과하지 않도록 조정

3) 납품 자료의 품질 보증 및 보완 체계 수립

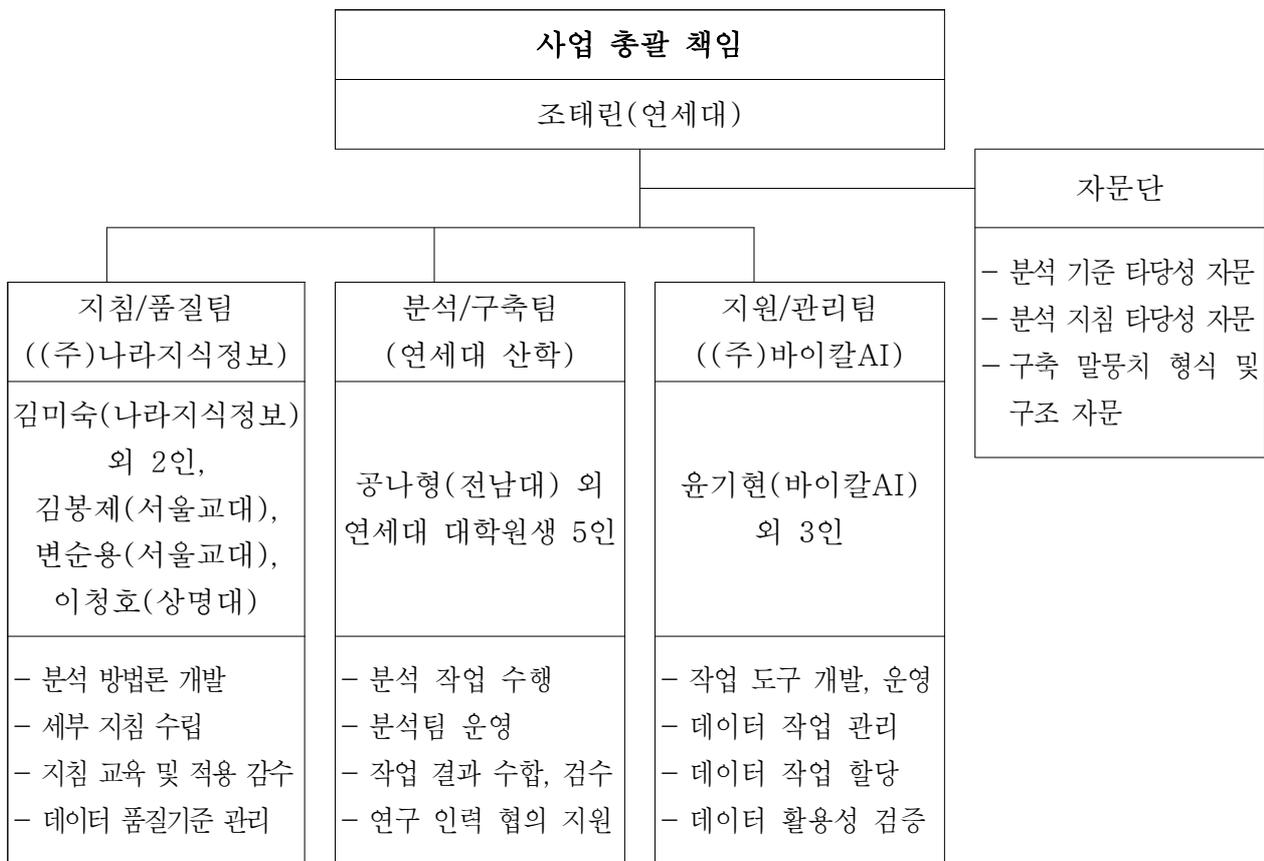
- 납품 자료에 대한 품질 보증 및 검증 체계 수립, 이행

- 최종 자료 납품 이후 보완 체계 수립·제안

1.3. 연구진의 구성

1.3.1. 과제 수행 조직 및 인원 현황

- 관련 연구 및 과제 수행 경험이 풍부한 전문가들로 다음과 같이 인력을 구성하고 업무 분장을 하였다.



[그림 1] 과제 수행 조직 및 인원 현황

1.3.2. 분석/구축팀 구성 및 운영

- 과제 수행 조직 중 분석/구축팀은 말뭉치 부적절성 분석과 부적절성 말뭉치 구축 작업을 실제로 수행한다는 점에서 효율적인 운영 및 감수 체계를 구축하는 것이 중요하였다. 이를 위해 팀별로 팀장 1명(대학원생)과 팀원(분

석 작업자, 학부생) 4명으로 구성하여 총 5개 팀을 구성하였고, 분석 작업자 1명당 최대 250,000어절을 배정하여 작업하도록 하였다. 또한 분석 작업자 1명이 추출한 비윤리적 문장은 팀원 간 상호 검토 및 팀장 검수를 통해 이견을 조정하고 합의를 도출하였다.

- 분석 및 구축 작업의 통일성과 정확성을 제고하기 위해 분석 작업자(크라우드 워커)는 연세대학교 국어국문학과 등 어문계열 학부생을 중심으로 언어학 관련 전공 기초 이수자 및 말뭉치 유경험자 우선 선발하였다.

<표 1> 분석/구축팀 구성 명단

총괄 공동연구원	팀	팀장	팀원 (분석 작업자)	학과/전공
공나형	1팀	심주희	남기웅	국어국문학과
			김나림	국어국문학과
			김예지	국어국문학과
			이연서	국어국문학과
	2팀	박예슬	김민수	국어국문학과
			김세현	교육학과
			김현진	국어국문학과
			허자인	국어국문학과
	3팀	김승래	최수연	국어국문학과
			어지영	국어국문학과
			남유림	국어국문학과
			김도아	국어국문학과
	4팀	배승희	이승아	국어국문학과
			고유빈	국어국문학과
			박지은	국어국문학과
			김보경	국어국문학과
	5팀	최민경	박영신	국어국문학과
			김나현	철학과
			신현	국어국문학과
			전현지	국어국문학과

- 팀장 선에서 결정하기 어려운 문제가 발생할 때에는 팀장 회의를 개최하여 분석/구축팀 총괄 공동연구원이 이견을 조정하고 합의를 도출하였으며, 팀장 회의를 통해서도 결정하기 어려운 문제가 발생할 때에는 공동연구원 회의를 개최하여 연구책임자가 이견을 조정하고 합의를 도출하였다.

1.4. 사업 수행 전략

- 본 사업의 과제 범위는 크게 두 가지로 구분할 수 있는데, 하나는 ‘부적절성 분석을 위한 분석 방법론 및 세부 지침을 수립하는 것이며, 다른 하나는 부적절성 분석 말뭉치를 구축하는 것이다. 이러한 두 가지 과제를 충실하게 수행하기 위해서는 관련 선행 연구 결과를 검토하고 분석하는 작업이 필요하였다. 이에 연구진은 먼저 비윤리성을 비롯한 부적절성 문제를 판단하고 분류하는 것과 관련된 기존의 다양한 연구(윤리학적 접근, 사회언어학적 접근 등)에 대한 면밀한 검토를 수행하였다. 다음으로 국립국어원에서 수행된 <말뭉치 언어의 사회적 인식 조사·분류(2021)>, <비윤리적 표현 말뭉치 연구 분석 및 시범 구축(2022)> 등 과제의 중간 및 최종 결과는 물론이고 한국지능정보사회진흥원(NIA) 등에서 수행된 관련 사업 결과를 검토하고 분석하였다. 이 밖에도 비윤리성 분석을 위해 국립국어원에서 제공하는 2020년 구축 ‘개체명 분석 말뭉치’ 관련 보고서와 기타 국립국어원 말뭉치 사업 관련 보고서들과 구축 자료 등을 참고하였다.
- 본 사업은 과제 관련 학문적, 실천적 관심과 경험이 풍부한 연구 인력과 과제 수행에 필요한 기술적 능력과 경험이 풍부한 사업 인력이 필요할 뿐만 아니라 참여 인력 간의 체계적인 업무 분담과 효율적 운영이 매우 중요하였다. 이에 본 참여 인력은 각자의 전문성과 경험을 살려 맡은 바 역할을 전담하는 분야별 책임제의 방식으로 업무를 분담하고, 주요 사안에 대해서는 연구책임자를 비롯한 공동연구원과 사업 참여 기관이 상호 협의하고 조정하는 방식으로 과제 수행의 효율성과 통일성을 확보하였다.
- 본 사업은 국립국어원에서 이루어진 기존 사업의 결과를 참조하고 역시 국립국어원에서 제공하는 기구축 말뭉치를 바탕으로 부적절성 분석 말뭉치를 구축해야 하며, 분석 범위와 기준, 유형 등에 관한 많은 조정이 필요하다는 점에서 발주 기관인 국립국어원과 긴밀하게 소통하고 협의하는 것이 필수적이었다. 본 사업의 참여 인력은 국립국어원의 각종 사업에 참여한 경험이 풍부하므로 발주 기관인 국립국어원이 요구하고 기대하는 바를 잘 이해할

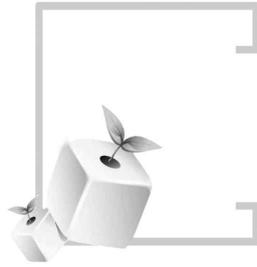
수 있었고, 사업 진행의 전반적인 과정에서 원활하고 효율적인 소통과 협의가 이루어졌다.

1.5. 기대 효과

- 본 사업의 결과물은 인공지능 상시 평가 체계 플랫폼의 평가 과제로 활용될 예정이다. 특히 국립국어원의 관련 과제인 ‘2022년 인공지능의 언어 능력 평가 체계 운영 및 말뭉치 정비’와 연계하여 인공지능의 언어 능력을 상시적으로 평가하고 기구축 말뭉치를 정비하는 데 활용할 수 있을 것으로 기대된다.
- 본 사업의 결과물은 대화 데이터를 다루는 기업의 비윤리성을 비롯한 부적절성 문제의 관리 역량을 강화하는 것에 기여할 수 있을 것이다. 온라인 공간에서 발생할 수 있는 비윤리적, 공격적, 편향적, 비하적 대화 내용을 대상으로 하는 기업들이 필연적으로 맞닥뜨리게 되는 대화체 콘텐츠의 부적절성의 정도를 측정할 수 있는 효율적인 검증 기준을 제공할 수 있을 것으로 기대된다. 또한 대화체 콘텐츠에 적용될 수 있는 부적절성 검증 틀을 제공함으로써 특정 기업이 내부 인력과 자원을 바탕으로 한 부적절성 검증 틀을 마련할 경우 발생할 수 있는 인력난과 기술적 한계를 해소하는 데 도움이 될 수 있을 것이다.
- 본 사업의 결과물은 범용적으로 활용 가능한 대화체 텍스트 부적절성 검증 틀을 개발하고 고도화하는 것에 기여할 수 있을 것이다. 대화 데이터에 대한 부적절성 검증 틀을 공개함으로써 특정 기업별로 구축하고 보유한 언어 자료가 내부 기밀로 처리되는 경우 발생할 수 있는 일관되지 않고 편중된 부적절성 검증의 비효율성을 방지할 수 있을 것으로 기대된다. 부적절성 검증의 틀을 공유하고 공개함으로써 대화체 콘텐츠를 다루는 기업들뿐 아니라 관련 분야의 각종 기관들, 관공서 및 다양한 민간 부문에서 활용할 수 있도록

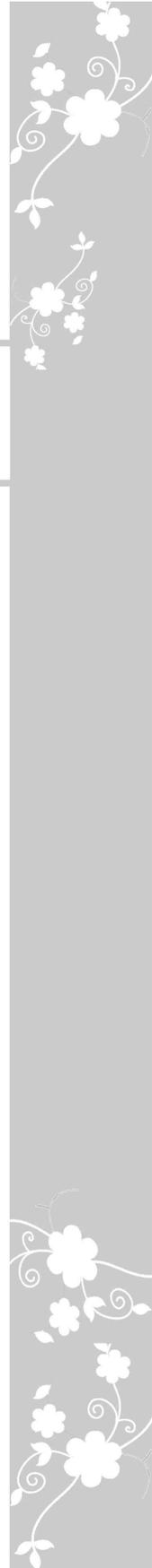
록 하여 지속적으로 발전하고 향상될 수 있는 한국형 부적절성 검증 표준을 정립하도록 하는 생태계가 구축될 수 있기를 기대한다.

- 본 사업의 결과로 구축되는 부적절성 검증 틀을 포함한 말뭉치는 다양한 연구 자료로도 활용이 가능하다. 온라인에서 실제 사용된 대화체 텍스트를 통해 우리 사회에 실존하는 다양한 부적절성 문제 양상을 파악하게 됨으로써 현실적이고 정교한 검증 방안을 논의할 수 있는 학술적 논의의 토대가 마련될 수 있을 것이다. 또한 현대의 일상적 의사소통 형식으로 자리 잡은 온라인 대화에서 드러나는 언어 형식과 혐오·차별적 어휘를 통해 시대·사회적 특성을 파악하는 연구 자료로도 활용할 수 있을 것이다.



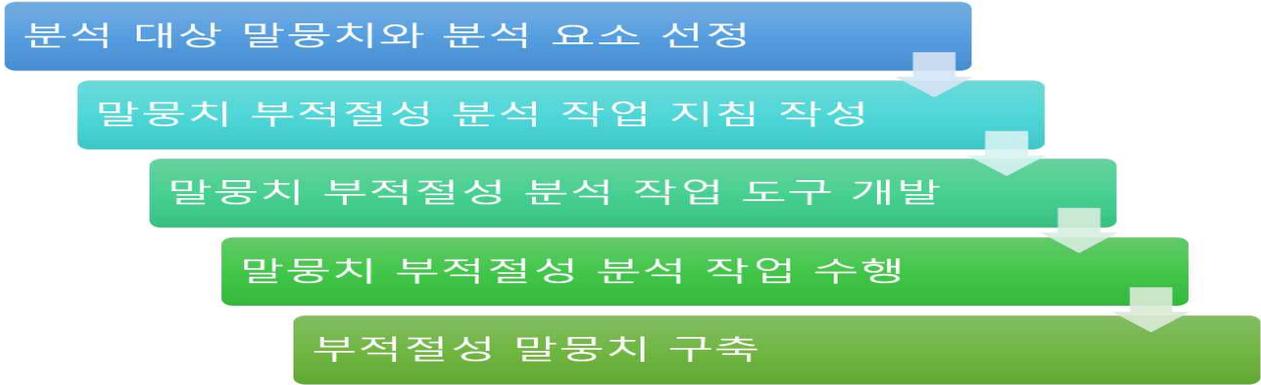
제 2 장

사업 추진 경과



2.1. 사업 수행 절차와 일정

- 본 사업의 전체적인 수행 절차는 다음과 같다.



[그림 2] 사업 수행 절차

- 본 사업은 다음과 같은 일정으로 수행되었다.

<표 2> 사업 수행 일정

일정 / 단계별 사업 내용	2022년					2023년	
	8월	9월	10월	11월	12월	1월	2월
계약 및 연구 착수	●						
지침 수립 및 공유		●	●	●	●		
1차 분석/주석 작업				●	●	●	
중간 보고						●	
지침 수정 및 2차 분석/주석 작업						●	●
작업 결과 검수						●	●
작업 결과 분석						●	●
결과보고서 및 구축 말뚎치 제출							●

2.2. 분석 대상 말뭉치와 분석 요소 선정

2.2.1. 분석 대상 말뭉치의 선정과 특성

- 본 사업에서는 국립국어원에서 제공하는 말뭉치를 부적절성 분석 대상으로 선정하였는데, 이는 <2020년 개체명 말뭉치 연구 분석> 사업의 결과물로 구축된 것이다. 해당 말뭉치는 2019년 국립국어원에서 구축한 웹 원시 말뭉치를 대상으로 4차 산업혁명에 대비한 인공지능 언어 처리 기술 개발 및 국어 연구의 전반으로 활용하기 위해 구축된 500만 어절 규모의 개체명 분석 웹 말뭉치이다.
- 개체명 분석 웹 말뭉치는 다양한 사회적 의사소통망 서비스(SNS)를 포함하는 개개인의 웹 문서를 대상으로 하였다는 점에서 공적 장르의 문어 말뭉치에 비해 본 사업의 분석 및 연구 대상인 부적절성을 파악하기에 적절한 말뭉치였지만, 배경 맥락을 파악하기 어려운 짧은 어절의 감탄사나 약어 표현 등이 많다는 점에서 본 사업에서 중요하게 고려하고자 한 비명시적 부적절성을 파악하는 데에는 다소의 어려움을 발생시키기도 했다.

2.2.2. 말뭉치 부적절성 분석 요소 제안

- 위의 사업 개요에서 밝힌 것처럼, 본 사업은 원래 말뭉치의 ‘비윤리성’을 분석하고 비윤리성 분석 말뭉치를 구축하는 것으로 발주되었으나, 착수 보고 이후 연구진 회의, 분석/구축팀 시범 작업, 주관 기관과의 협의 등을 거치면서 분석 및 구축 대상 관련 개념이 ‘비윤리성’에서 ‘부적절성’으로 변경되었다. ‘부적절성’은 기존의 ‘비윤리성’ 외에도 ‘공격성’, ‘편향성’, ‘비하성’ 등의 네 가지 특성 중 하나 이상이 나타나는 것을 의미한다. 부적절성의 분석 요소 역시 원래에는 비윤리성의 내용(혐오, 욕설/선정, 차별/편향 등)과 영역(성, 연령/세대, 출신, 신체, 문화 등)의 두 가지로 발주되었고, 연구진이 추가로 강도(상, 중, 하)와 대상(개인, 공동체, 문화, 자연 등)을 제안하였다.

2.2.3. 말뭉치 부적절성 분석 요소 1차 선정

- 본 사업의 연구진은 착수 보고 시 부적절성의 내용, 영역, 강도, 대상 등 네 가지 분석 요소를 제안하였는데, 착수 보고 이후 이루어진 분석 요소 1차 선정에서는 몇 가지 중요한 변화가 발생하였다.
- 첫째, 부적절성의 내용을 명시성(명시/비명시)으로 변경하였다. 이는 기존의 유사 연구에서 부적절성(비윤리성) 내용의 유형 구분을 위한 객관적인 기준을 마련하기 어렵고 이에 따라 분석 결과의 신뢰성이 부족하다는 비판이 많이 제기된 점을 고려한 것이다. 이에 본 사업에서는 부적절성이 특정 표현을 통해 명시적으로 나타나는지, 명시적인 표현은 없어도 맥락을 통해 나타나는지를 판정하는 명시성 분석으로 분석의 객관성과 신뢰성을 높이고자 했다.
- 둘째, 부적절성의 영역을 성, 연령/세대, 출신, 신체, 문화, 사회적 조건, 가족, 기타 등 8개 유형으로 분류하였다. 이러한 영역 유형 분류는 관련 선행 연구들과 국가인권위원회 법 제2조 제3항에서 제시된 차별 행위의 목록 등을 참고한 결과이다.
- 셋째, 부적절성의 대상을 분석 요소에서 제외하였다. 이는 부적절성의 영역이 8개 유형으로 상세화되어 새롭게 도입되면서 기존의 대상 유형 분류(국내에서의 개인/공동체/문화/자연 등 분류, 영미권에서의 individual, group, others 등 분류)가 잉여적인 측면이 커졌다는 점에서 부적절성의 대상을 분석해야 할 필요성이 현저하게 감소했기 때문이다.
- 넷째, 부적절성의 강도를 분석 요소에서 제외하였다. 이는 부적절성의 강도 판정을 위한 객관적 기준을 마련하기 어렵고 분석 작업자에 따라 강도 판정의 편차가 심해 분석 결과에 대한 신뢰성을 확보하는 데에도 한계가 있다는 점을 고려한 것이다.
- 결과적으로 1차 선정 시에는 부적절성의 ‘명시성’과 ‘영역’ 등 2개의 분석

요소가 선정되었다.

2.2.4. 부적절성 분석 대상 및 요소 선정을 위한 자문회의 개최

- 자문회의 일시: 2022. 11. 16. (수) 16:00~18:00
- 참석 자문위원: 카이스트 전산학부 오혜연 교수,
창원대 컴퓨터공학과 차정원 교수,
Softly AI 박성준 대표, KT 장두성 상무
- 자문회의에서 논의된 주요 의견과 본 사업에의 반영 결과는 다음과 같다.
 - 1) 자문위원 A
 - 결과 활용의 측면에서는 문장 단위로 태깅하는 것이 필요함. 단어, 어절 단위는 학습하기에 적절한 단위는 아님.
 - > 판정의 과정에서는 어절 단위로 판단하지만, 최종적으로는 문장 단위로 구축하였다.
 - 태깅 작업자 특성(성별, 연령 등)에 따라 태깅 결과가 다를 수 있음.
 - > 작업자, 화자의 변인에 따라 나타나는 차이가 중요하다는 점에 공감하고 다양한 변인을 고려한 작업자 구성을 하지 못한 한계를 인정하였다. 실제로 작업자들의 성별과 연령은 여성과 20대에 편중되어 있지만, 그러한 편중으로 인한 문제 발생 시에는 상대적으로 다양한 성별 및 연령이 포진된 팀장과 공동 연구진 차원의 검수와 논의를 통해 문제를 최소화하였다.
 - ‘문맥에 상관없이 부적절한 경우’ 그전에 나온 문장, 또는 어떤 기사에 대한 댓글이라든지 문맥에 따라 적절/부적절이 판단될 수 있음. 또한 사용되는 상황(ex. 일베, 학교 게시판에서 적절한 것이 다름)에 따라서도 판단 가능함.
 - > 맥락에 따라서 판정되는 경우는 작업 도구 상에 부적절한 문장 외에 선행 문장으로 5문장씩을 제시하였다.

2) 자문위원 B

- 부적절성의 기준이 시간, 상황, 대상 문서에 따라 변화할 수 있다는 점을 추후 과제 진행을 위하여 절차를 명확하게, 그 기준을 폭넓게 잡아 지침에 포함하면 좋을 것으로 보임.

-> 부적절성이 여러 변인에 따라 달라질 수 있다는 점에 공감하지만, 본 사업의 분석 작업 지침은 2020년대라는 현재의 시점에서 다양한 사회적 의사소통망 서비스(SNS)를 포함하는 개개인의 웹 문서를 대상으로 하였음을 밝혔다.

- 문자가 깨져 있는 경우(초성으로만 되어 있는 경우, 언어와 이미지가 결합된 경우)를 고려하였는지에 대한 문제

-> 선·후행 문장을 통해 맥락을 확인 가능한 경우에만 분석 대상으로 삼았고, 본 사업의 과제 범위에 따라 이미지는 고려하지 않았다.

- 어절 단위의 태깅 및 문장 단위의 말뭉치 구축 적절해 보임.

- 실현 가능성 문제: 실현 가능성이 있는 부분을 지침에 최대한 많이 포함하는 것을 고려할 필요가 있음.

-> 실현 가능성의 문제는 작업 지침의 일관된 적용 가능한지에 대한 문제이므로, 작업 지침에서 논란의 여지가 최소화될 수 있도록 객관성과 일관성에 초점을 두었다.

3) 자문위원 C

- 판정의 기준과 실제 언중들이 느끼는 부적절성 사이의 간극이 존재함. 특정한 사용자를 염두에 두고 진행하는 작업이 아닌, 일반적인 부적절성을 판정하기 때문에 나타나는 문제임. 너무 많은 것을 포괄하다 보면 자연스럽게 실효성이 떨어지는 문제가 있는데, 이에 대해서는 논의가 필요함.

-> 타깃이 분명하지 않고 일반적인 목적을 두고 있기 때문에 지적된 한계를 인정하였다.

- 명시적/비명시적 분류: 서비스를 만드는 과정에서 크게 중요한 기준은 아님. 공격적인 표현이 있을 때 그 표현이 가려졌을 때, 그 가려진 표현이 무엇인지 상상할 수 있는 패턴, 머릿속에서 자동 완성되도록 하는 방식이 활용됨.

- 부적절성의 영역: 영역 자체의 구분은 잘 되어 있고 이해가 쉬움. 그러나 사용자들이 중요하게 생각하는 것은 부적절성 표현에 대한 대상임. '누구를 향하는가?'가 중요함. 또한 이것은 부적절성의 강도와 밀접하게 관련됨. 대상이 없이 나타나는 부적절성은 강도가 약하다고 느끼는 경향성이 있음. 부적절성의 영역과 함께 타깃 정보(개인, 집단 등)가 있으면 활용도가 높을 것으로 보임.

-> 대상 유형은 본 과제의 영역에 수렴될 수 있으므로, 유효성 및 작업 단순화 차원에서 제외하였다.

- 부적절성 말뭉치를 현장에서 사용하기 위해서는 부적절성의 강도(정도)를 조절할 수 있는 것을 원함. 이에 대해서도 고려한다면 실용적 가치가 높아질 것.

-> 작업자들 간에 판정의 통일성과 일관성을 확보하기 어려워 분석 결과에 대한 객관성과 신뢰성이 부족하다는 점에서 분석 요소 1차 선정 시에는 강도를 제외하였지만, 2차 선정 시에는 부적절성 말뭉치를 활용할 학계와 산업계의 요구를 반영하여 강도를 다시 분석 요소에 포함하였다. 다만 본 사업에서는 강도 분석의 복잡성을 최소화하고 일관성을 최대화하고자 기존의 유사 연구에서 많이 사용한 3단계(상, 중, 하) 판정이 아닌 2단계(강, 약) 판정을 적용하였다.

4) 자문위원 D

- 현재 부적절성 영역은 추가 정보로 활용이 가능할 것으로 보임. 문장 단위의 부적절성 여부를 판단하고, 영역 정보는 메타 정보로 활용이 가능해 보임. 현재 분류 체계가 유용할지에 대한 고민이 필요함.

-> 분석 요소 2차 선정에서 영역 유형을 축소하는 방향으로 수정하였다.

- 소스나 원천 등 메타 정보들을 통해서 부적절성 여부를 판단할 수 있기 때문에 이를 표시하면 데이터 활용 측면에서 더 유용할 것으로 보임.

-> 분석 대상 말뭉치에 이미 제공되어 있으므로 반영하였다.

- 작업자들이 어떻게 태깅을 했는지에 대한 정보를 제공하는 것도 유용할 것으로 보임.

- 서비스별로 원하는 것이 다르기 때문에 강도가 필요한 부분이 있음.
- > 2차 선정 시에는 강도를 다시 분석 요소로 포함하였다.
- 부적절성 표현을 마스킹, 제외, 순화하는 서비스 또는 연구가 필요하기 때문에 이를 위한 작업도 추후에 이루어지면 좋겠음.
- > 향후 과제에는 반영이 가능하겠으나 본 과제의 범위를 넘어서는 부분이므로 제외하였다.

2.2.5. 말뭉치 부적절성 분석 요소 2차 및 최종 선정

- 자문회의 이후 연구진 회의, 분석/구축팀 시범 작업, 주관 기관과의 협의 등을 거치면서 분석 요소 2차 선정에서는 다시 몇 가지 중요한 변화가 발생하였다.
- 첫째, 1차 선정 시에는 제외하였던 부적절성의 ‘강도’를 2차 선정에서는 다시 분석 요소로 포함하였다. 이는 본 사업의 결과물로 구축되는 부적절성 말뭉치를 활용할 학계와 산업계의 현실적 요구를 반영한 것이다. 다만 본 사업에서는 강도 분석의 복잡성을 최소화하고 일관성을 최대화하고자 기존의 유사 연구에서 많이 사용한 3단계(상, 중, 하) 판정이 아닌 2단계(강, 약) 판정을 적용하였다.
- 둘째, 1차 선정 시에는 존재하지 않았던 부적절성의 ‘맥락’을 2차 선정에서는 새로운 분석 요소로 추가하였는데, 이는 부적절성의 맥락이 ‘부정적’인지, ‘긍정적’인지를 판정하는 것이다. 맥락을 부적절성의 새로운 분석 요소로 추가한 것은 판정된 부적절성이 화자의 태도(의도) 측면에서는 부정적이지 않더라도 그 맥락 내용 측면에서는 부적절한 경우, 부적절성으로 판정되었으나 화자의 태도(의도)와 맥락 내용 측면 모두에서 긍정적이거나 무표적인 것으로 판단되는 경우 등을 구별하여 분석하는 것이 필요하기 때문이다.
- 결과적으로 2차 선정 시에는 부적절성의 ‘명시성’, ‘영역’, ‘강도’, ‘맥락’ 등의 분석 요소가 선정되었고, 이러한 4가지가 최종적인 분석 요소가 되었다.

2.3. 말뭉치 부적절성 분석 작업 지침 작성

2.3.1. 작업 지침 초안 작성

- 착수 보고, 연구진 회의, 주관 기관과의 협의, 분석 요소 1차 선정 등을 통해 이루어진 논의와 분석/구축팀 시범 작업을 통해 발견한 문제와 해결 방안 등을 반영하여 작업 지침 초안을 작성하였다.
- 작업 지침 초안 작성 일정
 - 작업 지침 초안 작성: 2022. 9. 27. ~ 2022. 11. 14.

2.3.2. 작업 지침 수정안 작성

- 자문회의 의견과 부적절성 분석 요소 2차 선정, 분석/구축팀의 강도 태깅 샘플 작업 결과 등을 반영하여 작업 지침 수정안을 작성하였고, 이후 여러 차례에 걸친 주관 기관의 피드백과 연구진 논의 결과 등을 반영하여 작업 지침 최종안을 작성하였다. 작업 지침은 작업자들의 이해와 일관성 있는 작업, 그리고 최종보고서를 통해 공개되는 작업 지침의 독자들의 이해를 고려하여 작성하였고, 이를 위해 작업 지침 관련 주요 개념과 작업 절차, 예시 등을 최대한 상세하고 명확하게 제시하였다.
- 작업 지침 수정안 작성 일정
 - 강도 태깅 샘플 작업: 2022. 11. 17. ~ 2022. 11. 25.
 - 강도 태깅 샘플 작업에 대한 국립국어원 피드백: 2022. 11. 25. ~ 2022. 11. 30.
 - 작업 지침 수정 방향 작성: 2022. 11. 30. ~ 2022.12. 1.
 - 작업 지침 수정 방향에 대한 국립국어원 피드백: 2022. 12. 1.
 - 작업 지침 1차 수정안 작성: 2022. 12. 2. ~ 2022. 12. 8.
 - 작업 지침 1차 수정안에 대한 국립국어원 피드백: 2022. 12. 9. ~ 2022.

12. 14.

- 작업 지침 v1 작성: 2022. 12. 15. ~ 2022. 12. 23.
- 작업 지침 v1에 대한 국립국어원 피드백: 2022. 12. 24. ~ 2022. 12. 28.
- 작업 지침 v1 확정 및 연구진 공유: 2022. 12. 29.
- 작업 지침 v1 분석/구축팀 팀장 공유 및 수정, 보완: 2023. 1. 2.
- 작업 지침 v2 완성: 2023. 1. 2.

2.3.3. 작업 지침 교육

- 분석/구축팀 총괄 공동연구원이 전체 작업자(학부생 팀원) 대상으로 작업 지침 v2에 대한 교육을 실시하였다. 해당 교육에서는 특히 명시성, 맥락, 영역, 강도 등 부적절성 분석 요소에 대한 판정 기준과 예시에 대한 집중 교육이 이루어졌다(2023. 1. 5).

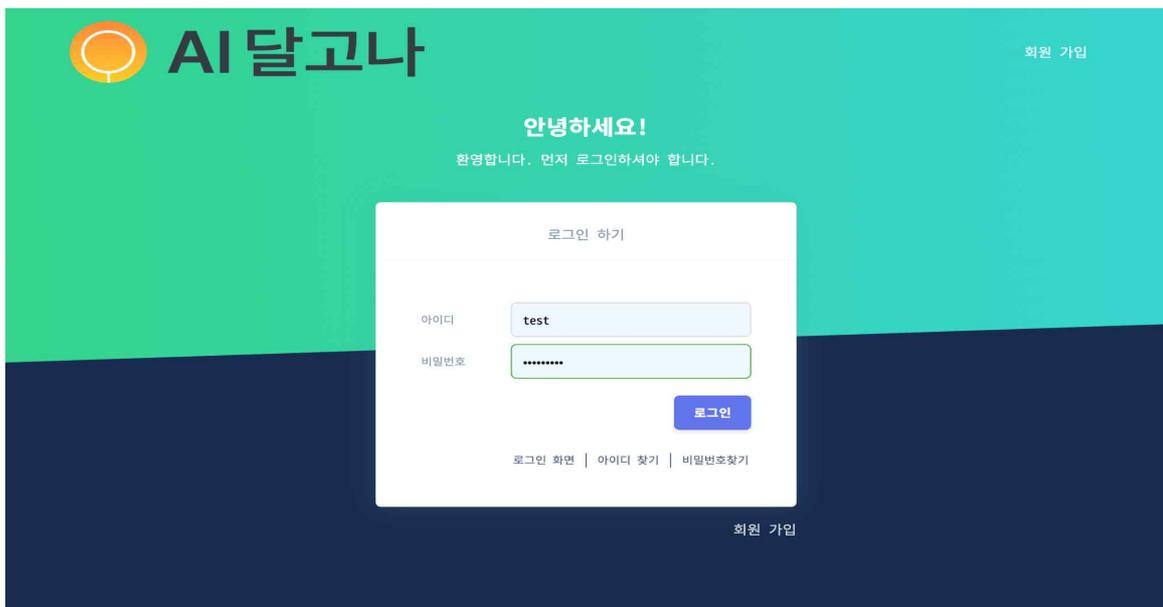
2.4. 말뭉치 부적절성 분석 작업 도구 개발

2.4.1. 작업 도구의 구성

- 말뭉치 부적절성 분석 작업 도구는 다음과 같은 2단계로 구성하였다.
 - 1단계: 부적절성 문장의 선별
 - 2단계: 부적절성이 나타나는 부분에 대한 시작/종료 표지 부착 및 분석 요소 주석

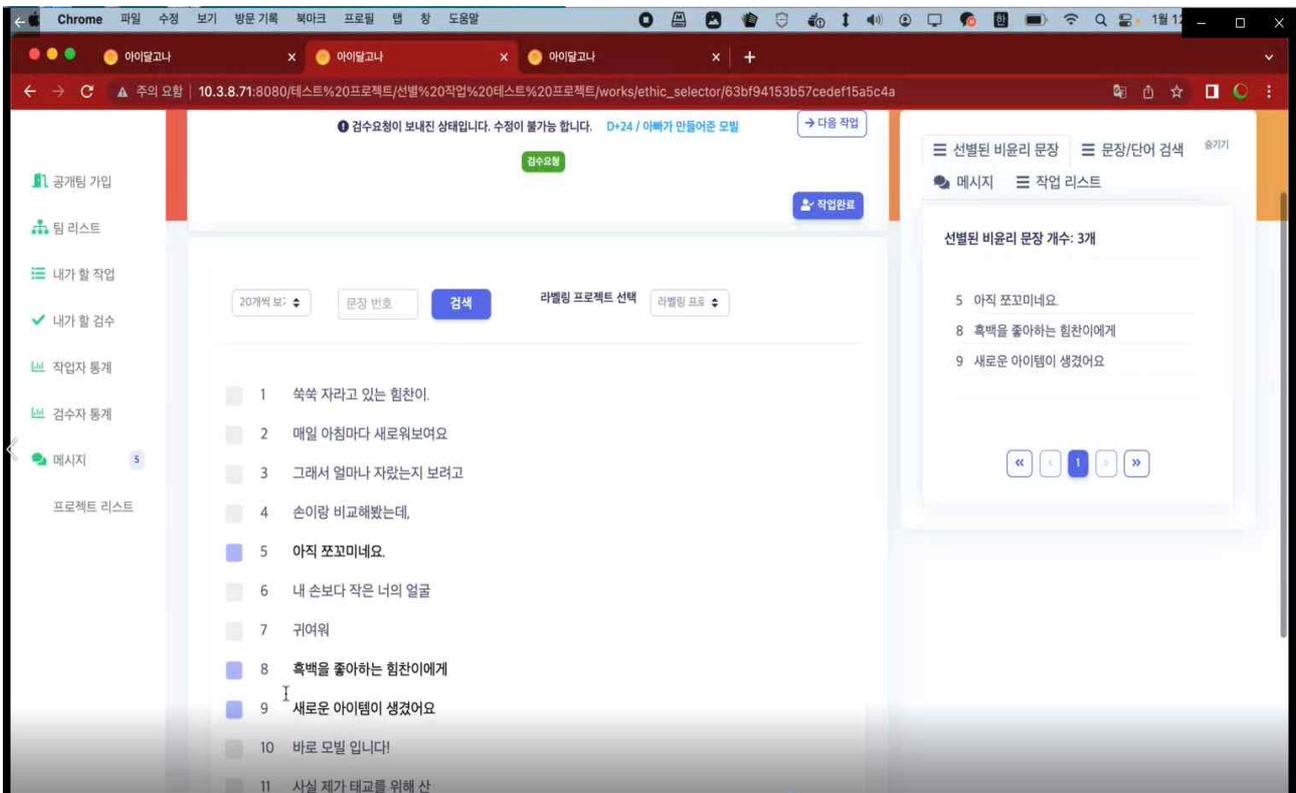
2.4.2. 작업 도구의 개발

- 말뭉치 부적절성 분석 작업 도구는 2021년 한국지능정보사회진흥원(NIA)에서 발주한 인공지능 학습용 데이터 구축 사업의 일환으로 진행된 ‘텍스트 윤리 검증 데이터’ 사업과 2021년 국립국어원에서 발주한 ‘비윤리성 말뭉치 연구 분석 및 시범 구축’ 사업 등 2020년 이후 7개 사업에서 작업 도구로 사용되어 그 성능과 기능을 확인했던 ‘아이달고나(AI달고나)를 기반으로 본 과제의 필요에 맞게 새롭게 개발하였다.



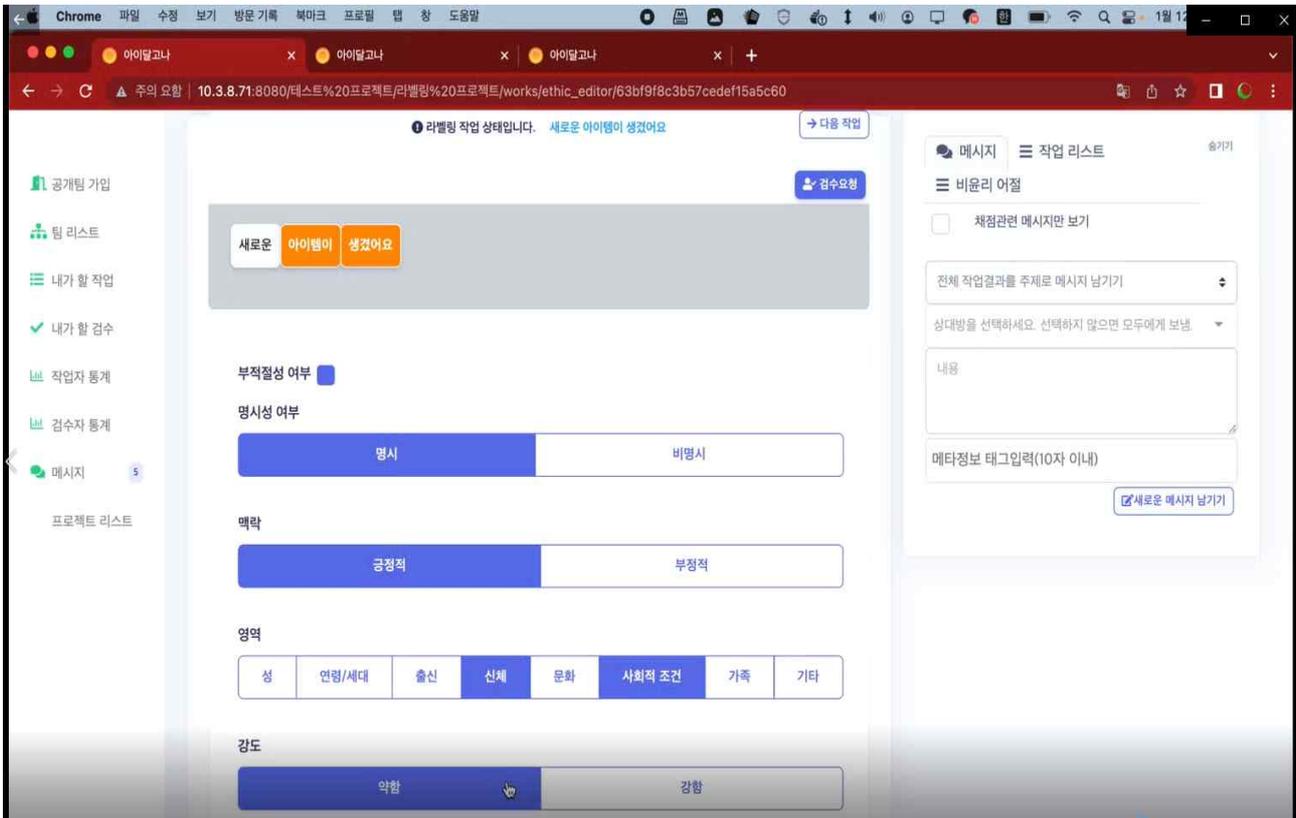
[그림 3] 아이달고나 로그인 화면

- 작업 도구에서 말뭉치 부적절성 분석 작업을 수행하기 위해서는 작업자별 25만 어절의 말뭉치를 할당하는 과정이 필요하였다. 그런데 주관 기관에서 제공한 분석 대상 말뭉치는 문장 단위로만 분할되어 있으므로 작업자 1인당 할당되는 말뭉치를 약 25만 어절 내외로 재분할하는 작업이 필요하여 별도의 말뭉치 선별기를 구축하여 활용하였다(2023. 1. 3. ~ 2023. 1. 9).



[그림 4] 아이달고나 1단계 부적절성 문장 선별 화면 예시

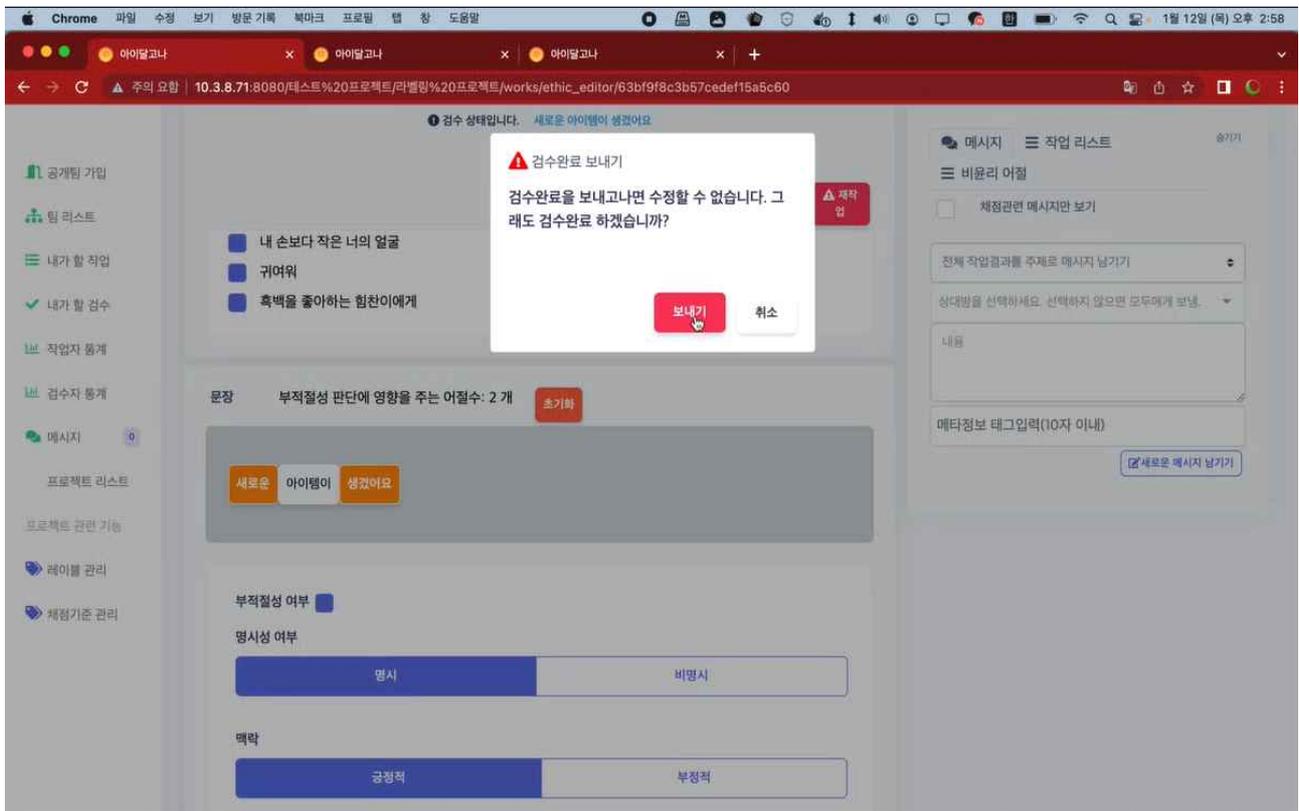
- 말뭉치 부적절성 분석 작업 도구는 착수 보고회 이후 바로 1차 개발을 완료하였으나(2022. 9), 자문회의 이후 수정된 작업 지침 v2에 따라 맥락, 강도 등 새롭게 추가된 분석 요소를 분석하는 작업이 가능하도록 작업 도구 개선이 이루어졌다(2023. 1. 3. ~ 2023. 1. 9).



[그림 5] 아이달고나 2단계 부적절성 분석 요소 주석 화면 예시

2.4.3. 작업 도구 사용 방법의 공유 및 교육

- 작업 도구 개발이 완료된 후에는 분석/구축팀을 중심으로 작업 도구 사용 방법을 공유하기 위한 교육이 이루어졌다. 먼저 공동연구원과 분석/구축팀 총괄 공동연구원(관리자), 팀장(검수자) 중심의 작업 도구 사용 방법 공유가 이루어졌다. 그 핵심 내용은 관리자가 작업 도구 상에서 작업자들에게 말뭉치를 할당하는 방법, 검수자가 작업 도구 상에서 검수 절차를 진행하는 방법 등이다. 다음으로 분석/구축팀 작업자를 대상으로 작업 도구 사용 방법에 대한 상세한 교육이 이루어졌다. 그 핵심 내용은 작업자가 작업 도구 상에서 부적절성 문장을 선별하고 분석 요소를 주석하는 방법이다(2023. 1. 10. ~ 2023. 1. 12.).

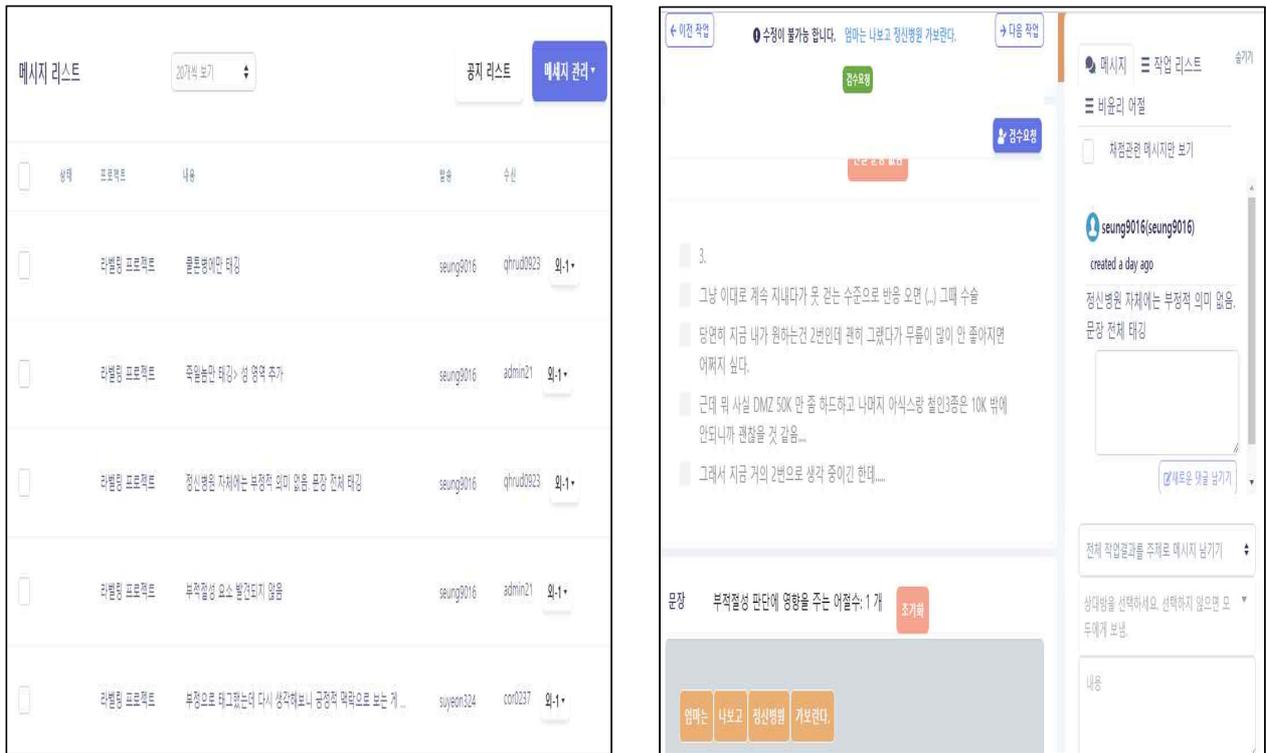


[그림 6] 아이달고나 2단계 검수 작업 화면 예시

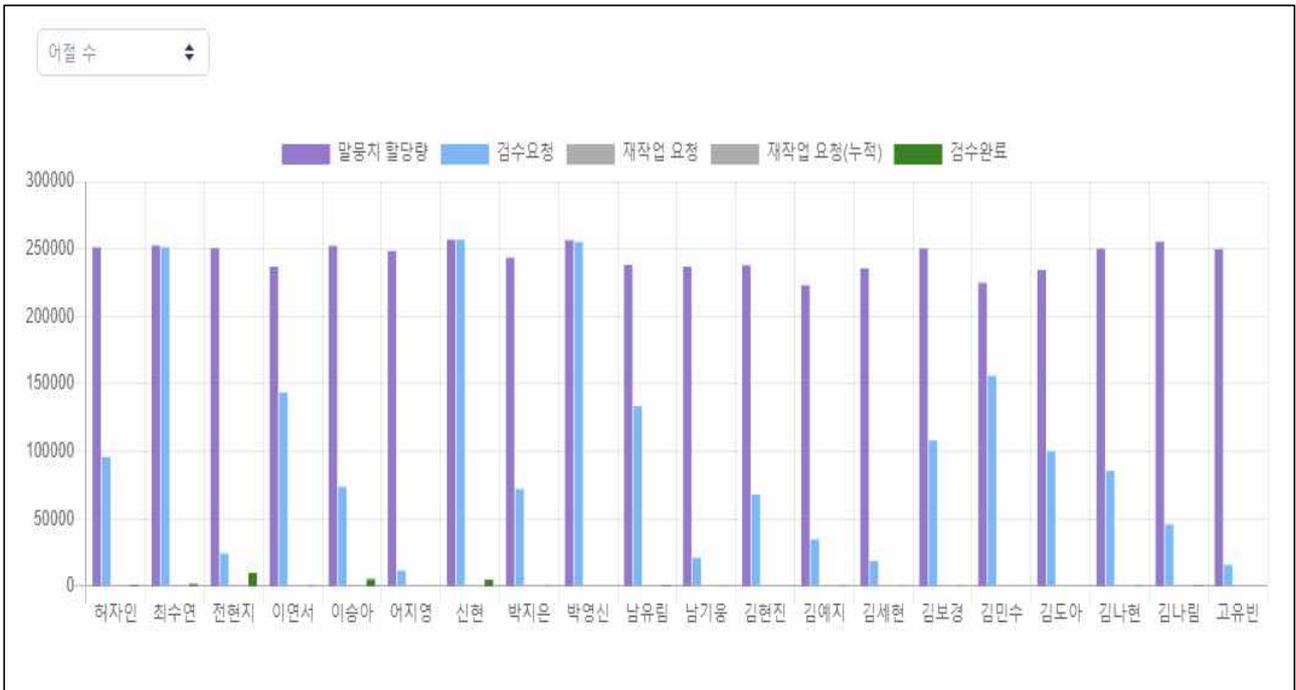
2.5. 말뭉치 부적절성 분석 작업 수행

2.5.1. 1차 분석 작업

- 작업 지침 v2 및 작업 도구의 공유와 교육을 바탕으로 중간보고용 1차 분석 작업과 검수 작업이 다음과 같이 본격적으로 진행되었다(2023. 1. 13. ~ 2023. 1. 25.).
 - 중간보고용 1차분 분석 작업 진행
 - : 2023년 1월 18일까지 3,000건 이상의 문장 선별 및 라벨링 작업 완료
 - 중간보고용 1차분 분석 결과물 검수 진행
 - : 2023년 1월 25일까지 선별 및 라벨링 작업 완료한 문장에 대한 검수와 재검수 진행
- ㉠ 작업자들의 오타깅에 대하여 팀장들이 메모 기재 후 표본 반려
- ㉡ 팀장들의 오검수에 대하여 총괄 공동연구원이 메모 기재



[그림 7] 검수 과정에서의 총괄 공동연구원 메모 예시



[그림 8] 작업자 분석 작업 현황(1월 25일 기준)

2.5.2. 중간보고회 개최

- 중간보고회 일시: 2023. 2. 1.(수) 16:00~18:00
- 중간보고회 주요 논의사항과 본 사업에의 반영 결과는 다음과 같다.
 - 1) 욕설의 어원을 고려하여 영역 태깅하는 문제
 - 어원 측면에서 ‘성’과 관련된 욕설(씨발, 존나, 졸라 등)에 ‘성’을 영역 태깅하는 것의 적절성에 대한 문제가 제기됨
 - > ‘성’ 관련 내용이나 맥락이 없이 단순한 감정의 배설로 사용되는 욕설은 가급적 ‘성’을 태깅하지 않았다.
 - 2) 영역(domain)별 비율 조정 문제
 - 제안서에서는 특정 영역의 비율이 35%를 넘지 않게 조정하기로 하였으나, 1차 작업 결과에서는 특정 영역이 50%를 상회하고 다른 영역은 5% 미만으로 나타나는 등 영역별 비율 편차가 매우 심하여 인공지능의 학습 자료로 적절할지에 대한 우려가 제기됨

-> 당초와 달리 영역을 복수 태깅하게 되었으므로 기계적 비율 조정의 필요성은 감소하였으나, 지나치게 비율이 낮게 나타난 영역은 유사 영역과 통합하여 영역별 비율 편차를 최대한 줄이는 방향으로 조정하였다.

3) 욕설이 맥락 없이 단독 실현되는 경우의 처리 문제(선정 및 영역 태깅)

- 욕설이 맥락 없이 단독으로 실현되는 경우를 ‘기타’ 영역으로 태깅하면서 ‘기타’ 영역의 비율이 지나치게 높아지고 있으므로 해당 경우를 분석 대상에서 제외하는 문제가 제기됨

-> 욕설만 단독으로 나오는 문장도 부적절성 문장이므로 이를 모두 분석 대상에서 제외하는 것은 적절하지 않으므로 포함시켰다. 다만 해당 욕설이 완전히 동일한 형태로 중복하여 나타나는 경우에는 하나만을 남기고 나머지는 제외하였다.

4) 부적절성 관련 어휘(표현) 목록 제시 문제

- 명시적 부적절성 문장을 판정하는 기준이 된 어휘(표현) 목록을 결과보고서에 부록으로 제시하는 방안이 제안됨

-> 국어사전에 제시된 욕설, 비어, 비하성 속어 등을 모두 제시하기는 어렵지만, 관련 보고서에 제시된 차별, 혐오, 선정 표현 등은 별도로 정리하여 결과보고서의 부록으로 제시하였다.

5) 맥락 및 강도 관련 작업 지침 적용 문제

- 맥락 및 강도 관련 작업 지침이 제대로 적용되지 않은 경우가 발견되었으며, 이에 대한 검수가 충분히 이루어진 것인지에 대한 의문이 제기됨

-> 긍정적 맥락인데, 강도가 ‘강’으로 태깅된 경우는 명백한 오류이므로 수정하였고, 시스템 반영 과정의 시간 문제로 검수가 이루어졌지만 아직 결과물에 반영되지 못한 경우는 최종 반영 여부를 확인하였다.

6) json 양식 문제

- json 양식을 정의한 적이 없어 21년 과제 기준으로 제출된 것으로 보임
- > 국립국어원에서 정하여 제공한 양식에 따라 json 파일을 출력하여 제출하였다.

2.5.3. 작업 지침 최종안 작성

- 작업 지침 v2에 따라 수행한 1차 분석 작업 과정에서 혼란과 오류가 많이 발생한 ‘강도’의 유형별 범위 기술을 명시성을 일차적 기준으로 제시하는 방향으로 다음과 같이 수정하였다.
- 성, 연령/세대, 출신, 신체, 문화, 사회적 조건, 가족, 기타 등 8개 유형으로 나누었던 부적절성 영역을 지나치게 비율이 낮게 나타난 영역은 유사 영역과 통합하여 관계/조건, 문화, 신체, 성, 연령/세대, 기타 등 6개 유형으로 축소하였다. 이는 사회적 조건, 출신, 가족 등 3개 영역을 ‘관계/조건’ 영역으로 통합한 결과이다.
- 비하성 속어로 혼돈할 수 있는 경우를 <주의>에 추가하였다. 예를 들어, 사전에 ‘낮잡아’와 유사한 ‘낮추어’로 표현되더라도 겸양의 의미를 지닌 경우에는 비하성이 없으므로 ‘명시’는 물론이고 ‘부적절성’으로도 판정하지 않도록 관련 내용을 <주의>를 통해 강조하였다.
- 차별적 측면이 있어도 부적절성으로 판정하지 않는 경우를 <주의>에 추가하였다. 관련 연구 및 조사 보고서에 차별 표현으로 지적되었고 차별적 측면이 있다 하더라도 다른 대안이 없어서 그대로 사용할 수밖에 없는 다음의 경우에는 ‘명시’는 물론이고 ‘부적절성’으로도 판정하지 않도록 관련 내용을 <주의>를 통해 강조하였다.
- 작업 지침 최종안 작성 일정
 - 작업 지침 v2에 대해 중간보고회에서 논의: 2023. 2. 1.
 - 작업 지침 v2를 중간보고회 논의를 반영하여 수정한 작업 지침 v3 작성:

2023. 2. 1. ~ 2023. 2. 3.

- 작업 지침 v2에 대한 국어원 피드백: 2023. 2. 1. ~ 2023. 2. 6.

- 작업 지침 최종안 작성 및 완성: 2023. 2. 6. ~ 2023. 2. 7.

※ 붙임 1(말뭉치 부적절성 분석 작업 지침_최종) 참조

2.5.4. 2차 분석 작업

○ 중간보고회 이후 분석 지침 최종안이 완성되기 이전에도 1차 분석 작업 결과에 대한 오류 재검수와 2차 분석 작업이 다음과 같이 진행되었다.

- 중간보고용 1차분 분석 결과에 대한 재검수

(2023. 2. 2. ~ 2023. 2. 6.)

: 1차분 분석 결과물의 오류, 특히 강도 주석 오류 중심으로 재검수를 진행하고 수정 반영

- 2차 분석 및 검수 작업

(2023. 2. 2. ~ 2023. 2. 7.)

: 일단 작업 지침 v2에 기반하여 2차분 작업자 분석 및 검수자 검수를 계속하여 진행

○ 작업 지침 최종안의 완성과 공유 및 교육을 바탕으로 2차 분석 작업과 검수 작업이 다음과 같이 본격적으로 진행되었다.

- 작업 지침 v4 검수자 교육 및 작업자 교육 실시

(2023. 02. 07. ~ 2023. 02. 08.)

- 기존의 2차 작업 결과 수정 및 신규 분석 진행

(2023. 02. 08. ~ 2023. 02. 18.)

: 2차분 분석 작업은 기본적으로 팀원 작업자들이 진행하였지만, 중간보고회 이후 작업 기간이 부족한 관계로 일부 작업 속도가 지나치게 더디거나 연락이 되지 않아 작업 진행을 확인할 수 없는 경우에는 해당 팀장과 분석 구축팀 공동연구원이 직접 작업

- 2차분 검수 작업 본격화

(2023. 2. 8. ~ 2023. 2. 18.)

: 팀장과 분석구축팀 공동연구원이 분석 작업 결과가 나오는 대로 검수 작업 동시 진행

○ 최종 단계에서는 2차 분석 작업 결과물에 대한 개인정보 비식별화 작업을 다음과 같이 진행하였다.

- 공동연구진에서 개발한 비식별화 편집기를 통해 비식별화 작업 진행
(2023. 2. 14. ~ 2023. 2. 18.)

: 최종 16,240문장을 작업자 1인당 800여 개씩 배정하여 비식별화 작업을 진행하고 팀장이 검수



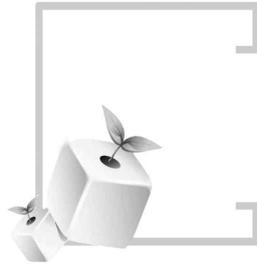
[그림 9] 작업자 분석 작업 현황(2월 18일 기준)

2.6. 부적절성 말뭉치 구축

- 위와 같은 분석 작업의 결과로 최종 16,240개 문장으로 구성된 부적절성 말뭉치를 구축하였고, 이를 국립국어원에서 제공한 json 양식에 맞추어 출력하여 제출하였다.

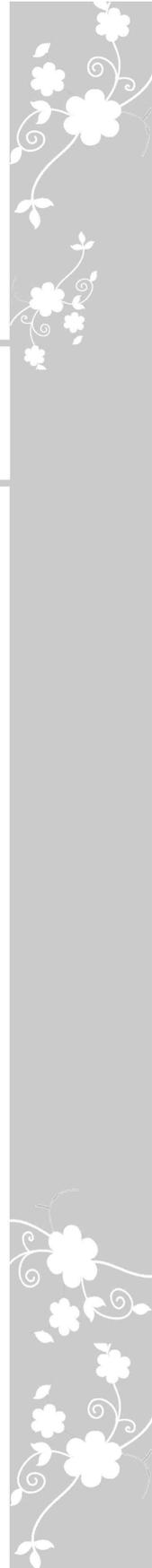
```
{
  "id": "EXAU2202302210",
  "metadata": {
    "title": "국립국어원 웹 말뭉치 추출 EXAU2202302210",
    "creator": "국립국어원",
    "distributor": "국립국어원",
    "year": "2022",
    "category": [
      "웹>블로그",
      "웹>누리소통망"
    ],
    "annotation_level": [
      "부적절성 분석"
    ],
    "sampling": "부분 추출 - 특정 부분 추출"
  },
  "document": [
    {
      "id": "ESRW1904002673.2335",
      "metadata": {
        "author": "나민규",
        "crawl_date": "20190911 22:56:13",
        "date": "20130108",
        "publisher": "facebook",
        "title": "NA",
        "url": "https://www.facebook.com/minkyu.rha/posts/300204870101239"
      },
      "sentence": {
        "sentence_id": "ESRW1904002673.2335.1.1",
        "sentence_form": "이 노래 듣고 배시시 웃었다면 소싯적 &location&애니 오덕이렸다.",
        "prev_sentence": [],
        "next_sentence": [],
        "inappropriate_sentence": {
          "context": "NEGATIVE",
          "is_explicit": false,
          "domain": "관계/조건,문화,기타",
          "intensity": "WEAK"
        },
        "inappropriate_expression": [
          {
            "begin": 0,
            "end": 40,
            "form": "이 노래 듣고 배시시 웃었다면 소싯적 &location&애니 오덕이렸다."
          }
        ]
      }
    }
  ]
}
```

[그림 10] 부적절성 말뭉치 json 파일 예시



제 3 장

사업 주요 내용



3.1. 부적절성의 개념 설정

- ‘부적절성’의 개념은 본 사업의 발주 시에 제시되었던 ‘비윤리성’의 개념이 지닌 한계를 극복하기 위해 설정된 것이다. 기존의 유사 연구에서 가장 많이 사용되었던 ‘비윤리성’의 기본적 개념은 사람으로서 마땅히 지키거나 행해야 할 도리에 어긋나거나 도덕적 규범을 위반하는 것을 의미하는 것이다. 이러한 비윤리성은 일반적으로 ‘공격성’, ‘편향성’, ‘비하성’ 등의 부정적 특성들을 포함하지만, 이러한 부정적 특성들은 엄격한 의미에서의 비윤리성으로 간주하기 어려운 경우에도 나타난다. 그런데 이처럼 비윤리성으로 간주하기 어려운 경우에 나타나는 공격성, 편향성, 비하성 등이 나타나는 발화는 그것을 접하는 누군가에게 다양한 유형의 부정적 감정을 유발할 수 있다는 점에서 문제가 있다.

- 본 사업에서는 비윤리성으로 간주하기 어려운 경우에도 부정적 감정을 유발하는 공격성, 편향성, 비하성 등을 말뭉치 분석 및 구축 대상으로 설정하기로 했다. 같은 문제의식에서 기존의 유사 과제에서는 ‘비윤리성’ 대신 ‘부정성’, ‘반사회성’ 등의 개념이 제안되기도 했지만, 본 사업에서는 기존의 용어 대신 ‘부적절성’이라는 용어를 사용하였다. 본 사업에서 사용하는 부적절성이라는 용어의 개념은 다음과 같은 네 가지 특성 중 하나 이상이 나타나는 것을 의미한다.
 - 사람으로서 마땅히 지키거나 행해야 할 도리에 어긋나거나 도덕적 규범을 위반하는 것을 의미하는 ‘비윤리성’
 - 비난, 저주, 모욕, 위협, 혐오, 폭력선동 등에서 나타나는 ‘공격성’
 - 차별, 편견, 배제, 불필요한 언급 등에서 나타나는 ‘편향성’
 - 멸시, 폄하, 무시, 조롱 등에서 나타나는 ‘비하성’

3.2. 부적절성의 분석 및 주석(태깅) 단위 설정

- 본 사업에서 부적절성 분석의 기본 단위는 ‘문장’이다. 이는 부적절성이 문장 내의 특정 어휘나 표현에 의해 발생하더라도 문장 단위로 분석한다는 것을 의미한다. 이에 따라 부적절성의 4가지 분석 요소인 명시성, 맥락, 영역, 강도 등은 문장 내의 특정 어휘나 표현과 관련된 것이더라도 전체 문장의 부적절성으로 제시하였다. 또한 분석 대상 문장은 분석 대상 말뭉치에서 한글 맞춤법이나 띄어쓰기 오류 등 어문 규범에 어긋나는 요소들 또는 오타 등을 포함하고 있더라도 수정 없이 그대로 분석하였다.

- 부적절성 분석의 결과를 태깅하는 기본 단위도 ‘문장’이다. 이는 부적절성의 명시성, 맥락, 영역, 강도가 문장 내의 특정 어휘나 표현에 의해 판정된 것이더라도 문장 단위로 주석한다는 것을 의미한다. 그리고 하나의 문장에 명시성, 맥락, 영역, 강도 등의 분석 요소가 각각 복수로 판정되는 경우에는 문장 단위에 다음의 기준에 따라 주석하였다.
 - 명시와 비명시가 함께 나타나는 문장은 ‘명시’로 주석
 - 부정적 맥락과 긍정적 맥락이 함께 나타나는 문장은 ‘부정적’ 맥락으로 주석
 - 두 개 이상의 영역에 해당하는 부적절성이 나타나는 문장은 복수 주석을 권장
 - 강한 부적절성과 약한 부적절성이 함께 나타나는 문장은 ‘강’한 부적절성으로 주석

3.3. 부적절성의 명시성 분석 및 주석

3.3.1. 명시성의 개념

- ‘명시성’은 부적절성이 구체적인 어휘나 표현을 통해 명시적으로 드러나는지(명시), 해당 문장의 맥락에서 드러나는지(비명시)를 판정하는 것을 의미한다. ‘명시성’의 분석은 부적절성 여부에 대한 논란이 적은 명시적 부적절성 문장을 우선적으로 선별할 수 있게 함으로써 부적절성의 판정을 최대한 객관적으로 할 수 있게 한다는 점에서 ‘명시성’은 부적절성의 일차적인 분석 요소이다.

<표 3> ‘명시성’의 유형과 범위

명시성 유형	범위
명시	- 사전에 근거하는 욕설, 비어, 비하성 속어 등이나 관련 보고서에 근거하는 차별 표현, 혐오 표현, 선정적 표현 등과 같은 명시적 부적절성 표현이 나타나는 문장
비명시	- 명시적 부적절성 표현이 나타나지는 않지만, : 해당 문장의 내외 맥락에서 비윤리성, 공격성, 편향성, 비하성 등의 부적절성이 확인되는 문장 : 대상에 대한 부정적 평가 혹은 태도를 드러내는 표현이 포함된 문장 : 현대적인 기준에서 윤리성을 크게 위배하였다고 판단되는 내용을 담은 표현이 나타나는 문장

- ‘명시성’은 본 사업의 발주 시에 주요한 분석 요소의 하나로 제시되었던 부적절성의 ‘내용(혐오, 욕설/선정, 차별/편향 등)’을 대체하는 개념이다. 기존의 유사 연구에서 부적절성(비윤리성) 내용은 유형 구분을 위한 객관적인 기준을 마련하기 어렵기 때문에 분석 및 주석 결과의 신뢰성도 부족하다는 비판이 많이 제기되었다. 이에 본 사업에서는 부적절성이 특정 표현을 통해 명시적으로 나타나는지, 명시적인 표현은 없어도 맥락을 통해 나타나는지를 판정하는 명시성 분석으로 분석의 객관성과 신뢰성을 높이고자 했다.

3.3.2. 명시성의 표현 범위

- 명시성은 문장 단위로 주석하지만, 명시인지, 비명시인지에 따라 부적절성이 표현되는 범위를 달리 표시하였다.
- ‘명시’의 경우에는 명시적 부적절성 표현이 나타나는 ‘어절’ 단위(체언+조사, 용언+어미 등)에 표현 범위를 표시하였고, 하나의 문장에 명시적 부적절성을 발생시키는 어절이 두 개 이상 나타나는 경우에는 각각의 어절에 표현 범위를 표시하였다. 그리고 명시적 부적절성을 발생시키는 어절을 포함하는 긴 문장이 띄어쓰기가 전혀 이루어지지 않은 채 하나의 어절로 제시되는 경우에는 분석 대상 말뭉치의 띄어쓰기를 수정하지 않는다는 원칙에 따라 어절이 나뉘지 않는 전체 문장에 표현 범위를 표시하였다. 마지막으로 명시적 부적절성을 발생시키는 관용구는 하나의 어절로 간주하여 구 단위에 표현 범위를 표시하였다.
- ‘비명시’의 경우에는 자동적으로 ‘문장’ 전체에 표현 범위를 표시하였다. 비명시로 주석된 문장만으로 비명시적 부적절성이 발생하는 맥락을 이해하기 어려운 경우에는 맥락 이해를 돕기 위해 선행 문장 또는 후행 문장을 각각 최대 5개까지 추가로 제시하였다. 그러나 해당 문장과 선행 또는 후행 문장 각각 5개까지를 포함하여 최대 11문장을 넘어선 범위까지 고려해야 하는 경우에는 부적절성 문장으로 선정하지 않았다. 또한 비명시적 부적절성이 나타나는 하나의 문장이 분석 대상 말뭉치에서는 둘 이상의 불완전 문장으로 나뉘어 있는 경우에는 부적절성의 직접적인 대상이나 주체가 포함된 문장에 표현 범위를 표시하였다.
- ‘명시’와 ‘비명시’가 함께 나타나는 경우에는 두 가지 표현 범위를 모두 표시하지 않고 명시적 표현 범위만 표시하였다.

3.3.3. 명시성 분석 및 주석의 기준

- 명시성 분석 및 주석의 구체적 기준과 상세한 예시는 부록 [붙임] 1의 작업 지침을 참조할 수 있다.

3.4. 부적절성의 맥락 주석

3.4.1. 맥락의 개념

- ‘맥락’은 해당 문장의 부적절성이 화자의 태도(의도)나 맥락 내용 측면에서 부정적인지, 긍정적인지를 판정하는 것을 의미한다. ‘맥락’의 분석은 판정된 부적절성이 긍정적 맥락에서 나타나는 경우의 특수성을 고려할 수 있게 하고, 후술하는 부적절성의 ‘강도’ 분석에서도 하나의 기준 또는 변수가 될 수 있다는 점에서 ‘맥락’은 부적절성의 중요한 분석 요소이다.

<표 4> ‘맥락’의 유형과 범위

맥락 유형	범위
부정적	<ul style="list-style-type: none"> - 부적절성이 화자의 태도(의도)나 맥락 내용 측면 모두에서 부정적으로 판단되는 문장 - 부적절성이 화자의 태도(의도) 측면에서는 부정적이지 않더라도 맥락 내용 측면에서 성 관련 폭력성, 선정성 등의 부적절성을 나타내는 문장 - 부적절성이 화자의 태도(의도) 측면에서 무표적으로 판단되거나 긍정성/부정성 판단이 불가능한 문장
긍정적	<ul style="list-style-type: none"> - 부적절성이 화자의 태도(의도)와 맥락 내용 측면 모두에서 긍정적인 것으로 판단되는 문장

- ‘맥락’은 본 사업의 발주 시나 제안 및 착수 보고 시에는 존재하지 않던 분석 요소인데, 자문회의 이후에 새로운 분석 요소로 추가되었다. 이는 판정된 부적절성이 화자의 태도(의도) 측면에서는 부정적이지 않더라도 그 맥락 내용 측면에서는 부적절한 경우, 부적절성으로 판정되었으나 화자의 태도(의도)와 맥락 내용 측면 모두에서 긍정적이거나 무표적인 것으로 판단되는 경우 등을 구별하여 분석하는 것이 필요하기 때문이다.

3.4.2. 맥락 분석 및 주석의 기준

- 맥락 분석 및 주석의 구체적 기준과 상세한 예시는 부록 [붙임] 1의 작업 지침을 참조할 수 있다.

3.5. 부적절성의 영역 주석

3.4.1. 영역의 개념

- ‘영역’은 해당 문장의 부적절성이 내용적 측면에서 어떤 영역과 관련되는지를 판정하는 것을 의미한다. ‘영역’의 분석은 부적절성이 발생하는 영역의 경향과 비중을 파악하고 예측할 수 있게 하며, 영역별 대응 또는 처리 방안을 모색할 수 있게 한다는 점에서 ‘영역’은 부적절성의 중요한 분석 요소이다.

<표 5> ‘영역’의 유형과 범위

영역 유형	범위
성	성별, 성적 지향, 성희롱 등
연령/세대	연령, 세대 등
신체	장애, 건강, 질병, 외모, 임신, 출산 등
문화	종교, 정치, 풍습, 예술, 행태, 사고방식 등
관계/조건	출신(인종, 국가, 지역 등) 사회적 조건(직업, 지위, 학력, 재산, 능력, 지력 등) 사적 관계(혼인, 가족 형태, 가족 관계, 친구(온라인 포함), 연인 등)
기타	사물(무정물), 감정의 배설 등 위의 유형으로 분류되지 않는 사례

- ‘영역’은 본 사업의 제안 및 착수 보고 시에는 관련 선행 연구들과 국가인권위원회 법 제2조 제3항에서 제시된 차별 행위의 목록 등을 참고하여 성, 연령/세대, 출신, 신체, 문화, 사회적 조건, 가족, 기타 등 8개 유형으로 분류하였다. 그런데 1차 작업 결과에서는 특정 영역이 50%를 상회하고 다른 영역은 5% 미만으로 나타나는 등 영역별 비율 편차가 매우 심하였다. 이러한 결과는 특정 영역의 비율이 35%를 넘지 않아야 한다는 제안요청서의 요건을 충족시키지 못한 것이며, 이렇게 구축되는 부적절성 말뭉치가 인공지능의 학습 자료로 적절할지에 대한 우려가 제기되었다. 이에 따라 최종 결과물에서 매우 낮은 비중으로 나타난 ‘출신(3%)’, ‘사회적 조건(14%)’, ‘가족

(2%)’ 등 세 개의 영역을 ‘관계/조건’ 영역으로 통합하였다.

3.5.2. 영역 분석 및 주석의 기준

- 영역 분석 및 주석의 구체적 기준과 상세한 예시는 부록 [붙임] 1의 작업 지침을 참조할 수 있다.

3.6. 부적절성의 강도 주석

3.6.1. 강도의 개념

- ‘강도’는 해당 문장의 부적절성이 그 심각성의 측면에서 어느 정도(강/약) 인지를 판정하는 것을 의미한다. 강도’의 분석은 부적절성의 강도 차이에 따라 그에 대한 대응 또는 처리 방안을 달리 마련할 수 있게 한다는 점에서 ‘강도’는 부적절성의 중요한 분석 요소이다.

<표 6> ‘강도’의 유형과 범위

강도 유형	범위
강	- 명시적 부적절성이 부정적 맥락에서 나타나는 문장 - 비명시적 부적절성이 성적 폭력성, 선정성 등 관련 부정적 맥락에서 나타나는 문장
약	- 명시적 부적절성이 긍정적 맥락에서 나타나는 문장 - 비명시적 부적절성이 성적 폭력성, 선정성 등 관련 부정적 맥락이 아닌 경우에 나타나는 문장

- ‘강도’는 작업자들 간에 판정의 통일성과 일관성을 확보하기 어려워서 분석 결과에 대한 객관성과 신뢰성이 부족하다는 점에서 분석 요소 1차 선정 시에는 제외하였지만, 2차 선정 시에는 부적절성 말뭉치를 활용할 학계와 산업계의 요구를 반영하여 다시 분석 요소에 포함하였다. 다만 본 사업에서는 강도 분석의 복잡성을 최소화하고 일관성을 최대화하고자 기존의 유사 연구에서 많이 사용한 3단계(상, 중, 하) 판정이 아닌 2단계(강, 약) 판정을 적용하였다.

3.6.2. 강도 분석 및 주석의 기준

- 강도 분석 및 주석의 구체적 기준과 상세한 예시는 부록 [붙임] 1의 작업 지침을 참조할 수 있다.

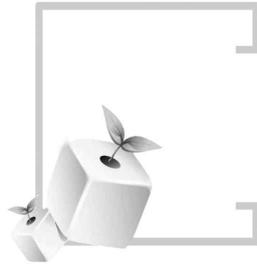
3.7. 부적절성 말뭉치의 비식별화

3.7.1. 개인정보 비식별화 기준

- 본 사업에서 구축되는 부적절성 말뭉치의 개인정보 비식별화 기준은 주관 기관과의 협의 하에 <2022년 인공지능의 언어 능력 평가 체계 운영 및 말뭉치 정비> 과제에서 정리된 ‘개인정보 판별 기준 및 비식별화 태그 세트’를 수정 없이 그대로 적용하였다. 이는 개인정보 포함 문장 판별 기준과 개인정보 비식별화 태그를 포함하고 있다.
- 개인정보 포함 문장을 판별하고 해당 문장 내 개인정보를 비식별화 기호로 처리하는 구체적 기준과 상세한 예시는 부록 [붙임] 1의 작업 지침을 참조할 수 있다.

3.7.2. 개인정보 비식별화 작업

- 본 사업의 최종 결과물에 대한 개인정보 비식별화 작업을 다음과 같이 진행하였다.
 - 공동연구진에서 개발한 비식별화 편집기를 통해 비식별화 작업 진행 (2023. 02. 14. ~ 2023. 02. 18.)
 - : 최종 16,240문장을 작업자 1인당 800여 개씩 배정하여 비식별화 작업을 진행하고 팀장이 검수



제 4 장

사업 결과와 논의 및 제안 사항



4.1. 사업 결과

- 본 사업의 결과로 구축된 부적절성 말뭉치의 문장 개수는 총 16,240개이다. 작업자들에 의해 검수 요청된 문장은 16,532개이었으나 검수 과정에서 작업 지침에 맞지 않는 문장, 비식별화 과정에서 부적절성 파악이 불가해지는 문장, 최종 확인 과정에서 동일 중복 출현 문장 등이 삭제되었다.

4.1.1. ‘명시성’, ‘강도’, ‘맥락’ 관련 분포

- 부적절성 말뭉치에서 나타나는 ‘명시성’, ‘강도’, ‘맥락’ 등 3개 분석 요소 관련 종합적 분포는 다음과 같다.

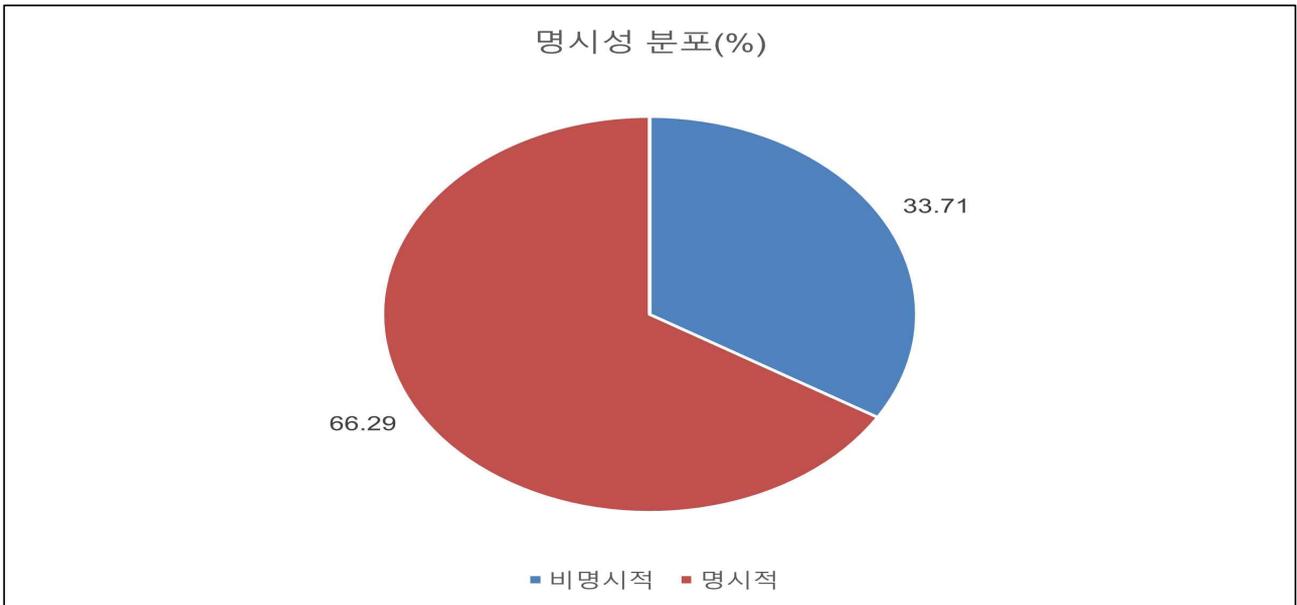
<표 7> 부적절성 문장의 ‘명시성’, ‘강도’, ‘맥락’ 관련 종합적 분포

(단위: 개수)

			강도		맥락	
			강	약	부정	긍정
명시성	명시적	10,765	7,769	2,996	7,799	2,966
	비명시적	5,475	848	4,627	4,729	746
합계		16,240	8,617	7,623	12,528	3,712
비중(%)			53.06%	46.94%	77.14%	22.86%

4.1.1.1. ‘명시성’ 관련 분포

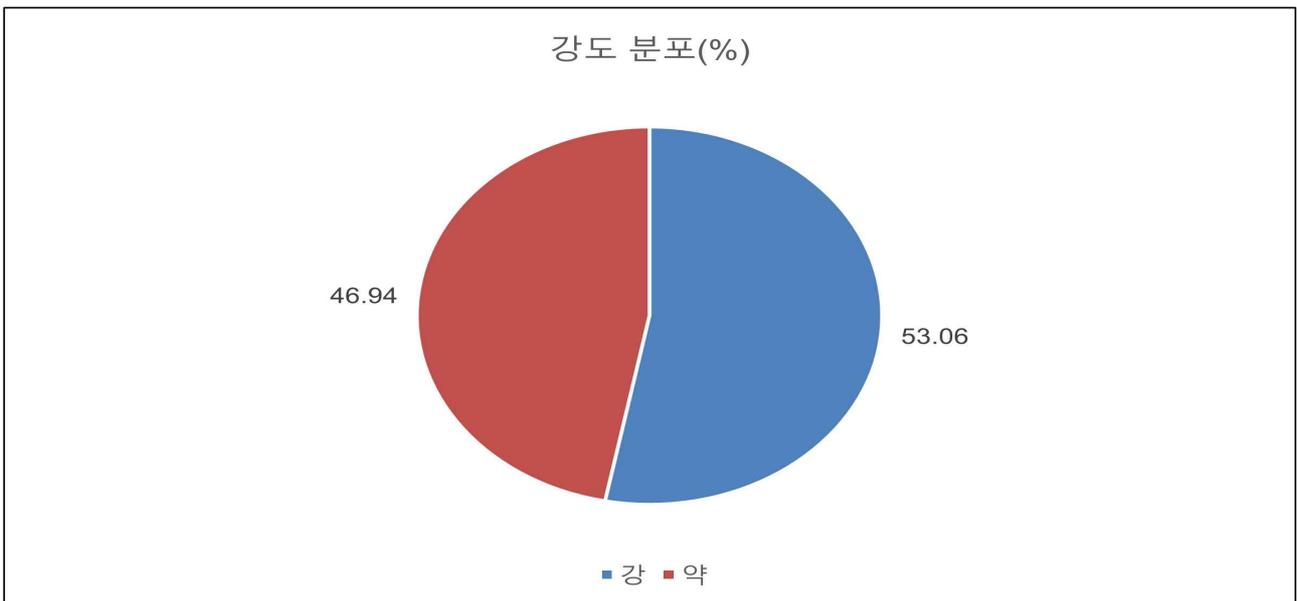
- 부적절성 말뭉치에서는 비명시적 부적절성 문장이 명시적 부적절성 문장보다 2배 가까이 많은 것으로 나타났다.



[그림 11] 부적절성 문장의 ‘명시성’ 관련 분포

4.1.1.2. ‘강도’ 관련 분포

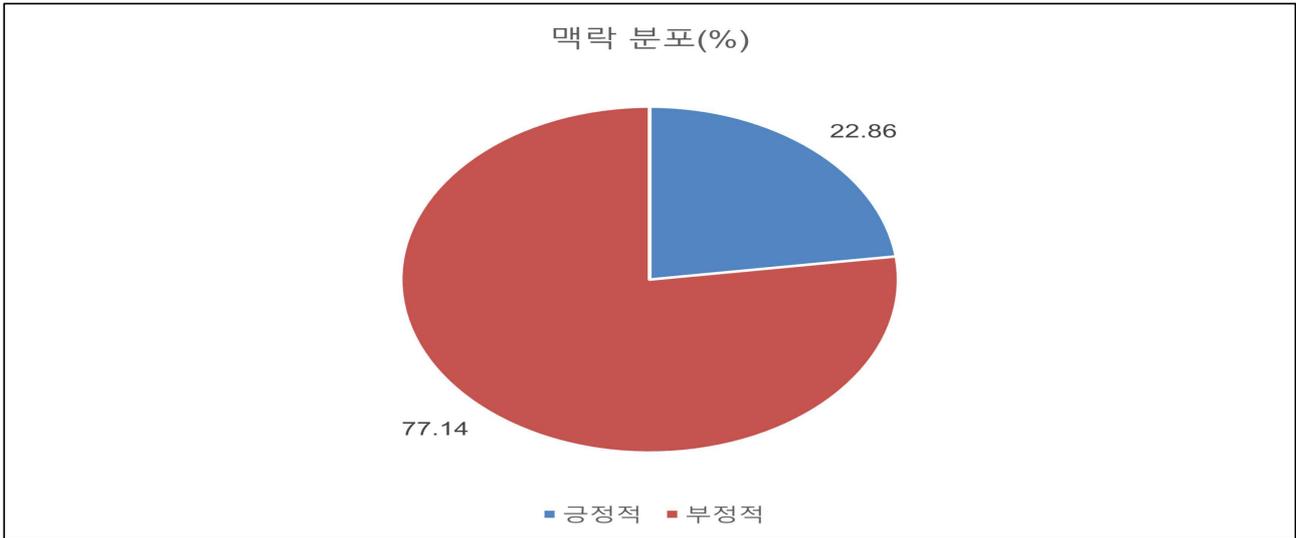
- 부적절성 말뭉치에서는 강도가 강한 부적절성 문장이 약한 부적절성 문장보다 조금 더 많은 것으로 나타났다.



[그림 12] 부적절성 문장의 ‘강도’ 관련 분포

4.1.1.3. ‘맥락’ 관련 분포

- 부적절성 말뭉치에서는 긍정적 맥락보다 부정적 맥락의 부적절성 문장이 3배 이상 많은 것으로 나타났다.



[그림 13] 부적절성 문장의 ‘맥락’ 관련 분포

4.1.2. ‘명시성’, ‘영역’ 관련 분포

- 부적절성 말뭉치에서 나타나는 ‘명시성’과 ‘영역’(단일 영역 기준)의 2개 분석 요소 관련 종합적 분포는 다음과 같다.

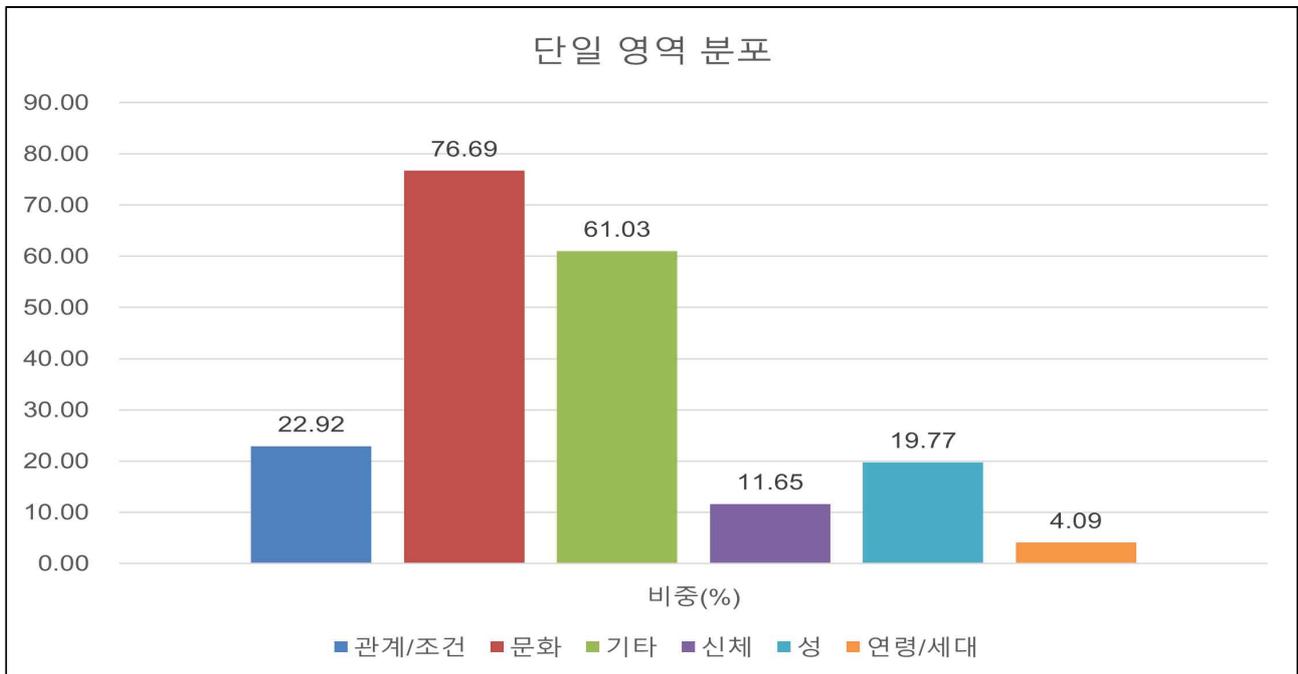
<표 8> 부적절성 문장의 ‘명시성’, ‘영역’ 관련 종합적 분포(단일 영역 기준)

(단위: 개수)

			영역					
			관계 /조건	문화	기타	신체	성	연령 /세대
명시성	명시적	10,765	1,773	4,429	2,413	814	1,372	242
	비명시적	5,475	1,950	8,026	7,498	1,078	1,839	422
합계		16,240	3,723	12,455	9,911	1,892	3,211	664
비중(%)			22.92%	76.69%	61.03%	11.65%	19.77%	4.09%

4.1.2.1. 단일 영역 기준 ‘영역’ 관련 분포

- 부적절성 말뭉치에 포함된 부적절성 문장의 영역 분포를 단일 영역 기준으로 보면, ‘문화’(76.69%)와 ‘기타’(61.03%) 영역의 비중이 지나치게 높게 나타났다.



[그림 14] 부적절성 문장의 ‘영역’ 관련 분포(단일 영역 기준)

4.1.2.1. 복합 영역 기준 ‘영역’ 관련 분포

- 부적절성 말뭉치에 포함된 부적절성 문장의 영역 분포를 복합 영역 기준으로 보면, 특정 하나의 영역이 전체 구축량의 35%를 초과하지 않는 것으로 나타났다. 예를 들어, 가장 많은 비중을 차지하는 ‘문화, 기타’ 영역의 비중이 31.35%이다. 하나의 문장에 둘 이상의 부적절성 영역이 관련되는 경우가 일반적임을 고려할 때, 복합 영역을 기준으로 영역 분포를 파악하는 것이 합리적일 것이다. 또한 부적절성 말뭉치를 인공지능의 언어 능력 평가 등에 활용하는 경우에도 복합 영역을 분류 기준으로 적용하는 것이 효율적일 것이다. (좀 더 상세한 내용은 4.2절의 논의 사항 참조)

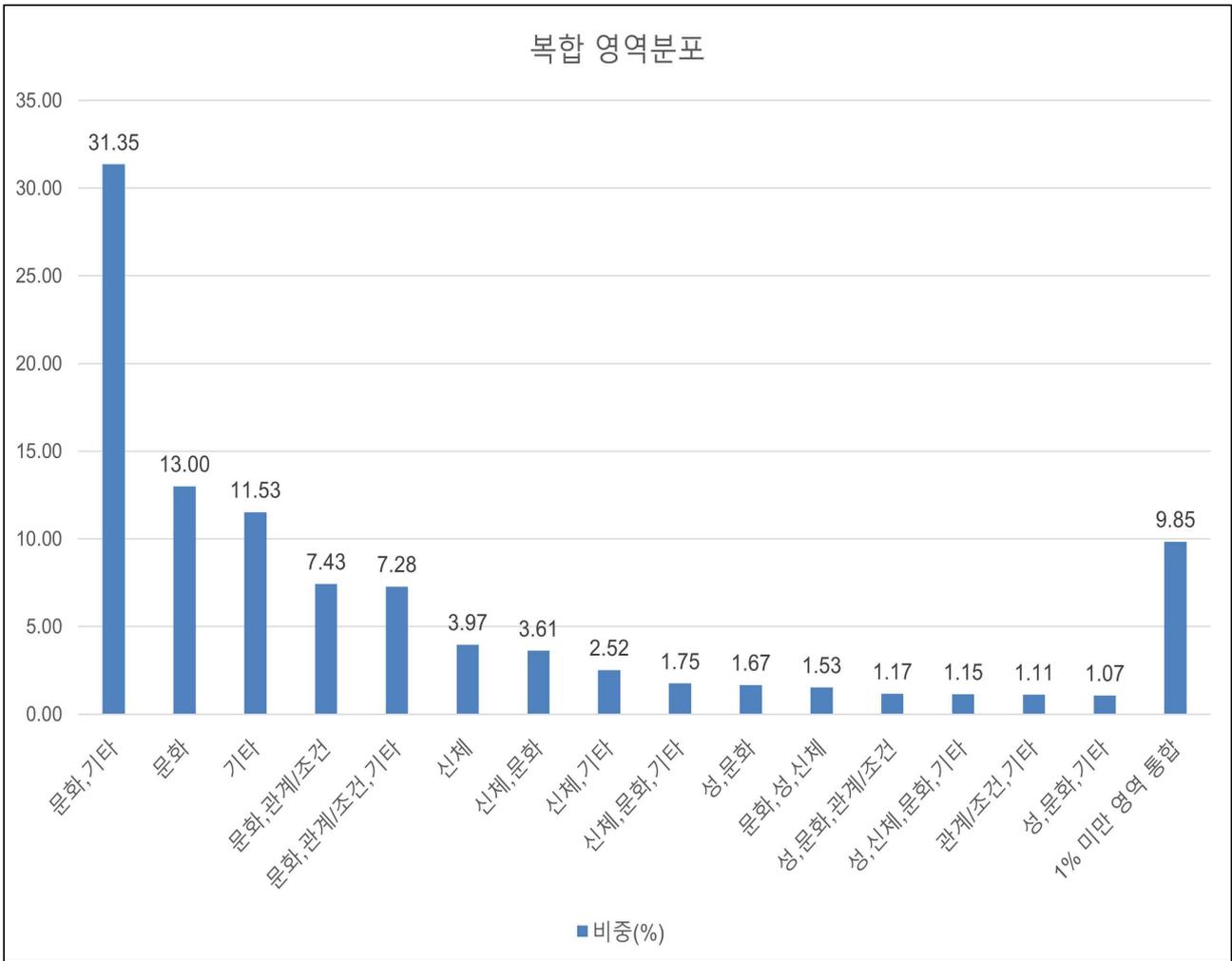
<표 9> 부적절성 문장의 '명시성', '영역' 관련 종합적 분포(복합 영역 기준)

(단위: 개수)

	영역					
	문화, 기타	문화	기타	문화, 관계/조건	문화, 관계/조건, 기타	신체
비명시	980	948	282	615	599	301
명시	4112	1164	1590	592	584	344
합계	5092	2112	1872	1207	1183	645
%	31.35	13.00	11.53	7.43	7.28	3.97

	영역					
	신체, 문화	신체, 기타	신체, 문화, 기타	성, 문화	문화, 성, 신체	성, 문화, 관계/조건
비명시	320	70	112	169	138	78
명시	267	339	172	102	110	112
합계	587	409	284	271	248	190
%	3.61	2.52	1.75	1.67	1.53	1.17

	영역				총합
	성, 신체, 문화, 기타	관계/조건, 기타	성, 문화, 기타	1% 미만 영역 통합	
비명시	30	82	81	670	5475
명시	156	99	93	929	10765
합계	186	181	174	1599	16240
%	1.15	1.11	1.07	9.85	100.00



[그림 15] 부적절성 문장의 '영역' 관련 분포(복합 영역 기준)

4.2. 논의 및 제언

4.2.1. 동일 부적절성 문장 처리 문제

- 동일한 부적절성 문장이 분석 대상 말뭉치의 다른 부분에서 중복하여 출현하는 경우가 있는데, 특히 ‘씨발’, ‘새끼’, ‘미친’ 등 욕설 표현 단독으로 하나의 문장이 형성된 부적절성 문장의 중복 출현이 잦은 편이었다. 이에 동일한 부적절성 문장은 한 번만 부적절성 분석 말뭉치로 구축하고 나머지 중복 출현 문장은 제외하는 것을 원칙으로 하였다.
- 다만, 동일한 부적절성 문장이 서로 다른 선후 맥락에서 나타나거나 해당 문장 전후로 문장부호, 이모티콘 등이 다르게 나타나는 경우에는 서로 다른 의미를 발생시키고 서로 다른 ‘영역’과 관련될 수 있다. 이에 동일한 부적절성 문장은 선후 맥락이 없거나 선후 맥락까지 동일한 경우에만 두 번째부터 출현하는 중복 문장을 제외하고, 형태적 변이가 있는 경우, 서로 다른 선후 맥락에서 나타나는 경우, 해당 문장 전후로 문장부호, 이모티콘 등이 다르게 나타나는 경우는 별도의 부적절성 문장으로 구축하였다.
- 본 사업은 매우 유사한 부적절성 문장이라도 약간의 차이로 인해 서로 다른 의미가 발생할 수 있고 서로 다른 영역과 관련될 수 있음을 확인하였다는 점에서 향후 관련 사업에서도 동일한 부적절성 문장의 판정 기준을 위와 같이 엄격히 설정할 것을 제언한다.

4.2.2. 분석 요소 ‘영역’별 비중 조정 문제

- 본 사업의 결과물로 구축된 부적절성 말뭉치에서 ‘출신’, ‘사회적 조건’, ‘가족’ 등 3개 영역은 부적절성 문장의 출현 비중이 매우 미미한 것으로 나타났다. 그런데 출신(인종, 국가, 지역 등), 사회적 조건(직업, 지위, 학력, 재산, 능력, 지력 등), 가족(혼인, 가족 형태, 가족 관계 등) 등 각 영역에 포함되는 하위 범위는 ‘관계/조건’이라는 개념으로 영역 통합이 가능하였다.

이에 최종 결과물에서는 출신, 사회적 조건, 가족 등 3개 영역을 ‘관계/조건’ 영역으로 통합하였다.

- 반면, ‘문화’와 ‘기타’ 영역은 해당 영역만 주석된 경우와 다른 영역과 함께 주석된 경우를 합산하면 부적절성 문장의 출현 비중이 지나치게 높은 것으로 나타났다. 그런데 이 두 영역은 해당 영역만 주석된 경우의 비중이 높지 않은 점을 고려할 때 대부분 다른 영역과 함께 주석되고 있음을 확인할 수 있다. 특정 영역이 지나치게 높은 비중을 차지하는 문제를 해결하기 위해 문화(종교, 정치, 풍습, 예술, 행태, 사고방식 등)와 기타 영역(사물/무정물, 감정의 배설 등)은 각 영역에 포함되는 하위 범위 중 일부를 또 하나의 다른 영역으로 설정하여 분리하는 방안을 고려해 볼 수 있으나, 이렇게 하면 새로운 영역과 기존 영역 간에 층위가 맞지 않는 문제가 발생한다. 이에 최종 결과물에서의 영역별 비중 제시에서는 하나의 영역만 주석된 경우와 다른 영역과 함께 주석된 경우를 합산하지 않고 별도로 분리하여 후자의 경우도 하나의 영역(복합 영역)으로 간주하는 방식을 취하였다.

-> 본 사업은 부적절성의 영역을 하나의 대표 영역만으로 배타적으로 주석하지 않고 여러 가지 영역을 복수로 주석하는 방식을 취했다는 점에서 위와 같은 영역 비중 제시가 적절하였으며, 향후 관련 사업에서도 같은 주석 방식을 취한다면 동일한 영역 비중 제시 방식을 취할 것을 제언한다.

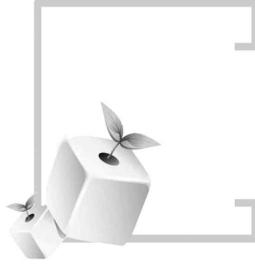
4.2.3. 사업 수행 관련 문제

- 본 사업의 과제 범위는 크게 부적절성 분석을 위한 분석 방법론 및 세부 지침을 수립하는 것과 수립된 방법론 및 지침을 바탕으로 부적절성 분석 말씀치를 구축하는 것 등의 두 가지로 나눌 수 있는데, 전자의 과제를 수행하는 데에 전체 사업 기간의 80% 이상이 소요되었다. 착수 이후 사업 초반에는 연구진 논의와 국어원 협의를 통해 부적절성 분석 작업의 객관성과 일관성을 중시하게 되면서 분석 요소를 ‘명시성’과 ‘영역’으로 한정하는 방향으로 분석 지침을 수립하고 표본 작업을 진행하면서 분석 지침을 보완하는 데에 이미 많은 시간이 소요되었다. 그런데 사업 기간의 절반을 넘긴 시점에

개최된 자문회의에서 제시된 인공지능(AI) 관련 학계와 업계의 의견 및 요구를 반영하게 되면서 분석 요소에 ‘맥락’과 ‘강도’를 추가하게 되었다.

- 추가된 분석 요소인 ‘맥락’과 ‘강도’는 기존 분석 요소인 ‘명시성’과 ‘영역’에 비해 작업자의 주관성이 개입될 여지가 상대적으로 크기에 주관성을 최소화하는 작업 지침 수정과 지침 최종안 공유 및 교육 등에 다시 추가적인 시간이 상당히 소요되었다. 결과적으로 분석 작업 방법론과 지침 수립 과정은 충실하게 수행된 편인 반면에 부적절성 분석 말뭉치 구축 과정은 매우 짧은 기간에 급박하게 수행되었다.

-> 본 사업은 말뭉치 부적절성 분석을 위한 분석 방법론 및 세부 지침을 수립했다는 점에 가장 큰 의의가 있으며, 본 사업에서 구축한 부적절성 분석 말뭉치는 인공지능 언어 능력 평가 체계 운영 등에 활용할 수 있을 것으로 기대된다. 다만 분석 요소 간 균형성과 인공지능의 언어 능력 평가 및 학습에의 활용성이 높은 부적절성 말뭉치의 구축을 위해서는 장기적인 계획 하에 다양한 출처에 기반한 대규모 말뭉치 분석 및 구축 작업을 후속 사업으로 추진할 것을 제언한다.



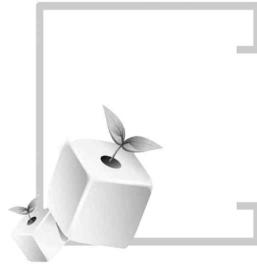
참고문헌



- 과학기술정보통신부(2020), <사람이 중심이 되는 인공지능 윤리기준>, 과학기술정보통신부.
- 과학기술정보통신부(2020), <윤리적 인공지능을 위한 국가정책 수립>, 과학기술정보통신부.
- 김수한·김성훈·김현철(2019), 해외 인공지능 교육동향과 학습도구 분석, <한국컴퓨터교육학회 학술발표대회논문집> 23(2), 25-28.
- 김봉제(2020), 인공지능 거짓말의 특성 이해 - 거짓말에 대한 윤리학적 담론을 중심으로, <인공지능인문학연구> 6, 79-97.
- 김봉제 외(2022), <비윤리적 표현 말뭉치 연구 분석 및 시범 구축>, 국립국어원.
- 김진웅(2021), 자연언어처리에서 윤리적 문제와 해결 방안: 연령 및 지역 편향성 극복의 출발점으로서 방언자료 수집, <연구방법논총> 6(1), 157-180.
- 김한성 외(2019), <모두를 위한 인공지능 윤리>, 한국교육학술정보원.
- 박재현·이승희(2009), <사회적 의사소통 연구: 지역·민족·인종에 대한 차별적 언어표현 개선 연구>, 국립국어원.
- 변순용(2020), 데이터 윤리에서 인공지능 편향성 문제에 대한 연구, <윤리연구> 1(128), 143-158.
- 변순용(2020), AI 시민성 교육에 대한 시론. <초등도덕교육> 67, 427-445.
- 변순용·김봉제·이청호(2021), 인공지능 윤리교육 내용체계 구성에 관한 연구, <한국초등도덕교육학회 학술대회 자료집>, 265-288.
- 윤상오(2018), 인공지능 기반 공공서비스의 주요 쟁점에 관한 연구: 챗봇(ChatBot) 서비스를 중심으로, <한국공공관리학보> 32(2), 83-104.
- 이찬규 외(2021), <말뭉치 언어의 사회적 인식 조사·분류>, 국립국어원.
- 이청호·김봉제·김형주·변순용·이찬규(2021). 윤리적 인공지능을 위한 비도덕 문장 판별 온톨로지 구축에 대한 연구, <인공지능인문학연구> 7, 149-170.
- 조태린(2019), 언어 사전의 정보적 기능과 윤리적 문제에 대한 소고, <한국사전학> 34, 105-126.
- 조태린(2021), 언어의 품격과 공공언어의 품격(성) 문제에 대한 비판적 고찰, <문법 교육> 41, 97-124.
- 조태린 외(2006), <차별적, 비객관적 언어 표현 개선을 위한 기초 연구>, 국립국어원.

- 조태린·김신각·유희재·김예지·이주희(2018), 대화형 인공지능의 윤리적 언어 표현을 위한 기초 연구 - 단어 단위의 비윤리적 언어 표현의 유형 분류를 중심으로, <어문학> 140, 65-96.
- 조태린·김신각·신유리·공나형·신아영(2018), 비윤리적 언어 표현의 의미 자질에 따른 하위 유형 연구 - 대화형 인공지능의 윤리적 언어 표현을 위한 기초 연구(2), <한민족 문화연구> 63(63), 147-184.
- 차정원 외(2020), <2020년 개체명 말뭉치 연구 분석>, 국립국어원.
- 최현철·변순용(2019), 인공적 도덕 행위자에 대한 융합접근의 철학적 기획, <윤리연구> 1(124), 1-16.
- Armstrong, H.(2015), Machines that learn in the wild: Machine learning capabilities, limitations and implications. *Technical report*, Nesta, London, England.
- Cutler, A., Pribić, M., & Humphrey, L. (2019). *Everyday ethics for artificial intelligence*, IBM Corporation.
- Dadvar, M., Trieschnigg, D., Ordelman, R. and de Jong, F. (2013). Improving cyberbullying detection with user context. In *European Conference on Information Retrieval*. Springer, 693-696.
- Davidson, T., Warmsley, D., Macy, M. and Weber, I. (2017). Automated hate speech detection and the problem of offensive language. In *Proceedings of the International AAAI Conference on Web and Social Media*, 11(1).
- Dinakar, K., Reichart, R. and Lieberman, H. (2011). Modeling the detection of textual cyberbullying. *The Social Mobile Web* 11(02).
- Hosseini, H, Kannan, S, Zhang, B, & Poovendran, R.(2017), Deceiving Google's perspective API built for detecting toxic comments, *arXiv* 1702:08138.
- Jobin, A., Ienca, M., & Vayena, E. (2019). The global landscape of AI ethics guidelines. *Nature Machine Intelligence*, 1(9), 389-399.
- Justo, R., Corcoran, Th., Lukin, S. M., Walker, M. and Torres, M. I. (2014). Extracting relevant knowledge for the detection of sarcasm and nastiness in the social web. *Knowledge-Based Systems* 69:124 - 133.

- Kumar, R., Ojha, A. K., Malmasi, S., & Zampieri, M.(2018). Benchmarking aggression identification in social media. In *Proceedings of the First Workshop on Trolling, Aggression and Cyberbullying* (TRAC-2018): 1-11.
- Lee, Geon-hwan, Park, Joo-chan, Choi, Dong-won, Lee, Yeon-gyeong, Choi, Ho-bin, and Han, Youn-hee (2019). Design and implementation of profanity filtering chat program based on deep learning. In *Korea Information Processing Society Conference Proceedings*, 26(2), 998-1001.
- Park, Kyohyeon and Lee, Jee Hyung (2006). Developing a Vulgarity Filtering System for Online Games using SVM, In *Korean Institute of Information Scientists and Engineers Academic Conference Proceedings*, 33(2B), 260-263.
- Roberts, H., Cowls, J., Morley, J., Taddeo, M., Wang, V., & Floridi, L. (2021). The Chinese approach to artificial intelligence: an analysis of policy, ethics, and regulation. *AI & SOCIETY*, 36(1), 59-77.
- Shahriari, K., & Shahriari, M. (2017, July). IEEE standard review—Ethically aligned design: A vision for prioritizing human wellbeing with artificial intelligence and autonomous systems. In *2017 IEEE Canada International Humanitarian Technology Conference* (IHTC) (pp. 197-201). IEEE.
- Vakkuri, V., Kemell, K., Jantunen, M., Halme, E., & Abrahamsson, P.(2021). ECCOLA – A method for implementing ethically aligned AI systems. *Journal of Systems and Software*, 182 doi:10.1016/j.jss.2021.111067.
- Waseem, Z., and Hovy, D. (2016). Hateful symbols or hateful people? predictive features for hate speech detection on twitter. In *Proceedings of the NAACL student research workshop*, 88-93.
- Wassem, Z., Davidson, Th., Warmesley, D. and Weber I. (2017). Understanding abuse: A typology of abusive language detection subtasks. In *Proceedings of the First Workshop on Abusive Language Online*, 78-84.
- Zampieri, M., Malmasi, S., Nakov, P., Rosenthal, S., Farra, N., & Kumar, R. (2019). Predicting the type and target of offensive posts in social media. *arXiv preprint arXiv:1902.09666*.



부록

[붙임 1]

말뭉치 부적절성

분석 작업 지침



말뭉치 부적절성 분석 작업 지침_최종

차 례

1. 부적절성의 개념	2
2. 부적절성의 분석 요소	3
2.1. 명시성	3
2.2. 맥락	3
2.3. 영역	4
2.4. 강도	4
3. 부적절성의 분석 단위와 주석(태깅) 단위	5
3.1. 분석 단위	5
3.2. 주석(태깅) 단위	5
4. 부적절성의 주석 방법	7
4.1. 명시성 주석	7
4.1.1. 명시성의 표현 범위	7
4.1.2. 명시	9
4.1.3. 비명시	14
4.1.4. 명시와 비명시가 함께 나타나는 경우	19
4.2. 맥락 주석	20
4.2.1. 부정적 맥락	20
4.2.2. 긍정적 맥락	21
4.3. 영역 주석	23
4.3.1. 복수 영역의 처리	23
4.3.2. 성	25
4.3.3. 연령/세대	26
4.3.4. 신체	27
4.3.5. 문화	28
4.3.6. 관계/조건	29
4.3.7. 기타	33
4.4. 강도 주석	35
4.4.1. 강	35
4.4.2. 약	37
5. 개인정보 판별 기준과 비식별화 태그 세트	39
5.1. 개인정보 포함 문장 판별 기준	39
5.2. 개인정보 비식별화 태그	41

1. 부적절성의 개념

- ‘부적절성’은 ‘비윤리성’은 물론이고 엄격한 의미에서의 ‘비윤리성’으로 간주하기 어려운 부정적 특성까지 포함하는 폭넓은 개념
 - ‘부적절성’을 구성하는 기본적 개념인 ‘비윤리성’은 사람으로서 마땅히 지키거나 행해야 할 도리에 어긋나거나 도덕적 규범을 위반하는 것을 의미
 - 비윤리성은 일반적으로 ‘공격성’, ‘편향성’, ‘비하성’ 등의 부정적 특성들을 포함
 - 비난, 저주, 모욕, 위협, 혐오, 폭력선동 등에서 나타나는 ‘공격성’
 - 차별, 편견, 배제, 불필요한 언급 등에서 나타나는 ‘편향성’
 - 멸시, 폄하, 무시, 조롱 등에서 나타나는 ‘비하성’
 - 이상의 부정적 특성들은 엄격한 의미에서의 비윤리성으로 간주되기 어려운 경우에도 나타나며, 관련 발화를 접하는 누군가에게 다양한 유형의 부정적 감정 유발 가능
- => 이에 따라 본 과제에서는 기존의 유사 과제에서 사용해온 ‘비윤리성’, ‘부정성’, ‘반사회성’ 등의 용어 대신 ‘부적절성’이라는 용어를 사용

2. 부적절성의 분석 요소

2.1. 명시성

○ ‘명시성’은 부적절성이 구체적인 어휘나 표현을 통해 명시적으로 드러나는지(명시), 해당 문장의 맥락에서 드러나는지(비명시)를 판정하는 것을 의미

명시성 유형	범위
명시	- 사전에 근거하는 욕설, 비어, 비하성 속어 등이나 관련 보고서에 근거하는 차별 표현, 혐오 표현, 선정적 표현 등과 같은 명시적 부적절성 표현이 나타나는 문장
비명시	- 명시적 부적절성 표현이 나타나지는 않지만, : 해당 문장의 내외 맥락에서 비윤리성, 공격성, 편향성, 비하성 등의 부적절성이 확인되는 문장 : 대상에 대한 부정적 평가 혹은 태도를 드러내는 표현이 포함된 문장 : 현대적인 기준에서 윤리성을 크게 위배하였다고 판단되는 내용을 담은 표현이 나타나는 문장

- ‘명시성’의 분석은 부적절성 여부에 대한 논란이 적은 명시적 부적절성 문장을 우선적으로 선별할 수 있게 함으로써 부적절성의 판정을 최대한 객관적으로 할 수 있게 한다는 점에서 부적절성의 일차적인 분석 요소

2.2. 맥락

○ ‘맥락’은 해당 문장의 부적절성이 화자의 태도(의도)나 맥락 내용 측면에서 부정적인지, 긍정적인지를 판정하는 것을 의미

맥락 유형	범위
부정적	- 부적절성이 화자의 태도(의도)나 맥락 내용 측면 모두에서 부정적으로 판단되는 문장 - 부적절성이 화자의 태도(의도) 측면에서는 부정적이지 않더라도 맥락 내용 측면에서 성 관련 폭력성, 선정성 등의 부적절성을 나타내는 문장 - 부적절성이 화자의 태도(의도) 측면에서 무표적으로 판단되거나 긍정성/부정성 판단이 불가능한 문장
긍정적	- 부적절성이 화자의 태도(의도)와 맥락 내용 측면 모두에서 긍정적인 것으로 판단되는 문장

- ‘맥락’의 분석은 판정된 부적절성이 긍정적 맥락에서 나타나는 경우의 특수성을 고려할 수 있게 하고, 후술하는 부적절성의 ‘강도’ 분석에서도 하나의 기준 또는 변수가 될 수 있다

는 점에서 부적절성의 중요한 분석 요소

2.3. 영역

○ ‘영역’은 해당 문장의 부적절성이 내용적 측면에서 어떤 영역과 관련되는지를 판정하는 것을 의미

영역 유형	범위
성	성별, 성적 지향, 성희롱 등
연령/세대	연령, 세대 등
신체	장애, 건강, 질병 ²⁾ , 외모, 임신, 출산 등
문화	종교, 정치, 풍습, 예술, 행태, 사고방식 등
관계/조건 ³⁾	출신(인종, 국가, 지역 등) 사회적 조건(직업, 지위, 학력, 재산, 능력, 지력 등) 사적 관계(혼인, 가족 형태, 가족 관계, 친구(온라인 포함), 연인 등)
기타	사물(무정물), 감정의 배설 등 위의 유형으로 분류되지 않는 사례

- ‘영역’의 분석은 부적절성이 발생하는 영역의 경향과 비중을 파악하고 예측할 수 있게 하며, 영역별 대응 또는 처리 방안을 모색할 수 있게 한다는 점에서 부적절성의 중요한 분석 요소

2.4. 강도

○ ‘강도’는 해당 문장의 부적절성이 그 심각성의 측면에서 어느 정도(강/약)인지를 판정하는 것을 의미

강도 유형	범위
강	- 명시적 부적절성이 부정적 맥락에서 나타나는 문장 - 비명시적 부적절성이 성적 폭력성, 선정성 등 관련 부정적 맥락에서 나타나는 문장
약	- 명시적 부적절성이 긍정적 맥락에서 나타나는 문장 - 비명시적 부적절성이 성적 폭력성, 선정성 등 관련 부정적 맥락이 아닌 경우에 나타나는 문장

- ‘강도’의 분석은 부적절성의 강도 차이에 따라 그에 대한 대응 또는 처리 방안을 달리 마련할 수 있게 한다는 점에서 부적절성의 중요한 분석 요소

2) 장애, 건강, 질병은 신체적인 것과 정신적인 것을 모두 포함하여 ‘신체’ 영역으로 판정한다.

3) 작업 과정에서는 ‘출신’, ‘사회적 조건’, ‘가족’을 각각 독립적인 영역 유형으로 구별하였으나, 최종 결과물에서 ‘출신(3%)’, ‘사회적 조건(14%)’, ‘가족(2%)’의 영역이 매우 낮은 비중으로 실현되었기 때문에 이 세 가지 영역을 ‘관계/조건’ 영역으로 통합하여 제시하고, ‘가족’은 ‘사적 관계’로 명칭을 변경하였다.

3. 부적절성의 분석 단위와 주석(태깅) 단위

3.1. 분석 단위

○ 부적절성 분석의 기본 단위는 ‘문장’

- 부적절성이 문장 내의 특정 어휘나 표현에 의해 발생하더라도 부적절성 분석의 기본 단위는 ‘문장’
- 문장은 분석 대상 말뭉치에서의 문장이 한글 맞춤법이나 띄어쓰기 오류 등 어문 규범에 어긋나는 요소들 또는 오타 등을 포함하고 있더라도 수정 없이 그대로 분석

¶ 환경이 그지(√저지) 같아.

¶ 블로그짓한다고(√블로그짓 한다고) 열심히 사진찍고 있습(√있음)

¶ 진짜 저 사진 찍어놓고 너무 멍청해보여서(√멍청해보여서) 한동안 웃었다.ㅋㅋ

¶ 내 장동건은 지금처럼 망가지지 않을 때가 있었어!! 쥘압(√제발) 이상한 cf좀 고만 찍어라! 멍청해뵈다(√멍청해 뵈다) T T

- 부적절성의 4가지 분석 요소인 명시성, 맥락, 영역, 강도 등은 문장 내의 특정 어휘나 표현과 관련된 것이라도 전체 문장의 부적절성으로 제시
- 최종 결과물로서의 부적절성 말뭉치 구축의 단위도 ‘문장’ (15,000 문장 이상)

3.2. 주석(태깅) 단위

○ 부적절성 분석의 결과를 태깅하는 기본 단위도 ‘문장’

- 부적절성의 명시성, 맥락, 영역, 강도는 문장 내의 특정 어휘나 표현에 의해 판정된 것이더라도 문장 단위로 태깅

○ 하나의 문장에서 명시성, 맥락, 영역, 강도 등의 분석 요소가 각각 복수로 판정되는 경우는 문장 단위에 다음의 기준에 따라 태깅(분석 요소별 상세한 태깅 방식은 후술 참조)

- 명시와 비명시가 함께 나타나는 문장은 복수 태깅을 하지 않고 ‘명시’로 태깅
- 부정적 맥락과 긍정적 맥락이 함께 나타나는 문장은 복수 태깅을 하지 않고 ‘부정적’ 맥락으로 태깅
- 두 개 이상의 영역에 해당하는 부적절성이 나타나는 문장은 복수 태깅을 권장
- 강한 부적절성과 약한 부적절성이 함께 나타나는 문장은 ‘강’한 부적절성으로 태깅

¶ 머리 검은짐승은 거두는게 아니랬는데 재들은 뭐.. 그냥 **깜둥이들은** 거두면 안되는 듯.

- ▶ [명시성] 명시
- ▶ [맥락] 부정적
- ▶ [영역] 관계/조건, 문화
- ▶ [강도] 강

☞ 이 문장은 명시적 부정절성(깜둥이들은)과 비명시적 부정절성(타인을 믿으면 배신당한다는 사고방식)을 모두 포함하지만, '명시'와 '비명시'가 함께 나타나는 문장은 '명시'로 태깅한다는 원칙에 따라 '명시'로 태깅하고, 부정절성의 강도도 명시는 '강'인 반면에 비명시는 '약'으로 판정되지만, '강'과 '약'이 함께 나타나는 문장은 '강'으로 태깅한다는 원칙에 따라 '강'으로 태깅하되, 영역은 복수 태깅을 권장한다는 원칙에 따라 명시와 관련된 '관계/조건'과 비명시와 관련된 '문화'를 모두 태깅한다.

4. 부적절성의 주석 방법

4.1. 명시성 주석

4.1.1. 명시성의 표현 범위(시작-종료, begin-end)

○ 명시성은 문장 단위로 태깅하지만, 명시인지, 비명시인지에 따라 부적절성이 표현되는 범위를 달리 표시
(작업 도구에서는 클릭으로 표시하지만, 본 지침에서는 표현 범위를 붉은색으로 표시함)

4.1.1.1. 명시성의 표현 범위

- ‘명시’의 경우에는 명시적 부적절성 표현이 나타나는 ‘어절’ 단위(체언+조사, 용언+어미 등)에 표현 범위를 표시

¶ 여기 **들딱들이** 대거 나온 듯.

☞ ‘들딱’이 명시적 부적절성 표현이지만 ‘들딱들이’라는 어절 전체에 표현 범위를 표시한다.

¶ 오늘도 **구라치다** 하루가 다 갔어~

☞ ‘구라’가 명시적 부적절성 표현이지만 ‘구라치다’라는 어절 전체에 표현 범위를 표시한다.

- 하나의 문장에 명시적 부적절성을 발생시키는 어절이 두 개 이상 나타나는 경우에는 각각의 어절에 표현 범위를 표시
(단, 해당 문장의 부적절성 맥락, 영역, 강도 등은 문장 단위로 태깅, 후술 참조)

¶ 이번의 똥 준 **눈(남자라면 너무 무섭다..:)**을 보면서 또 한 번 느꼈는데 빠순이는 악의는 없지만 정말 정말 이해가 안 간다는 거.

- 명시적 부적절성을 발생시키는 어절을 포함하는 긴 문장이 띄어쓰기가 전혀 이루어지지 않은 채 하나의 어절로 제시되는 경우에는 분석 대상 말뭉치의 띄어쓰기를 수정하지 않는다는 원칙에 따라 어절이 나뉘지 않는 전체 문장에 표현 범위를 표시

¶ **시발진짜내가쪽팔림을무릅쓰고학원쌤한테잠깐만폰해도되냐고물어봤는데**

☞ ‘시발’이 명시적 부적절성 표현이지만 띄어쓰기가 전혀 이루어지지 않았으므로 문장 전체에 표현 범위를 표시한다.

- 명시적 부적절성을 발생시키는 관용구⁴⁾는 하나의 어절로 간주하여 구 단위에 표현 범위를 표시

¶ 그러나 나는 이미 **눈이 빠였으니** 선생님이 웃을 때 보이는 눈가의 잔주름까지 마냥 귀여워보였다.

☞ ‘눈(이) 빠다’는 <표준>에 ‘뻥한 것을 잘못 보고 있을 때 비난조로 이르는 말’이라는 뜻의 관용구로 등재되어 있으므로 구를 구성하는 두 어절에 표현 범위를 표시한다.

4.1.1.2. 비명시의 표현 범위

4) 관용구의 판정은 대사전류 및 관련 보고서에 근거한다.

- '비명시'의 경우에는 자동적으로 '문장' 전체에 표현 범위를 표시

¶ 하지만 전문가 말만큼 믿을 게 못 되는 것도 없다.

¶ 외국 아가씨들이랑 놀기 위해..

¶ 원작판에 비해 키도 커지고 몸매도 좋아졌다.

- 비명시로 태깅된 문장만으로 비명시적 부적절성이 발생하는 맥락을 이해하기 어려운 경우에는 맥락 이해를 돕기 위해 선행 문장 또는 후행 문장을 각각 최대 5개까지 추가 제시 가능

¶ 기껏 국물이라고 준 오뎅국물?은
미지근해서 따뜻하게 땀혀달라고 해도
고대로였어요 ㅋㅋㅋㅋ
상무지구 오징어나라는
그냥 다음부터 안갈 것 같아요

왜 사람들은 초심을 잃는거죠

☞ "왜 사람들은 초심을 잃는 거죠"라는 문장이 불친절한 서비스에 대한 비난과 조롱에 해당하는 비명시적 부적절성 문장임을 확인할 수 있도록 선행 문장 5개를 함께 제시한다.

<주의>

- 비명시적 부적절성 문장의 판정과 분석은 해당 문장 내 맥락만이 아니라 선행 및 후행 문장의 맥락을 고려하되, 해당 문장과 선행 또는 후행 문장 각각 5개까지를 포함하여 최대 11문장을 넘어선 범위까지 고려해야 하는 경우에는 부적절성 문장으로 판정하지 않음

¶ 몸에서 옥수가 나오나여?

☞ 이 문장은 '남들 흑시 옥수생은'이라는 제목의 글에 포함된 문장으로서 제목을 함께 고려하였을 경우 해당 문장은 부적절성 문장으로 판정하는 것이 타당하다. 그러나 기구축 말뭉치 언어 자료에 제목이 포함되어 있지 않아 선·후행 문장이 제시되어 있지 않아 해당 문장을 시험 6수생에 대한 비하로 볼 수 있는 근거가 부족하므로 이를 부적절성 문장으로 판정하지 않는다.

¶ 아 떡방아 찢고싶다

☞ 위 문장은 선·후행 문장이 제시되지 않은 채 실현된 것으로, 축자적 의미 그대로 '방아를 찢고 싶다'인 것인지 성적인 행위를 비속하게 표현한 것인지 알 수 없으므로 맥락 의미를 정확히 추측하기 어려운 경우 이를 부적절성 문장으로 판정하지 않는다.

- 비명시적 부적절성이 나타나는 하나의 문장이 분석 대상 말뭉치에서는 둘 이상의 불안전 문장으로 나뉘어 있는 경우에는 부적절성의 직접적인 대상이나 주체가 포함된 문장에 표현 범위를 표시

(단, 나뉘어 있는 선행 또는 후행 불안전 문장도 추가 제시 필요)

¶ 정말이지 중국 사람들..

엄청 시끄럽다

☞ 위의 두 문장은 하나의 문장일 수 있지만 말뭉치에서 두 개로 나누어 제시하고 있으므로 부적절성의 직접

적인 대상이나 주체가 포함된 선행 문장 “정말이지 중국 사람들...”만을 비명시적 부적절성 문장으로 태깅하되 그 맥락을 이해할 수 있도록 후행 문장 “엄청 시끄럽다”를 추가 제시한다.

4.1.1.3. 명시와 비명시가 함께 나타나는 경우의 표현 범위

- 명시와 비명시의 표현 범위를 모두 표시하는 것은 불가하고 명시의 표현 범위만 표시

¶ 왜 키우냐 묻거든 강 웃지요 근데 **꼴에** 머리는 겁나 좋다 그냥 똥똥한 사람같다 ㅋㅋ 오래 살아라,

☞ ‘꼴’은 <표준> 등의 대사전류에 ‘사람의 모양새나 행태를 낮잡아 이르는 말’로 뜻풀이되어 있으므로 신체를 비하하는 ‘명시적 부적절성’으로 판정하고, 비명시적 부적절성으로는 고양이를 키우는 행태에 대한 자조적 내용과 함께 사람의 체형에 대한 (대체로) 부정적 태도가 함께 드러나지만, 명시에 해당하는 ‘꼴’에만 표현 범위를 표시한다.

¶ 이번의 똥 준 **논(남자라면 너무 무섭다...)**을 보면서 또 한 번 느꼈는데 빠순이는 악의는 없지만 정말 정말 이해가 안 간다는 거.

☞ 두 개의 명시적 부적절성 외에 행동을 분노에 비유하는 비명시적 부적절성이 나타나지만, 명시에 해당하는 두 어절에만 표현 범위를 표시한다.

4.1.2. 명시

○ ‘명시’는 다음의 조건들을 엄격하게 충족하는 문장들에 한정하여 태깅

4.1.2.1. 사전에 욕설, 비어, 비하성 속어로 기술된 어휘나 표현이 나타나는 문장

- <표준국어대사전(이하 표준)>과 <고려대한국어대사전(이하 <고려대>)> 등의 대사전류 중 어느 하나에서라도 욕설, 비어⁵⁾, 비하성 속어⁶⁾ 등으로 기술되는 어휘나 표현이 나타나는 문장은 ‘명시’로 태깅

¶ (가족의 상황에 대한 묘사) 돈벌러 가신 **남정네분들은** 늦게 퇴근하고.

▶ [명시성] 명시 [맥락] 긍정적 [영역] 성, 관계/조건 [강도] 약

☞ ‘남정네’는 <고려대>에서 ‘여자가 사내를 조금 낮추어 이르는 말’로 뜻풀이되어 있어 비어에 해당하므로 이를 포함하는 문장은 ‘명시’로 태깅한다.

¶ 오늘도 **구라치다** 하루가 다 갔어~

▶ [명시성] 명시 [맥락] 부정적 [영역] 문화 [강도] 강

☞ <고려대>에서 ‘구라치다’는 ‘비속하게’로 제시되어 있어 비하성 속어에 해당하므로 이를 포함하는 문장은 ‘명시’로 태깅한다.

¶ 대학생이나 되는데, 생각은 참 **초딩** 같아.

▶ [명시성] 명시 [맥락] 부정적 [영역] 연령/세대, 문화, 관계/조건 [강도] 강

☞ ‘초딩’은 <고려대>에 ‘얕잡아 이르는 말.’로 제시되어 있어 비어에 해당하므로 이를 포함하는 문장은 ‘명시’로 태깅한다.

5) 주로 ‘낮잡아’, ‘낮추어’, ‘홀하게’, ‘비하하는’ 등의 메타 술어와 함께 기술됨

6) 주로 ‘비속하게’ 등의 메타 술어와 함께 기술됨

- 대사전류에서 다의어 중 욕설, 비어, 비하성 속어에 해당하는 의미로 사용된 어휘나 표현이 나타나는 문장은 '명시'로 태깅

¶ 환경이 **그지(√거지)** 같아.

- ▶ [명시성] 명시 [맥락] 부정적 [영역] 문화, 관계/조건, 기타 [강도] 강
- ☞ '거지'는 <고려대>에서 다의어 의미 중 '행색이 지저분하고 초라하여 볼품없거나 남에게 빌붙어 사는 사람 등을 욕하여 이르는 말'에 해당하는 의미로 사용되어 욕설에 해당하므로, 이를 포함하는 문장은 '명시'로 태깅한다.

¶ -_-; 꼬질꼬질한 골목길 사이로 할머니 얼굴이 그려진 간판이 보이고, 가게에 들어서자마자 피죽도 못 먹은 **노인네** 둘이 비틀거리며 나오는 중.

- ▶ [명시성] 명시 [맥락] 부정적 [영역] 연령/세대, 신체, 문화, 기타 [강도] 강
- ☞ '노인네'는 <고려대>에 '나이든 사람을 얕잡아 이르는 말.'로 풀이되어 있어 비어에 해당하므로 이를 포함하는 문장은 '명시'로 태깅한다.

¶ **망할**_스마트폰을 집어던지든지 해야지 이메일을 너무 자주 확인해서 이런 걸 보게 되나 봐.

- ▶ [명시성] 명시 [맥락] 부정적 [영역] 문화, 기타 [강도] 강
- ☞ '망하다'는 <표준>에서 다의어 의미 중 "(주로 '망할' 꼴로 쓰여) 못마땅한 사람이나 대상에 대하여 저주의 뜻으로 이르는 말"에 해당하는 의미로 사용되어 욕설에 해당하므로 이를 포함하는 문장을 '명시'로 태깅한다.

- 대사전류에 욕설, 비어, 비하성 속어 등으로 표현된 어휘나 표현 중 고빈도로 사용되는 경우(주로 '미친', '시발', '존나', '망할' 및 그것들의 변이형)에는 초성 등으로 (단독) 실현되더라도 그것이 나타나는 문장은 '명시'로 태깅

¶ **ㅁㅈ**

¶ **ㅈㄴ웃겨** 알겠어요 **ㅈㄴ별거** 아닌데

¶ **ㅅㅂ너무싫음** 웅성우말고 웃이;

¶ **ㅁㅈ거야나 씨발좀나무서워**

- ☞ 위의 문장들은 모두 '존나', '시발', '미친' 등의 부적절성 표현을 유추할 수 있으므로 명시적 부적절성으로 판정한다.

<주의>

- 사전에 욕설, 비어, 비하성 속어로 기술되지 않거나 '낯잡아', '비속하게' 등으로 표현되는 비하성이 없고 '속되게' 등으로만 표현되는 단순 속어가 나타나는 문장은 '명시적 부적절성'으로 판정하지 않음

¶ 크라상은 참 고급스런 빵이자 부드러운 감촉과 버터향이 아주 **죽인다**.

- ☞ '죽이다'는 다의어 중 단순 속어의 의미로 사용되었으므로 '명시적 부적절성'으로 판정하지 않는다.

¶ 오늘 참 **빡세다**.

- ☞ '빡세다'는 단순 속어이고 해당 표현이 포함된 문장에서 비하성이나 조롱성이 표상되지 않으므로 '명시적 부적절성'으로 판정하지 않는다.

- 사전에 단순 속어로 기술된 표현이 나타나는 문장이라도 맥락을 통해서 비하나 조롱 등이 나타나는 경우에는 '비명시적 부적절성 문장'으로 판정

¶ 트위터일안하냐진짜빡치네 실트좀올려주지

☞ ‘빡치다’는 단순 속어이지만 해당 표현이 포함된 문장이 회사 혹은 회사의 행태에 대한 비난에 해당하므로 이를 ‘비명시적 부적절성’으로 판정한다.

¶ 검색어몇개지웠다고 나라가망할듯이 지랄들이야 ~~

☞ ‘지랄(하다)’는 단순 속어이지만 해당 표현이 포함된 문장이 행위에 대한 객관적 근거를 결여한 비난에 해당하므로 이를 ‘비명시적 부적절성’으로 판정한다.

¶ 불매 강 다 들은거갸다

☞ ‘들다’는 단순 속어이지만 해당 표현이 포함된 문장이 불매라는 행태에 대한 객관적 근거를 결여한 비난에 해당하므로 이를 ‘비명시적 부적절성’으로 판정한다.

- 사전에 ‘낮잡아’와 유사한 ‘낮추어’로 표현되더라도 겸양의 의미를 지닌 경우에는 비하성이 없으므로 ‘명시’는 물론이고 ‘부적절성’으로도 판정하지 않음

¶ 저는 빠질래요.

☞ ‘저’는 대사전류에 자기를 낮추어 가리키는 일인칭 대명사로 설명되어 있으므로 ‘부적절성’으로 판정하지 않는다.

4.1.2.2. 연구 보고서에 차별, 혐오, 선정 등의 문제로 지적된 표현이 나타나는 문장

- 공공기관의 관련 연구 및 조사 보고서(국립국어원의 <사회적 의사소통 연구> 시리즈 등)에서 차별 표현, 혐오 표현, 선정적 표현 등으로 지적된 어휘나 표현이 나타나는 문장은 ‘명시’로 태깅

¶ 여배우들은 피부가 재산이고....

▶ [명시성] 명시 [맥락] 부정적 [영역] 신체, 관계/조건 [강도] 강

☞ 팔호 안의 ‘여’와 같은 표현은 관련 보고서(안상수 외, 2007)에서 불필요한 성별 언급을 특정 성에 대한 차별적 표현으로 지적했으므로, 이를 포함하는 문장은 ‘명시’로 태깅한다.

¶ XXX XX에서 털털한 여배우의 매력을 뽐내고, 얼마전 종영한 <X XXX>에서

▶ [명시성] 명시 [맥락] 긍정적 [영역] 성, 문화, 관계/조건 [강도] 약

¶ 완벽한 S라인을 자랑하는 아름다운 몸매가 그 첫 번째.

▶ [명시성] 명시 [맥락] 부정적 [영역] 신체, 성 [강도] 강

☞ ‘S라인’은 관련 보고서(조태린 외, 2006)에서 불필요하게 외모를 강조하거나 묘사하는 선정적이고 자극적인 표현으로 지적했으므로, 이를 포함하는 문장은 ‘명시’로 태깅하며 부정적 맥락에서 성적 대상화를 하는 것에 해당하므로 ‘강’으로 태깅한다.

¶ 서울시 초,중,고 학생과 학부형, 선생님 여러분께 마음의 상처를 더 이상 주지 않고 반성하는 마음에서 결자해지하는 모습을 보여주시기를 기대해 봅니다.

▶ [명시성] 명시 [맥락] 긍정적 [영역] 성, 관계/조건 [강도] 강

☞ ‘학부형’은 남성형으로 여성까지 포괄한 예로 서울시여성가족재단에서 펴낸 <서울시 성평등 언어 사전>에서 ‘학부모’라는 순화어를 제시한 바 있으므로 이를 ‘명시’로 태깅한다.

¶ 넌 조선족 사는데 같이 살고 싶어?

▶ [명시성] 명시 [맥락] 부정적 [영역] 관계/조건 [강도] 강

☞ 중국이 아닌 한국에서 사용되는 '조선족'은 관련 보고서(박재현 외, 2009)에서 차별적, 비하적으로 사용됨을 지적하고 '재중동포' 또는 '한국계 중국인' 등으로 대체하는 것을 제안하고, 국립국어원 행정용어순화어(2018. 3. 27.)에서도 동일한 순화어를 제안하므로 이를 포함하는 문장은 '명시'로 태깅한다.

¶ 예로부터 **백인** > **흑인** > 개 > **황인** 이더라고

▶ [명시성] 명시 [맥락] 부정적 [영역] 관계/조건 [강도] 강

☞ '백인', '흑인', '황인' 등은 사전적으로는 비하적 의미가 없으나 관련 보고서(박재현 외, 2009)에서는 백인 중심의 인종 차별적 표현이며 불필요한 인종 언급이라고 지적한 바 있으므로 이들을 포함하는 문장은 '명시'로 태깅한다.

¶ **좌좀 전라디언인 흥어들** 제일 큰 문제가 뭔지 아냐 ㅋㅋㅋ

▶ [명시성] 명시 [맥락] 부정적 [영역] 관계/조건, 문화 [강도] 강

☞ '전라디언'과 '흥어'는 관련 보고서(박미숙 외, 2017)에서 지역 혐오 표현으로 사용됨을 지적했으므로, 이들을 포함하는 문장은 '명시'로 태깅한다. 또한 '좌좀'은 '좌익 좀비'의 줄임말로써 관련 보고서에서는 다루지 않은 용레이지만 정치 사상에 대한 혐오적 신어로 간주할 수 있으므로 '명시'로 태깅한다.

¶ XXX의 **미망인** XXX **여사** 의기투합했고..

▶ [명시성] 명시 [맥락] 부정적 [영역] 성, 관계/조건 [강도] 강

☞ '미망인'은 관련 보고서(조태린 외, 2006)에서 결혼한 여성이 남편이 사망했음에도 '아직 따라 죽지 못한 사람'이라는 봉건적 의미를 담고 있는 차별적 표현으로 지적했으므로, 이를 포함하는 문장은 '명시'로 태깅한다.

- 대사전류에 등재되어 있지 않고, 공공기관의 관련 연구 및 조사 보고서에서 차별 표현, 혐오 표현, 선정적 표현 등으로 제시되지 않았다 하더라도 같은 방식과 내용으로 만들어진 유사 표현이 나타나는 문장은 '명시'로 태깅

¶ 여기 **틀딱들이** 대거 나온 듯.

▶ [명시성] 명시 [맥락] 부정적 [영역] 연령/세대, 문화 [강도] 강

☞ '틀딱'은 노인 세대에 대한 차별 및 혐오 표현이므로 이를 포함하는 문장은 '명시'로 태깅한다.

¶ 솔직히 **기균층이나 지균층**은 동기란 생각 안 해.

▶ [명시성] 명시 [맥락] 부정적 [영역] 문화, 관계/조건, 기타 [강도] 강

☞ '기균층'과 '지균층'은 특정 입시 제도(기회균형선발전형, 지역균형선발전형)를 통해 들어온 대학생에 대한 차별 및 혐오 표현이므로 이를 포함하는 문장은 '명시'로 태깅한다.

¶ 나 수영강사 바꿨는데 설명을 진짜 **존나** 못하고 말귀더 잘 멧알아듣거든!

전에 **수제비는** 잘 가르치는데 세세한거 정확하게 안해도 넘어가줬거든 할 수만 있으면..

▶ [명시성] 명시 [맥락] 부정적 [영역] 문화, 관계/조건 [강도] 강

☞ 위의 문장에서 '수제비'는 '수영 강사'와 '제비'의 합성어로 '남성 수영 강사'에 대한 차별 표현에 해당하고, 음식을 지칭하는 '수제비'로 오해할 가능성이 없다는 점에서 이를 포함하는 문장은 '명시'로 태깅한다.

¶ 여긴 **감자국**과 다르게 **존나뽕뽕해**

▶ [명시성] 명시 [맥락] 부정적 [영역] 관계/조건, 기타 [강도] 강

☞ 위의 문장에서 '감자국'은 강원도에 대한 지역 비파 표현에 해당하고, 음식을 지칭하는 '감자국'으로 오해할 가능성이 없다는 점에서 이를 포함하는 문장은 '명시'로 태깅한다.

¶ **미소빌런**

▶ [명시성] 명시 [맥락] 긍정적 [영역] 문화 [강도] 약

☞ ‘빌런’은 악당을 의미하는 영어 단어를 음차한 것으로 국어 사전에 등재되어 있지 않지만, 누군가를 악당으로 표현하기 위한 자극적이고 선정적인 표현이므로 이를 포함하는 문장은 ‘명시’로 태깅한다.

<주의>

- 관련 연구 및 조사 보고서에 차별 표현으로 지적되었고 차별적 측면이 있다 하더라도 특히 성별 제시 순서와 같이 다른 대안이 없어서 그대로 사용할 수밖에 없는 경우에는 ‘명시’는 물론이고 ‘부적절성’으로도 판정하지 않음

¶ 부모, 자녀, 남녀, 형제자매, 학부모

☞ 위의 단어들은 관련 보고서에 남성을 표준으로 여성에 앞서 호명되는 단어로 지적되어 호명 순서가 성차별적인 측면이 있지만 다른 대안이 없으므로 ‘부적절성’으로 판정하지 않는다.

¶ 모부자가정, 편모편부, 엄마아빠

☞ 위의 단어들은 관련 보고서에 여성이 남성보다 앞서는 단어로 지적되어 호명 순서가 성차별적인 측면이 있지만 다른 대안이 없으므로 ‘부적절성’으로 판정하지 않는다.

- 차별 표현, 혐오 표현, 선정적 표현이 나타나더라도 해당 표현을 비판적으로 언급하는 경우에는 ‘명시’는 물론이고 ‘부적절성’으로도 판정하지 않음

¶ 미망인 의미 진짜 개빡친다

☞ 위의 문장은 맥락상 ‘미망인’이라는 어휘에 대한 부정적 감정을 표출한 것이며 명시적으로 부적절한 단어를 사용한 것이 아니며 ‘개빡친다’ 또한 단순 속어로 볼 수 있으므로 맥락을 고려하여 ‘부적절성’으로 판정하지 않는다.

4.1.2.3. 기타 ‘명시’로 태깅하는 문장

- 비식별화되어 있거나 그림 문자, 초성 등으로 되어 있더라도 명시적 부적절성 표현임이 명백한 표현이 나타나는 문장은 ‘명시’로 태깅

¶ 젠장 새로 이사 온 것들은 코빼기도 안 보이고 만날 망할 개새끼만 2마리씩 짚어대니 죽겠네(선행 문장)

그나저나 오늘 날씨 왠케 좋아!!! 이런 X!! 욕 나오게 하네.

▶ [명시성] 명시 [맥락] 부정적 [영역] 기타 [강도] 강

☞ 원 말씀치에서는 경우에 따라 X로 문자열이 가려진 경우가 존재한다. 이 문장에서는 ‘욕 나오게 하다’라는 표현으로 인하여 X의 문자열이 육설임이 명백하므로 이를 포함하는 문장은 ‘명시’로 태깅한다.

¶ 시발

▶ [명시성] 명시 [맥락] 부정적 [영역] 기타 [강도] 강

¶ 산책못가면 생기는 일



▶ [명시성] 명시 [맥락] 부정적 [영역] 문화, 기타 [강도] 강

☞ 강아지의 행동에 대한 문장으로 ‘지랄 발광’이 그림 문자 및 영문자가 포함된 형태로 형상화되어 있고 ‘발광’은 비어의 의미를 지닌 표현이므로 ‘명시’로 태깅한다.

¶ **초심잡아서뭐하노ssibal**

▶ [명시성] 명시 [맥락] 부정적 [영역] 문화, 기타 [강도] 강

☞ 욕설 '시발'이 영문(ssibal)로 표기되어 있으나 맥락 및 표기상 욕설임이 명백하므로 '명시'로 태깅한다.

<주의>

- 개인정보 등이 비식별화되거나 오타가 발생한 경우 등으로 인해 명시적인 부적절 표현인지 여부를 정확하게 판단하기 어려운 표현이 나타나는 문장은 '명시'는 물론이고 '부적절성'으로도 판정하지 않음

¶ 내가 물주가 아니니 바로 xx에게 “들었죠? 맛난이 사주셔야 됩니다.” 한 번 찢러놓고..

☞ 원 말씀치에서는 경우에 따라 X로 문자열이 가려진 경우가 존재한다. X의 문자열이 명시적 부적절성 표현으로 사용되었는지 판단이 불가능하므로 '부적절성'으로 판정하지 않는다.

¶ 한국작가 중 사인받고 싶은 넘버 원이 스마트교이고 투가 황모씨인데 황작가 사인은 예전에 노받은게 있다 그림은 부탁도 안했는데 그려줘서 완전 반했음

☞ '노'를 단순 오타로 판단할 수도 있고 명시적 부적절성 표현으로 사용되었는지 판단하기 어려우므로 '부적절성'으로 판정하지 않는다.

4.1.3. 비명시

- '비명시'는 '명시'의 조건들을 엄격하게 충족하지 않음에도 부적절성이 나타나는 문장들을 대상으로 태깅

4.1.3.1. 명시적 부적절성 표현이 나타나지 않지만, 맥락에서 부적절성이 확인되는 문장

- 명시적 부적절성 표현이 나타나지는 않지만, 문장 내 맥락에서 부적절성이 확인되는 문장은 '비명시'로 태깅

¶ **하지만 전문가 말만큼 믿을 게 못 되는 것도 없다.**

▶ [명시성] 비명시 [맥락] 부정적 [영역] 관계/조건, 문화 [강도] 약

☞ 문장 내 맥락에서 전문가에 대한 비하와 비난 등 관련 부적절성이 확인되므로 '비명시'로 태깅한다.

¶ **외국 아가씨들이랑 놀기 위해..**

▶ [명시성] 비명시 [맥락] 부정적 [영역] 성, 연령/세대, 문화, 관계/조건 [강도] 강

☞ 문장 내 맥락에서 외국의 젊은 여성에 대한 성적 대상화, 편견과 비하 등 관련 부적절성이 확인되므로 '비명시'로 태깅한다.

¶ **원작판에 비해 키도 커지고 몸매도 좋아졌다.**

▶ [명시성] 비명시 [맥락] 부정적 [영역] 신체 [강도] 강

☞ 문장 내 맥락에서 불필요한 외모 언급 관련 부적절성이 확인되므로 '비명시'로 태깅한다.

- 명시적 부적절성 표현이 나타나지는 않지만, 선후 맥락에서 부적절성이 확인되는 문장은 '비명시'로 태깅

: 맥락을 확인할 수 있도록 선행 및 후행 문장을 각각 최대 5개까지 추가 제시 필요

¶ 기껏 국물이라고 준 오뎅국물?은
 미지근해서 따뜻하게 덤혀달라고 해도
 고대로였어요 ㅋㅋㅋㅋ
 상무지구 오징어나라는
 그냥 다음부턴 안갈 것 같아요
왜 사람들은 초심을 잃는거죠
 ▶ [명시성] 비명시 [맥락] 부정적 [영역] 문화, 기타 [강도] 약
 ☞ “왜 사람들은 초심을 잃는 거죠”가 비난에 해당하는 비명시적 부적절성 문장임을 확인할 수 있도록 위와 같이 최대 5개의 선행 문장을 함께 제시한다.

4.1.3.2. 명시적이지 않지만, 주로 부적절성 관련 의미로 사용되는 표현이 나타나는 문장

- 사전의 용례나 일상 발화에서 주로 비하성이나 선정성 등의 부적절성 관련 의미로 사용되더라도 사전이나 보고서에 근거하는 명시적 부적절성은 아닌 표현이 나타나는 문장은 ‘비명시’로 태깅

¶ (영화배우 키아누 리브스Keanu Reeves에 대한 이야기) **오늘의 팬서비스는 평소 키아누 하는 짓에 비해 너무 파격적이라 진짜 놀랐다.**
 ▶ [명시성] 비명시 [맥락] 긍정적 [영역] 문화 [강도] 약
 ☞ ‘짓’은 <표준> 등의 대사전류에 욕설, 비어, 비하성 속어 등으로 명시되어 있지 않아도 맥락 등을 고려하였을 때 비하성 등의 부적절성 등을 추측할 수 있으므로, 이를 포함하는 문장은 ‘비명시’로 태깅한다.

¶ **니시끼들은** 미워할수가없다 **이쁜짓만** 골라서 하네 ㅎ
 ▶ [명시성] 명시 [맥락] 긍정적 [영역] 문화 [강도] 약
 ☞ 위의 맥락에서 ‘짓’은 비하성 등의 부적절성을 추측할 수 없으므로 ‘짓’으로 인한 부적절성을 태깅하지 않는다.

¶ **오늘 인턴들이 불러내서 밥 먹었는데 인턴 주제에 왜 이렇게 밥은 비싼 거 먹어서 내 일주일 용돈을 날리는 거냐.**
 ▶ [명시성] 비명시 [맥락] 부정적 [영역] 관계/조건, 문화 [강도] 약
 ☞ ‘주제’는 <표준> 등의 대사전류에 욕설, 비어, 비하성 속어 등으로 명시되어 있지 않아도 “(흔히 ‘주제에’ 꼴로 쓰여) 변변하지 못한 처지”라는 뜻풀이와 사용례 등을 고려하면 주로 비하성 등의 부적절성을 나타내지만, 이를 포함하는 문장은 ‘비명시’로 태깅한다.

¶ **오늘 결국 싸구려 면바지 한 장 사고 10시간은 걸은 거 같다.**
 ▶ [명시성] 비명시 [맥락] 부정적 [영역] 기타, 관계/조건 [강도] 약
 ☞ ‘싸구려’는 <표준> 등의 대사전류에 욕설, 비어, 비하성 속어 등으로 명시되어 있지 않아도 사용례 등을 고려하면 주로 비하성 등의 부적절성을 나타내지만, 이를 포함하는 문장은 ‘비명시’로 태깅한다.

¶ **츄리닝에 잠바때기 입고 광화문 가는 길입니다.**
 ▶ [명시성] 비명시 [맥락] 부정적 [영역] 기타, 관계/조건 [강도] 약
 ☞ ‘잠바때기(잠바때기)’는 <고려대>에서 “‘잠바’를 속되게 이르는 말.”이라는 단순 속어로 제시되어 있어도 사용례 등을 고려하면 주로 비하성 등의 부적절성을 나타내지만, 이를 포함하는 문장은 ‘비명시’로 태깅한다.

¶ 너무나 **굵벵이같은** 속도 어쩔

- ▶ [명시성] 비명시 [맥락] 부정적 [영역] 문화, 관계/조건 [강도] 약
- ☞ '굵벵이'는 <표준> 등의 대사전류에 욕설, 비어, 비하성 속어 등으로 명시되어 있지 않아도 "동작이 굵뜨고 느린 사물이나 사람을 비유적으로 이르는 말"이라는 뜻풀이와 사용례 등을 고려하면 주로 비하성 등의 부적절성을 나타내지만, 이를 포함하는 문장은 '비명시'로 태깅한다.

¶ **제발** 애 **짜지**려면 기본상식 시험 치고 통과해서 **짜질렀으면** 좋겠다.

- ▶ [명시성] 비명시 [맥락] 부정적 [영역] 문화, 신체 [강도] 약
- ☞ '짜지르다'는 <표준> 등의 대사전류에 단순 속어로 제시되어 있어도 뜻풀이나 실제 사용례 등을 고려하면 비하성 등의 부적절성을 나타내지만, 이를 포함하는 문장은 '비명시'로 태깅한다.

¶ **여자는** 기력지가 긴 애들이 **먹기도** 좋지

- ▶ [명시성] 비명시 [맥락] 부정적 [영역] 문화, 성, 신체 [강도] 강
- ☞ '먹다'는 <표준> 등의 대사전류에 단순 속어로 제시되어 있어도 사전의 다의어 뜻풀이((남자가 여자를) 성적으로 침해하여 짓밟다) 및 사용례 등을 고려하면, 선정성(성폭력, 성희롱) 등의 부적절성을 나타내지만, 이를 포함하는 문장은 '비명시'로 태깅한다.

¶ 나랑 **떡칠래?**

- ▶ [명시성] 비명시 [맥락] 부정적 [영역] 문화, 성 [강도] 강
- ☞ '떡치다'는 <고려대>에서 "성적으로 관계를 맺는 일을 하다."라는 의미의 단순 속어로 제시되어 있어도 실제 사용례를 고려하면 선정성, 성희롱 등의 부적절성을 나타내지만, 이를 포함하는 문장은 '비명시'로 태깅한다.

¶ 저 중에 누가 제일 **따먹고** 싶게 생김?

- ▶ [명시성] 비명시 [맥락] 부정적 [영역] 문화, 성, 신체 [강도] 강
- ☞ '따먹다'는 <표준> 등의 대사전류에 단순 속어로 제시되어 있어도 '여자의 정조를 빼앗다'라는 뜻풀이와 사용례 등을 고려하면 선정성(성폭력, 성희롱) 등의 부적절성을 나타내지만, 이를 포함하는 문장은 '비명시'로 태깅한다.

- 사전의 용례나 일상 발화에서 주로 대상에 대한 부정적 평가나 묘사에 사용되더라도 사전이나 보고서에 근거하는 명시적 부적절성은 아닌 서술표현(주로 형용사)이 나타나는 문장은 '비명시'로 태깅

¶ (공정적 맥락에서 조카를 귀엽게 묘사) **시아머니와 며느리 버릇없는 조카 세연양과.**

- ▶ [명시성] 비명시 [맥락] 긍정적 [영역] 관계/조건, 문화, 연령/세대 [강도] 약
- ☞ '버릇없다'는 <표준> 등의 대사전류에 욕설, 비어, 비하성 속어 등으로 명시되어 있지 않아도 대상을 부정적으로 묘사하는 의미를 내재하고 있지만, 이를 포함하는 문장은 '비명시'로 태깅한다.

¶ **내 장동건은** 지금처럼 **망가지지 않을** 때가 있었어!! **젤압** 이상한 cf**좀** **고만** **찍어라!** **멍청해**빈다 ㄱ ㄱ)

- ▶ [명시성] 비명시 [맥락] 부정적 [영역] 문화, 신체, 관계/조건, 기타 [강도] 약
- ☞ 해당 문장에서 '망가지다', '이상하다', '멍청하다' 등의 단어는 <표준> 등의 대사전류에 욕설, 비어, 비하성 속어 등으로 명시되어 있지 않아도 대상을 부정적으로 묘사하는 의미를 내재하고 있지만, 이를 포함하는 문장은 '비명시'로 태깅한다.

- 사전의 용례나 일상 발화에서 주로 비하성 등의 부적절성 관련 의미를 표현하는 데 사용되더라도 사전이나 보고서에 근거하는 명시적 부적절성은 아닌 형식 형태소('따위' 같은 의

¶ 이런 거 어디 없냐여? ㄴ ㄴ 최저가라고 해서 클릭하면 핑크 **딱뽀**만 팔거나(이 세상 컬러 중 가장 싫어하는 색), 이름만 걸고 다른 모델을 팔거나 하는 낚시가 많네.

▶ [명시성] 비명시 [맥락] 부정적 [영역] 문화, 기타 [강도] 약

¶ 장사 그 **딱뽀**로 하지 마라..

▶ [명시성] 비명시 [맥락] 부정적 [영역] 관계/조건, 문화 [강도] 약

¶ 생각하는 수준이 그러니까 아파트 경비나 하지

▶ [명시성] 비명시 [맥락] 부정적 [영역] 관계/조건 [강도] 약

¶ 말하는 것 보니 평생 **고시원이나** 전전하게 생겼지

▶ [명시성] 비명시 [맥락] 부정적 [영역] 관계/조건, 문화, 기타 [강도] 약

존 명사, '(이)나' 같은 조사 등이 나타나는 문장은 '비명시'로 태깅

4.1.3.3. 기타 '비명시'로 태깅하는 문장

- 발화자 자신에 대한 내용이라도 비명시적 부적절성이 나타나는 문장은 '비명시'로 태깅

¶ 이 나이 먹도록 취업도 못하고 **엄마한테 빌붙어 살고..**

▶ [명시성] 비명시 [맥락] 부정적 [영역] 관계/조건, 연령/세대, 문화 [강도] 약

☞ '빌붙다'는 <표준> 등의 대사전류에 욕설, 비어, 비하성 속어 등으로 명시되어 있지 않아도 사용례 등을 고려하면 비하성 등의 부적절성을 나타내지만, 이를 포함하는 문장은 '비명시'로 태깅한다. 또한 이 문장은 발화자 자신에 대한 내용이라도 부적절성이 비명시적으로 나타나는 것으로 판정한다.

¶ 그럼 무식한 티 안 내려면 나도 붙여야겠지.

▶ [명시성] 비명시 [맥락] 부정적 [영역] 관계/조건, 문화 [강도] 약

☞ 이 문장은 발화자 자신에 대한 내용이라도 부적절성이 비명시적으로 나타나는 것으로 판정한다.

¶ 없는 돈에 **편드하는 짓은 이제 하지 말아야지.**

▶ [명시성] 비명시 [맥락] 부정적 [영역] 관계/조건, 문화 [강도] 약

☞ 이 문장은 발화자 자신에 대한 내용이라도 부적절성이 비명시적으로 나타나는 것으로 판정한다.

¶ 한약 먹는 보람도 없이 이게 무슨 짓이람..

▶ [명시성] 비명시 [맥락] 부정적 [영역] 문화 [강도] 약

☞ 이 문장은 발화자 자신에 대한 내용이라도 부적절성이 비명시적으로 나타나는 것으로 판정한다.

- 정치, 사회, 종교 등에 대한 지지나 반대 의견에서도 그 맥락에서 객관성을 상실한 비하성(조롱, 비난 등)이나 편향성(차별 등), 공격성(혐오) 등의 부적절성이 확인되는 문장은 '비명시'로 태깅

¶ '**국민안전처**'는 '**국민이 안전한 척**' 하고'

▶ [명시성] 비명시 [맥락] 부정적 [영역] 문화, 관계/조건, 기타 [강도] 약

¶ 국방부'는 '**국방 방심부**' 같구나...

▶ [명시성] 비명시 [맥락] 부정적 [영역] 문화, 관계/조건, 기타 [강도] 약

¶ **착실하게 독재웨이** 걷는 중

▶ [명시성] 비명시 [맥락] 부정적 [영역] 문화, 관계/조건 [강도] 약

¶ **난 박근혜가 뭘해도 싫다.**

▶ [명시성] 비명시 [맥락] 부정적 [영역] 문화, 관계/조건 [강도] 약

¶ (사드 도입에 새누리당 의원 일부가 찬성한 것에 대해) **쌍수 들고 환영했으니 당연히
너님들이 유치하셔야죠.**

▶ [명시성] 비명시 [맥락] 부정적 [영역] 문화, 관계/조건 [강도] 약

¶ **오늘 새벽같이 일어나서 XX청 가서 교육 받고 그 후에 YY서로 이동, 개 끌려다니듯
계속 높으신 분들의 방에서 말씀을 들었다.**

▶ [명시성] 비명시 [맥락] 부정적 [영역] 문화, 관계/조건, 기타 [강도] 약

¶ **가 아니고 고급 노동자 화이트컬러.**

▶ [명시성] 비명시 [맥락] 부정적 [영역] 문화, 관계/조건 [강도] 약

¶ **동성애는 나라를 망칩니다**

▶ [명시성] 비명시 [맥락] 부정적 [영역] 성, 문화 [강도] 강

¶ **조두순은 죽여도 무죄. 제발**

▶ [명시성] 비명시 [맥락] 부정적 [영역] 문화 [강도] 강

- 특정 생물을 가리키는 이름이더라도 부적절성이 드러나는 경우 ‘비명시’로 태깅

¶ **머느리밀씻개**

▶ [명시성] 비명시 [맥락] 부정적 [영역] 성, 관계/조건, 기타 [강도] 강

<주의>

- 객관성을 상실한 부적절성이 나타나지 않고 정치, 사회, 종교 등에 대한 단순 지지나 반대 의견을 나타내는 문장은 ‘비명시’는 물론이고 ‘부적절성’으로도 판정하지 않음

¶ **요새 정치인 가족이라면 더 중형을 때린다**

☞ 부적절성이 나타나지 않는 개인의 정치적 의견이므로 ‘부적절성’으로 판정하지 않는다.

¶ **새누리당의 압승을 막아야 할 텐데 걱정입니다.**

☞ 부적절성이 나타나지 않는 개인의 정치적 의견이므로 ‘부적절성’으로 판정하지 않는다.

¶ **오보를 하고도 인정하지 않는 후안무치의 자세를 누가 바로잡아야 할 텐데..**

지난 주말 6차 촛불집회는 1,500여 단체가 주최했다.

시위 참가자 숫자는 전국적으로 232만명을 넘어섰다는 언론 보도는 경찰 추산 42만9천 명 보다 5배 이상 큰 차이가 나지만 모든 언론사들이 경쟁적으로 부풀리기를 되풀이하면서 국민 여론이라고 포장한다.

☞ 부적절성이 나타나지 않는 개인의 정치적 의견이므로 ‘부적절성’으로 판정하지 않는다.

¶ **검은색 쫄티에 짝끼는 흰색 스키니입은 사람을 봤어요 ㅋㅋ**

딱 달라붙는 흰색 스키니.....

보자마자 눈이 찌푸려지더라구요 ~~

☞ 타인의 옷차림에 대한 과도한 묘사로 보아 '부적절성'으로 판정하되 '눈이 찌푸려진다'는 개인의 취향으로 보아 '부적절성'으로 판정하지 않는다.

4.1.4. 명시와 비명시가 함께 나타나는 경우

○ 하나의 문장에 '명시'와 '비명시'가 함께 나타나더라도 복수 태깅은 불가

- 명시와 비명시가 함께 나타나는 문장은 '명시'로만 태깅
(단, 부적절성의 영역은 명시는 물론이고 비명시와 관련된 영역까지 모두 복수 태깅)

¶ 왜 키우냐 묻거든 강 웃지요 근데 **꼴에** 머리는 겁나 좋다 그냥 똥똥한 사람같다 ㅋㅋ 오래 살아라.

▶ [명시성] 명시 [맥락] 긍정적 [영역] 신체, 관계/조건, 문화 [강도] 약

☞ '꼴'은 <표준> 등의 대사전류에 '사람의 모양새나 행태를 낮잡아 이르는 말'로 뜻풀이되어 있어 비어에 해당하므로 이를 포함하는 문장은 '명시'로 태깅할 수 있고, 문장 전체에서는 고양이를 키우는 행태와 특정한 신체적 조건을 가진 사람을 비하하는 부적절성이 확인된다는 점에서 '비명시'로도 태깅할 수 있는데, 명시와 비명시가 함께 나타나는 문장은 '명시'로만 태깅한다.

¶ 이번의 똥 준 **넌(남자라면 너무 무섭다..:)**을 보면서 또 한 번 느꼈는데 빠순이는 악의는 없지만 정말 정말 이해가 안 간다는 거.

▶ [명시성] 명시 [맥락] 부정적 [영역] [영역] 성, 관계/조건, 문화, 기타 [강도] 강

☞ '넌'은 <표준> 등의 대사전류에 '여자를 낮잡아 이르는 말'로 뜻풀이되어 있으며 맥락상 '여성'이라는 것이 명백히 드러나는 비어에 해당하고 '빠순이'는 '연예인이나 운동선수 등을 맹목적으로 추종하고 따라다니는 극성팬 중 여자를 속되게 이르는 말'로 제시되어 있으므로 부적절성 표현으로 판정하지 않는다. 다만 이들에 대한 비하적 시선이 내재된 문장이므로 '비명시 무적절성'으로 판정한다. '비명시'와 '명시'가 중복되어 표상될 경우 '명시'를 우선하여 판정하므로 이들을 포함하는 문장은 '명시'로 태깅한다.

4.2. 맥락 주석

4.2.1. 부정적 맥락

4.2.1.1. 화자의 태도(의도)나 맥락 내용 측면 모두에서 부정적으로 판단되는 문장

- 판정된 부적절성이 화자의 태도(의도)나 맥락 내용 측면 모두에서 부정적으로 판단되는 문장은 ‘부정적’으로 태깅

- ¶ 대학생이나 되는데, 생각은 참 **초딩** 같아.
▶ [명시성] 명시 [맥락] 부정적 [영역] 연령/세대, 문화, 관계/조건 [강도] 강
- ¶ 오늘도 **구라치다** 하루가 다 갔어~
▶ [명시성] 명시 [맥락] 부정적 [영역] 문화 [강도] 강
- ¶ 제발 애 **싸지려면** 기본상식 시험 치고 통과해서 **싸질렀으면** 좋겠다.
▶ [명시성] 비명시 [맥락] 부정적 [영역] 문화, 신체 [강도] 약
- ¶ **정말이지 중국 사람들..**
엄청 시끄럽다
▶ [명시성] 비명시 [맥락] 부정적 [영역] 관계/조건, 문화 [강도] 약
- ¶ **착실하게 독재웨이 걷는 중**
▶ [명시성] 비명시 [맥락] 부정적 [영역] 문화, 관계/조건 [강도] 약

4.2.1.2. 화자의 태도(의도)가 부정적이지 않더라도 맥락 내용 측면에서 부적절한 문장

- 판정된 부적절성이 화자의 태도(의도) 측면에서는 부정적이지 않더라도 그 맥락 내용 측면에서 성폭력, 성추행, 성희롱 등의 성 관련 폭력적인 내용이나 성관계, 성적 대상화(상품화) 등의 선정적인 내용을 포함하는 문장은 ‘부정적’으로 태깅

- ¶ 완벽한 **S라인**을 사랑하는 아름다운 몸매가 그 첫 번째.
▶ [명시성] 명시 [맥락] 부정적 [영역] 신체, 성 [강도] 강
☞ ‘S라인’은 화자가 대상에 대한 칭찬을 의도로 사용한 것일지라도 관련 보고서(조태린 외, 2006)에서 불필요하게 외모를 강조하거나 묘사하는 선정적이고 자극적인 표현으로 지적했으므로, 이를 포함하는 문장은 ‘부정적’으로 태깅한다.
- ¶ **나랑 떡칠래?**
▶ [명시성] 비명시 [맥락] 부정적 [영역] 문화, 성 [강도] 강
☞ ‘떡치다’는 <고려대>에서 “성적으로 관계를 맺는 일을 하다.”라는 의미의 단순 속어로 제시되어 있고, 화자도 단지 그런 의미로만 사용한 것일지라도 실제 사용례 등을 고려하면 성관계와 관련한 선정적이고 자극적인 표현이므로, 이를 포함하는 문장은 ‘부정적’으로 태깅한다.
- ¶ **여자는 기력지가 긴 애들이 먹기도 좋지**
▶ [명시성] 비명시 [맥락] 부정적 [영역] 문화, 성, 신체 [강도] 강
☞ ‘먹다’는 <표준> 등의 대사전류에 단순 속어로 제시되어 있고 화자도 단지 그런 의미로만 사용한

것일지라도 사전의 다의어 뜻풀이((남자가 여자를) 성적으로 침해하여 짓밟다) 및 사용례 등을 고려하면, 선정성(성폭력, 성희롱) 등의 부적절성을 나타내는 표현이므로, 이를 포함하는 문장은 '부정적'으로 태깅한다.

¶ **저 중에 누가 제일 따먹고 싶게 생김?**

- ▶ [명시성] 비명시 [맥락] 부정적 [영역] 문화, 성, 신체 [강도] 강
- ☞ '따먹다'는 <표준> 등의 대사전류에 단순 속어로 제시되어 있고, 화자도 단지 그런 의미로만 사용한 것일지라도 '여자의 정조를 빼앗다'라는 뜻풀이와 사용례 등을 고려하면 선정성(성폭력, 성희롱) 등의 부적절성을 나타내는 표현이므로, 이를 포함하는 문장은 '부정적'으로 태깅한다.

4.2.1.3. 화자의 태도(의도)가 무표적이거나 긍정성/부정성 판단이 불가능한 문장

- 판정된 부적절성이 화자의 태도(의도) 측면에서 무표적으로 판단되거나 긍정성/부정성 판단이 불가능한 문장은 '부정적'으로 태깅

¶ **헐 시발 황희찬**

- ▶ [명시성] 명시 [맥락] 부정적 [영역] 기타 [강도] 강
- ☞ 위의 맥락에서 '시발'은 축구선수 황희찬에 대한 감탄의 의도인지, 실망의 의도인지 정확히 알 수 없으므로 '부정적'으로 태깅한다. 또한 '황희찬'은 추후 비식별화 작업을 거친 후에는 공개되지 않으므로 해당 인명을 검색하여 영역을 추측하여 태깅하지 않는다.

¶ **시발:::**

¶ **시발ㅋㅋㅋ**

- ▶ [명시성] 명시 [맥락] 부정적 [영역] 기타 [강도] 강
- ☞ '::, ㅋㅋㅋ, ㅎㅎ' 등의 이모티콘이 붙은 경우에도 화자의 의도를 정확히 알 수 없으므로 이를 '부정적'으로 태깅한다.

¶ **와 좀 말안되게생기셨**

- ▶ [명시성] 비명시 [맥락] 부정적 [영역] 신체 [강도] 약
- ☞ 위의 문장은 사람의 외모에 대한 불필요한 언급으로 부적절 문장으로 판정할 수 있으나 맥락을 정확히 파악할 수 없는데, 이러한 경우에도 '부정적'으로 태깅한다.

4.2.2. 긍정적 맥락

○ 화자의 태도(의도)와 맥락 내용 측면 모두에서 긍정적으로 판단되는 문장

- 명시적 또는 비명시적 부적절성 문장으로 판정되었지만, 화자의 태도(의도)와 맥락 내용 측면 모두에서 긍정적이거나 무표적인(부정적이지 않은) 것으로 판단되는 문장은 '긍정적'으로 태깅

¶ **기장도 살짝 다듬어주시고, **거지존** 극복 할 수 있게**

- ▶ [명시성] 명시 [맥락] 긍정적 [영역] 문화, 관계/조건, 기타 [강도] 약
- ☞ '거지'는 <고려대>에서 다의어 의미 중 '행색이 지저분하고 초라하여 볼품없거나 남에게 빌붙어 사는 사람' 등을 욕하여 이르는 말'에 해당하는 의미로 사용되어 욕설에 해당하므로, 이를 포함하는 문장은 '명시'로 태깅하되, 머리카락의 지저분한 부분을 깔끔하게 다듬는다는 맥락에서 사용된 점을 고려하여 '긍정적'으로 태깅한다.

¶ **미소빌런**

- ▶ [명시성] 명시 [맥락] 긍정적 [영역] 문화 [강도] 약

☞ '빌런'은 누군가를 악당으로 표현하기 위한 자극적이고 선정적인 표현이므로 이를 포함하는 문장은 '명시'로 태깅하되, 미소가 매우 큰 매력을 발산하는 사람이라는 맥락에서 사용된 점을 고려하여 '긍정적'으로 태깅한다.

¶ **진짜 저 사진 찍어놓고 너무 명충해보여서 한동안 웃었다.ㅋㅋ**

▶ [명시성] 비명시 [맥락] 긍정적 [영역] 관계/조건, 신체 [강도] 약

☞ 해당 문장에서 '명충하다(✓명칭하다)'를 포함하는 문장은 '비명시'로 태깅하되, 고양이 사진이 재미있음을 표현하는 맥락에서 사용된 점을 고려하여 '긍정적'으로 태깅한다.

¶ (긍정적 맥락에서 조카를 귀엽게 묘사) **시어머니와 며느리 버릇없는 조카 세연양과.**

▶ [명시성] 비명시 [맥락] 긍정적 [영역] 관계/조건, 문화, 연령/세대 [강도] 약

☞ '버릇없다'를 포함하는 문장은 '비명시'로 태깅하되, 조카의 행동을 귀엽게 묘사하는 맥락에서 사용된 점을 고려하여 '긍정적'으로 태깅한다.

¶ (영화배우 키아누 리브스Keanu Reeves에 대한 이야기) **오늘의 팬서비스는 평소 키아누 하는 짓에 비해 너무 파격적이라 진짜 놀랐다.**

▶ [명시성] 비명시 [맥락] 긍정적 [영역] 문화 [강도] 약

☞ '짓'을 포함하는 문장은 '비명시'로 태깅하되, 배우의 팬서비스에 대한 놀라움을 표현하는 맥락에서 사용된 점을 고려하여 '긍정적'으로 태깅한다.

¶ (긍정적 맥락에서 조카를 귀엽게 묘사) **시어머니와 며느리 버릇없는 조카 세연양과.**

- 상황에 대한 정확한 맥락 판정이 어렵더라도 주로 명시적 부적절성 표현이 긍정적 슬어와 공기하여 긍정적 의미가 표상되는 문장은 '긍정적'으로 태깅

¶ **시발** 존잼/존맛/개웃겨/개짤어

▶ [명시성] 명시 [맥락] 긍정적 [영역] 문화, 기타 [강도] 약

¶ **시발** 존멋/존예

▶ [명시성] 명시 [맥락] 긍정적 [영역] 신체, 기타 [강도] 약

¶ **시발** 개짤어

▶ [명시성] 명시 [맥락] 긍정적 [영역] 문화, 기타 [강도] 약

☞ 위의 맥락에서 앞뒤 문장을 알 수 없는 경우 해당 문장의 맥락을 정확히 파악할 수는 없으나 통상 긍정적 의미를 내재한 '재밌다', '멋있다', '예쁘다', '짤다(속어)'의 의미에 근거하여 '긍정적'으로 태깅한다.

¶ **개웃기다시발** 나 왜 속음

▶ [명시성] 명시 [맥락] 부정적 [영역] 문화, 기타 [강도] 강

☞ 그러나 경우에 따라 위의 예문처럼 문장의 맥락이 부정적인 것이(위의 문장은 '웃기다'의 의미가 '어이없다'의 부정적 의미로 사용되었음을 문장을 통해 알 수 있음) 확실한 경우가 존재하는데, 이 경우 슬어의 의미가 긍정적이라 하더라도 문장 전체의 맥락을 파악하여 '부정적'으로 태깅한다.

4.3. 영역 주석

4.3.1. 복수 영역의 처리

4.3.1.1. 부적절성의 영역이 두 개 이상 확인되는 문장

- 부적절성의 영역이 두 개 이상 확인되는 문장은 해당 영역들을 개수 제한 없이 모두 복수 태깅 허용 및 권장

¶ **고등학교 2학년 때만 해도 친구가 돼지같이 또 떡볶이 먹으러 가냐?`**라고 놀렸는데..

▶ [명시성] 명시 [맥락] 부정적 [영역] 신체, 문화 [강도] 강

☞ '돼지'는 <표준> 등의 대사전류의 다의어 의미 중에 '몹시 뚱뚱한 사람을 놀림조로 이르는 말'이라는 의미로 사용되어 비어에 해당하므로 이를 포함하는 문장은 외모 관련 영역인 '신체'를 태깅한다. 그리고 이 문장에서는 잘 먹는 행태를 비하하는 비하성이 확인된다는 점에서 '문화'를 복수 태깅한다.

¶ **-_-; 꼬질꼬질한 골목길 사이로 할머니 얼굴이 그려진 간판이 보이고, 가게에 들어서자마자 피죽도 못 먹은 노인네** 둘이 비틀거리며 나오는 중.

▶ [명시성] 명시 [맥락] 부정적 [영역] 연령/세대, 신체, 문화, 기타 [강도] 강

☞ '노인네'는 <고려대>에 '나이든 사람을 얕잡아 이르는 말'로 뜻풀이되는 비어이므로 '연령/세대'를 태깅한다. 그리고 이 문장에서는 노인의 신체적 허약함과 행태를 비하하는 비하성이 확인된다는 점에서 '신체'와 '문화'를 추가로 복수 태깅한다.

¶ **외국 아가씨들이랑 놀기 위해..**

▶ [명시성] 비명시 [맥락] 부정적 [영역] 성, 연령/세대, 문화, 관계/조건 [강도] 강

☞ 이 문장은 맥락에서 젊은 여성을 성적 대상화하는 부적절성이 나타난다는 점에서 '성', '연령/세대'를 복수 태깅한다. 그리고 특히 외국의 젊은 여성에 대한 편향성과 비하성 관련 부적절성이 확인되므로 '관계/조건 (출신)'도 복수 태깅한다. 또한 특정 행태에 대한 비하성이 확인된다는 점에서 '문화'를 추가적으로 복수 태깅한다.

¶ **(내가 좀::) 1박 2일에서 찬물에 들어갔다 나온 찌이 있는데 난 박찬호 허벅지랑 찬물에 젖은 고탄력 흰...**

▶ [명시성] 비명시 [맥락] 부정적 [영역] 신체, 성 [강도] 강

☞ 이 문장은 남성의 특정 신체 부위를 성적 대상화하여 불필요하게 언급하는 편향성이 확인되므로 '신체'와 '성'을 복수 태깅한다.

- 의미적 고정성이 낮은 표현('망할', '짓', '따위', '(이)나' 등)을 포함하는 부적절성 문장은 해당 표현의 수식 또는 첨가 대상이나 맥락 내용에 근거하여 해당 영역들을 복수 태깅

¶ **망할** 스마트폰을 집어던지든지 해야지 이메일을 너무 자주 확인해서 이런 걸 보게 되나봐.

▶ [명시성] 명시 [맥락] 부정적 [영역] 문화, 기타 [강도] 강

☞ '망하다'는 <표준>에서 다의어 의미 중 "(주로 '망할' 꼴로 쓰여) 못마땅한 사람이나 대상에 대하여 저주의 뜻으로 이르는 말"에 해당하는 의미로 사용되어 욕설에 해당하므로 이를 포함하는 문장은 감정의 배설 관련 영역인 '기타'로 태깅하고, 이 문장에서는 스마트폰을 너무 자주 사용하는 행태에 대한 비하성이 확인되므로 '문화'를 복수 태깅한다.

¶ 이 **망할** 병은 식욕이 없어지고 소화효소가 안나오는 게 증상이라 심해지면 밥을

못먹는다.

▶ [명시성] 명시 [맥락] 부정적 [영역] 신체, 기타 [강도] 강

☞ '망할'은 욕설에 해당하므로 이를 포함하는 문장은 감정의 배설 관 영역인 '기타'로 태깅하고, 질병에 대한 저주 등의 공격성이 확인된다는 점에서 '신체'를 복수 태깅한다.

¶ 이런 거 어디 없냐? ㄴ ㄴ 최저가라고 해서 클릭하면 핑크 따위만 팔거나(이 세상 켄리 중 가장 싫어하는 색), 이름만 걸고 다른 모델을 팔거나 하는 남시가 많네.

▶ [명시성] 비명시 [맥락] 부정적 [영역] 관계/조건, 문화, 기타 [강도] 약

☞ 의미적 고정성이 낮은 표현인 '따위'를 포함하는 이 문장은 특정 직업 종사자에 대한 비하성이 확인된다는 점에서 '관계/조건'을 태깅하고, 특정 색깔을 선호하는 행태에 대한 비하성이 확인된다는 점에서 '문화'를 복수 태깅하며, 특정 색깔을 저주하는 공격성이 확인된다는 점에서 '기타'를 복수 태깅한다.

¶ 장사 그 따위로 하지 마라..

▶ [명시성] 비명시 [맥락] 부정적 [영역] 관계/조건, 문화 [강도] 약

☞ '따위'를 포함하는 이 문장은 특정 직업 종사자에 대한 비하성이 확인된다는 점에서 '관계/조건'을 태깅하고, 해당 직업 종사자의 행태에 대한 비하성이 확인된다는 점에서 '문화'를 복수 태깅한다. 복수 태깅하며, 특정 색깔을 저주하는 공격성이 확인된다는 점에서 '기타'를 복수 태깅한다..

¶ 생각하는 수준이 그러니까 아파트 경비나 하지

▶ [명시성] 비명시 [맥락] 부정적 [영역] 관계/조건 [강도] 약

☞ 의미적 고정성이 낮은 표현인 '(이)나'를 포함하는 이 문장은 특정 직업 종사자에 대한 비하성이 확인된다는 점에서 '관계/조건'을 태깅하고, 특정 사고방식에 대한 비하성이 확인된다는 점에서 '문화'를 복수 태깅한다.

¶ 말하는 것 보니 평생 고시원이나 전전하게 생겼지

▶ [명시성] 비명시 [맥락] 부정적 [영역] 관계/조건, 문화, 기타 [강도] 약

☞ '(이)나'를 포함하는 이 문장은 말하는 태도와 행태에 대한 비하성이 확인된다는 점에서 '문화'를 태깅하고, 고시원이라는 주거 형태와 그곳에 거주하는 사람에 대한 비하성이 확인되므로 '기타'와 '관계/조건'을 복수 태깅한다.

4.3.1.2. 명시와 비명시가 함께 나타나는 문장

- 하나의 문장에 명시와 비명시가 함께 나타나는 문장은 각각의 부적절성 관련 영역을 모두 복수 태깅

¶ 오늘 지하철에서 웬 변태 새끼를 만나서 좀 거시기했지만 별로 대단한 찌질이는 아니었기에 금방 잊었다.“

▶ [명시성] 명시 [맥락] 부정적 [영역] 문화, 기타 [강도] 강

☞ '새끼'는 <고려대>에 '어떤 사람을 욕하여 이르는 말. 주로 남자에게 쓴다'로 뜻풀이되어 욕설에 해당하므로 이를 포함하는 문장은 감정의 배설 관련 영역인 '기타'와 성별 관련 영역인 '성'을 복수로 태깅할 수 있으나 맥락상 남성이라는 것이 명확하지 않으므로 '성'을 태깅하지 않는다. 그리고 이 문장에서는 '변태'가 명시적 부적절성 표현은 아니지만 성적 지향과 특정 행태를 비정상적으로 간주하는 비하성이 확인된다는 점에서 '성'과 '문화'를 복수로 태깅한다. 또한 '찌질이'도 명시적 부적절성 표현은 아니지만, 특정 행태를 비하하는 비하성이 확인된다는 점에서 '문화'를 태깅할 수 있다.

¶ 줄리 인터뷰 내용만 잠깐 기사로 읽었는데 어찌면 하나같이 말도 푹부러지게 잘하는지! 대개 배우하면 빈머리가 생각나지만 이 여인은 다른 것 같네.

▶ [명시성] 명시 [맥락] 긍정적 [영역] 성, 관계/조건, 문화 [강도] 약

☞ '빈머리'는 사고 능력이 부족한 사람에 대한 차별적이고 비하적인 표현이므로 이를 포함하는 문장은 '관계/조건'을 태깅하고, 이 문장에서는 특정 직업 종사자의 사고 능력과 성별에 대한 비하성이 확인된다는 점에서 '관계/조건'과 '성'을 복수 태깅한다.

¶ 요즘 젊은 **넌놈들은** 뇌가 없는 것 같아.

▶ [명시성] 명시 [맥락] 긍정적 [영역] 연령/세대, 관계/조건, 문화 [강도] 강

☞ 이 문장에서는 특정 연령층의 사고 능력과 행태에 대한 비하성이 확인된다는 점에서 '연령/세대', '관계/조건', '문화'를 복수 태깅한다.

4.3.2. 성

- 성별, 성적 지향, 성희롱, 성폭력 성 편견 등과 관련한 부적절성이 나타나는 문장은 '성'으로 태깅

¶ 이번의 똥 준 **논(남자라면 너무 무섭다..)**을 보면서 또 한 번 느꼈는데 **빠순이**는 악의는 없지만 정말 정말 이해가 안 간다는 거.

▶ [명시성] 명시 [맥락] 부정적 [영역] 성, 관계/조건, 문화, 기타 [강도] 강

☞ '넌'은 <표준> 등의 대사전류에 '여자를 낮잡아 이르는 말'로 뜻풀이되는 비어이고, '논'이라는 명사가 여성을 가리키는 것이 분명하므로 '성'을 태깅한다. 그리고 이 문장에서는 특정 행태를 배설물이라는 '기타'에 비유하여 비하하는 부적절성이 확인된다는 점에서 '문화'를 복수 태깅한다.

¶ (가족의 상황에 대한 묘사) 돈벌러 가신 **남정네분들은** 늦게 퇴근하고.

▶ [명시성] 명시 [맥락] 긍정적 [영역] 성, 관계/조건 [강도] 약

☞ '남정네'는 <고려대>에서 '여자가 사내를 조금 낮추어 이르는 말'로 뜻풀이되어 있고 맥락상 남성임을 유추할 수 있으므로 이를 포함하는 문장은 '성'을 태깅한다. 그리고 가족 구성원의 일부에 대한 비하적 표현이라는 점에서 '관계/조건(사적 관계)'을 복수 태깅한다.

¶ (특정 정당의 지지자들을 비난하는 맥락) 며칠 전까지만 해도 방송에서 돈 받은 걸로 죽일 **놈** 살릴 **놈** 하던 걸 똑똑히 봤는데..

▶ [명시성] 명시 [맥락] 부정적 [영역] 문화 [강도] 강

☞ '놈'은 <표준> 등의 대사전류에 '남자를 낮잡아 이르는 말'로 뜻풀이되어 있으나 맥락상 '남성'임을 확인할 수 없으므로 '성'을 태깅하지 않는다. 그리고 이 문장에서는 특정 정당에 대한 사고방식을 비하하는 부적절성이 확인된다는 점에서 '문화'를 복수 태깅한다.

¶ 다음 주에는 **덩치 큰 아저씨 하나 안 나오니 조용할 거 같네.**

▶ [명시성] 비명시 [맥락] 부정적 [영역] 성, 신체, 연령/세대, 문화 [강도] 약

☞ 이 문장에서는 남성의 외모에 대한 편향성과 비하성이 확인된다는 점에서 '성'과 '신체'를 복수 태깅한다. 그리고 특정 연령대와 행태에 대한 비하성도 확인된다는 점에서 '연령/세대'와 '문화'를 복수 태깅한다.

¶ **말을 잘 못하는 남자 특유의 어법이다.**

▶ [명시성] 비명시 [맥락] 부정적 [영역] 성, 문화 [강도] 약

☞ 이 문장에서는 남성의 특정 행태에 대한 비하성이 확인된다는 점에서 '성'과 '문화'를 복수 태깅한다.

¶ **본품도 귀엽죠? 약간 로즈골드에 화이트조합이라 엄청 여성스럽고 우아해요 ㅋㅋ**

▶ [명시성] 비명시 [맥락] 부정적 [영역] 성, 문화, 기타 [강도] 약

☞ 이 문장에서는 여성의 성향이나 행태를 일반화하고 고정하는 편향성이 나타난다는 점에서 '성'과 '문화'를

복수 태깅한다. 그리고 특정 색깔의 조합에 대한 비하성이 확인된다는 점에서 '기타'를 복수 태깅한다.

¶ **외국 아가씨들이랑 놀기 위해..**

▶ [명시성] 비명시 [맥락] 부정적 [영역] 성, 연령/세대, 문화, 관계/조건 [강도] 강

☞ 이 문장에서는 젊은 여성을 성적 대상화하는 부적절성이 나타난다는 점에서 '성', '연령/세대'를 복수 태깅한다. 그리고 특히 외국의 젊은 여성에 대한 편향성과 비하성 관련 부적절성이 확인되므로 '관계/조건(출신)'도 복수 태깅한다. 또한 특정 행태에 대한 비하성이 확인된다는 점에서 '문화'를, 결혼하지 않은 여성에 대한 지칭/호칭이라는 점에서 '관계/조건(사적 관계)'을 추가적으로 복수 태깅한다.

4.3.3. 연령/세대

- 연령이나 세대와 관련한 부적절성이 나타나는 문장은 '연령/세대'로 태깅

¶ **여기 틀딱들이 대거 나온 듯.**

▶ [명시성] 명시 [맥락] 부정적 [영역] 연령/세대, 문화 [강도] 강

☞ '틀딱'은 특정 연령대의 사람들과 그들의 행태에 대한 차별 및 혐오 표현이므로 이를 포함하는 문장은 '연령/세대'와 '문화'를 복수 태깅한다.

¶ **짤방은 도통 철없는 노인네 사진으로 마무리..**

▶ [명시성] 명시 [맥락] 부정적 [영역] 연령/세대, 문화 [강도] 강

☞ '노인네'는 <고려대>에 "나이든 사람을 얕잡아 이르는 말"로 뜻풀이되는 비어이므로 '연령/세대'를 태깅하고, 이 문장에서는 해당 세대의 행태에 대한 비하성이 확인된다는 점에서 '문화'를 복수 태깅한다.

¶ **아마도 이 늙은 고양이가 퐁달고 와서 실수로 한 개 떨궜나보다.**

▶ [명시성] 비명시 [맥락] 부정적 [영역] 연령/세대, 문화 [강도] 약

☞ 이 문장에서는 연령이 높은 세대의 비위생적 행태에 대한 비하성이 확인된다는 점에서 '연령/세대'와 '문화'를 복수 태깅한다.

¶ **폴햄에서 겨우 고른건 사이즈가 없고 난닝구라는 집에서 정말 아저씨 목 늘어난 난닝구같은 티를 찾았다! 이거다, 베이지색에 유넥, 약간 넉넉한 오버핏, 면백, 만원.**

▶ [명시성] 비명시 [맥락] 부정적 [영역] 연령/세대, 성, 문화, 기타 [강도] 약

☞ 이 문장에서는 특정 연령대의 남성에 대한 비하성이 확인된다는 점에서 '연령/세대'와 '성'을 복수 태깅한다. 그리고 해당 남성의 일상적 행태와 특정 의류 제품에 대한 비하성이 확인된다는 점에서 '문화'와 '기타'를 추가로 복수 태깅한다.

¶ **할배랑 사는 애들은 생각하는 것도 딱 할배야**

▶ [명시성] 비명시 [맥락] 부정적 [영역] 연령/세대, 관계/조건, 성, 문화 [강도] 약

☞ '할배'는 '할아버지'의 방언형으로 명시적 부적절성 표현이 아니지만, 이 문장에서는 할아버지라는 가족 구성원에 대한 비하성이 확인된다는 점에서 '연령/세대', '관계/조건(사적 관계)', '성' 등을 복수 태깅한다. 그리고 그들의 사고방식과 사고 능력에 대한 비하성이 확인된다는 점에서 '문화'와 '관계/조건'을 추가로 복수 태깅한다.

¶ **아스트레스!!!!!!!!!!!!!! 애비때문에 자살한녀가있다? 접니다..,**

▶ [명시성] 명시 [맥락] 부정적 [영역] 관계/조건, 성, 문화, 기타 [강도] 약

¶ **아부지식단뽕☺☺**

▶ [명시성] 명시 [맥락] 긍정적 [영역] 성, 연령/세대, 문화, 관계/조건 [강도] 약

☞ ‘아버비’는 ‘아버지’와 동일한 의미의 어휘이지만 위의 맥락에서는 가족으로서의 부친이 표상되어 있고, 아래는 다이어트를 하며 본인의 식이가 본인의 나이에 먹을 법한 음식이 아니라는 의미 또한 표상되므로 아래의 경우에는 ‘연령/세대’를 포함하되 위의 경우는 맥락을 고려하여 ‘연령/세대’를 태깅하지 않는다.

4.3.4. 신체

- 신체적인 것과 정신적인 것을 막론하는 장애, 건강, 질병 등과 관련한 부적절성이나 외모, 임신, 출산 등과 관련한 부적절성이 나타나는 문장은 모두 ‘신체’로 태깅

¶ 관종인스타홍내정병리

▶ [명시성] 명시 [맥락] 부정적 [영역] 신체, 문화 [강도] 강

☞ ‘관종’은 관심을 받기 위해 여러 가지 행동을 하는 사람을 정신적으로 이상한 사람으로 간주하고 비하하는 데 사용되는 차별적 표현이고, ‘정병리’는 정신병을 앓는 사람을 비하하는 데 사용되는 차별적 표현이므로 이를 포함하는 문장은 ‘신체’를 태깅한다. 그리고 이 문장에서는 스타를 흉내 내는 행태를 비하하는 부적절성이 확인된다는 점에서 ‘문화’를 복수 태깅한다.

¶ 꼬마니콜라에 보면 **동보** 친구가 매일 크림을 문히고 크라상을 먹는데 내가 그 심정 알지.

▶ [명시성] 명시 [맥락] 부정적 [영역] 신체, 문화 [강도] 강

☞ ‘동보’는 <표준> 등의 대사전류에 ‘살이 찌서 뚱뚱한 사람을 놀림조로 이르는 말’로 뜻풀이되는 비어이므로 이를 포함하는 문장은 ‘신체’를 태깅한다. 그리고 이 문장에서는 음식을 먹는 행태에 대한 비하성이 확인된다는 점에서 ‘문화’를 복수 태깅한다.

¶ 완벽한 **S라인**을 자랑하는 아름다운 몸매가 그 첫 번째.

▶ [명시성] 명시 [맥락] 부정적 [영역] 신체, 성 [강도] 강

☞ ‘S라인’은 관련 보고서(조태린 외, 2006)에서 불필요하게 외모를 강조하거나 묘사하는 선정적이고 자극적인 표현으로 지적했으므로, 이를 포함하는 문장은 ‘신체’를 태깅한다. 그리고 이는 여성의 신체를 가리키는 데 주로 사용되므로 ‘성’을 복수 태깅한다.

¶ 키스 사진은 충격파가 약간 있었지만 이 정도 **꽃미모니** 니가 뭘 하든 난 다 용서할 수 있어.

▶ [명시성] 명시 [맥락] 부정적 [영역] 신체, 문화 [강도] 강

☞ ‘꽃미모(미모)’는 관련 보고서(조태린 외, 2006)에서 불필요하게 외모를 강조하거나 묘사하는 부적절한 표현으로 지적했으므로, 이를 포함하는 문장은 ‘신체’를 태깅한다. 그리고 외모를 기준으로 대상을 평가하는 행태 관련 편향성이 확인된다는 점에서 ‘문화’를 추가로 복수 태깅한다.

¶ 왜 키우냐 묻거든 강 웃지요 근데 **꿀애** 머리는 겁나 좋다 그냥 뚱뚱한 사람같다 ㅋㅋ 오래 살아라,

▶ [명시성] 명시 [맥락] 부정적 [영역] 신체, 관계/조건, 문화 [강도] 강

☞ ‘꿀’은 <표준> 등의 대사전류에 사람의 모양새나 행태를 낮잡아 이르는 말로 뜻풀이되는 비어이므로, 이를 포함하는 문장은 ‘신체’와 ‘관계/조건’을 복수 태깅한다. 그리고 이 문장에서는 고양이를 키우는 행태에 대한 자조와 사람의 체형에 대한 비하 관련 부적절성이 확인된다는 점에서 ‘문화’와 ‘신체’를 추가로 복수 태깅한다.

¶ 카메라 **좆같아요**

▶ [명시성] 명시 [맥락] 부정적 [영역] 성, 신체, 기타 [강도] 강도

¶ **좨나** 빠쳐

▶ [명시성] 명시 [맥락] 부정적 [영역] 성, 신체, 기타 [강도] 강도

☞ 이 문장에서는 남성의 성기를 비속하게 이르는 '좨'의 철자가 그대로 실현되었기 때문에 의미적으로 표상된다고 보아 '성'과 '신체'를 태깅한다.

¶ **결정장애자인 저에게~ 참 좋은 선택이었던거 같아요!**

▶ [명시성] 비명시 [맥락] 긍정적 [영역] 신체, 문화 [강도] 약

☞ '결정장애자'는 결정을 망설이는 행태 또는 신중한 사고방식을 장애로 비유하는 차별적 표현으로 비하성이 나타난다는 점에서 이를 포함하는 문장은 '신체'와 '문화'를 복수 태깅한다.

¶ **원작판에 비해 키도 커지고 몸매도 좋아졌다.**

▶ [명시성] 비명시 [맥락] 긍정적 [영역] 신체 [강도] 약

☞ 이 문장에서는 외모에 대한 불필요한 언급 관련 편향성이 확인된다는 점에서 '신체'를 태깅한다.

¶ **시크릿 가든도 끝나고 무슨 낙으로 사나 했는데 돌아오는구나!!! 머리술이 좀 부족한 닐의 친구가 죽느냐 사느냐에 기로에 선 채 끝났는데 어떻게 될까.**

▶ [명시성] 비명시 [맥락] 부정적 [영역] 신체 [강도] 약

☞ 이 문장에서는 외모에 대한 불필요한 언급 관련 편향성이 확인된다는 점에서 '신체'를 태깅한다.

<주의>

- '좨' 파생어인 '좨나'의 변이형이라 하더라도 '죤나', '죤라' 등은 그 철자에서 '성'이나 '신체'를 추측할 수 없으므로 맥락에 근거하여 영역을 태깅
- 유사 부사 '열라', '겁나' 등은 '좨'파생 부사로 보기 어렵고 동일한 철자의 방언형 또한 존재하므로 이를 단순 속어로 처리

¶ **죤나웃겜늘네**

▶ [명시성] 명시 [맥락] 긍정적 [영역] 문화 [강도] 약

¶ **겜 죤넌 웃기네 ㅋㅋㅋㅋ**

▶ [명시성] 명시 [맥락] 긍정적 [영역] 문화, 기타 [강도] 약

¶ **몸 죤라 으슬으슬**

▶ [명시성] 명시 [맥락] 부정적 [영역] 문화 [강도] 강

¶ **아침부터출금이백만원은겜나아프고공허하다**

4.3.5. 문화

- 종교, 정치, 풍습, 예술, 행태, 사고방식 등과 관련한 부적절성이 나타나는 문장은 '문화'로 태깅

¶ **오늘도 구라치다** 하루가 다 갔어~

▶ [명시성] 명시 [맥락] 부정적 [영역] 문화 [강도] 강

☞ '구라'는 <고려대>에서 '거짓말을 비속하게 이르는 말'로 뜻풀이되는 비하성 속어이고, '구라치다'는 거짓말을 하는 행태에 대한 비하 표현이 되므로, 이를 포함하는 문장은 '문화'를 태깅한다

¶ 정말 이 일에는 **거지같은** 녀들이 꼬일 수밖에 없는 건가.

▶ [명시성] 명시 [맥락] 부정적 [영역] 문화, 관계/조건 [강도] 강

☞ '거지'는 <고려대>에서 다의어 의미 중 '행색이 지저분하고 초라하여 볼품없거나 남에게 빌붙어 사는 사람 등을 욕하여 이르는 말'로 뜻풀이되는 욕설로 관련 행태와 능력에 대한 비하성이 나타나므로, 이를 포함하는 문장은 '문화'와 '관계/조건'을 복수 태깅한다.

¶ 기장도 살짝 다듬어주시고, **거지존** 극복 할 수 있게

▶ [명시성] 명시 [맥락] 긍정적 [영역] 문화, 관계/조건, 기타 [강도] 약

☞ '거지'는 <고려대>에서 다의어 의미 중 '행색이 지저분하고 초라하여 볼품없거나 남에게 빌붙어 사는 사람 등을 욕하여 이르는 말'로 뜻풀이되는 욕설로 관련 행태와 능력에 대한 비하성이 나타나므로, 이를 포함하는 문장은 '문화'와 '관계/조건'을 복수 태깅한다 그리고 '거지존'은 지저분한 머리카락이 있는 부분이라는 부정물에 대한 비하적 의미로 사용되고 있으므로, 이를 포함하는 문장은 '기타'를 추가로 복수 태깅한다.

¶ 그러나 나는 이미 **눈이 빠졌으니** 선생님이 웃을 때 보이는 눈가의 잔주름까지 마냥 귀여워보였다.

▶ [명시성] 명시 [맥락] 부정적 [영역] 문화 [강도] 강

☞ '눈(이) 빠다'는 <표준> 등의 대사전류에 '뻥한 것을 잘못 보고 있을 때 비난조로 이르는 말'로 뜻풀이되는 관용구로 특정한 사고방식에 대한 비하성이 나타나므로, 이를 포함하는 문장은 '문화'를 태깅한다.

¶ (영화배우 키아누 리브스Keanu Reeves에 대한 이야기) **오늘의 팬서비스는 평소 키아누 하는 짓에 비해 너무 파격적이라 진짜 놀랐다.**

▶ [명시성] 비명시 [맥락] 긍정적 [영역] 문화 [강도] 약

☞ 이 문장에서는 팬서비스 관련 특정 행태를 비하하는 비하성이 확인된다는 점에서 '문화'를 태깅한다.

¶ **착실하게 독재웨이 걷는 중**

▶ [명시성] 비명시 [맥락] 부정적 [영역] 문화, 관계/조건 [강도] 약

☞ 이 문장에서는 특정한 정치 상황을 독재라고 비난하는 공격성이 확인된다는 점에서 '문화'를 태깅한다. 그리고 정치인으로서의 지위나 능력을 비하하는 비하성이 확인된다는 점에서 '관계/조건'을 복수 태깅한다.

¶ **나 같은 건 자살이 답이다**

▶ [명시성] 비명시 [맥락] 부정적 [영역] 문화, 관계/조건 [강도] 강

☞ '자살'은 대사전류에 따른 명시적 부적절성 표현은 아니지만, 현대적 기준에서 보았을 때 윤리적 기준을 크게 위배하는 행태이므로, 이를 포함하는 문장은 '문화'를 태깅한다. 그리고 이 문장에서는 화자 자신의 지위나 능력을 비하하는 비하성이 확인된다는 점에서 '관계/조건'을 복수 태깅한다.

4.3.6. 관계/조건

- 인종, 국가, 지역 등과 같은 '출신'과 관련한 부적절성이 나타나는 문장은 '관계/조건'으로 태깅

¶ 머리 검은짐승은 거두는게 아니했는데 재들은 뭐.. 그냥 **깜둥이들은** 거두면 안되는 듯.

▶ [명시성] 명시 [맥락] 부정적 [영역] 관계/조건, 문화 [강도] 강

☞ '깜둥이'는 특정 인종에 대한 차별적 표현이므로 이를 포함하는 문장은 '관계/조건(출신)'을 복수 태깅하고, 이 문장에서는 타인을 믿으면 배신당한다는 사고방식이 확인된다는 점에서 '문화'를 복수 태깅한다.

¶ 지구의 암 같은 존재 라니까 **짱깨 새끼들은**

▶ [명시성] 명시 [맥락] 부정적 [영역] 관계/조건, 신체, 기타 [강도] 강

☞ '짱깨'는 관련 보고서(조태린 외, 2006)에 중국인을 비하하는 차별적 표현으로 지적되어 있으므로 이를 포함하는 문장은 '관계/조건(출신)'을 태깅하며 감정의 배설 관련 영역인 '기타'를 복수 태깅한다. 또한 이 문장에서는 특정 국적 사람들을 암이라는 질병에 비유하는 부적절성이 확인된다는 점에서 '신체'를 추가로 복수 태깅한다.

¶ 여자애들은 왜 이렇게 징징거리냐. **종특이나**

▶ [명시성] 명시 [맥락] 부정적 [영역] 관계/조건, 성, 문화 [강도] 강

☞ '종특'은 '종족특성'의 준말로 특정 종족이나 인종에 대한 비하성과 편향성이 나타나는 차별 또는 혐오 표현이므로 '관계/조건(출신)'을 태깅한다. 그리고 이 문장에서는 '종특'의 개념을 여성으로 확장하여 여성의 특정 행태를 일반화, 고정화하는 부적절성이 확인된다는 점에서 '성'과 '문화'를 추가로 복수 태깅한다.

¶ 부산국제공항은 인천국제공항에 비교하면 **시골** 공항에 불과하다.

▶ [명시성] 명시 [맥락] 부정적 [영역] 관계/조건, 기타 [강도] 강

☞ '시골'은 관련 보고서(박재현 외, 2009)에서 서울이 아닌 지역에 대한 비하성이 나타나는 차별적 표현으로 지적하였으므로 이를 포함하는 문장은 '관계/조건(출신)'을 태깅한다. 그리고 이 문장에서는 특정 공항의 규모나 시설에 대한 비하성이 확인되므로 '기타'를 추가로 복수 태깅한다.

¶ **외국인 여자분이** 일하시는데

뭔가 대답을 쓰시면서 해서

말하기도 싫어서

물 달라고 안했어요

▶ [명시성] 명시 [맥락] 부정적 [영역] 관계/조건, 성, 문화 [강도] 강

☞ '외국인'과 '여자분'은 종업원을 굳이 외국인과 여성으로 불필요하게 언급한 것이고, 관련 보고서(안상수 외, 2007)에서 불필요한 성별 언급을 특정 성에 대한 차별적 표현으로 지적하기도 했으므로, '관계/조건(출신)'과 '성'을 복수 태깅한다. 그리고 해당 종업원의 서비스 행태를 비난하는 공격성이 확인된다는 점에서 '문화'를 추가로 복수 태깅한다.

¶ **정말이지 중국 사람들..**

엄청 시끄럽다

▶ [명시성] 비명시 [맥락] 부정적 [영역] 관계/조건, 문화 [강도] 약

☞ 이 문장에서는 중국인을 시끄러운 사람들로 일반화하고 고정화하는 편향성이 확인된다는 점에서 '관계/조건(출신)'을 태깅한다. 그리고 중국인의 시끄러운 행태를 비난하는 공격성이 확인된다는 점에서 '문화'를 복수 태깅한다.

¶ **일본하고 중국은 계속 이렇게 유치원생 수준에 머물러 주면 된다**

▶ [명시성] 비명시 [맥락] 부정적 [영역] 관계/조건, 문화, 연령/세대 [강도] 약

☞ 이 문장에서는 특정 국가를 비하하는 비하성이 확인된다는 점에서 '관계/조건(출신)'을 태깅한다. 그리고 특정 국가의 행태를 조롱하는 비하성이 확인된다는 점에서 '문화'를 복수 태깅한다. 또한 유치원생 연령대의 지적 수준을 비하하는 비하성이 확인된다는 점에서 '관계/조건'과 '연령/세대'를 복수 태깅한다.

¶ **말을 섞다보면 외국인이지만 이거 한국사람 아닌가 싶게 뻗속까지 비슷해서 좀 불쾌하다.**

▶ [명시성] 비명시 [맥락] 부정적 [영역] 관계/조건, 문화 [강도] 약

☞ 이 문장에서는 외국인에게 기대한 모습에 대한 편향성이 확인된다는 점에서 '관계/조건(출신)'을 태깅한다. 그리고 외국인의 특정 행태를 불쾌하게 여기고 비난하는 공격성이 확인된다는 점에서 '문화'를 추가로 복수 태깅한다.

- 직업, 지위, 학력, 재산, 능력, 지력 등과 같은 '사회적 조건'과 관련한 부정절성이 나타나는 문장은 '관계/조건'으로 태깅

¶ **망할** 미래에셋 어찌구 편드로 구렁이 알 같은 내 돈 7만원 날리고서 완전히 손 털었다.

▶ [명시성] 명시 [맥락] 부정적 [영역] 관계/조건, 기타, 문화 [강도] 강

☞ '망하다'는 <표준>에서 다의어 의미 중 "(주로 '망할' 꼴로 쓰여) 못마땅한 사람이나 대상에 대하여 저주의 뜻으로 이르는 말"로 뜻풀이되는 욕설로 저주의 대상을 고려하여 영역을 태깅한다. 이 문장에서는 미래에셋이라는 기업 종사자나 기업 자체(무정물)를 비난하는 공격성이 나타나므로 '관계/조건'과 '기타'를 태깅한다. 그리고 미래에셋이라는 기업의 편드 운용 행태를 비난하는 공격성이 나타난다는 점에서 '문화'를 추가로 복수 태깅한다.

¶ 소리가 나올 정도로 생2의 미래가 창창(?)했는데

어느새 이 곳 생2는 **미친** 곳이야 어서 여기를 나가야...

수능을 정지합니다. 어 안되지않아?

틀렸어! 꿈도 희망도 표본도... 다 없어!

생2에서 고통받은 영혼을 힐링해줄 다른 과탐으로 ㅠㅠ

▶ [명시성] 명시 [맥락] 부정적 [영역] 관계/조건, 기타, 문화 [강도] 강

☞ '생2'라는 과목과 '생2'를 선택한 집단에 대한 부정절성이 표상된 문장으로 '과하게 공부를 열심히 해야 하는 과목 등의 지력과 연관될 수 있는 의미가 표상되며 해당 의미에 대한 조롱으로 볼 수 있으므로 '관계/조건'을 태깅한다.

¶ **미친** 점프 개쩨어

▶ [명시성] 명시 [맥락] 긍정적 [영역] 관계/조건, 문화, 신체, 기타 [강도] 약

☞ 이 문장은 점프의 능력과 관계있는 맥락으로 볼 수 있으므로 '문화'와 '신체' 그리고 능력과 관계하는 '관계/조건'을 태깅한다.

¶ 오늘 인턴들이 불러내서 밥 먹었는데 인턴 **주제**에 왜 이렇게 밥은 비싼 거 먹어서 내 일주일 용돈을 날리는 거냐.

▶ [명시성] 비명시 [맥락] 부정적 [영역] 관계/조건, 문화 [강도] 약

☞ 이 문장에서는 인턴이라는 지위를 가진 대상을 비하하는 비하성이 확인된다는 점에서 '관계/조건'을 태깅한다. 그리고 인턴들의 식사 관련 행태를 비난하는 공격성이 확인된다는 점에서 '문화'를 복수 태깅한다.

¶ 그런데 대학 졸업자라는 이 사람은 대부분의 보통 학력 혹은 공부 쪽과는 거래가 먼 상인들보다 훨씬 부정확하고 정직하지 못한 점이 내 분노를 샀다.

▶ [명시성] 비명시 [맥락] 부정적 [영역] 관계/조건, 문화 [강도] 약

☞ 이 문장에서는 학력이나 지력, 직업 등을 차별하는 편향성이 확인된다는 점에서 '관계/조건'을 태깅한다. 그리고 특정 대상에 대한 기대와 실망, 분노 등의 사고방식 관련 부정절성이 확인된다는 점에서 '문화'를 복수 태깅한다.

¶ 하지만 전문가 말만큼 믿을 게 못 되는 것도 없다.

▶ [명시성] 비명시 [맥락] 부정적 [영역] 관계/조건, 문화 [강도] 약

☞ 이 문장에서는 전문가라는 특정 집단의 능력과 지위를 비하하는 비하성이 확인된다는 점에서 '관계/조건'을 태깅한다. 그리고 전문가의 말하기 행태를 비난하는 공격성이 확인된다는 점에서 '문화'를 복수 태깅한다.

¶ 스타벅스야, 너희 커피 안 산다고 빈도 안 알아주냐? - - 치사해서 안 간다.

▶ [명시성] 비명시 [맥락] 부정적 [영역] 관계/조건, 기타, 문화 [강도] 약

☞ 이 문장에서는 '스타벅스'라는 기업 종사자와 기업 자체를 비난하는 공격성이 확인된다는 점에서 '관계/

조건'과 '기타'를 복수 태깅한다. 그리고 원두를 갈아주지 않는 행태를 비난하는 공격성이 확인된다는 점에서 '문화'를 추가로 복수 태깅한다.

¶ **이런 거 어디 없냐? ㅏ ㅏ 최저가라고 해서 클릭하면 핑크 따위만 팔거나(이 세상 컬러 중 가장 싫어하는 색), 이름만 걸고 다른 모델을 팔거나 하는 낚시가 많네.**

▶ [명시성] 비명시 [맥락] 부정적 [영역] 관계/조건, 문화, 기타 [강도] 약

☞ 의미적 고정성이 낮은 표현인 '따위'를 포함하는 이 문장은 특정 직업 종사자에 대한 비하성이 확인된다는 점에서 '관계/조건'을 태깅한다. 그리고 특정 색깔을 선호하는 행태에 대한 비하성과 특정 색깔을 저주하는 공격성이 확인된다는 점에서 '문화'와 '기타'를 복수 태깅한다.

¶ **장사 그 따위로 하지 마라..**

▶ [명시성] 비명시 [맥락] 부정적 [영역] 관계/조건, 문화 [강도] 약

☞ '따위'를 포함하는 이 문장은 특정 직업 종사자에 대한 비하성이 확인된다는 점에서 '관계/조건'을 태깅한다. 그리고 해당 직업 종사자의 행태에 대한 비하성이 확인된다는 점에서 '문화'를 복수 태깅한다.

¶ **생각하는 수준이 그러니까 아파트 경비나 하지**

▶ [명시성] 비명시 [맥락] 부정적 [영역] 관계/조건 [강도] 약

☞ 의미적 고정성이 낮은 표현인 '(이)나'를 포함하는 이 문장은 특정 직업 종사자에 대한 비하성이 확인된다는 점에서 '관계/조건'을 태깅한다. 그리고 타인의 특정 사고방식에 대한 비하성이 확인된다는 점에서 '문화'를 복수 태깅한다.

¶ **말하는 것 보니 평생 고시원이나 전전하게 생겼지**

▶ [명시성] 비명시 [맥락] 부정적 [영역] 관계/조건, 기타, 문화 [강도] 약

☞ '(이)나'를 포함하는 이 문장은 고시원이라는 특정 주거 형태에 거주하는 사람과 주거 형태 자체를 비하하는 비하성이 확인된다는 점에서 '관계/조건'과 '기타'를 복수 태깅한다. 그리고 말하는 태도나 방식 등의 특정 행태를 비하하는 비하성이 확인된다는 점에서 '문화'를 추가로 복수 태깅한다.

- 혼인, 가족 형태, 가족 관계 등과 같은 '사적 관계' 관련한 부적절성이 나타나는 문장은 '관계/조건'으로 태깅

¶ **말본새를 보면 에미도 없는 것들이지**

▶ [명시성] 명시 [맥락] 부정적 [영역] 관계/조건, 성, 문화 [강도] 강

☞ '에미(어미)'는 <표준> 등의 대사전류에 '어머니의 낮춤말'로 뜻풀이되는 비어로 이를 포함하는 문장은 '관계/조건(사적 관계)'과 '성'을 태깅한다. 또한 이 문장에서는 말하기의 특정 방식이나 태도 관련 행태를 비하하는 비하성이 확인된다는 점에서 '문화'를 복수 태깅한다.

¶ **가만보니 이상하게 느끼한 쪽은 유부남이더구만.**

▶ [명시성] 비명시 [맥락] 부정적 [영역] 관계/조건, 성, 문화 [강도] 약

☞ 이 문장에서는 아내가 있는 남성에게 대한 비하성이 된다는 점에서 '관계/조건(사적 관계)'과 '성'을 복수 태깅한다. 그리고 그들의 행태에 대한 비하성이 확인된다는 점에서 '문화'를 추가로 복수 태깅한다.

¶ (긍정적 맥락에서 조카를 귀엽게 묘사) **시어머니와 며느리 버릇없는 조카 세연양과.**

▶ [명시성] 비명시 [맥락] 긍정적 [영역] 관계/조건, 문화, 연령/세대 [강도] 약

☞ 이 문장에서는 조카의 버릇없는 행태 관련 부적절성이 확인된다는 점에서 '관계/조건(사적 관계)'과 '문화'를 복수 태깅한다. 그리고 나이 어린 조카의 버릇없음은 귀엽게 간주된다는 편향성이 확인된다는 점에서 '연령/세대'를 추가로 복수 태깅한다.

<주의>

- ‘돌다’, ‘미치다’, ‘미친놈’, ‘돌(똥)아이’ 등의 표현은 사회적 지력과 의미적 관계성을 지닌 표현이라 하더라도 맥락상 그러한 관련성을 파악할 수 없다면 ‘관계/조건’을 태깅하지 않음

¶ **미친 쇼나 예쁘자나** ㄱ ㅏ ㅏ

- ▶ [명시성] 명시 [맥락] 긍정적 [영역] 신체 [강도] 약
- ☞ 이 문장에서는 아이돌 외모에 대한 감탄으로 ‘미친’이 사용되어 지력과 크게 관계없으므로 ‘관계/조건’을 태깅하지 않는다.

¶ **야깁씨발** ㅋㅋㅋㅋ**미친놈들인가** ㅋㅋㅋㅋㅋㅋㅋㅋㅋㅋㅋㅋㅋㅋㅋㅋ

- ▶ [명시성] 명시 [맥락] 부정적 [영역] 문화, 기타 [강도] 강
- ☞ 인터넷 게시판 사용자들에 대한 문장으로 의도를 정확히 파악할 수 없지만 욕설 등의 부적절성 표현이 나타나므로 ‘부정적’으로 맥락을 태깅하고, 게시판 사용과 관련한 ‘문화’, 게시판에 대한 ‘기타’를 태깅하되 지력과는 크게 관계없으므로 ‘관계/조건’을 태깅하지 않는다.

¶ **월요일같아서 그런가, 단거 미친듯이 땡기네**

- ☞ ‘미친 듯이’는 고려대 대사전에 ‘무언가에 몰입하여 매우 열심히.’로 제시되어 있으므로 특별한 부적절이 관찰되지 않는다면 부적절성으로 판정하지 않는다.

4.3.7. 기타

- 사물(무정물)이나 감정의 배설 등과 같이 위의 유형으로 분류되지 않는 사례들과 관련한 부적절성이 나타나는 문장은 ‘기타’로 태깅

¶ **특히 수박 잘라오라는 거!! 수박은 잘 안 잘리는 과일이거든! 씨팔** 내 팔뚝 0배는 굵으면 서! 몸 안 움직이는 그대여 계속 똥똥해지세요!!! - -

- ▶ [명시성] 명시 [맥락] 부정적 [영역] 기타, 신체, 문화 [강도] 강
- ☞ ‘씨팔’은 감정의 배설과 관련한 욕설이므로, 이를 포함하는 문장은 ‘기타’를 태깅한다. 그리고 이 문장에서는 ‘똥똥하다’라는 부정적 서술어를 통해 특정 신체 조건을 조롱하는 비하성이 확인된다는 점에서 ‘신체’를 복수 태깅한다. 또한 잘 움직이지 않으려 하는 행태를 저주하는 공격성이 확인된다는 점에서 ‘문화’를 추가로 복수 태깅한다.

¶ **그나저나 오늘 날씨 왠케 좋아!!! 이런 X!!** 욕 나오게 하네.

- ▶ [명시성] 명시 [맥락] 부정적 [영역] 기타 [강도] 강
- ☞ ‘X’는 감정의 배설과 관련한 욕설임이 분명하므로, 이를 포함하는 문장은 ‘기타’로 태깅한다.

¶ **현피뜨자 개새끼들아**

- ▶ [명시성] 명시 [맥락] 부정적 [영역] 기타, 문화 [강도] 강
- ☞ 위의 문장이 단독으로 실현된 경우 ‘개새끼’는 지칭하는 바가 명확하지 않지만 맥락상 사람임이 비교적 정확하고, ‘하는 짓이 알밋거나 더럽고 됴됨이가 좋지 아니한’으로 주로 행동과 결부되어 사용되는 욕설임을 고려하여 ‘문화’로 태깅하고 감정 배설로 볼 수 있으므로 ‘기타’ 또한 복수 태깅한다.

¶ **소방관을 왜 여자 뽑음? 민원서들임?**

- ▶ [명시성] 명시 [맥락] 부정적 [영역] 기타, 관계/조건, 문화, 성 [강도] 강
- ☞ ‘민원서들’은 민원 관련 업무와 업무 담당자를 비하하는 차별적 표현이므로, 이를 포함하는 문장은 ‘기타’와 ‘관계/조건’을 복수 태깅한다. 그리고 이 문장에서는 특정 성별이 소방관을 선발하는 행태를 비난하는 공격성과 특정 성별의 소방관을 비하하는 비하성이 확인된다는 점에서 ‘문화’와 ‘성’을 추가로 복수

태깅한다.

¶ 오늘 결국 **싸구려** 면바지 한 장 사고 10시간은 걸은 거 같다.

▶ [명시성] 비명시 [맥락] 부정적 [영역] 기타, 관계/조건 [강도] 약

☞ 이 문장에서는 가격이 저렴한 옷(무정물)을 비하하는 비하성이 확인된다는 점에서 '기타'를 태깅한다. 그리고 이러한 옷을 입은 사람의 지위와 능력이 낮을 것이라고 일반화하여 생각하는 편향성이 확인된다는 점에서 '관계/조건'을 복수 태깅한다.

¶ **츄리닝에 잠바때기** 입고 광화문 가는 길입니다.

▶ [명시성] 비명시 [맥락] 부정적 [영역] 기타, 관계/조건 [강도] 약

☞ 이 문장에서는 잠바라는 옷(무정물)을 비하하는 비하성이 확인된다는 점에서 '기타'를 태깅한다. 그리고 이러한 옷을 입은 사람의 지위와 능력이 낮을 것이라고 일반화하여 생각하는 편향성이 확인된다는 점에서 '관계/조건'을 복수 태깅한다.

¶ **북한에 최소한 달이나 명왕성만한 소행성 하나 떨어지면 좋겠다**

▶ [명시성] 비명시 [맥락] 부정적 [영역] 기타, 문화 [강도] 약

☞ 이 문장에서는 북한이라는 국가(무정물)를 저주하는 공격성이 확인된다는 점에서 '기타'를 태깅한다. 그리고 북한이라는 정치 체제 또는 권력을 혐오하는 공격성이 확인된다는 점에서 '문화'를 복수 태깅한다.

- 사물(무정물)이나 감정의 배설 등과 같이 위의 유형으로 분류되지 않는 사례들과 관련한 부적절성이 나타나는 문장은 '기타'로 태깅

4.4. 강도 주석

4.4.1. 강

4.4.1.1. 명시적 부적절성이 부정적 맥락에서 나타나는 문장

- <표준> 등의 대사전류에 욕설, 비어, 비하성 속어로 직접 기술된 명시적 부적절성 표현이 부정적 맥락에서 나타나는 문장은 ‘강’으로 태깅

¶ 말본새를 보면 **에미도** 없는 것들이지

▶ [명시성] 명시 [맥락] 부정적 [영역] 관계/조건, 성, 문화 [강도] 강

☞ ‘에미(어미)’는 <표준> 등의 대사전류에 ‘어머니의 낮춤말’로 뜻풀이되는 비어로 부정적 맥락에서 이를 포함하는 문장은 ‘강’으로 태깅한다.

¶ -_-; 꼬질꼬질한 골목길 사이로 할머니 얼굴이 그려진 간판이 보이고, 가게에 들어서자마자 피죽도 못 먹은 **노인네** 둘이 비틀거리며 나오는 중.

▶ [명시성] 명시 [맥락] 부정적 [영역] 연령/세대, 신체, 문화, 기타 [강도] 강

☞ ‘노인네’는 <고려대>에 ‘나이든 사람을 얕잡아 이르는 말’로 뜻풀이되는 비어로 부정적 맥락에서 이를 포함하는 문장은 ‘강’으로 태깅한다.

- 공공기관의 관련 보고서 등에서 지적된 차별/혐오/선정 표현 등의 명시적 부적절성 표현이 부정적 맥락에서 나타나는 문장은 ‘강’으로 태깅

¶ XXX의 **미망인** XXX 여사가 의기투합했고..

▶ [명시성] 명시 [맥락] 부정적 [영역] 성, 관계/조건 [강도] 강

☞ ‘미망인’은 관련 보고서(조태린 외, 2006)에서 결혼한 여성이 남편이 사망했음에도 ‘아직 따라 죽지 못한 사람’이라는 봉건적 의미를 담고 있는 차별적 표현으로 지적했으므로, 부정적 맥락에서 이를 포함하는 문장은 ‘강’으로 태깅한다.

¶ 넌 **조선족** 사는데 같이 살고 싶어?

▶ [명시성] 명시 [맥락] 부정적 [영역] 관계/조건 [강도] 강

☞ 중국이 아닌 한국에서 사용되는 ‘조선족’은 관련 보고서(박재현 외, 2009)에서 차별적, 비하적으로 사용됨을 지적하고 ‘재충동포’ 또는 ‘한국계 중국인’ 등으로 대체하는 것을 제안하고, 국립국어원 행정용어순화어(2018)에서도 동일한 순화어를 제안하므로 부정적 맥락에서 이를 포함하는 문장은 ‘강’으로 태깅한다.

¶ 부산국제공항은 인천국제공항에 비교하면 **시골** 공항에 불과하다.

▶ [명시성] 명시 [맥락] 부정적 [영역] 관계/조건, 기타 [강도] 강

☞ ‘시골’은 관련 보고서(박재현 외, 2009)에서 서울이 아닌 지역에 대한 비하성이 나타나는 차별적 표현으로 지적하였으므로, 부정적 맥락에서 이를 포함하는 문장은 ‘강’으로 판정한다.

- 공공기관의 관련 보고서 등에서 차별/혐오/선정 표현 등으로 지적되지 않았더라도 같은 방식과 내용으로 만들어진 유사 표현이 부정적 맥락에서 나타나는 문장은 ‘강’으로 태깅

¶ 여기 **틀딱들이** 대거 나온 듯.

▶ [명시성] 명시 [맥락] 부정적 [영역] 연령/세대, 문화 [강도] 강

☞ ‘틀딱’은 특정 연령대의 사람들과 그들의 행태에 대한 차별 및 혐오 표현이므로 부정적 맥락에서 이를 포함하는 문장은 ‘강’으로 태깅한다.

¶ 솔직히 **기균충이나 지균충**은 동기란 생각 안 해.

▶ [명시성] 명시 [맥락] 부정적 [영역] 문화, 관계/조건, 기타 [강도] 강

☞ '기균충'과 '지균충'은 특정 입시 제도를 통해 들어온 대학생에 대한 차별 및 혐오 표현이므로 부정적 맥락에서 이를 포함하는 문장은 '강'으로 태깅한다

4.4.1.2. 비명시적 부적절성이 성적 폭력성, 선정성 등 관련 부정적 맥락에서 나타나는 문장

- 비명시적 부적절성 문장이라도 성폭력, 성추행, 성희롱 등의 성 관련 폭력적인 내용이나 성 관계, 성적 대상화(상품화) 등의 선정적인 내용을 포함하는 부정적 맥락에서 나타나는 문장은 '강'으로 태깅

¶ **나랑 떡칠래?**

▶ [명시성] 비명시 [맥락] 부정적 [영역] 문화, 성 [강도] 강

☞ '떡치다'는 <고려대>에서 "성적으로 관계를 맺는 일을 하다."라는 의미의 단순 속어로 제시되어 있어도 실제 사용례를 고려하면 선정성, 성희롱 등의 부적절성을 나타내므로, 이 문장은 '강'으로 태깅한다.

¶ **여자는 기력지가 긴 애들이 먹기도 좋지**

▶ [명시성] 비명시 [맥락] 부정적 [영역] 문화, 성, 신체 [강도] 강

☞ '먹다'는 <표준> 등의 대사전류에 단순 속어로 제시되어 있고 화자도 단지 그런 의미로만 사용한 것일지라도 사전의 다의어 뜻풀이((남자가 여자를) 성적으로 침해하여 짓밟다) 및 사용례 등을 고려하면, 선정성(성폭력, 성희롱) 등의 부적절성을 나타내므로, 이 문장은 '강'으로 태깅한다.

¶ **저 중에 누가 제일 따먹고 싶게 생김?**

▶ [명시성] 비명시 [맥락] 부정적 [영역] 문화, 성, 신체 [강도] 강

☞ '따먹다'는 <표준> 등의 대사전류에 단순 속어로 제시되어 있고, 화자도 단지 그런 의미로만 사용한 것일지라도 '여자의 정조를 빼앗다'라는 뜻풀이와 사용례 등을 고려하면 선정성(성폭력, 성희롱) 등의 부적절성을 나타내므로, 이 문장은 '강'으로 태깅한다.

¶ **외국 아가씨들이랑 놀기 위해..**

▶ [명시성] 비명시 [맥락] 부정적 [영역] 성, 연령/세대, 문화, 관계/조건 [강도] 강

☞ 이 문장에서는 젊은 여성을 성적 대상화하는 부적절성이 나타난다는 점에서 '강'으로 태깅한다.

¶ **(내가 좀;;) 1박 2일에서 찬물에 들어갔다 나온 찌이 있는데 난 박찬호 허벅지랑 찬물에 젖은 고탄력 흰...**

▶ [명시성] 비명시 [맥락] 부정적 [영역] 신체, 성 [강도] 강

☞ 이 문장은 남성의 특정 신체 부위를 성적 대상화하여 불필요하게 언급하는 편향성이 확인되므로 '강'으로 태깅한다.

- 비명시적 부적절성 문장이라도 현대의 윤리적 기준을 크게 위배하는 내용을 포함하는 부정적 맥락에서 나타나는 문장은 '강'으로 태깅

¶ **솔까 남자가 저런 책 썼으면 출판은커녕 화형 당할 듯**

▶ [명시성] 비명시 [맥락] 부정적 [영역] 문화, 기타, 성 [강도] 강

☞ '화형'은 사전적 의미에서는 특정 형벌을 의미하는 객관적 단어이지만, 역사적으로 해당 형벌이 행해졌다고 하더라도 현대적 기준에서는 윤리적 기준을 크게 위배하는 내용을 담고 있으므로, 이를 포함하는 문장은 '강'으로 태깅한다.

¶ 나 같은 건 자살이 답이다

▶ [명시성] 비명시 [맥락] 부정적 [영역] 관계/조건, 문화 [강도] 강

☞ '자살'은 대사전류에 따른 명시적 부적절성 표현은 아니지만, 현대적 기준에서 보았을 때 윤리적 기준을 크게 위배하는 내용을 담고 있으므로, 이를 포함하는 문장은 '강'으로 태깅한다.

¶ 동성애는 나라를 망칩니다

▶ [명시성] 비명시 [맥락] 부정적 [영역] 성, 문화 [강도] 강

☞ 이 문장에서는 동성애를 차별하고 혐오하는 편향성과 공격성이 확인되므로 '강'으로 태깅한다.

¶ 조두순은 죽여도 무죄. 제발

▶ [명시성] 비명시 [맥락] 부정적 [영역] 문화 [강도] 강

☞ 이 문장에서는 현대의 윤리적 기준을 크게 위배하는 범죄자는 죽여도 무죄라는 내용이 확인되므로 '강'으로 태깅한다.

¶ 뭐 x 리메이크버전 ㅇㅈ?ㅠㅠㅠㅠ나죽여ㅠㅠㅠㅠ나죽이려고작정한건가

▶ [명시성] 명시 [맥락] 긍정적 [영역] 문화, 기타 [강도] 약

☞ 이 문장에서는 '죽이려고 작정하다'만을 보면 현대의 윤리적 기준을 크게 위배하지만 맥락을 고려하였을 때 그것이 부정물에 대한 감탄의 비유적 표현으로 확인되므로 '명시', '긍정적', '약'을 태깅한다.

4.4.2. 약

4.4.2.1. 명시적 부적절성이 긍정적 맥락에서 나타나는 문장

- <표준> 등의 대사전류에 욕설, 비어, 비하성 속어로 직접 기술된 명시적 부적절성 표현을 포함하는 문장이라도 긍정적 맥락에서 나타나는 문장은 '약'으로 태깅

¶ (가족의 상황에 대한 묘사) 돈벌러 가신 남정네분들은 늦게 퇴근하고.

▶ [명시성] 명시 [맥락] 긍정적 [영역] 성, 관계/조건 [강도] 약

☞ '남정네'는 <고려대>에서 '여자가 사내를 조금 낮추어 이르는 말'로 뜻풀이되는 비어이므로 명시적 부적절성 표현이지만, 이를 포함하는 문장의 맥락이 긍정적이므로 '약'으로 태깅한다.

¶ 왜 키우냐 묻거든 강 웃지요 근데 꼴에 머리는 겁나 좋다 그냥 똥똥한 사람같다 ㅋㅋ 오래 살아라.

▶ [명시성] 명시 [맥락] 긍정적 [영역] 신체, 관계/조건, 문화 [강도] 약

☞ '꼴'은 <표준> 등의 대사전류에 사람의 모양새나 행태를 낮잡아 이르는 말로 뜻풀이되는 비어이므로 명시적 부적절성 표현이지만, 이를 포함하는 문장의 맥락이 긍정적이므로 '약'으로 태깅한다.

¶ 그나저나 오늘 날씨 웰케 좋아!!! 이런 X!! 욕 나오게 하네.

▶ [명시성] 명시 [맥락] 긍정적 [영역] 기타 [강도] 약

☞ 'X'는 뒤에 이어지는 '욕 나오게 하네'로 인해 욕설임이 명백하므로 명시적 부적절성 표현이지만, 이를 포함하는 문장의 맥락이 긍정적이므로 '약'으로 태깅한다

4.4.2.2. 비명시적 부적절성이 성적 폭력성, 선정성 등 관련 부정적 맥락이 아닌 경우에 나타나는 문장

- 성 관련 폭력성, 선정성 등의 내용이 아닌 비명시적 부적절성 문장은 맥락이 부정적인지,

긍정적인지 상관없이 ‘약’으로 태깅

¶ 이 나이 먹도록 취업도 못하고 엄마한테 빌붙어 살고..

▶ [명시성] 비명시 [맥락] 부정적 [영역] 관계/조건, 연령/세대, 문화 [강도] 약

¶ ‘국민안전처’는 ‘국민이 안전한 척’ 하고`

▶ [명시성] 비명시 [맥락] 부정적 [영역] 문화, 관계/조건, 기타 [강도] 약

¶ 국방부`는 `국방 방심부` 같구나...

▶ [명시성] 비명시 [맥락] 부정적 [영역] 문화, 관계/조건, 기타 [강도] 약

¶ 오늘 새벽같이 일어나서 XX청 가서 교육 받고 그 후에 YY서로 이동, 개 끌려다니듯 계속 높이신 분들의 방에서 말씀을 들었다.

▶ [명시성] 비명시 [맥락] 부정적 [영역] 문화, 관계/조건, 기타 [강도] 약

¶ 하지만 전문가 말만큼 믿을 게 못 되는 것도 없다.

▶ [명시성] 비명시 [맥락] 부정적 [영역] 관계/조건, 문화 [강도] 약

¶ 생각하는 수준이 그러니까 아파트 경비나 하지

▶ [명시성] 비명시 [맥락] 부정적 [영역] 관계/조건, 문화 [강도] 약

¶ 한약 먹는 보람도 없이 이게 무슨 짓이람..

▶ [명시성] 비명시 [맥락] 부정적 [영역] 문화 [강도] 약

¶ 오늘 결국 싸구려 면바지 한 장 사고 10시간은 걸은 거 같다.

▶ [명시성] 비명시 [맥락] 부정적 [영역] 기타, 관계/조건 [강도] 약

¶ (긍정적 맥락에서 조카를 귀엽게 묘사) 시어머니와 며느리 버릇없는 조카 세연양과.

▶ [명시성] 비명시 [맥락] 긍정적 [영역] 관계/조건, 문화, 연령/세대 [강도] 약

¶ 진짜 저 사진 찍어놓고 너무 명충해보여서 한동안 웃었다.ㅋㅋ

▶ [명시성] 비명시 [맥락] 긍정적 [영역] 관계/조건, 신체 [강도] 약

5. 개인정보 판별 기준과 비식별화 태그 세트

5.1. 개인정보 포함 문장 판별 기준

- 이름, 출신/소속 번호, 온라인 계정, 주소, 상호명, 상표명은 부적절성 표현의 대상 여부와 관계없이 모두 비식별화
- 그 외 장소명, 창작물명은 부적절성 표현의 대상일 경우만 비식별화

1) 부적절성 표현 대상 여부와 관계없이 이름, 출신/소속, 번호, 온라인 계정, 주소, 상호명, 상표명 등에 대해서 모두 비식별화

예) 아까 소라 언냐 찾아와서 야기했는데 앙골가고 나서 (이름)

&account& 팔로우하셈 (온라인 계정)

미림여고 다녀, (출신/소속)

&tel-num&이야 (전화번호)

주민번호 &social-security-num&(주민등록번호)

카드번호 먼저 &card-num&(카드번호)

저런 짜치는 도구들 다이소에 다 있잖아(상호명)

맛동산이 얼마나 맛있는데 엄마가 끈대입맛이래(상표명)

2) 이름의 경우, 실존 인물만 비식별화 대상으로 삼는다. 게임/소설/영화/드라마의 가상 인물/캐릭터는 제외

예) 봉준호가 자신안의 한남과 강남좌파를 드러내는 영화요

예) 강호동이랑 유재석은 각자 다른이유이긴 하지만 똑같은 레벨로 혐오스러움

예) 트럼프는 문재인보다 더한 빨갱이야...!

예) 어떤 미친 한남새끼가 팀버튼 죽었다고 한거

→ ‘팀버튼’은 부적절성 표현 대상이 아니지만 이름이므로 비식별화

예) 그래서 김지영이 한남 죽이기라도 했냐고ㅋㅋㅋㅋ 또 또 한남들 읽지도 않고 지랄해대지

→ 소설/영화 ‘82년생 김지영’의 주인공. 실존 인물이 아니므로 비식별화하지 않음

- 실명 외에도 인명 별칭, 대화명, 변형된 형태의 인명에 대해서도 맥락상 누구인지 유추가 가능하다면 비식별화

예) 문제의 맥심 어찌고 동님이 알려해서 1번 봤는데 그것은 그냥.. 지난주에 싸운 거래처 직원 얼굴이랑 똑같이 생긴 평범한 한남 얼굴이라 타격 0이었던

예) 박근혜, 박근혜, 안지호, 박대통령

예) 차녹얼빠 한남 감독 중 최고 클래스 아닌가요 → 박찬욱 감독

- 부적절성 표현 뒤에 이름이 나타나는 경우 이름만 비식별화

예) 한남노엘, 한남춘배, 노콘준상

3) 장소명, 창작물 등에 대해서는 부적절성 표현의 대상이 되는 경우에만 비식별화

예) 후쿠시마 원료를 사용하는데다가 여혐 마케팅 → 지역

예) 미우새 존나 싫어하는데 저걸 보면서 자칭 영포티한테 희망을 주면서 키덜트에 대한 환상이나 미화같은거만 보여주는 것 같아서 싫고 또 거기 나오는 출연진 어머니란 사람들이 간간이 하는 그 세대 특유의 빵은 말때문에 싫어함 → 방송 프로그램 이름

예) 이마무라 쇼헤이의 우나기 선생, 지금 시점에서는 이 할아버지가 개저라는 인상은 지을수가 없는건데 그럼에도 재밌는것은 뒷골목으로 표현되는 당시의 인간군상들에 대한 묘사들.

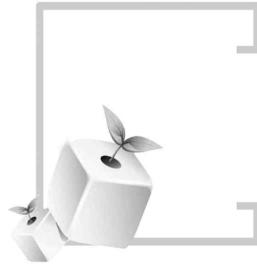
→ 책 이름. 직접적인 비윤리적 표현의 대상이 아니므로 비식별화하지 않음

5.2. 개인정보 비식별화 태그

○ 개인정보가 포함된 문장으로 판별되면, 해당 문장 내 개인정보를 비식별화 기호로 처리

- 이름, 출신/소속 번호, 온라인 계정, 주소, 상호명, 상표명을 제외하고는 비윤리적 표현의 대상일 경우만 비식별화함에 주의

분류		태그	항목
이름		&name&	실명, 특수 애칭, 별명, 대화명, 필명 가수 그룹명도 포함 예) 황교안, 장성규, 문재인, 안재현, 김구라, 김정은, 김희철, 개동수, 러브리즈 예) 박근혜, 안소희, 차녹얼빠, 박대통령 예) 한남노엘, 한남춘배, 노콘준상 예) 문제의 핵심 어찌고 똘님, 김롯데
출신 소속	출신 학교, 지역	&affiliation&	출신 학교, 지역 예) 나 OO 중학교 나왔어. 출신 지역이 아닌 지역명은 장소로 주석
	온라인 커뮤니티		예) 메갈, 클리앙, 디씨갤, 워마드, 일베, 허갤, 베캠
	정당		예) 자한당, 자유한국당, 애국당 예) 민중당, 자살당 → 변형된 형태이지만 맥락에서 어떤 정당인지 유추가 가능하면 주석
	팬클럽		예) 아재리너스, 리브리너스, 엑소엘
	기타		예) 그러나 음악시간은 노는시간이었고 합창부는 짚따나 가는 동아리 취급이었음
번호	고유 식별 번호	&social-security-num&	주민등록번호
	전화번호	&tel-num&	
	카드번호	&card-num&	
	계좌번호	&bank-account&	
	기타 번호	&num&	일련번호, (구매자) 식별 번호, 사업자 등록 번 호, 비밀번호
온라인 계정		&online-account&	아이디, 이메일 주소
주소		&address&	상세 주소, 아파트 및 거주 건물명
상호명		&company&	기업/회사/상점 이름 예) 롯데, YG, 스킵, 트위터, 넷플릭스, 조선 일보, 다이소
장소명		&location&	나라, 도시 이름 예) 폴란드, 후쿠시마
상표명		&brand&	제품명, 브랜드명 예) 오레오, 맛동산, 체리마루
창작물명		&art&	소설, 영화, 드라마, 만화 등의 작품명 예) 키리시마가 동아리 활동 그만둔대, 가우스 전자, 김지영, 우모페, 보랩
기타		&other&	위에서 언급하지 않은 항목



부록

[붙임 2] 부적절성 관련 어휘(표현) 목록



기존 관련 보고서에서 차별적, 혐오적 표현 등으로 지적된 어휘(표현)

대범주	소범주	언어 표현	유사 예시	비고
지역	비(非)-	비강남, 비서울(권), 비수도권, 비서구		
	촌(村)-류	촌티, 촌사람 등		
		시골	시골 사람, 시골 출신	
		지방으로 내려가다	서울로 올라가다	상경, 낙향도 동일하게 처리
		지방대		수도권 소재와 지방 소재의 지리적 구분을 명시해야 하는 경우 제외
		수도권대학		수도권 소재와 지방 소재의 지리적 구분을 명시해야 하는 경우 제외
		낙후 지역		지역명이 공개된 경우
		판자촌	비닐하우스촌	
인종		백인(종)	항인(종), 유색인(종), 흑인(종) 등	불필요하게 인종을 강조하는 경우
		살색		
		조선족		
		에스키모		
		인디언		
혈통	-계(系) 류	한국계	아시아계, 유태계, 아프리카계	불필요하게 혈통을 강조하는 경우
		한국인 피	토종 한국인, 하프 코리아	불필요하게 혈통을 강조하는 경우
		교포	동포, 한국계	불필요하게 혈통을 강조하는 경우
		혼혈(아, 인)	잡종	불필요하게 혈통을 강조하는 경우
		다문화		호칭으로 사용되는 경우
		라이파이한	신(新라이파이한, 라이베리아판(版) 라이파이한	
이주		동남아 노동자		
		외래종	토종	사람에 대해 사용하는 경우
		귀순자	탈북자, 새터민	
		여성 결혼 이민자.	결혼 이민자	불필요하게 이주 배경을 강조하는

		결혼 이주 여성		경우 불필요하게 성별을 강조하는 경우
장애	시력	장님	애꾸눈, 외눈박이, 오대박이	관련한 속담(장님 코끼리 다리 만지듯 등) 및 관용구 표현 포함(눈 먼 돈, 눈 먼 횡포 등)
	청력	귀머거리	농아(인), (자)	사전에 비어로 올라 있음
	언어	언청이	언청새님	
		병어리		관련한 속담(병어리 냉가슴 앓는다, 귀머거리 3년 병어리 3년 등) 및 관용구 표현 포함
	신체	얕은뱅이		
		절름발이		
		곶추	곶사등이, 곶사	
		병신	얕은뱅이, 반신불구	고빈도 욕설인 '병신'의 경우, 장애와 관계되지 않은 경우 존재
		난쟁이	땅딸보, 딸보, 땅개, 작다리	
		곰배	곰배팔이, 곰배팔, 외팔이, 외팔뚝이, 조막손	
		절름발이	절뚝발이, 절뚝이	관용구 표현 포함(절름발이 행정 등)
		곶장다리	안짱다리	
		뺨장다리		
		외다리	외짱다리	
	정신	무뇌아		
		정신박약자		
		바보		
		등신		
		칠뜨기		
		팔푼이		
얼간이				
떨떨이				
기타	장애우, 장애자			
	장애에도 불구하고	정상인 못지않게		
	불구	불구자		

직업	잡(雜)-류	잡상인	잡역부, 잡역꾼	특정 직업을 비하할 목적으로 사용된 경우
	-쟁이 류	그림쟁이	구두쟁이, 도배쟁이	
		간호원		
		청소부		
		가정부	식모, 파출부	
		철밥통		
		노가다		
		기레기		
연령		틀딱		
		개저씨		
		잼민이		
	-딩 류	초딩	중딩, 고딩	
가족, 출생		결손가정	편모가정, 편부가정	
		사생아		
		팔삭둥이	미숙아	
성	불필요한 성별 강조	여성 혹은 남성임을 불필요하게 표현하거나 고정관념에 근거한 성역할이나 성별 속성 강조하는 어휘 표현	이00(여, 43세), 주부000, 여류000, 여직원, 남자 간호사, 기센 여주인 등	
		여성의 고정관념적 속성을 불필요하게 강조하는 어휘 표현	옛되어 보이는, 꼬리친다, 양탈부리다, 양갈지다, 야들야들, 여우, 여성미, 여성스러운, 질투, 가녀린, 청순 등	
	고정관념적 속성 강조	남성의 고정관념적 속성을 불필요하게 강조하는 어휘 표현	과감한, 스케일이 큰, 능글능글, 무뚝뚝, 대장부 등	
		성차별적 이데올로기를 포함하는 어휘 표현	미망인, 출가외인, 집사람, 안사람, 기동서방, 도련님, 집안의 대들보, 마누라, 여시, 여편네, 처녀 출전, 처녀 생식	

		특정 성을 자극적으로 표현하는 어휘 표현	흑진주, 신데렐라, 레이싱걸, 엽기녀, 쿼카, 매력녀, 슛총각, 영계, 꽃미남, 완소남, 킹카, 김치녀 등	
		특정 성의 성적/신체적 측면을 이용하는 어휘 표현	쭈쭈빵빵, 섹시 가슴, S라인, 환상의 바디라인, 울끈불끈 가슴근육, 조각 같은, 탄탄한 근육, 로린이	
		특정 성을 비하하는 어휘 표현	여편네, 부엌데기, 암캐, 접대부, 계집애, 머슴, 기생오라비, 제비족, 수컷, 마마보이, 탕아, 한남	
	성별 언어 구조	하나의 성으로 다른 성까지 대표하는 경우	스포츠맨십, 샐러리맨, 학부형	
기타	-충(蟲) 류	맘충	기균충, 급식충	
		꼰대		
	-병(病) 류	중2병		
		흥어	전라디언	흥어'의 경우 특정 지역을 차별하는 의미로 사용될 때에 한정
		신용불량자		

<기획·연구>

국립국어원 강미영 언어정보과장

국립국어원 유희정 학예연구사

국립국어원 이민주 연구원

국립국어원 박미은 연구원

국립국어원 정영은 연구원

<연구 참여자>

연구 책임자 조태린(연세대학교)

공동 연구원 공나형(전남대학교)

김미숙(나라지식정보), 한용운(나라지식정보)

박승희(나라지식정보), 변순용(서울교육대학교)

김봉제(서울교육대학교), 이청호(상명대학교)

윤기현(바이칼AI), 김성진(바이칼AI)

연구 보조원 박예슬(연세대학교), 심주희(연세대학교)

김승래(연세대학교), 배승희(연세대학교)

최민경(연세대학교), 이재엽(바이칼AI), 김미영(바이칼AI)

발행인: 국립국어원장

발행처: 국립국어원

서울시 강서구 금남화로 154

전화 02-2669-9775, 전송 02-2669-9727

인쇄일: 2023년 2월 21일

발행일: 2023년 2월 21일

인 쇄: 연세대학교 POD센터

※ 이 보고서는 국립국어원의 용역비로 수행한 ‘2022년 말뭉치 비윤리성 분석 및 연구’ 사업의 결과물을 발간한 것입니다.