

KAIST

카이스트 어휘망과 명사의미 연구

배선미, 시정곤
KORTERM/BOLA

차례

□ 카이스트 어휘의미망 (CoreNet)

- ◆ 목적, 특징, 구축과정, 구성

□ 카이스트 어휘망 브라우저 데모

- ◆ 단어, 개념명, 코어넷 번호 검색
- ◆ 격틀 검색

□ 한국어 명사 개념지도

- ◆ 계층별 전체 개념분포 분석
- ◆ 명사의 계층별 개념분포 분석
- ◆ 개념명에 따른 명사 어휘 분포

□ 맺음말

- ◆ 보완점 및 향후 과제

Korea Advanced Institute of Science and Technology

카이스트 어휘망의 목적

□ 카이스트 한국어 단어망은 자연언어저리를 함에 있어 의미의 애매성 애소를 위하여 개념체계를 기반으로 구축된 어휘 의미망이다.

□ 중의성 애소 활용 예:

Korean All Words Sense Tagging

[Te15912] [Source-Channel] [MEL] [MEM] [SYM] [ALL]

[Correct Incomert] [Correct after update, Press on Update loading] [ends to correct. Press on Print out.]

```

215 알자 enc 121120410000
766 up
471 알자 enc 121160000000
472 알자
473 알자 enc 123420000000 112201188000
474 알자
475 알자
476 알자 enc 123420000000 112212000000
477 알자 enc 123420000000 112212000000
478 알자 enc 123420000000 112212000000
479 알자 enc 123420000000 112212000000
480 알자
481 알자 enc 123420000000 112212000000
482 알자
483 알자 enc 123420000000 112212000000
484 알자
485 알자 enc 123420000000 112212000000
486 알자
487 알자 enc 123420000000 112212000000
488 알자
489 알자 enc 123420000000 112212000000
490 알자
491 알자 enc 123420000000 112212000000
492 알자
493 알자
494 알자 enc 123420000000 112212000000
495 알자
496 알자
497 알자 enc 123420000000 112212000000
498 알자
499 알자 enc 123420000000 112212000000
500 알자

```

3

Korea Advanced Institute of Science and Technology

카이스트 어휘망의 특징

□ **종합적 체계**

- ◆ 명사가 먼저 구축된 다음, 명사 어휘망을 활용하여 동사와 형용사 단어망을 만들(하나의 개념체계 안에 명사/동사/형용사 연결)
- ◆ 동사와 형용사 어휘망에서는 격들과 명사 어휘망까지 연결하여 명사, 동사, 형용사, 격들을 아우르는 입체적이고 종합적인 체계로 구축

□ **코퍼스 기반**

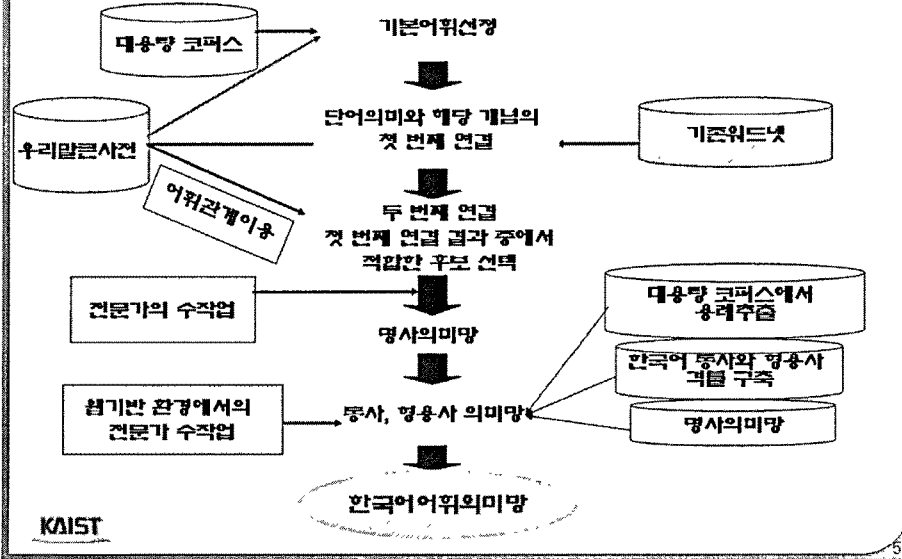
- ◆ 한국어 어휘망은 어휘 선정에서부터 서술어 논항 목록까지 철저하게 코퍼스를 기반으로 구축되었다는 특징을 가진다

□ **다국어로의 확장 가능성**

- ◆ 한-중-일 다국어 개념체계 공유
- ◆ 한국어 단어망은 한국어를 중심으로 한/중/일/영 다국어 어휘망으로 확장(중국어는 첫 번째 버전 나눔)
- ◆ 뜻풀이 부여
 - 구축된 모든 어휘에 우리말본사전의 뜻풀이 부여

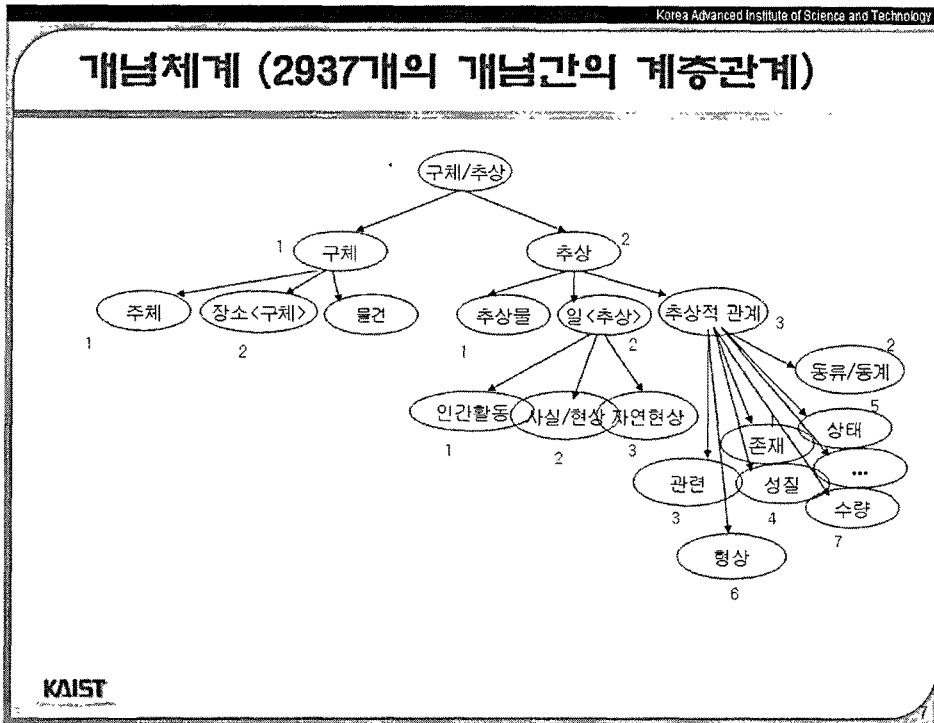
4

카이스트 어휘망의 구축과정



어휘망의 구성

- 2937개 개념 노드 (한-중-일)
- 한국어
 - ◆ 21368개 명사 (51172개 의미)
 - ◆ 1758개 동사 (5290개 의미)
 - ◆ 813개 형용사 (2081개 의미)
 - ◆ 989개 동사 격틀 연결
 - ◆ 1289개 형용사 격틀 연결
- 중국어
 - ◆ 20647개 명사 (30601개 의미)
 - ◆ 288개 동사 (911개 의미)
 - ◆ 80개 형용사 (129개 의미)
 - ◆ 1205개 동사 격틀 연결
 - ◆ 205개 동사 격틀 연결



Korea Advanced Institute of Science and Technology

명사의의미망

- 개념 노드와 단어 의미와의 연결
- 고빈도 단어 조사하여 목록 선정
- 자동으로 해당 개념번호 후보 제시
 - ◆ 기존 일본어 명사 어휘체계와 한국어 단일어 사전 이용
- 전문가의 수동 선택
- 한글학회 우리말큰사전 기준으로 의미 구분

KAIST

동사, 형용사 의미망

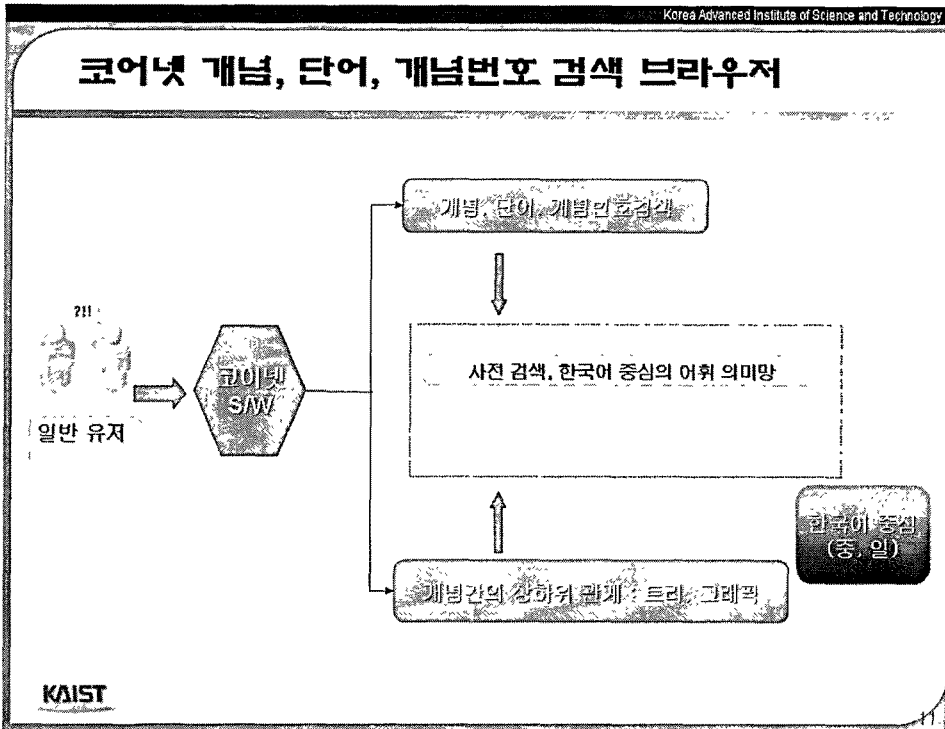
□작업 설명

- ◆ 동사와 형용사는 문장에서의 역할이 명사와 다르며, 의미 또한 동반하는 논항의 의미적, 통사적 특성에 의해 분화되므로 논항과 함께 고려하여야 한다.
- ◆ 한국어 동사 격을 있으며 형용사 격도 구축할 계획이므로 동사 자체의 의미를 체계화한다면 격들과 연동할 수 있을 것이다.
- ◆ “서술어가 상태와 동작을 기술하기 때문에 이를 실제화 하면 명사적 개념으로 표현 가능” (Ikehara S. et al., 1997)

KAIST

KAIST

코어넷 브라우저 데모



Korea Advanced Institute of Science and Technology

브라우저를 통한 카이스트 어휘망의 검색

□ 코어넷 브라우저의 기능

- ◆ 단어, 개념명, 코어넷 번호 검색
 - 상하위 관계 및 관련 어휘를 보여주는 그래픽 에디터 지원
 - 개념체계를 나무구조로 보여주는 에디터 지원
 - 우리말 큰사전과의 매핑 지원 (어휘번호 및 의미번호 매핑)
- ◆ 동사와 형용사 격을 검색
 - 동사의 개념번호에 따른 논항의 개념명 및 논항의 단어 검색 지원

□ 브라우저 데모

KAIST

Korea Advanced Institute of Science and Technology

코어넷 시스템(CoreNet System) 1

코어넷 메인

제법 트리창

그래프 창

검색 부분

검색 결과창

KAIST

13

Korea Advanced Institute of Science and Technology

코어넷 시스템(CoreNet System) 2

코어넷 단어 검색 (예: 사과)

제법 트리창

그래프 창

번호	의미번호	의미명	의미종류	의미원	의미어
1	0	사과 (사과)	명사	사과	사과
2	0	사과 (사과)	명사	사과	사과
3	0	사과 (사과)	명사	사과	사과
4	0	사과 (사과)	명사	사과	사과
5	0	사과 (사과)	명사	사과	사과
6	1	사과 (사과)	명사	사과	사과

검색 부분

검색 결과창

제법명에 해당되는
어휘어, 중국어 어휘

KAIST

사전검색 : 풀력

14

Korea Advanced Institute of Science and Technology

코어넷 시스템(CoreNet System) 3

개념 트리창

코어넷 단어 검색 (예 : 사과)

- 1 개념 트리 탐색이 가능
- 2 해당 개념의 여휘를 결과창에서 확인 가능
- 3 개념 상호간의 관계를 트리형태로 확인 가능

상위 노드

중위 노드

15

Korea Advanced Institute of Science and Technology

코어넷 시스템(CoreNet System) 4

그래프 창

코어넷 단어 검색 (예 : 사과)

- 1 기준노드 - 노란색
- 2 부모노드 - 주황색
- 3 형제노드 - 녹색
- 4 자식노드 - 파란색
- 5 관련어휘 - 회색(한국어), 분홍색(중국어)

16

코어넷 시스템(CoreNet System) 5

“에브다” 가질 수 있는 문형과 해당 논항

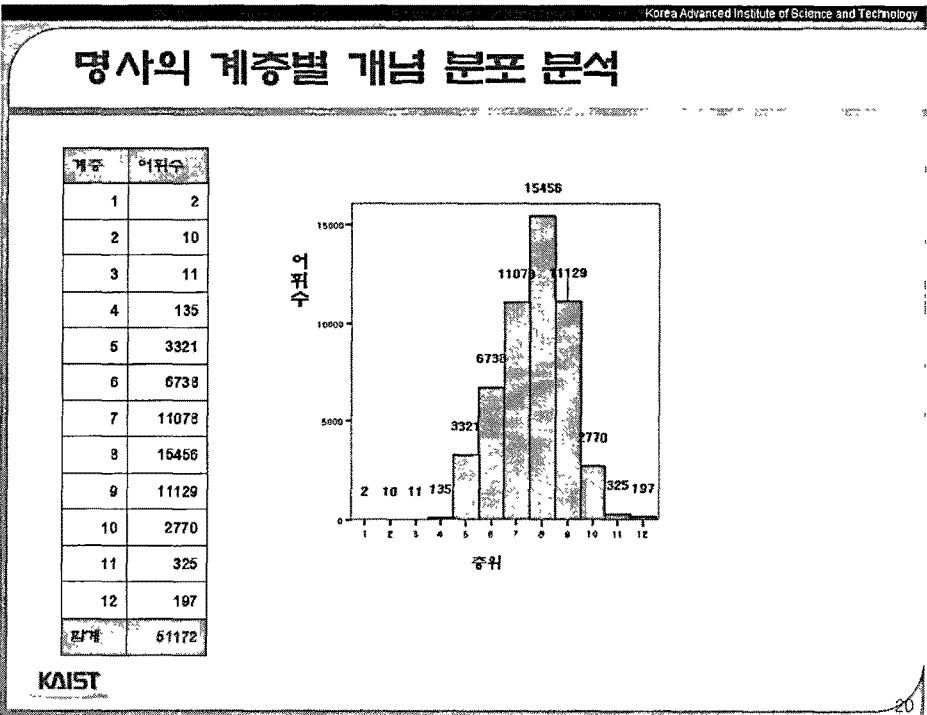
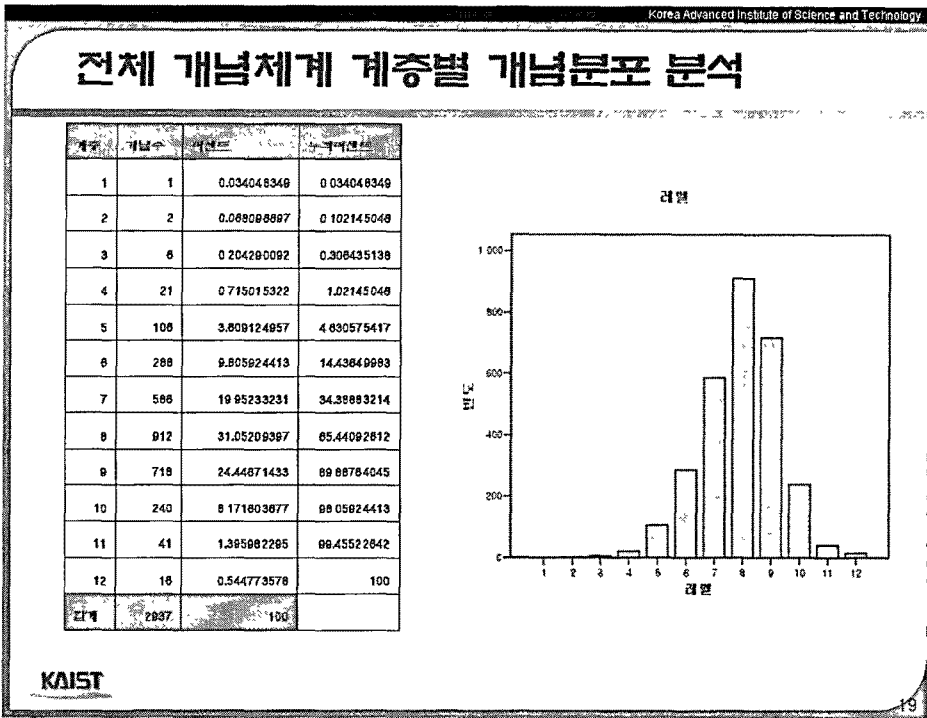
문형사	문형번호	문형명	N1	N2
에브다	12	주어	0	0
에브다	12	목적어	0	0
에브다	12	주어	0	0
에브다	12	목적어	0	0
에브다	12	주어	0	0
에브다	12	목적어	0	0
에브다	12	주어	0	0
에브다	12	목적어	0	0
에브다	12	주어	0	0
에브다	12	목적어	0	0

KAIST

KAIST

한국어 명사 개념지도

- 계층별 전체 개념분포 분석
- 명사의 개념분포 분석
- 개념명에 따른 명사 어휘 분포

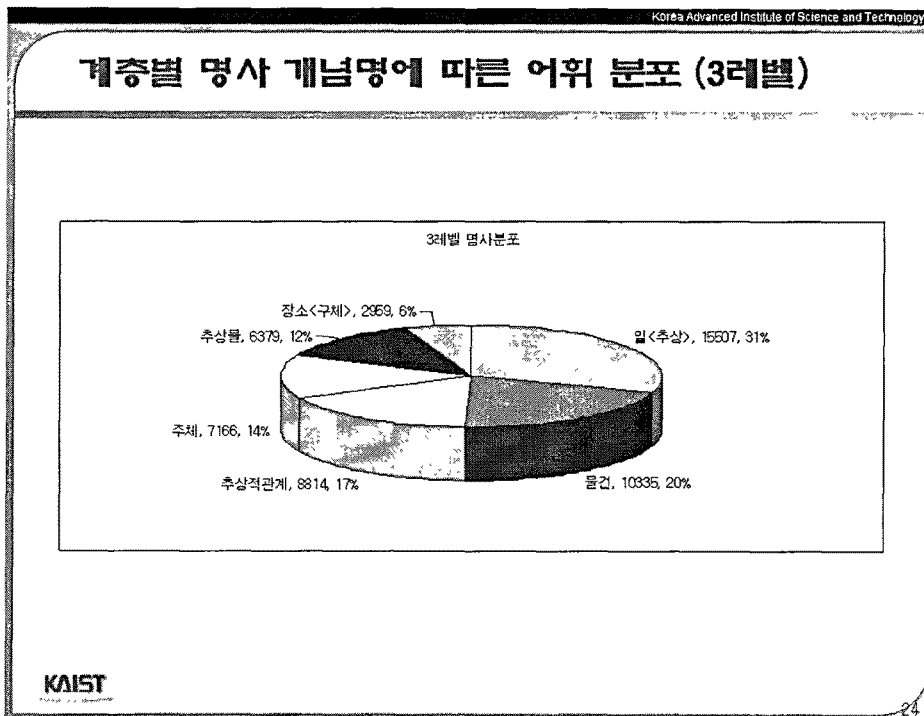
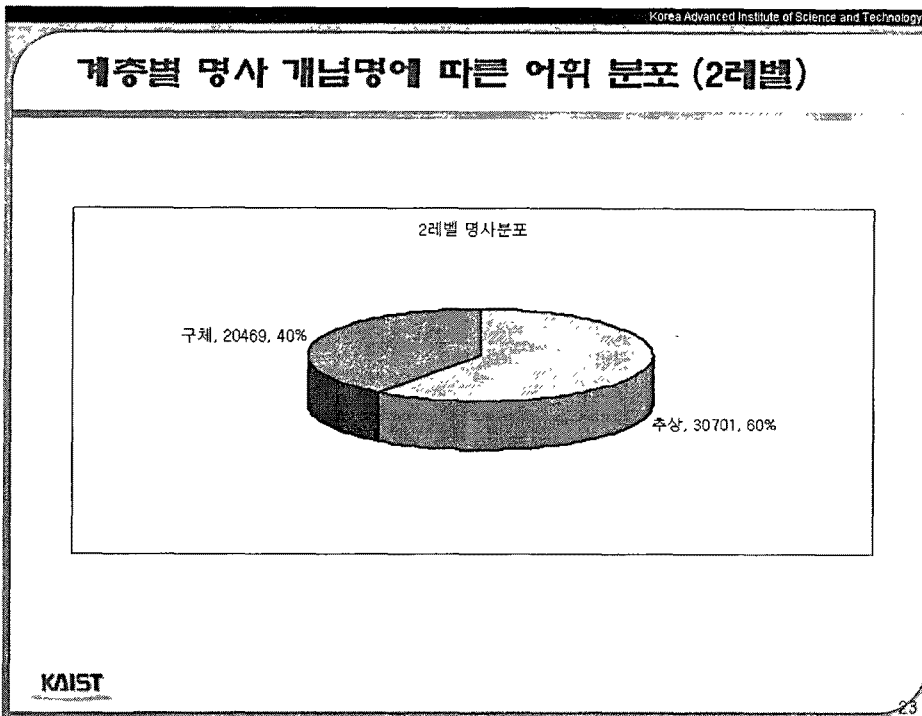


개념명에 따른 명사 어휘 분포 (상위 30위)

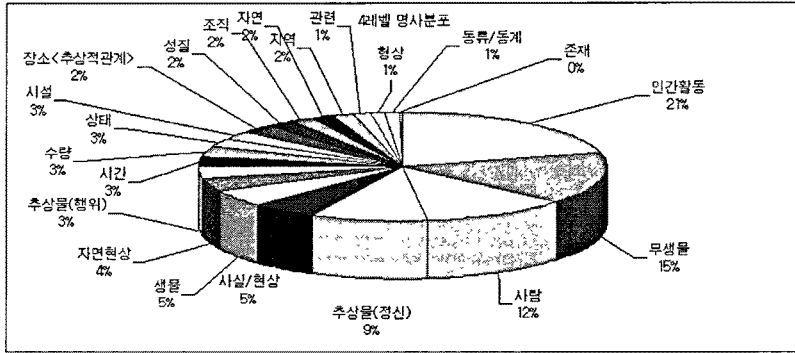
코딩번호	명사 개수	개념명	레벨	백분율	어휘
122323213	429	명의 종류	9	0.6363%	간담, 간담, 간담, 간담법, 번두름, 폐결핵, 폐결
12113121	382	인명	8	0.7465%	가명, 경, 경명, 경계, 경향, 정호, 최호
11113336	330	하는 사람	7	0.6449%	가담자, 가담자, 감독, 감독자, 감시원, 감청기, 검사
121142	317	문체류	6	0.6078%	가계부, 가계, 가축, 가력, 각호, 감독, 개성, 팔리러, 건물
122111	284	마음	5	0.5550%	가슴, 간담, 간담, 감동, 감성, 각의, 각의, 각의, 감동, 감동, 감동
12325	264	능관	5	0.5159%	각체각물, 각물, 간담, 감동, 감, 계급, 계층, 총, 총위
11322611	264	가족(논제)	8	0.5159%	가전물, 가계, 가족, 가력, 각호, 감독, 개성, 팔리러, 건물
111132A	260	사람<기타 지휘>	7	0.5081%	가관, 감찰, 감찰, 감독, 감시, 감청, 고문, 고문관, 감사관, 공인, 관철사
111131224	254	기타 관료	9	0.4964%	감사, 감사장, 경, 장관, 고관, 공판, 관직, 교장
1239221	245	기관(지휘/인간활동/etc)	7	0.4789%	간발, 개재, 개회기, 개회기, 청년기, 공소시효, 과도기
11111212	244	여자	8	0.4769%	계보, 계단, 계급, 계집녀, 계집애, 공주, 과부, 관기, 관기, 국모, 공녀, 공인, 귀부
121145	233	책(내용)	6	0.4553%	가이드, 가정집, 간행물, 공선, 교과서, 교본, 교편, 구약
11222	230	지역(인간활동)	5	0.4495%	각차, 액자, 거주, 거점, 거주, 거주지, 거주, 감문소, 감시장
1131111	227	포유류	7	0.4436%	감미사, 감, 감마, 감부, 감주마, 감찰원, 교활라, 교양도구, 교양미
113227222	225	출판물	9	0.4397%	가이드, 가정집, 경전, 고서, 고전, 고집, 관모, 교본, 교지, 국서, 국서, 그림책
12373	220	단위	5	0.4299%	가구, 가력, 가리, 가마, 가마나, 가호, 각, 각, 갈래, 감, 감, 개미, 거리, 건, 건물
113221	214	물품	6	0.4182%	가스레인지, 고기통, 고물, 팔통, 팔통, 팔통, 팔통, 팔통, 팔통
1221245	203	스포츠	7	0.3967%	감독구, 개민선, 계급, 감찰, 감청선, 감선, 감청, 감마, 감도, 감주, 팔, 팔프
121113E	200	방법	7	0.3909%	감정력, 교법, 구법, 구단법, 요법, 부법, 민간요법, 발정법, 향도, 향법, 향주, 향
11113191	200	군인	8	0.3909%	감시원, 감시장, 감비법, 계집사형관, 공명, 관관, 관경, 관관, 관연성, 군사, 군사
12113133	196	말	8	0.3830%	각설, 감연대정, 감연, 감환, 경구, 경어, 계사, 교명, 교머, 공치사, 과연, 구
1132275	186	부기	7	0.3635%	개마리판, 간, 감기, 감선, 감선, 감연, 감기통, 감포면, 관공, 기관총, 기포
111121	182	진술	6	0.3575%	관음모양, 구제주, 구주, 군신, 군신, 군신, 대장군, 도신, 독신, 동신
11322421	181	역록(논제)	8	0.3537%	가사, 감찰, 감찰, 감찰, 감주, 감주, 감주, 감주, 감주, 감주, 감주, 감주
1132231	179	약물류(의약품)	7	0.3499%	각삼제, 감주, 감삼제, 감삼제, 구단약, 구단, 구동제, 구동제, 기일, 가제,
121112	178	학문분야/학과	7	0.3478%	가사, 가정과, 경영학과, 경제, 경제학, 고고학, 고대사, 고학, 과학, 발력
112132	174	기타 거주지급	6	0.3400%	가전물, 가계, 가족, 가력, 각물, 고물, 고물, 고물, 고물, 고물, 고물, 고물
1235113	169	장소<사람><지휘>	7	0.3322%	가상, 감, 감독, 감시, 감시원장, 감관장, 건설부장관, 감사장, 감찰원장
1235113	169	장소/미상	7	0.3303%	가관, 가스, 격물, 경제, 경직, 교수명단, 교육, 구검실, 군용, 공산, 긴급

카리스트 명사 어휘망의 계층별 어휘분포

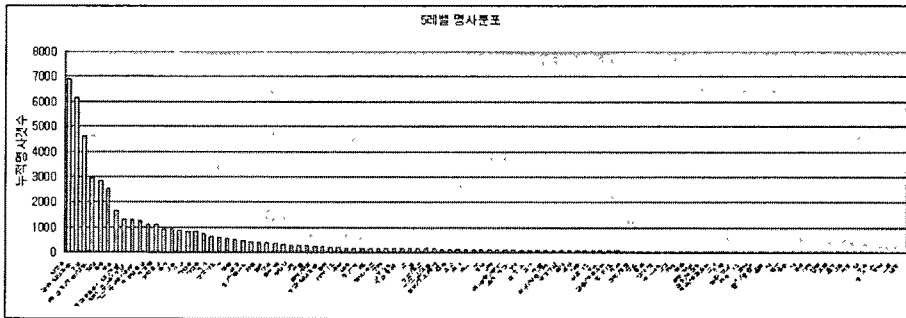
코딩번호	레벨	한국어	누적명사개수	누적명사비율
1	1	구체/주상	51172	100.00%
12	2	주상	30701	60.00%
11	2	구체	20469	40.00%
122	3	입<주상>	15507	30.30%
113	3	출간	10335	20.20%
123	3	주상적관계	8814	17.22%
111	3	주제	7165	14.00%
121	3	주상들	6379	12.47%
112	3	장소<구체>	2959	5.78%
1221	4	인간활동	10785	21.08%
1132	4	부상들	7769	15.18%
1111	4	사람	6092	11.90%
1211	4	주상들(정신)	4648	9.08%
1222	4	사실/현상	2671	5.22%
1131	4	상들	2563	5.01%
1223	4	자연현상	2050	4.01%
1212	4	주상들(행위)	1731	3.38%
1239	4	시간	1625	3.18%
1237	4	수항	1522	2.97%
1235	4	상태	1330	2.60%
1121	4	시점	1320	2.58%
1238	4	장소<주상적관계>	1217	2.38%
1234	4	상들	1067	2.09%
1112	4	조직	1067	2.09%
1123	4	자연	833	1.63%
1122	4	지역	806	1.58%
1233	4	관련	737	1.44%
1236	4	항상	570	1.11%
1231	4	동위/동계	542	1.06%
1231	4	현재	203	0.40%



계층별 명사 개념명에 따른 어휘 분포 (4레벨)



계층별 명사 개념명에 따른 어휘 분포 (5레벨)



맺음말

□ 계층별 전체 분포 및 명사 분포 8>9>7>6... 레벨 순

□ 명사의 레벨별 어휘 분포

- ◆ 2레벨: 추상(12)>구상(11)
- ◆ 3레벨: 일<추상>(122) > 물건(113)> 추상적 관계(123)> 주체(111)> 추상물(121)> 장소<구체>(112)
- ◆ 4레벨: 인간활동(1221) > 무생물(1132) > 사람(1111) > 추상물(정신)(1211) > 사실/현상(1222)> 생물(1131)> 자연현상(1223)> 추상물(행위)(1212)> 시간(1239)> 수량(1237)...
- ◆ 5레벨: 인공물 (11322)> 행위<인간활동>(12212)> 정신(12211)> 사람<직업/지위/역할>(11113)> 인간(11111)> 변동(12222)> 동물(11311)> 지적생산물(사고/학습)(12111)> 언어<추상물(정신)>(12113)> 제도<추상물(행위)>(12121)...

현 명사 의미체계의 문제점 및 보완점

□ 개념명과 어휘관계 미흡

- ◆ 개념명 자체의 수정
- ◆ 개념과 어휘와의 관계를 고려하여 수정
- ◆ 하위 개념 노드의 분화 미흡

□ 반의어 처리 미흡

□ 어휘 수의 보완 필요

향후 과제

- 카이스트 어휘망 2차 버전-개념명과 어휘 수정
- 영어판의 확장 및 중국어판 보완
- 동사 및 형용사의 어휘 보완
- 카이스트 기본어휘 등급 부여
- 카이스트 어휘망의 활용 가능성 및 사례 연구