



인공지능에서의 언어 처리: 자연어처리

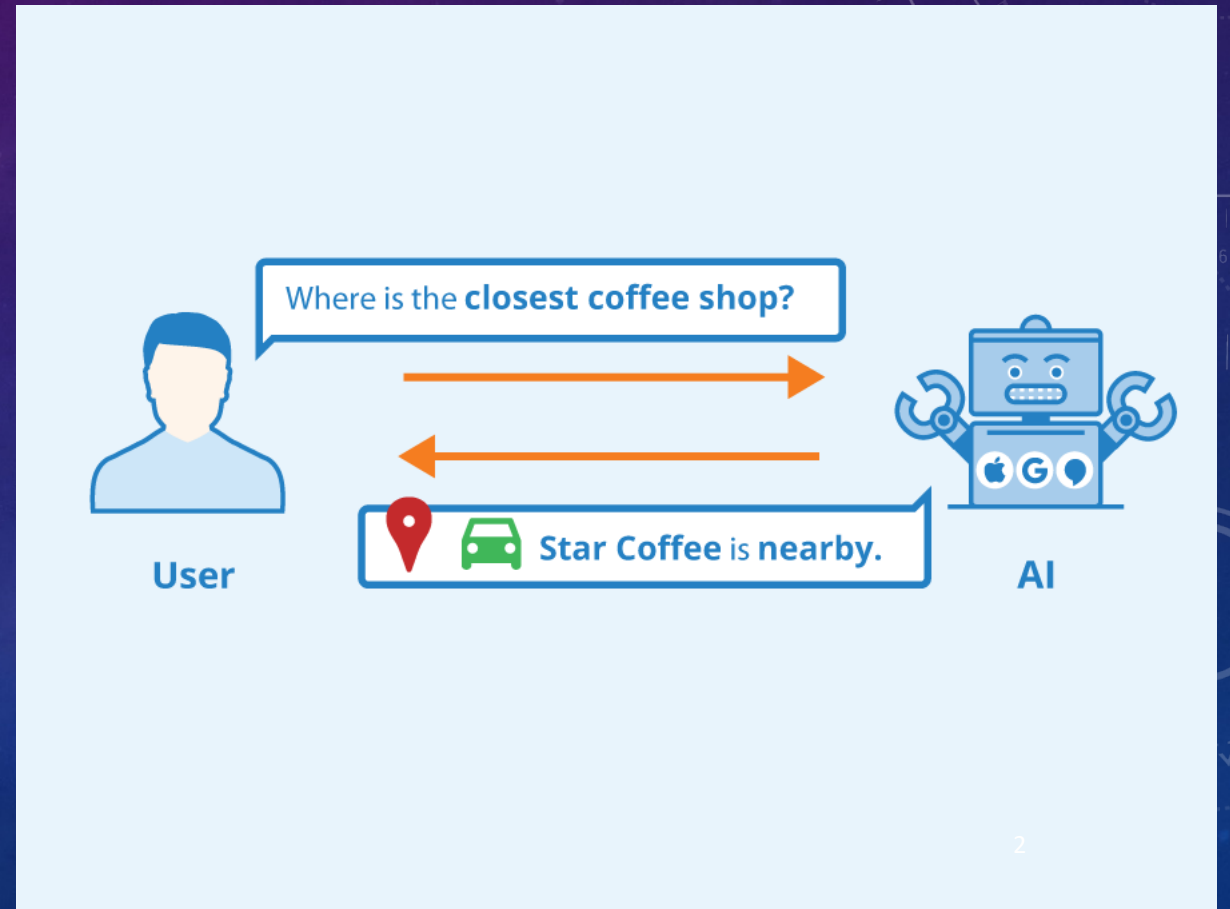
서울대학교 언어학과

신호필

HPSHIN@SNU.AC.KR

01. 자연어처리란?

- 자연어처리(Natural Language Processing)
또는컴퓨터언어학(Computational Linguistics)
 - *Natural Language Processing*, or NLP, is the sub-field of AI that is focused on enabling computers to understand and process human languages.
 - https://www.seobility.net/en/wiki/Natural_Language_Processing
- 인간의 언어와 관련되는 여러 분야-언어학, 컴퓨터학, 전기공학, 인지과학, 심리학, 통계학 등-에 걸치는 학제적인 분야



01. 자연어처리란?

- 컴퓨터는 인간의 언어를 이해할 수 있는가?
 - 인간의 언어를 이해할 수 있도록 많은 노력의 결과로 상당한 진전
 - 그러나 여전히 한계
 - 인간이 언어를 인식하는 것과 완전히 동일하지는 않지만 상당한 수준의 진전을 보임
 - 언어를 이해하기 위한 여러 Tool들이 존재함

02. 자연어처리의 응용 분야는?

- 텍스트분류
 - Spam detection
 - Sentiment Analysis
- 기계번역
- 정보검색
- 챗봇
- 문서요약/문서생성
- 질의-응답시스템

02. 자연어처리의 응용 분야

- Question Answering: IBM's Watson
 - 2011년 2월 우승
- WILLIAM WILKINSON'S "AN ACCOUNT OF THE PRINCIPALITIES OF WALLACHIA AND MOLDOVIA" INSPIRED THIS AUTHOR'S MOST FAMOUS NOVEL
 - → Bram Stoker



02. 자연어처리 응용 분야

- Information Extraction
 - 행사: 창의교육프로젝트설명회
 - 날짜: 2021년 11월 24일
 - 시작: 12시 30분
 - 끝: 13:30분
 - 장소:기초교육원 320호

전체 답장 | 삭제 | 보관 | 이동 | 업무로 등록 | ...

보낸 사람: 교무과 (11.23 13:25)
받는 사람: 신효필님

2021학년도부터 추진 중인 <Inno-Edu: 2031 서울대 창의교육프로젝트> 사업 설명회를 개최하오니
학내 구성원들의 많은 참석을 부탁드립니다.

○ 일시: 2021. 11. 24.(수) 12:30~13:30
○ 장소: 기초교육원(61동) 320호
* 코로나19 기본 방역수칙을 준수하여 100명 미만으로 참석 인원을 제한하고, 행사장 내 취식을 금지합니다.
* 행사 종료 후 샌드위치를 제공할 예정입니다.

Inno-Edu 2031: 서울대 창의교육프로젝트 설명회 개최

일시 | 2021. 11. 24.(수) 12:30~13:30
장소 | 기초교육원(61동) 320호
주관 | 교무처 교무과
문의 | 02-880-2077, sekim33@snu.ac.kr
* 현장 참석 인원은 100명 미만으로 제한되며, 샌드위치를 제공할 예정입니다.

본 메일은 서울대학교 대량메일시스템을 통해 발송된 메일입니다. 메일 수신을 원치 않으시면 [수신거부](#) 를 클릭하십시오.
본인의 수신거부 목록을 확인하려면 [수신거부목록 확인](#) 을 클릭하십시오.
This email has been sent through the SNU mass mailing system.
If you do not wish to receive emails of this category, please click [here](#).
To check the emailing lists you have unsubscribed from, please click [here](#).

서울대학교
SEOUL NATIONAL UNIVERSITY

08826 서울시 관악구 관악로 1 서울대학교
1 Gwanak-ro, Gwanak-gu, Seoul 08826, Korea
Copyright 2012 Seoul National University All Rights Reserved

02. 자연어처리 응용 분야

- Sentiment Analysis
 - Attributes: 가격, 상품상태, 매직스페이스

★★★★★ 5 신세계물 · we***** · 19.11.27.

지정일에 배송됐고 설치도 완벽◆

지정일에 배송됐고 설치도 완벽했습니다. 1인 가구이지만 수많은 날의 고민 끝에 이 모델로 결정했습니다. 양문형이나 4도어냐의 고민에서부터 세미빌트인이나 빌트인이나 까지... 결국 저의 모든 것을 만족시켜주는 건 lg 디오스 양문형 냉장고더라구요^^ 저는 1인가구라 냉동실 사용이 주를 이룰 수 밖에 없고 냉장고의 주요 사용 용도는 마실 것의 보관이었기에 냉동실 사용이 편리한 양문형으로 했고 음료 보관 및 이용이 편리한 매직스페이스가 있는 제품으로 선택하였습니다! 용량에 대한 고민도 많았는데.. 1인 가구가 쓰기에 800리터대의 제품은 너무 낭비이지 않을까 싶었지만.. 냉장고와 티비는 거거익선이라는 말은 정말 빕인인 것 같습니다~ 작아서 불편한 것 보단 넉넉히 사용하는 게 훨씬 나은 것 같습니다~ 냉장고를 딱딱 채워서 쓰면 오히려 냉장고가 제 기능을 발휘하기 힘들다고 하더라구요~ 솔직히 이 정도 크기와 성능에 이름 있는 브랜드 제품을 리뷰펼치기 ▾

★★★★★ 5 11번가 · pj***** · 20.08.09.

최고예요

가격이 좋아요 2등급이고색상시 고급지고 튼튼해보이는 스타일 매직스페이스가있어 음료수 물 편하게 이용하고 에너지절감도 되고부오밍게 선물했는데 아주 만족 좋아합니다다만 문쪽 선반이 작은거라도 맨아랫까지있음 더알찰텐데 조금아습네요냉동실 얼음얼리는곳시 크진않지만 두분이 쓰기만좋을만큼되네요 감추~ 리뷰펼치기 ^

★★★★★ 5 신세계물 · cc***** · 19.10.11.

친정엄마 냉장고가 10년 다되어◆

친정엄마 냉장고가 10년 다되어가는 기존 양문형 냉장고였는데 고장은 없으니 냉동실이 좀아 냉기순환이 안되는지 공공 얼지않아 선물 드렸는데 정말 마음에 들어하시네요 4도어 5도어가 대세이긴 하나 기능대비 가격 차이가 너무 많이 나서 선택했는데 대만족입니다 양문형이긴 하나 외관상으론 5도어처럼 보입니다 ㅋㅋ 특히 매직스페이스가 정말 좋아요 예전 홈바보다 사이즈도 크고 냉기보존도 잘 되는듯요 리뷰펼치기 ▾

★★★★★ 5 신세계물 · sa***** · 20.05.09.

*배송 배송은 엄청 빨라요. 전국

*배송

배송은 엄청 빨라요. 전국품질이라 오래 걸릴 줄 알았는데 5일 이내에 발송되었습니다. 배송기사님들 2명 오셔서 설치해주고 가셨어요. 폐가전수거는 하루 전이나 당일 아침에 전화하시면 이야기하시면 됩니다. 폐가전수거도 깔끔하게 해주셨어요.

*상품

상품상태는 아주 좋아요. 다른 사이트에서 살까하다가 ssg를 선택한 것은 경품, 서비스 면에서 신뢰가 갔기 때문이에요. 인증샷에서 배송된다고 해서 더욱 만족해요. 리뷰펼치기 ▾

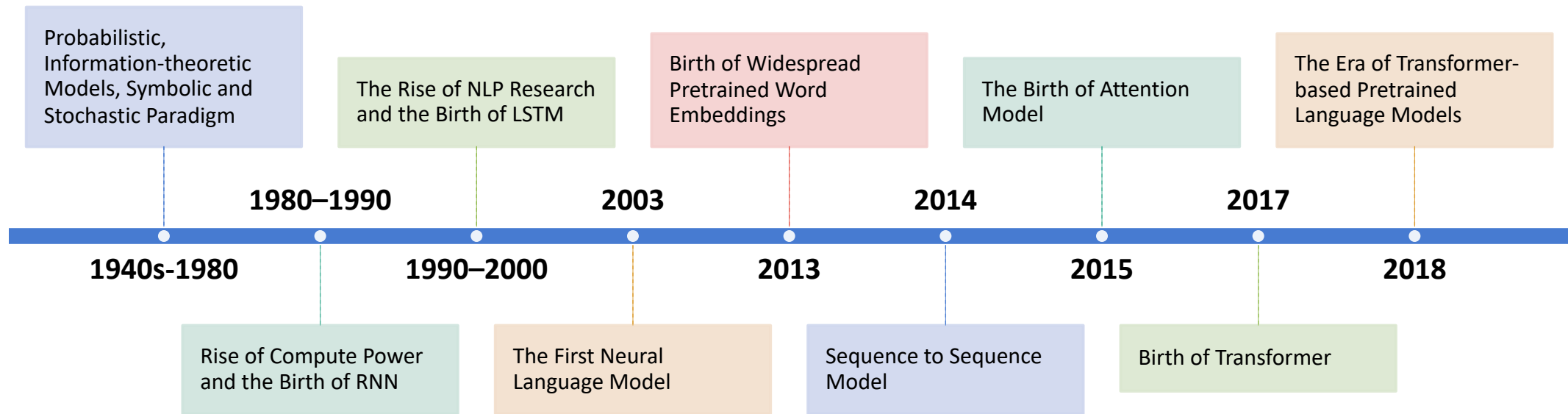
02. 자연어처리 응용 분야

- Machine Translation

The screenshot shows the Papago web interface. The input text is: "올해 주택분 종합부동산세 고지인원이 94만7000명이라고 기획재정부가 22일 밝혔다. 지난해(66만7000명)보다 42% 늘었다." The output text is: "The Ministry of Strategy and Finance said on the 22nd that the number of notices for the comprehensive real estate tax for housing this year was 947,000. It increased 42% from last year (667,000)." Below the translation, a glossary lists key terms with their English definitions and examples.

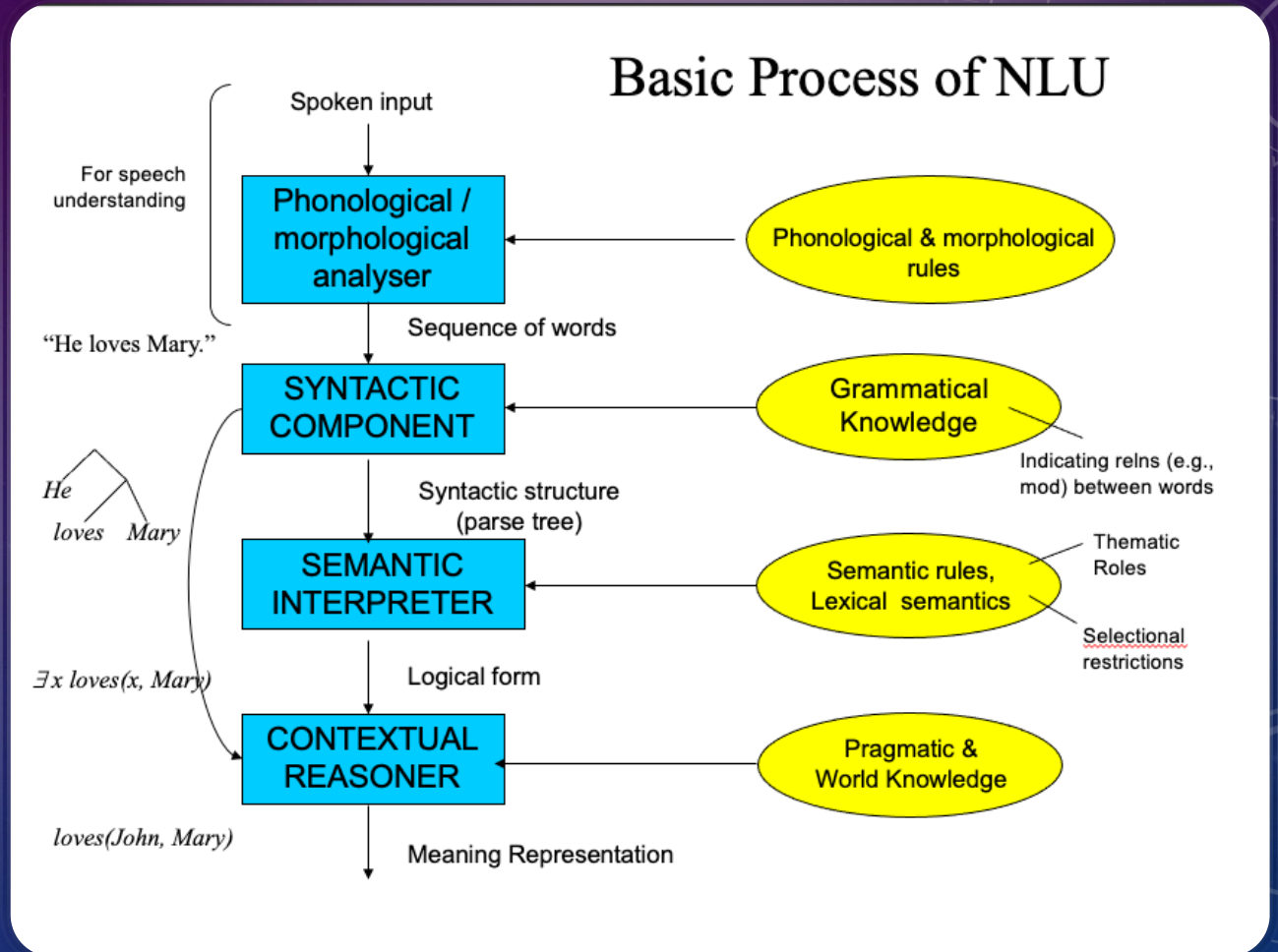
Term	Definition
고지 [高地]	1. (땅보다 높은 땅) highlands, heights, high ground, uplands; (넓고 평평한) tableland 2. (목표) goal 3. (유리한 조건) 처지
올해	1. this year, the present(current) year
주택 [住宅]	1. house, housing
종합 [綜合]	1. synthesize, put together, piece together
고지 [告知]	1. [명사] (알림) notice, notification, [동사] notify (sb of sth), inform
인원 [人員]	1. the number of people[persons]
기획 [企劃]	1. [명사] plan, planning, project [동사] plan, design
AND	1. [컴퓨터] 앤드, 논리곱 ((논리곱을 만드는 논리 연산자(論理演算子); cf. OR))
strategy 미국·영국, US·UK[ˈstrætədʒi]	1. (특정 목표를 위한) 계획[전략] 2. 계획[전략] 수립[집행] 3. (군사적인) 전략 (→tactic)
finance 미국·영국, US·UK[ˈfaɪnæns; ˈfɑːnæns; ˈfɛːnæns]	1. (사업·프로젝트 등의) 재원[자금] 2. (특히 정부나 기업의) 재정[재무] 3. 자금[재원]을 대다 (=fund)
ministry 미국·영국, US·UK[ˈmɪnɪstri]	1. (정부의 각) 부처 2. (집합적으로) 목사[성직자] 3. 목사[성직자]의 직책[임기]
say 미국·영국, US·UK[sɛɪ]	1. 말하다, ...라고 (말)하다 2. (특정한 어구를 반복해서) (말)하다[올조리다] 3. 발언권, 결정권 4. (놀람기쁨을 나타내어) 아[와] 5. (말을 저을 꺼낼 때) 저
notice 미국식, US[nɒtɪs]	1. 시경순 주로 알아채

03. 자연어처리는 어떻게 발전되어 왔나?



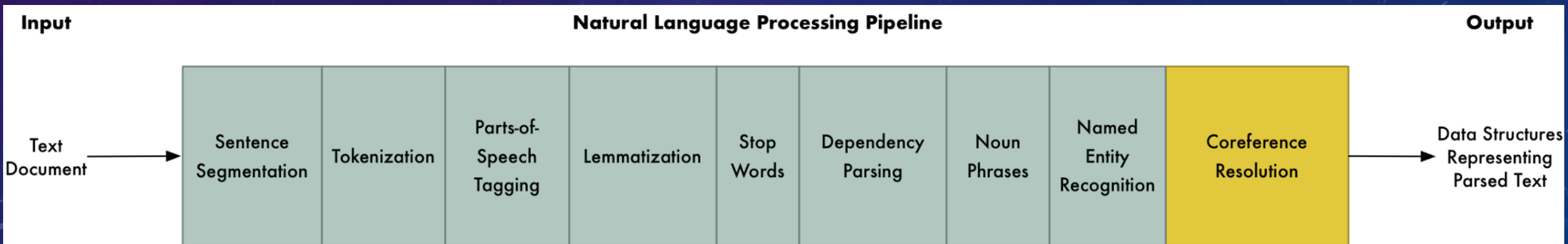
04. 자연어처리는 어떻게 이루어지나

- Top Down, Linguistic Analysis



04. 자연어처리는 어떻게 이루어지나

- Top Down, Linguistic Analysis
 - Image source (https://miro.medium.com/max/2000/1*zHLs87sp8R61ehUoXepWHA.png)



04. 자연어처리는 어떻게 이루어지나

- Bottom up, Data Driven, Classification

The screenshot shows the NAVER movie review page for the category 'Netizen Recommended Movies'. The page features a sidebar with navigation options like 'Movie Rankings' and 'Reviews'. The main content area displays a table of movies with their titles, ratings (represented by stars), and brief user comments. A pagination bar at the bottom indicates 13 pages of results, with the current page being 9. A search bar is also visible at the bottom of the page.

번호	강상평	글쓴이/날짜
17819835	둘 ★★★★★ 10 끝나가는게 너무 아쉽습니다. 치밀하게 워하나 부족하게 없는 영화.. 음향에 압도당하고 오래색 영상에 빠지고 살라메의 눈빛과 얼굴과 목소리에 취해.. 신고	lko5**** 21.11.27
17819834	너는 내 운명 ★★★★★ 8 아 눈 겁나부었어... 신고	alsq**** 21.11.27
17819833	메이드 인 이태리 ★★★★★ 8 아름다운 토스카나에서 평화로운 이야기에 힐링을 신고	sson**** 21.11.27
17819832	유체이탈자 ★★★★★ 10 막 재미있어 지려고 했는데 갑자기 끝난느낌 소재가 신선했고 액션연기 불만했음 넷플릭스 장편드라마로 보면 재미있을것같은 소재 단 러브스토리는 살짝만 가미되면 좋을듯 신고	minh**** 21.11.27
17819831	강릉 ★★★★★ 7 개연성은 별로없는데 그냥.. 재미음으로는 불만함 신고	yjsy**** 21.11.27
17819830	루카 ★★★★★ 10 가족이랑같이봤는데정말재밌고감동적이었어요 신고	dooi**** 21.11.27
17819829	현형 인 살렘: 악령의 마을 ★★★★★ 1 당신의 86분은 소중한 시간입니다 신고	adly**** 21.11.27
17819828	서터 ★★★★★ 8 거북목인 사람은 보지마세요 신고	thdw**** 21.11.27
17819827	자산어보 ★★★★★ 10 아직도 바뀌지 않은 씬섬비의 나라.. 신고	wowl**** 21.11.27
17819826	베놈 2: 렛 데어 비 카니지 ★★★★★ 10 베놈기업따 ㅋㅋㅋ 그렇게 재밌진않지만 신고	coco**** 21.11.27

현재 상영작 평점보기
[현재 상영작]

전체 리스트 총 13,761,953개의 평점이 있습니다.

영화 인기검색어 더보기

- 유체이탈자 ↑ 1
- 장르만 로맨스 ↓ 1
- 이타닐스 - 0
- 연애백진 로맨스 ↑ 3
- 둘 ↓ 1

2021.11.26

네티즌 최고 평점 더보기

현재 상영영화 모든 영화

- 코다 9.24
- 기적 9.18
- 파이판 9.08
- 타오르는 여인의 초상 9.06
- 꽃다발 같은 사람들. 9.0

2021.11.26까지의 한림루 누락 평점

가장 많이 추천된 리뷰

- [메메리]-장르만 로맨스. 장르만 로맨스 acts****
- 장르만 로맨스를 보고스. 장르만 로맨스 film****
- <이타닐스> 이제는 MC. 이타닐스 demo****
- 프랜치 디스패치 (메거진의. 프랜치 디스패치 reno****
- 영화[돈 록업]메인 예고. 돈 록업 shin****

2021.11.19-2021.11.26

평점은 영화별로 1인당 1회만 등록이 가능하며 하루에 최대 3개의 영화까지만 등록이 가능합니다. 게시판 용도에 맞지 않는 글은 운영자에 의해 삭제될 수 있습니다. ?

04. 자연어처리는 어떻게 이루어지나

BOTTOM UP, DATA DRIVEN, CLASSIFICATION

https://miro.medium.com/max/1000/1*x93sfurvt_bmomlzhgeglw.png

https://miro.medium.com/max/1000/1*hv4shrjhhk4j9zw27qb_r-g.png

Input: Email Text

"Dir Sir, I have 18 Million dollars that my uncle..."

Text Classification Model

Output: Spam score

Class: "Spam"

Input: Review Text

"This used to be a giant parking lot where government..."

Text Classification Model

Output: Predicted Stars

Class: "5"

05. 자연어처리는 왜 어려운가

- Ambiguity (중의성): 언어의 모든 층위의 중의성
 - 어휘적 층위: ‘감기’?, ‘Apple’?
 - 통사구조는 의미에 영향을 미친다
 - Flying planes is dangerous
 - Flying Planes are dangerous
 - Teacher Strikes Idle Kids
 - 의미와 세상지식(World Knowledge)는 통사구조에 영향을 미친다
 - *Flying insects is dangerous
 - Flying insects are dangerous
 - I saw the Grand Canyon flying to LA
 - I saw a condor flying to LA

05. 자연어처리는 왜 어려운가

비정형데이터

- accomplished! U taught us 2
- should never give up either❤️

Segmentation Issue

- The New York-New Haven Railroad

Idioms

- Dark horse
- 손이 크다

신조어

- Unfriend/Retweet/bromance/

World Knowledge

Tricky Entity Name

- Let it Be was recorded....

06. 언어모델(LANGUAGE MODEL)

- Predict Next Word
 - $P(\text{새빨간 거짓말}) > P(\text{새빨간 희망})$
 - $P(\text{이 강의는 참 재미있어요})$
- Assign a Probability to a sentence
 - N-gram based : Unigram, Bigram, Trigram...
 - Word Embedding
 - Transformer based Language Modeling

06. 언어모델(LANGUAGE MODEL)

- 문장의 확률을 어떻게?
 - 조건부 확률: $P(B|A) = P(A, B)/P(A) \rightarrow P(A, B) = P(A)p(B|A)$
 - $P(A, B, C, D) = P(A)P(B|A)P(C|A, B)P(D|A, B, C)$
 - Chain Rule:
 - $P(X_1, X_2, X_3, \dots, X_n) = P(X_1)P(X_2 | X_1)P(X_3 | X_1, X_2) \dots P(X_n | X_1, \dots, X_{n-1})$

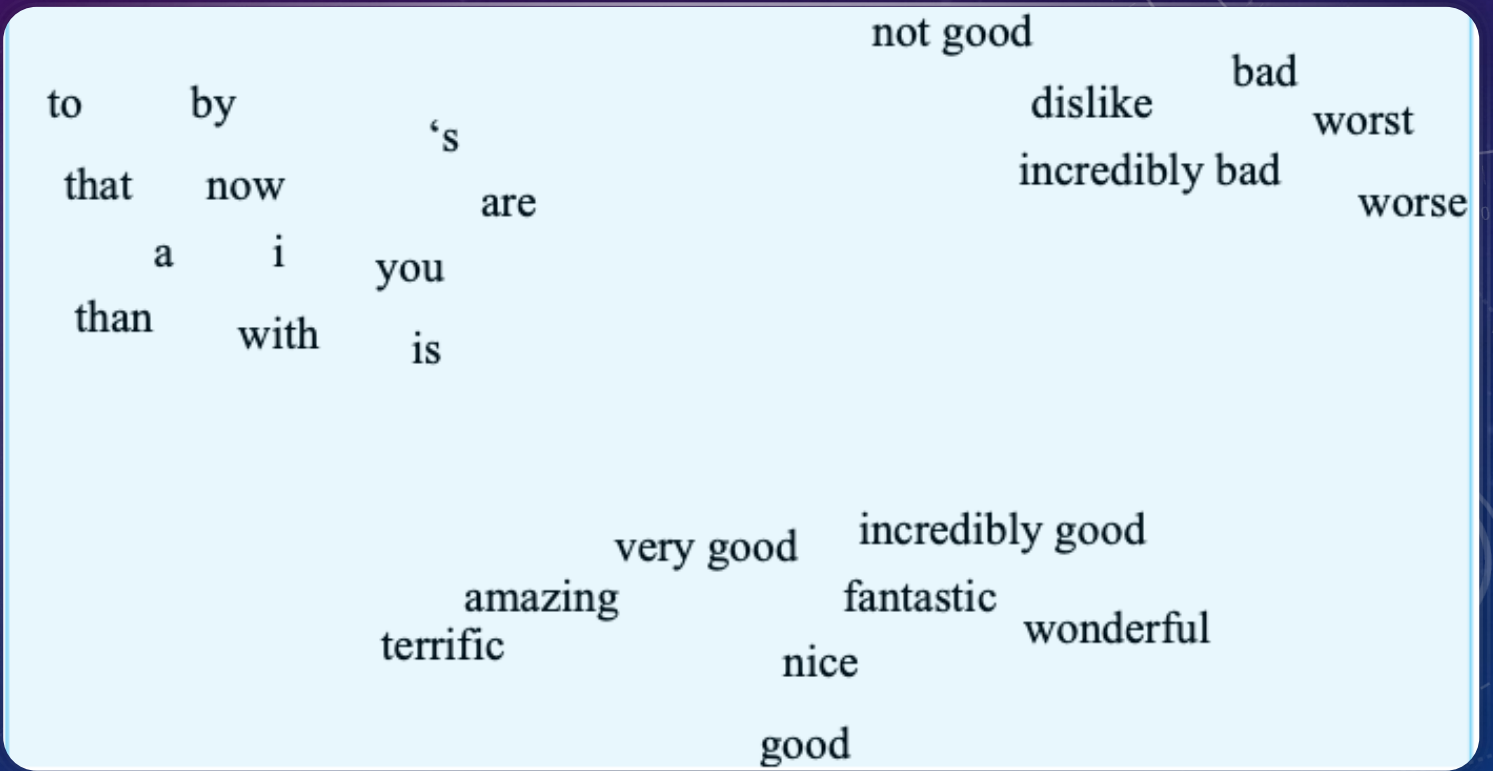
$$P(w_1 w_2 \dots w_n) = \prod_i P(w_i | w_1 w_2 \dots w_{i-1})$$

- $P(\text{“이 강의는 참 재미 있어요”}) = P(\text{이}) \times P(\text{강의는} | \text{이}) \times P(\text{참} | \text{이 강의는}) \times P(\text{재미} | \text{이 강의는 참}) \times P(\text{있어요} | \text{이 강의는 참 재미})$
 - $P(\text{있어요} | \text{이 강의는 참 재미}) = \text{Count}(\text{이 강의는 참 재미 있어요}) / \text{Count}(\text{이 강의는 참 재미})$
- Markov Assumption
 - $P(\text{있어요} | \text{이 강의는 참 재미}) \approx P(\text{있어요} | \text{재미})$ 또는 $P(\text{있어요} | \text{참 재미})$
 - $P(\text{“이 강의는 참 재미 있어요”}) = P(\text{이}) \times P(\text{강의는} | \text{이}) \times P(\text{참} | \text{강의는}) \times P(\text{재미} | \text{참}) \times P(\text{있어요} | \text{재미})$

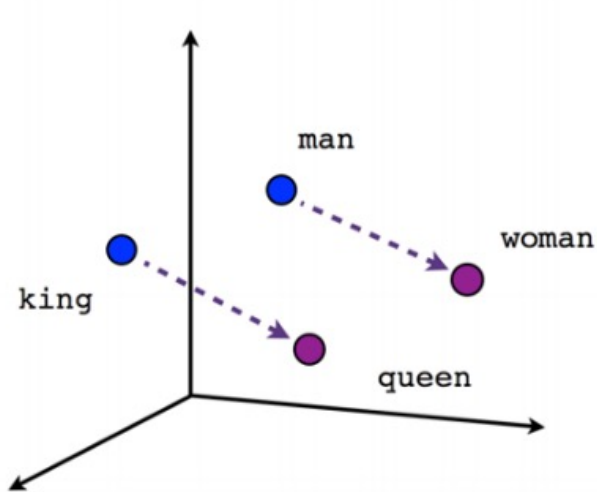
$$P(w_1 w_2 \dots w_n) \approx \prod_i P(w_i | w_{i-k} \dots w_{i-1})$$

07. 단어 임베딩(WORD EMBEDDING)

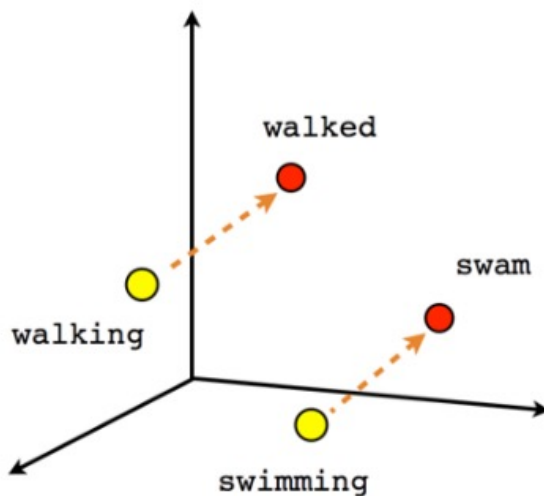
- Numericalization : 텍스트를 숫자로 바꾸는 방법
 - One-hot encoding, N-Gram 언어 모델, Bag of Words 언어 모델...
- Vector 의미론(Semantics)
 - 단어의 의미를 어떻게 표상할 수 있는가?
 - 동의어, 유사어, 다의어, 관련어...
- Word Embedding
 - Distributional Hypothesis
 - Word2Vec, Glove, FastText..
 - Static Word Embedding (다의어, 동의어 구별이 되지 않음)



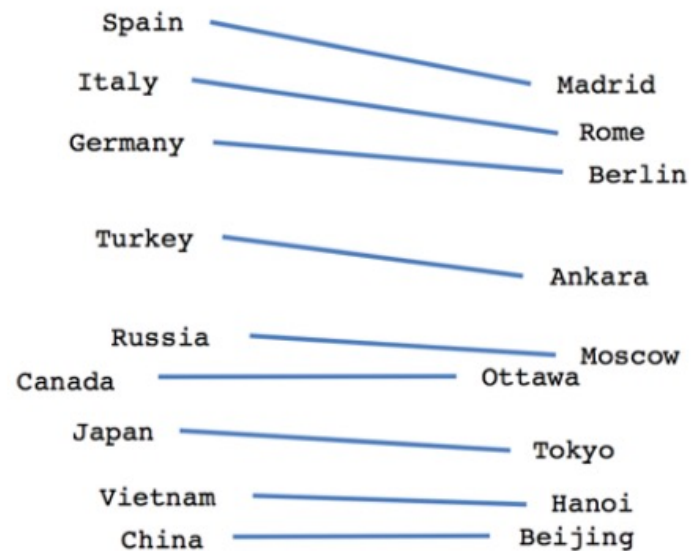
07. 단어 임베딩



Male-Female



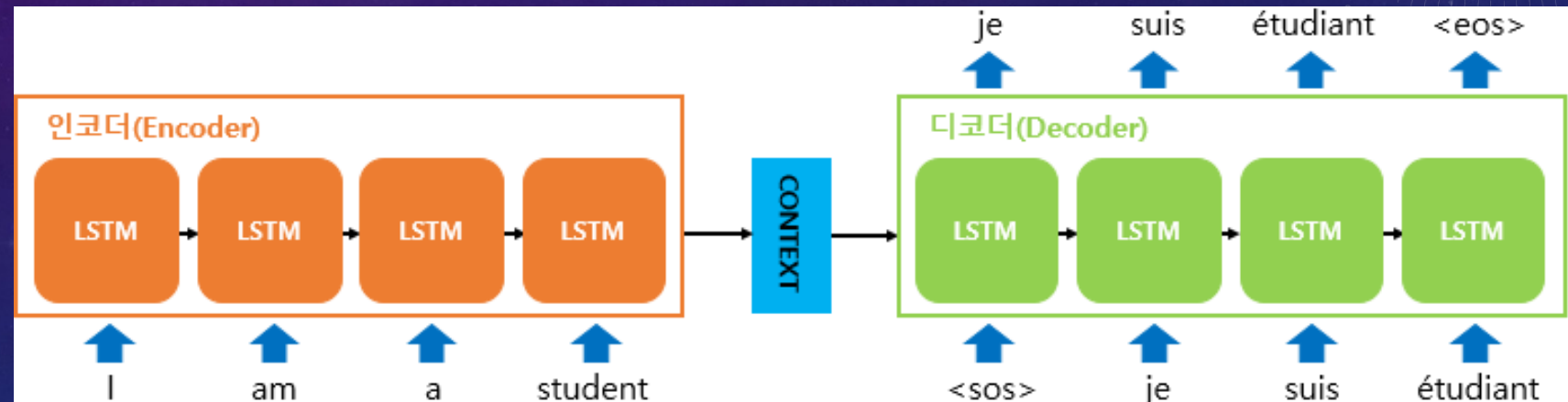
Verb tense

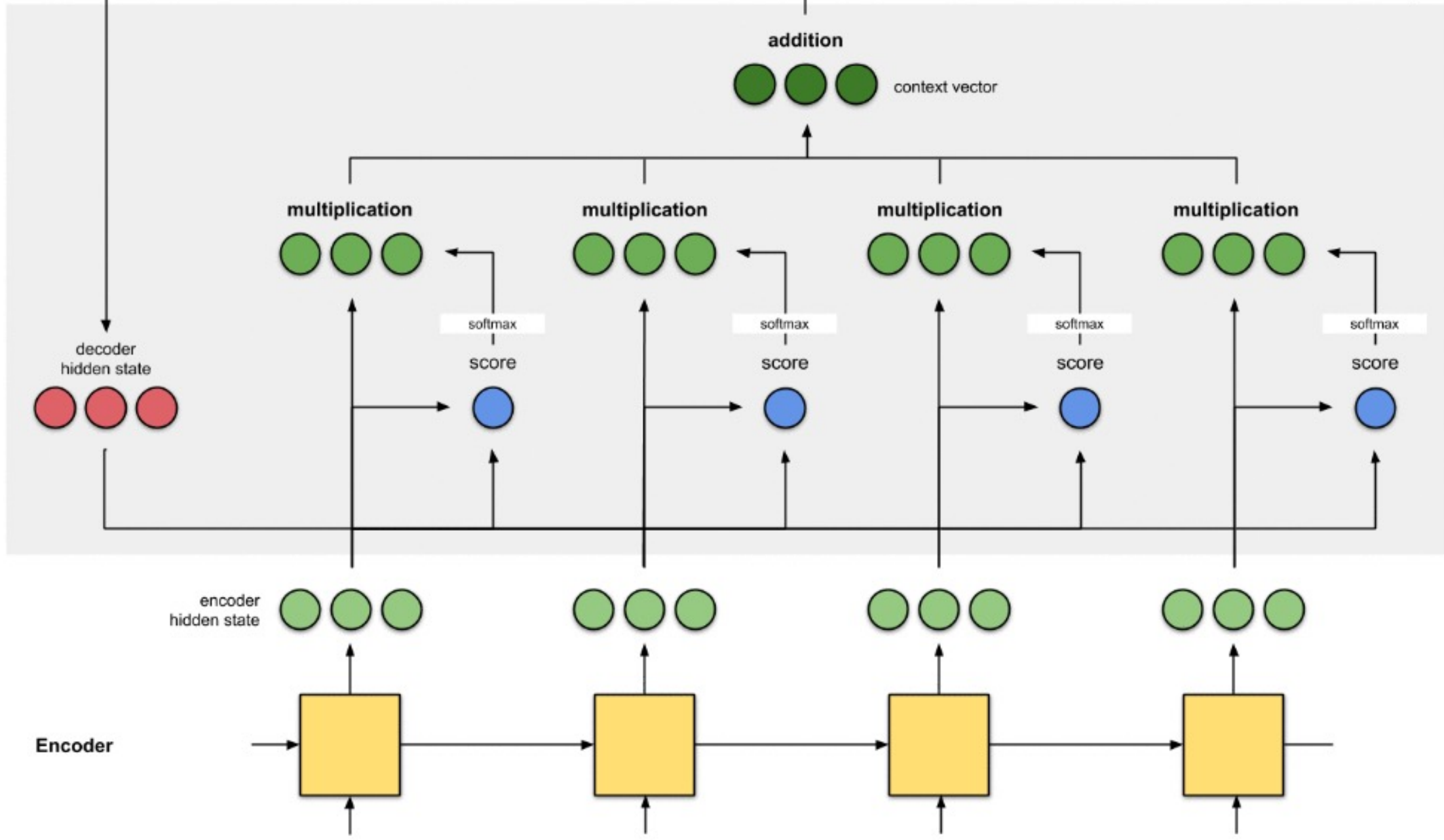


Country-Capital

08. 시퀀스투시퀀스 /어텐션 (SEQUENCE TO SEQUENCE/ATTENTION)

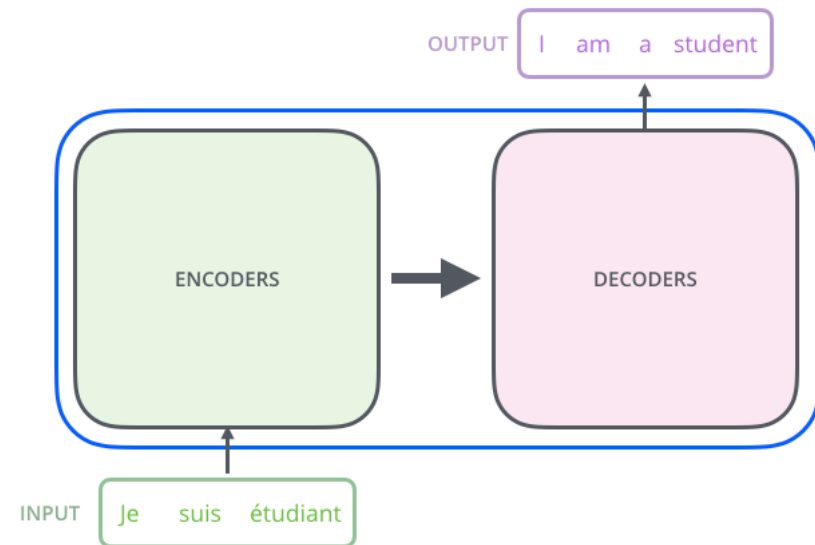
- Sequence to sequence Model
- Encoder-Decoder





09. 트랜스포머

- [https://jalammar.github.io/images/t/the transformer 3.png](https://jalammar.github.io/images/t/the%20transformer%203.png)
- [https://jalammar.github.io/images/t/The transformer encoders decoders.png](https://jalammar.github.io/images/t/The%20transformer%20encoders%20decoders.png)



09. 트랜스포머

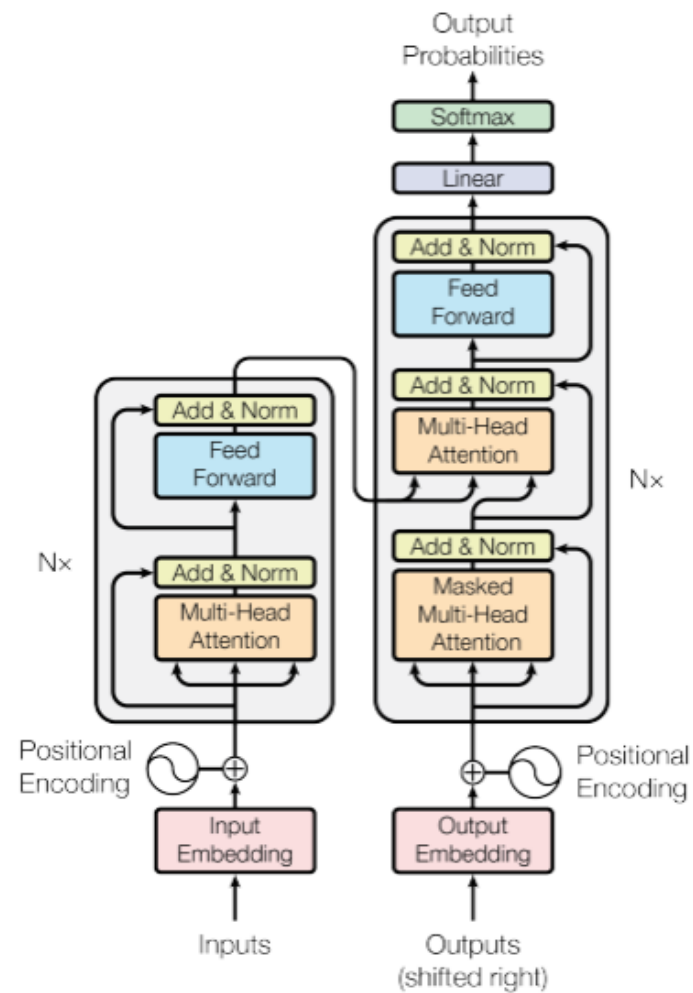
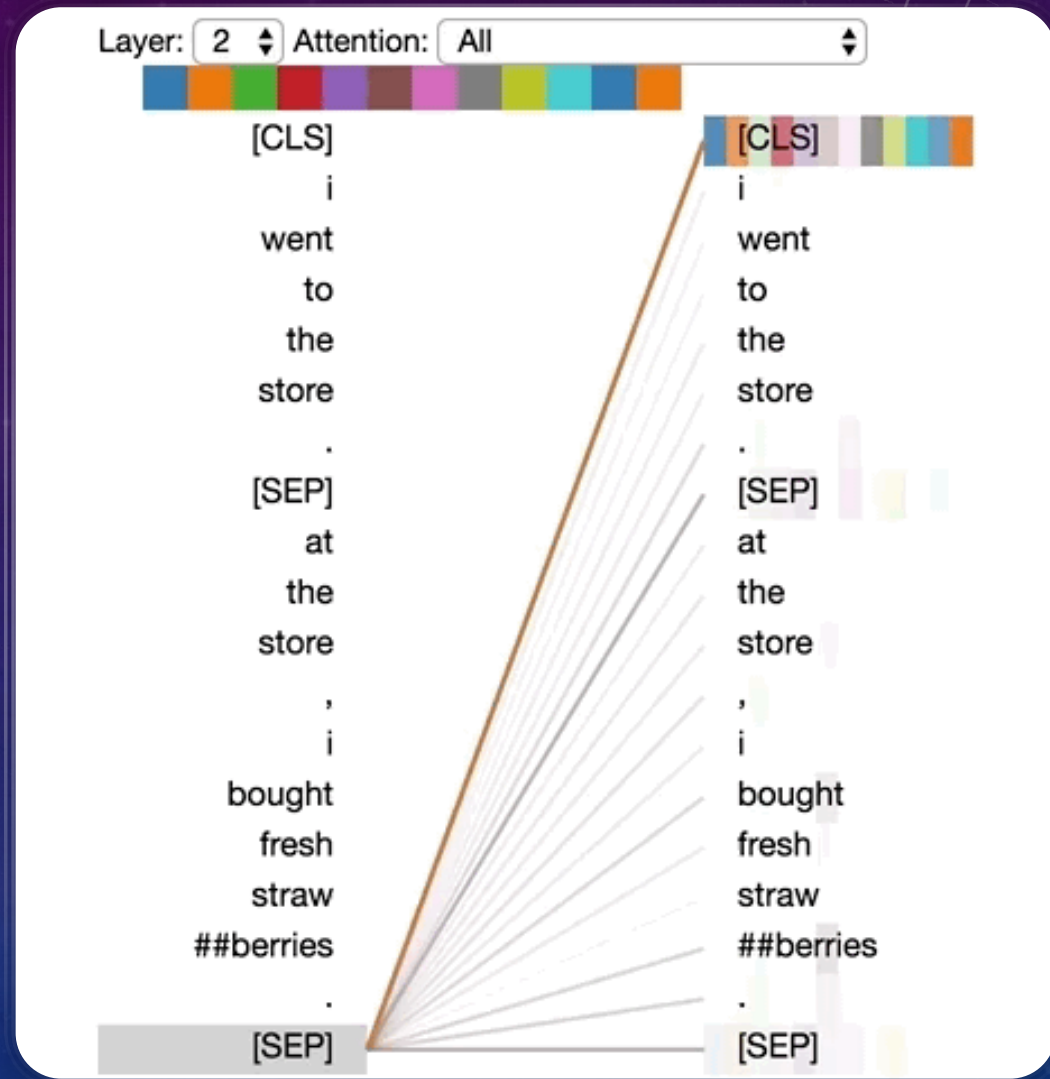


Figure 1: The Transformer - model architecture.

09. 트랜스포머

- Self Attention
 - <https://towardsdatascience.com/deconstructing-bert-distilling-6-patterns-from-100-million-parameters-b49113672f77>

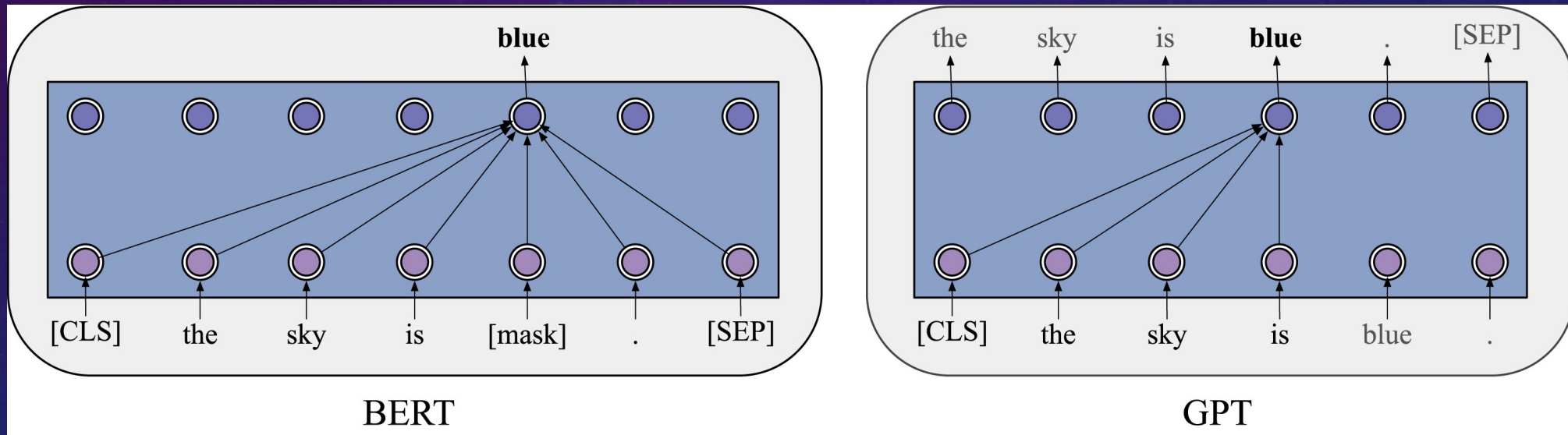


10. 사전학습모델

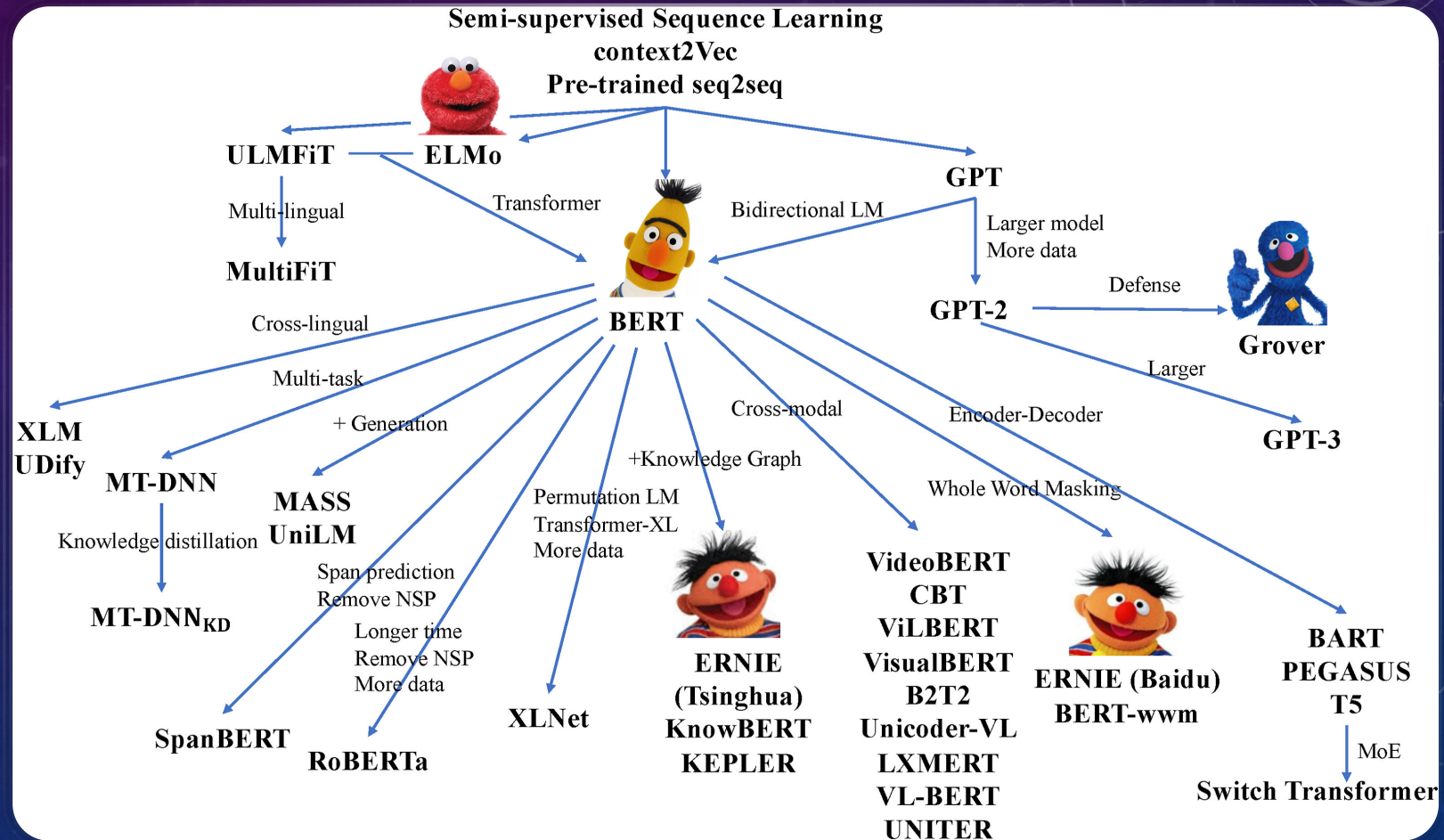
- Decoders or autoregressive models:
 - GPT, GPT-2, GPT3. CTRL(Conditional Transformer Language Model for Controllable Generation, Transformer-XL, Reformer, XLNet
- Encoders or autoencoding models
 - BERT, ALBERT, RoBERTa, DistilBERT, ConvBERT, XLM, XLM-RoBERTa, FlauBERT, ELECTRA, Funnel Transformer, Longformer
- Sequence-to-sequence models
 - BART, Pegasus, MarianMT, T5, MT5, Mbart, ProphetNet, XLM-ProphetNet
- Multimodal models
 - MMBT
- Retrieval-based models
 - DPR, RAG

10. 사전학습모델

- BERT and GPT (Generative Pre-Training)
- Auto Encoder Model vs. Auto Regressive Model



10. 사전학습모델



- <https://ars.els-cdn.com/content/image/1-s2.0-S2666651021000231-gr8.jpg>

거대언어모델(LARGE LANGUAGE MODEL)

사전학습모델(Pre-trained Language Model)의 모델크기나 데이터 크기를 확장하면 downstream task에서 모델 용량이 향상됨

- LLM은 기존 언어모델에서 볼 수 없던 새로운 능력을 보여줌
- LLM은 인간이 AI를 개발하고 사용하는 방식에 변화를 가져옴
 - LLM에 접근하는 방식은 prompt를 통해서 주로 이루어짐
 - 따라서 LLM이 이해할 수 있는 방식으로 지시문(instruction)을 만들어야 함
- chatGPT와 GPT-4의 출현으로 Artificial General Intelligence의 가능성 모색됨

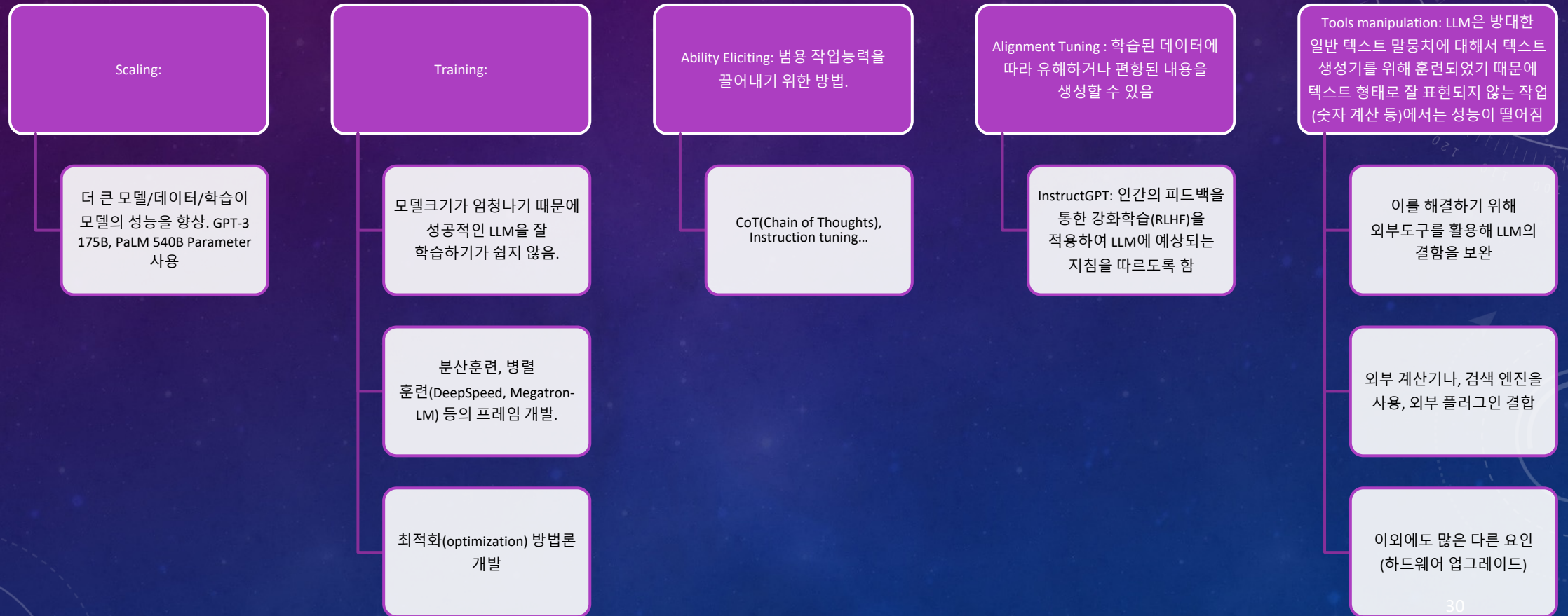
그럼에도 불구하고 LLM의 기본원리는 잘 밝혀지지 않음

- 작은 PLM이 아닌 LLM에서 새로운 기능이 발생하는지 의문
- LLM을 구축하기 위해서는 엄청난 양의 데이터와 고성능의 GPU가 대량으로 필요하기 때문에 훈련 비용이 많이 듦
 - 몇 거대 기업들만 훈련 가능한 상황
 - 구축된 모델 및 사용된 데이터셋의 비공개

LLM의 배경

- 일반적으로 거대언어모델은 트랜스포머(Transformer) 아키텍처를 기반으로 수천억 또는 그 이상의 parameter를 포함하는 언어모델을 지칭
 - Scaling: 거대언어모델의 능력향상은 모델의 크기에 따라 성능이 대략적으로 증가하는 scaling 법칙으로 부분적으로 설명가능.
 - Emergent Ability: “작은 모델에는 존재하지 않지만 큰 모델에서 발생하는 능력”. 규모가 일정수준에 도달하면 성능이 무작위보다 크게 상승함
 1. In-Context Learning(ICL): GPT-3에서 처음으로 도입. 언어모델에 자연어 지시 또는 여러 작업의 데모가 제공되었다면, 추가 학습이나 gradient 업데이트 없이 입력 테스트의 순서를 완성하여 예상 출력을 생성함
 2. Instruction-following: 자연어 설명을 통해 포맷된 다중 작업 데이터셋을 혼합하여 미세조정(instruction tuning)하면, LLM은 명령어 형태로 설명되는 보이지 않는 작업에서도 잘 수행됨
 3. Step-by-step Reasoning: 수학 연산, 논리적 추론 등을 위해 사고의 연쇄(Chain-of-Thoughts)를 사용하여, 최종 답을 도출하기 위한 중간 추론 단계가 포함된prompt를 활용하여 적용

LLM의 주요 기술(KEY TECHNICS)



LIST OF LARGE LANGUAGE MODELS ([HTTPS://EN.WIKIPEDIA.ORG/WIKI/LARGE_LANGUAGE_MODEL](https://en.wikipedia.org/wiki/Large_language_model))

Name	Release date ^[a]	Developer	Number of parameters ^[b]	Corpus size	License ^[c]	Notes
BERT	2018	Google	340 million ^[55]	3.3 billion words ^[55]	Apache 2.0 ^[56]	An early and influential language model, ^[2] but encoder-only and thus not built to be prompted or generative ^[57]
XLNet	2019	Google	~340 million ^[58]	33 billion words		An alternative to BERT; designed as encoder-only ^{[59][60]}
GPT-2	2019	OpenAI	1.5 billion ^[61]	40GB ^[62] (~10 billion tokens) ^[63]	MIT ^[64]	general-purpose model based on transformer architecture
GPT-3	2020	OpenAI	175 billion ^[24]	300 billion tokens ^[63]	public web API	A fine-tuned variant of GPT-3, termed GPT-3.5, was made available to the public through a web interface called ChatGPT in 2022. ^[65]
GPT-Neo	March 2021	EleutherAI	2.7 billion ^[66]	825 GiB ^[67]	MIT ^[68]	The first of a <i>series of free GPT-3 alternatives</i> released by EleutherAI. GPT-Neo outperformed an equivalent-size GPT-3 model on some benchmarks, but was significantly worse than the largest GPT-3. ^[68]
GPT-J	June 2021	EleutherAI	6 billion ^[69]	825 GiB ^[67]	Apache 2.0	GPT-3-style language model
Megatron-Turing NLG	October 2021 ^[70]	Microsoft and Nvidia	530 billion ^[71]	338.6 billion tokens ^[71]	Restricted web access	Standard architecture but trained on a supercomputing cluster.
Ernie 3.0 Titan	December 2021	Baidu	260 billion ^[72]	4 Tb	Proprietary	Chinese-language LLM. <i>Ernie Bot</i> is based on this model.
Claude ^[73]	December 2021	Anthropic	52 billion ^[74]	400 billion tokens ^[74]	Closed beta	Fine-tuned for desirable behavior in conversations. ^[75]
GLaM (Generalist Language Model)	December 2021	Google	1.2 trillion ^[76]	1.6 trillion tokens ^[76]	Proprietary	Sparse mixture-of-experts model, making it more expensive to train but cheaper to run inference compared to GPT-3.
Gopher	December 2021	DeepMind	280 billion ^[77]	300 billion tokens ^[78]	Proprietary	
LaMDA (Language Models for Dialog Applications)	January 2022	Google	137 billion ^[79]	1.56T words, ^[79] 168 billion tokens ^[78]	Proprietary	Specialized for response generation in conversations.
GPT-NeoX	February 2022	EleutherAI	20 billion ^[80]	825 GiB ^[67]	Apache 2.0	based on the Megatron architecture
Chinchilla	March 2022	DeepMind	70 billion ^[81]	1.4 trillion tokens ^{[81][78]}	Proprietary	Reduced-parameter model trained on more data. Used in the <i>Sparrow</i> bot.
PaLM (Pathways Language Model)	April 2022	Google	540 billion ^[82]	768 billion tokens ^[81]	Proprietary	aimed to reach the practical limits of model scale
OPT (Open Pretrained Transformer)	May 2022	Meta	175 billion ^[83]	180 billion tokens ^[84]	Non-commercial research ^[d]	GPT-3 architecture with some adaptations from Megatron
YaLM 100B	June 2022	Yandex	100 billion ^[85]	1.7TB ^[85]	Apache 2.0	English-Russian model based on Microsoft's Megatron-LM.

LIST OF LARGE LANGUAGE MODELS

Minerva	June 2022	Google	540 billion ^[86]	38.5B tokens from webpages filtered for mathematical content and from papers submitted to the arXiv preprint server ^[86]	Proprietary	LLM trained for solving "mathematical and scientific questions using step-by-step reasoning". ^[87] Minerva is based on PaLM model, further trained on mathematical and scientific data.
BLOOM	July 2022	Large collaboration led by Hugging Face	175 billion ^[88]	350 billion tokens (1.6TB) ^[89]	Responsible AI	Essentially GPT-3 but trained on a multi-lingual corpus (30% English excluding programming languages)
Galactica	November 2022	Meta	120 billion	106 billion tokens ^[90]	CC-BY-NC-4.0	Trained on scientific text and modalities.
AlexaTM (Teacher Models)	November 2022	Amazon	20 billion ^[91]	1.3 trillion ^[92]	public web API ^[93]	bidirectional sequence-to-sequence architecture
LLaMA (Large Language Model Meta AI)	February 2023	Meta	65 billion ^[94]	1.4 trillion ^[94]	Non-commercial research ^[6]	Trained on a large 20-language corpus to aim for better performance with fewer parameters. ^[94] Researchers from Stanford University trained a fine-tuned model based on LLaMA weights, called Alpaca. ^[95]
GPT-4	March 2023	OpenAI	Exact number unknown, approximately 1 trillion ^[f]	Unknown	public web API	Available for ChatGPT Plus users and used in several products .
Cerebras-GPT	March 2023	Cerebras	13 billion ^[97]		Apache 2.0	Trained with Chinchilla formula.
Falcon	March 2023	Technology Innovation Institute	40 billion ^[98]	1 Trillion tokens (1TB) ^[98]	Apache 2.0 ^[99]	The model is claimed to use only 75% of GPT-3's training compute, 40% of Chinchilla's, and 80% of PaLM-62B's.
BloombergGPT	March 2023	Bloomberg L.P.	50 billion	363 billion token dataset based on Bloomberg's data sources, plus 345 billion tokens from general purpose datasets ^[100]	Proprietary	LLM trained on financial data from proprietary sources, that "outperforms existing models on financial tasks by significant margins without sacrificing performance on general LLM benchmarks"
PanGu-Σ	March 2023	Huawei	1.085 trillion	329 billion tokens ^[101]	Proprietary	
OpenAssistant ^[102]	March 2023	LAION	17 billion	1.5 trillion tokens	Apache 2.0	Trained on crowdsourced open data
PaLM 2 (Pathways Language Model 2)	May 2023	Google	340 billion ^[103]	3.6 trillion tokens ^[103]	Proprietary	Used in Bard chatbot . ^[104]

상업적으로 사용 가능한 공개 언어 모델

	License	Data	Architecture	Weights	Size	Checkpoints	Language
Meta Llama2	Llama license	Open	Open	Open	7, 13, 70	Yes	English / Multilingual
EleutherAI Pythia	Apache 2.0	Open	Open	Open	7, 12	Yes	English
EleutherAI Polyglot	GPL-2.0	Open	Open	Open		Yes	English / Multilingual
GPT-J	MIT	Open	Open	Open	6	Yes	English
Databricks Dolly 2	Apache 2.0	Open	Open	Open	7, 12	Yes	English
Cerebras-GPT	Apache 2.0	Open	Open	Open	7, 13	Yes	English / Multilingual
StableLM	CC BY-SA-4.0	Open	Open	Open	3, 7, (15, 30, 65, 175)	Yes	English
Mosaic MPT	Apache 2.0	Open	Open	Open	7, 30	Yes	English
Falcon GPT	Apache 2.0	Open	Open	Open	7, 40	Yes	English

TECHNICAL EVOLUTION OF GPT-SERIES MODELS

A timeline of Existing Large Language Models(Zhao et al. (2023))

