

ChatGPT 기술 분석 백서

SAMSUNG SDS
데이터분석서비스팀

Last Updated: Apr. 5, 2023

Revision History		
Version	Date	Description
v1.0	'23.3.22	백서 초안 완성
v1.1	'23.3.31	GPT-4 등 최신 동향 반영
V1.2	'23.4.5	오탈자 수정

"본 문서는 급변하는 GPT 기술에 대한 이해를 돕기 위해 신속하게 작성되었으나,
일부 내용이 최신 정보와 상이하거나 오탈자 등을 포함할 수 있습니다.
본 문서는 삼성SDS의 공식 입장을 대표하지 않으며,
문서를 통해 제공되는 정보로 인해 발생하는 법적 문제나 책임을 지지 않습니다."

목차

1. General AI vs Narrow AI.....	3
2. ChatGPT란	5
2-1. 언어 모델의 발전	7
2-2. Transformer 기반 언어 모델.....	9
2-3. GPT(Generative Pre-trained Transformer)	10
2-4. ChatGPT	15
3. ChatGPT 활용	16
3-1. ChatGPT Web UI	16
3-2. ChatGPT API.....	16
3-3. Fine-tuning	17
3-4. GPT 기반 어플리케이션	20
3-5. 개인 업무 활용.....	21
3-6. 한계	22
3-7. 유의점.....	23
4. ChatGPT, 그리고 생성형 AI가 바꿀 미래.....	24
5. 맺음말	26

※ 본 문서는 최신 정보를 바탕으로 지속적으로 온라인 업데이트 될 예정입니다.

최신 버전의 문서는 아래 웹사이트를 참고해 주세요.

▶ <https://www.brightics.ai/>

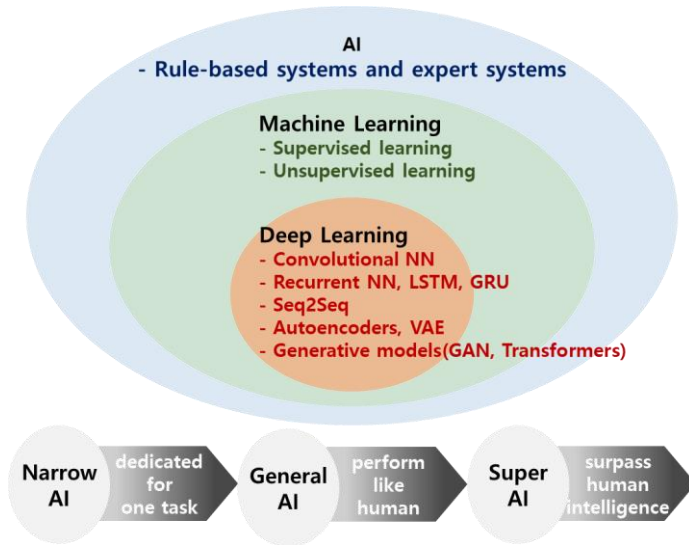
※ 본 백서 인용 시에는 아래와 같이 출처를 표기해 주시기 바랍니다.

“삼성SDS 『ChatGPT 기술 분석 백서(Mar. 2023)』”

1. General AI vs Narrow AI

- **인공 지능**(Artificial Intelligence, AI)은 컴퓨터 시스템을 사용하여 인간의 지능적인 작업을 모방하거나 수행하는 것을 말합니다. 초기의 인공 지능은 인간의 전문 지식을 기반으로 미리 정의된 규칙과 논리에 따라 단순 반복적인 작업을 주로 수행했습니다. 인공 지능의 하위 분야로서 등장한 **머신 러닝**(machine learning)은 인간이 쉽게 이해하거나 분별하기 어려운 패턴과 확률적 분포를 학습하여 인사이트를 제공함과 동시에 미래를 예측하여 의사 결정에 도움을 주었습니다. 머신 러닝에서 한 단계 더 나아가 발달한 **딥 러닝**(deep learning)은 인간의 신경계 구조에 착안하여 만들어진 인공 신경망을 기반으로 더 복잡하고 심층적인 정보를 처리할 수 있습니다. 이렇게 인공 지능의 성능은 지속적으로 향상되어 왔습니다.
- 2016년 3월(7년 전 이맘때), 전세계가 주목한, 인공 지능 분야의 혁신적인 사건이 있었습니다. **Google DeepMind**에서 개발한 인공 지능 **AlphaGo**가 바둑 세계 챔피언 이세돌 선수를 이긴 사건은 인공 지능 기술의 혁신과 발전을 대중들에게 크게 알린, 역사적인 순간이었습니다.
- 이 대국은, 제한된 범위 내에서 특정 작업을 수행하고 특화된 문제를 해결할 수 있도록 설계된 **Narrow AI**(좁은 인공지능; Weak AI, 약 인공지능)의 능력을 명확하게 보여 주는 예시입니다. 바둑을 학습한 Narrow AI는 특정 영역에서 매우 전문화 되어 있기 때문에, AlphaGo와 같이 바둑 게임을 놀라운 수준으로 수행할 수 있지만, 이와 관련 없는 다른 작업, 즉, 언어 이해, 이미지 분류, 음성 인식, 자율 주행, 수요 예측 등은 수행할 수 없습니다. 이와 달리, **General AI**(일반 인공지능; Strong AI, 강 인공지능)는 다양한 영역에서 인간이 할 수 있는 모든 지적 작업에 대해 인지/학습/수행할 수 있도록 설계된 것을 가리키며, 아직은 이론적 개념에 머물러 있습니다.
- 딥 러닝의 발전과 맞물려 등장한 **생성형 AI**(Generative AI)는 기존 패턴을 기반으로 오디오, 비디오, 이미지, 텍스트, 코드, 시뮬레이션 등의 새로운 콘텐츠를 생성하는 데 사용할 수 있는 알고리즘입니다. **OpenAI**에서 개발한 **ChatGPT**는 생성형 AI 분야에서 인공 지능 기술이 오랜 시간 지속적으로 진화해 오면서 만들어 낸 걸작으로서, Narrow AI에서 General AI로 향해 가는 변곡점이 되고 있습니다. 궁극적으로 인공 지능은 인간의 지능을 능가하는 **Super AI**를 달성하는 것이 목표입니다.
- ChatGPT는 주로 자연어 처리 작업(Natural Language Processing, NLP)을 위해 설계되었지만, Narrow AI에 비해 다양한 작업에 적용될 수 있습니다. 그러나 여전히 진정한 의미의 General AI는 아닙니다. ChatGPT의 지식과 능력이 사전에 정의된 학습 데이터의 범위에 의해 제한되며 언어 영역 대비 그 외의 작업은 비교적 약한 수행 능력을 가지고 있기 때문입니다.
- General AI를 향해 가는 여정에 있어서 ChatGPT 출시는, 인공지능 기술의 발전으로 인해 인간의 작업을 지원할 수 있는 영역 확장이 가속화되고 있으며, 현재까지 출시된 타 AI 모델들보다 큰 잠재력을 보여 준다는 점에서 의미 있습니다.

【그림1. 인공지능이 포함하는 개념】



(설명 출처: ChatGPT)

인공 지능(AI)는 지능형 시스템을 위한 방법론을 연구/개발하는 분야입니다.

초기 AI는 규칙 기반 또는 전문가 시스템으로서, 사전에 정의된 규칙 또는 인간의 전문 지식을 기반으로 작동했습니다.

머신러닝은 AI의 하위 분야로 등장하여 그간 발전해 왔으며, 딥러닝은 기계 학습의 하위 개념입니다. 딥러닝의 발전으로 인해 AI 성능이 크게 향상되었습니다.

AI 발전 단계 중 현재는 Narrow AI와 General AI 사이에 위치하고 있습니다. AI 연구의 미래는 General AI 뿐만 아니라 Super AI를 달성하여, 인간-기계 상호 작용에서 전례 없는 가능성을 창출하는 것을 목표로 합니다.

▪ OpenAI에서 ChatGPT를 출시('22.11.30)한 데 이어 Google에서는 **Bard** 론칭 발표('23.2.6)에서 데모를 시연하였습니다. 그리고 뒤이어 Google이 Workspace 제품군에 대한 AI 전면 적용 계획을 발표('23.3.15)하였으며, 이를 뒤 OpenAI와 파트너십을 맺고 있는 **Microsoft**에서 역시 자사 업무 생산성 도구 전반에 AI를 적용한 소프트웨어 **Microsoft 365 Copilot** 출시 계획을 발표('23.3.16)하였습니다. 불과 4개월도 채 되지 않는 기간에 GPT 기술을 둘러싼 글로벌 빅테크 기업들의 경쟁이 격화되었고, 전세계적으로 IT 업계에서는 **GPT**(Generative Pre-trained Transformer) 기술을 접목한 비즈니스 모델 개발에 열을 올리고 있습니다. ChatGPT 출시는 인공 지능의 발전 단계에 큰 획을 긋는 역사적인 사건으로 기록될 것입니다.

【표1. ChatGPT 관련 타임라인】

일자	상세 내용
'22.11.30	OpenAI, ChatGPT(GPT-3.5) 론칭
'22.12.5	ChatGPT 일간 활성화 사용자 수(DAU) 100만 명 돌파
'22.12.15	OpenAI, ChatGPT에 워터마크 도입
'22.12.25	ChatGPT 일간 활성화 사용자 수(DAU) 1000만 명 돌파
'22.12.26	Google, ChatGPT에 '코드 레드' 발령
'22.12.27	베스핀글로벌, '헬프나우 AI'에 GPT 기술 도입
'22.12.31	12월 월간 활성화 사용자 수(MAU) 5,700만 명 돌파
'23.1.23	Microsoft, OpenAI에 100억 달러 추가 투자
'23.1.31	12월 월간 활성화 사용자 수(MAU) 1억 명 돌파
'23.2.1	투블럭AI, 생성형 AI 활용 방법 관련 특허 등록
'23.2.3	네이버, 상반기 '서치GPT' 출시 계획 발표
'23.2.6	Google, Bard 론칭 발표 및 데모 시연
'23.2.7	Microsoft, ChatGPT 탑재 검색 엔진 '빙(Bing)' 발표
'23.3.9	업스테이지, OCR 기술에 ChatGPT를 결합한 '아숙업(AskUp)' 론칭
'23.3.14	OpenAI, ChatGPT(GPT-4) 론칭
'23.3.15	Google, Workspace 제품군에 대한 AI 전면 적용 계획 발표
'23.3.16	Microsoft, 업무 생산성 도구 전반에 AI를 적용한 Microsoft 365 Copilot 출시 계획 발표
'23.3.23	OpenAI, ChatGPT plugin 지원 발표

2. ChatGPT란

▪ ChatGPT는 OpenAI(Sam Altman이 '15.12월 설립하였으며, 인공지능에 의한 전 인류적 이익을 추구하는 연구 개발 회사)가 개발한 대화형 AI(챗봇) 서비스로서, GPT 아키텍처(본문 2-3. GPT 참고)를 기반으로 합니다.

▪ ChatGPT는 출시와 동시에 폭발적인 관심을 받으며 5일 만에 100만 명, 1개월 만에 약 1,000만 명의 사용자를 모은 데 이어 2개월 만에 월간 활성 사용자 수(Monthly Active Users, MAU) 1억 명을 돌파하였습니다. 타 IT 업체들의 경우, 100만 명 사용자 달성까지 걸린 시간이 **Netflix**(3.5년), **Airbnb**(2.5년), **Facebook**(10개월), **Spotify**(5개월), **Instagram**(2.5개월), **iPhone**(74일)인 것을 감안하면 엄청난 인기입니다.

【그림2. 플랫폼별 사용자 100만명 도달에 걸린 시간, 출처: Katharina Buchholz, statista, '23.1.24, 「ChatGPT Sprints to One Million Users」(<https://www.statista.com/chart/29174/time-to-one-million-users/>)】



▪ 기존의 챗봇은 질문에서 추출한 키워드를 DB에 저장된 정보와 단순히 패턴 매칭을 하는 방식 등의 규칙 기반 시스템으로 작동하여, 한정된 범위 내에서만 대화가 가능합니다. 반면에, ChatGPT는 표면적인 정보 뒤에 숨어 있는 맥락을 이해하고, 과거 대화 기록을 기억함으로써, 직접 사람과 대화하는 것과 유사한 수준의 정보 전달이 가능합니다.

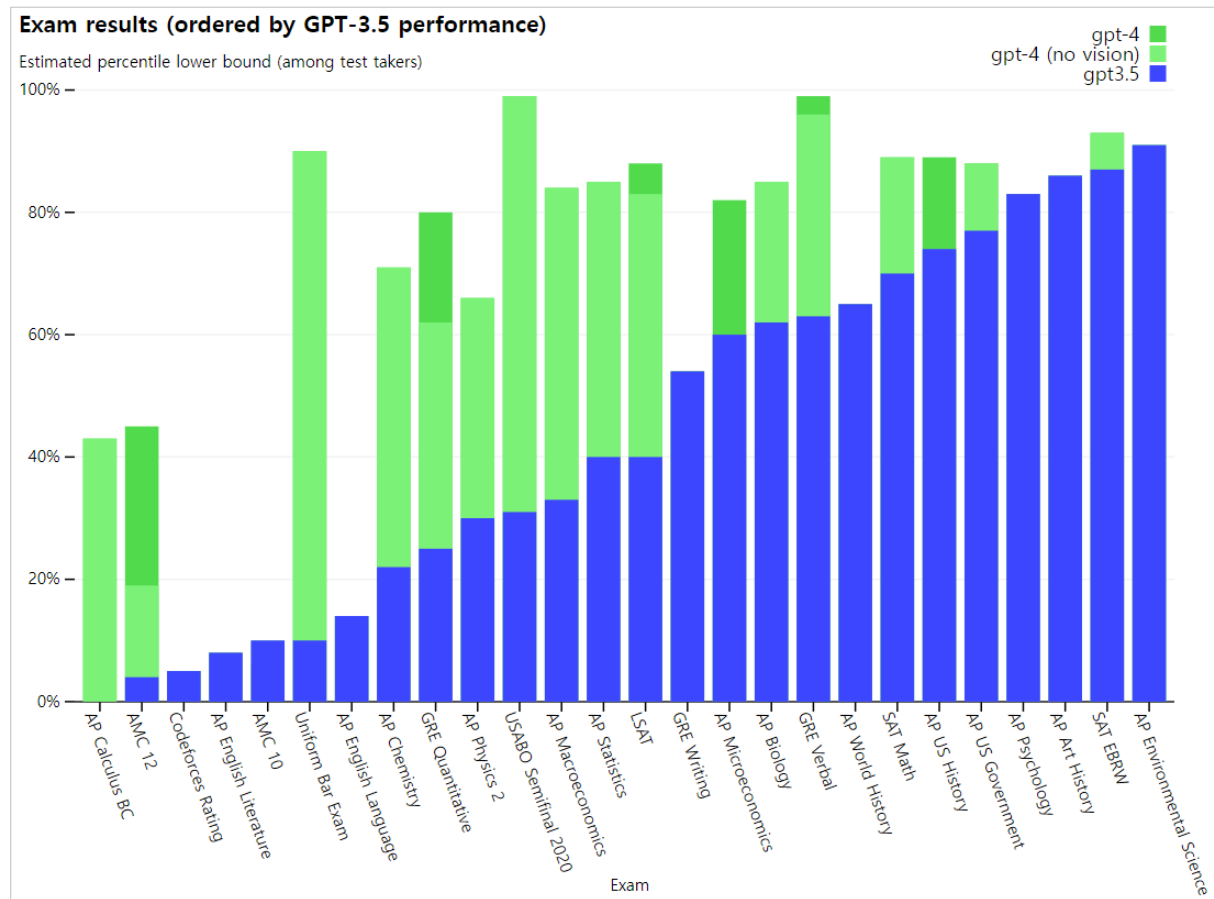
▪ 조금 더 구체적인 목적을 가진 예시들을 살펴 본다면, 텍스트 생성(예: 가사, 소설, 이메일, 기획안), 감성 분석(예: 고객 후기 긍정/부정 구분), 문서 분류(예: 뉴스 기사 주제 분류, VoC 분류), 문서 요약, 키워드 추출, 번역, Q&A, 문법 교정, 프로그래밍 코드 작성, 표 작성 등의 기능을 수행할 수 있습니다.

▪ 출시 후 여러 성능 테스트를 거친 결과, 미국 와튼스쿨 MBA, 미국 의사면허시험, 로스쿨 시험 등을 모두 무난하게 통과함으로써, 인간의 지적 능력과 관련된 업무 수행 역량을 입증하였습니다.

【표2. ChatGPT(3.5) 성능】

성능	상세 내용
미국 와튼스쿨 MBA 통과	필수 교과목인 '운영관리' 기말시험에서 B-에서 B 사이 점수 획득
미국 미네소타 로스쿨 시험 통과	객관식 문항 95개, 에세이 문항 12개로 이루어진 시험에서 C+ 점수 획득
미국 의사면허시험(USMLE) 통과	모든 시험에서 50% 이상의 정확도
학술 논문의 공동 저자로 등재	영국 맨체스터 교수 Siobhan O'Connor, 논문에 ChatGPT를 공동 저자로 등재
광고 대본 작성	영화배우 Ryan Reynolds의 스타일로 휴대폰 통신사 Mint Mobile 광고 대본 작성

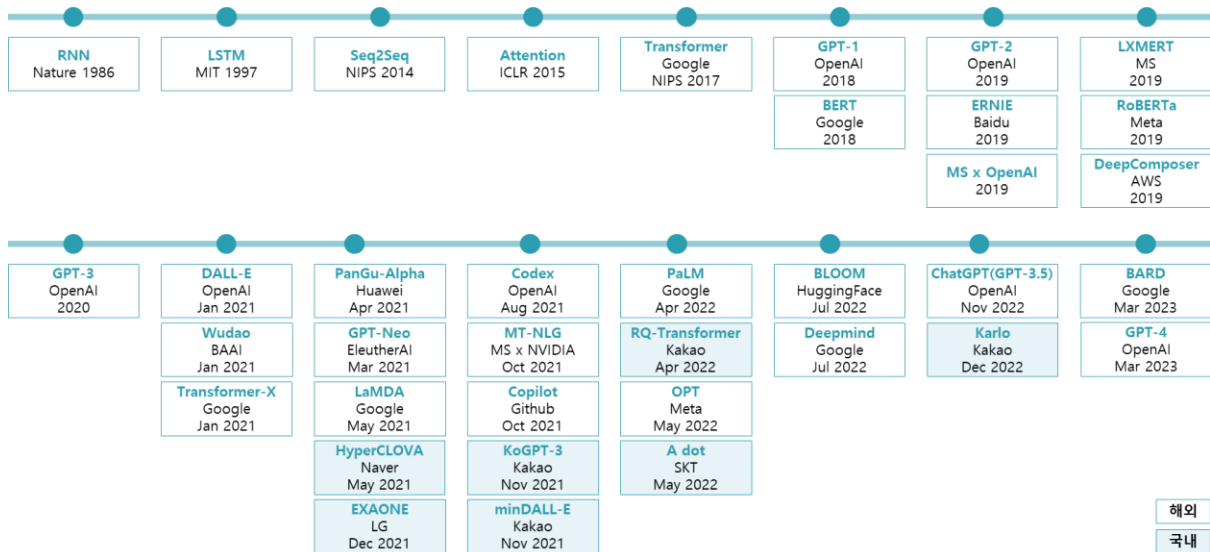
【그림3. ChatGPT 성능, 출처: OpenAI 공식 홈페이지(<https://openai.com/research/gpt-4>)】



2-1. 언어 모델의 발전

▪ ChatGPT는 **생성형 AI**(Generative AI) 분야에서 인공 지능 및 기계 학습(machine learning) 기술이 오랜 시간 지속적으로 진화해 오면서 만들어낸 걸작입니다. 생성형 AI는 기존 패턴을 기반으로 오디오, 비디오, 이미지, 텍스트, 코드, 시뮬레이션 등의 새로운 콘텐츠를 생성하는 데 사용할 수 있는 알고리즘으로서, 최근 언어 및 이미지 분석에서 크게 주목을 받고 있습니다. 생성형 AI 분야가 현재 수준으로 발전하기까지 크고 작은 공헌들이 있었지만, 언어 모델에 있어 큰 이정표 역할을 한 모델들 위주로 설명해 보겠습니다.

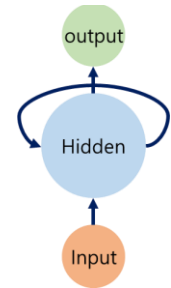
【그림4. 언어 모델의 주요 발전 타임라인】



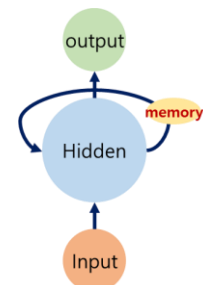
▪ 인간의 두뇌로는 한 눈에 쉽게 파악하기 어렵고 복잡한 규칙을 가진 【그림5. FNN 모식도】 현상을 기계에 학습시킴으로써, 현상 속에 내재된 패턴을 파악하고 미래를 예측하고자 합니다. 이를 위해서 현상에 대한 정보를 압축시킨 데이터를 인간의 신경계와 유사한 구조를 가지고 있는 인공 신경망에 넣어 줍니다. 데이터가 신경망 사이를 거쳐서 나아가는 동안 기계는 데이터의 대세적인 흐름, 유형화 할 수 있는 요소들, 패턴 파악에 필수적인 정보들을 학습하고 노이즈는 제거해 가면서 데이터에 대해 이해하는 과정을 거칩니다. 그리고 데이터가 신경망의 최종 지점에 도달했을 때 기계가 학습한 정보를 바탕으로 예측값을 내어 놓는데, 만약 정답과 다를 경우 데이터(오차)를 뒤로 돌려 보내서 다시 학습하는 방식을 조정합니다. 정답을 제대로 맞추고 목표한 바를 이룰 때까지 이 과정을 반복합니다. 이렇게 기계가 데이터를 학습할 수 있도록 앞으로 흘러 보내 주고 뒤로 돌려 보내는 것을 순방향 신경망과 역전파(Feedforward Neural Networks and **Backpropagation**, **FNNs**, 1980s)라고 합니다.



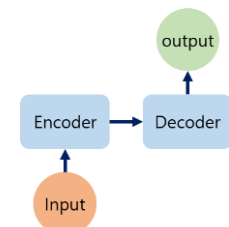
▪ 음성이나 텍스트처럼 데이터가 앞뒤 순서를 가지고서 순차적으로 들어올 때, 순방향 신경망에서는 데이터가 계속 한 방향으로만 흘러 가도록 넣어 주며, 이 입력되는 데이터에 의존해서 기계도 예측값을 산출합니다. 하지만 이번에 새로 들어온 데이터에 대한 정보를 기반으로 예측값을 산출할 때 바로 직전에 이미 흘러가 버린 정보도 함께 참고하는 것이 더 정확할 수 있습니다. 이전에 들어 와서 한 번 흘러가 버린 정보를 다시 망 내부로 유입시킬 수 있는 순환 구조를 가질 수 있도록 고안된 것이 순환 신경망(Recurrent Neural Networks, **RNNs**, 1980s-1990s)입니다. 과거 정보에 대한 기억과 새로 들어온 정보를 함께 고려하여 앞뒤 맥락에 대한 이해를 바탕으로 신경망을 학습함으로써, 과거와 현재를 모두 아우른 예측이 가능합니다. 시계열 데이터, 음성 인식, 자연어 처리와 같이 시간이 흐름에 따라 변화하는 순차적 데이터 처리에 적합합니다.



▪ 이제 RNN 구조를 기반으로 과거에 들어온 정보도 같이 참고해서 현재값을 예측할 수 있습니다. 하지만 직전에 들어 온 정보뿐만 아니라 한참 전에 들어 온 정보도 함께 참고를 해야 하는데, 시간이 점점 흐르고 유입되는 데이터가 누적되는 상황에서는 기억이 소실되고 학습 능력이 떨어지게 됩니다. 이를 해결하기 위해서 장기적으로 기억해야 할 중요한 정보와 잊어 버려도 되는 정보를 구분하여 제어함으로써, 길이가 긴 순차적 데이터도 효과적으로 학습할 수 있도록 고안된 것이 장단기 메모리(Long Short-Term Memory, **LSTM**, 1997)이며, RNN의 한 유형입니다.

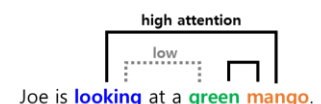


▪ 두 개의 RNN 구조, 즉, (1) 입력 데이터를 학습하고 필요한 정보와 맥락만 응축해서 처리하는 RNN 구조 '**인코더(encoder)**', 그리고 (2) 인코더로부터 과거 및 현재에 대한 응축된 정보를 전달 받아 예측값을 산출하는 RNN 구조 '**디코더(decoder)**'를 결합하여 구성한 신경망 유형이 **Seq2Seq**(Sequence to Sequence, 2014) 입니다. 이처럼 인코더와 디코더를 각각 구성하면 입력-출력 데이터 간 복잡한 관계를 더 효과적으로 파악할 수 있으므로 기계 번역과 같은 작업에 강점을 가집니다. 물론 이때 인코더 및 디코더를 구성하는 RNN을 LSTM 구조로 적용할 수 있으며, 여러 겹의 블록 형태로도 쌓을 수 있습니다.



▪ Seq2Seq 모델에서 정보가 압축될 때 손실이 발생하기 때문에, 입력 정보가 길어지면 모델 성능이 떨어집니다. 이를 해결하기 위해서, 순차적으로 입력되는 정보 중 필요한 부분만 선택적으로 주의를 집중할 수 있도록 하여 Seq2Seq 모델의 성능을 높인 것이

【그림9. Attention 원리】



Attention(주의) 메커니즘(2014) 입니다. 즉, 하나의 단어를 출력할 때마다 입력된 전체 문장을 참고하는데, 이때 입력 문장 내 모든 단어들을 동일한 비중으로 살피는 것이 아니라, 이번에 출력할 단어와 관련성이 높은 단어들에 조금 더 비중을 두는 것으로, 출력이 이루어 질 때마다 입력 문장 내에서 큰 비중으로 다뤄지는 단어들의 집합은 계속 달라지게 됩니다.

2-2. Transformer 기반 언어 모델

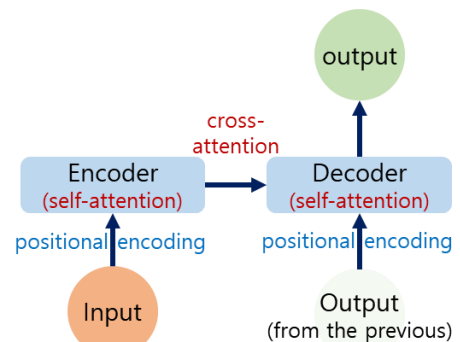
▪ 2017년, Google에서 제목 'Attention is all you need'로 논문을 발표하면서 Attention 메커니즘을 적용한 Transformer 모델을 소개했습니다. Transformer는 이전 모델을 훨씬 능가하는 성능을 보여주었으며, 이때부터 생성형 AI분야에서 언어 모델이 빠른 속도로 발전하게 됩니다. 현재 활용되고 있는 언어 모델들의 대부분이 Transformer를 기반으로 확장해 가고 있는 것을 보면 더 혁신적인 방법론 제안이 있기까지는 계속해서 중심 역할을 할 것으로 보입니다.

【표3. Transformer 기반의 모델】

BERT 계열	BERT, RoBERTa, ALBERT, DistilBERT, ELECTRA, XLM, XLM-R (XLM-RoBERTa), SpanBERT, MobileBERT, TinyBERT, DeBERTa
GPT 계열	GPT-1, GPT-2, GPT-3, GPT-Neo, GPT-NeoX, Megatron-GPT
Seq2Seq Transformer	T5(Text-to-Text Transfer Transformer), BART(Bidirectional and Auto-Regressive Transformers), mBART(Multilingual BART), PEGASUS(Pre-training with Extracted Gap-sentences for Abstractive Summarization Sequence-to-sequence)
Extended context	Longformer, BigBird, Transformer-XL, Universal Transformers

▪ **Transformer**(2017)도 Seq2Seq 모델처럼 입력 정보와 출력 정보 간 **【그림10. Transformer 모식도】**

관계를 매핑하는 인코더-디코더 구조를 가지긴 하지만, 두 모델의 가장 큰 차이점은 정보를 받아 들이는 방식입니다. RNN은 정보를 순차적으로 하나씩 입력 받으면서 순환 구조를 통해 과거에 대한 기억을 되살려 활용하는 반면에, Transformer는 정보를 한꺼번에 입력 받으며 Attention 메커니즘을 통해 집중이 필요한 정보들을 위주로 스캔하여 살펴 봅니다.



먼저 **인코더 Self-attention** 메커니즘을 통해 입력된 전체

문장에서 각 단어마다 다른 단어들과의 관계 및 중요도를 파악하고 함축된 정보로 인코딩

합니다. 이 정보를 받은 디코더는 **인코더-디코더 attention** 메커니즘을 통해 현재 출력할 단어가 입력 문장 내 각 단어들과 얼마나 관련성이 있는지 측정할 뿐 아니라, **디코더 Self-attention** 메커니즘을 통해 지금까지 출력된 단어들 간 관련성을 측정합니다.

즉, 입력 문장 내 단어들 간 관련성, 입력 문장과 출력 단어 간 관련성, 출력 문장 내 단어들 간 관련성을 두루 살피게 됩니다. 이 세 종류의 Attention 메커니즘으로 인해 입력 정보와 관련성이 높은 결과를 출력할 수 있으며, 이 Attention 메커니즘들은 병렬 처리되기 때문에 대량의 데이터를 효율적으로 다룰 수 있습니다.

참고로, RNN에서는 정보를 순차적으로 하나씩 입력 받기 때문에 문장 내 단어 순서 정보를 따로 처리하여 저장할 필요가 없으며 순환 구조를 통해 지나간 정보에 대한 기억을 되살리기만 하면 됩니다. 이와 달리, Transformer는 정보를 한꺼번에 입력 받으며, 문장 내 단어 위치 정보를 별도로 처리하여 결합시킴으로써 순서 정보를 보존(positional encoding)합니다.

▪ **BERT**(Bidirectional Encoder Representations from Transformers, 2018)는 Google이 Transformer 모델을 기반으로 언어 표현을 사전 학습시키기 위해 고안한 방법론이자, 대량의 텍스트 데이터로 사전 학습된 모델입니다. 이름에서도 알 수 있듯이 양방향 문맥 파악이 가능하다는 점이 주요 특징이고, 이로부터 마스크 언어 모델링, 다음 문장 예측 등에 활용하기 용이하며, 특정 작업을 위한 Fine-tuning(미세조정)이 가능합니다.

- 1) 양방향 문맥(Bidirectional context): 기존 언어 모델과 달리 앞뒤로 놓인 단어들을 모두 고려하여 문맥을 파악함으로써 언어 독해력을 높일 수 있습니다.
- 2) 마스크 언어 모델링(Masked Language Modeling, MLM): 입력된 문장 내에서 일정량의 단어들을 가린 뒤, 문맥을 고려하여 원래 단어가 무엇이었는지 예측하는 과정을 통해, 정확한 표현을 생성해내는 방법을 학습합니다.
- 3) 다음 문장 예측(Next Sentence Prediction, NSP): 짝지어진 문장 입력을 통해 연달아 나타날 수 있는 문장들을 학습시킴으로써, 모델이 단어 수준을 넘어서 문장 수준에서 문맥을 이해할 수 있습니다.
- 4) 작업 특화 미세조정(Task-specific fine-tuning): 사전 학습된 이후에는 특정 작업 관련 데이터로 미세조정해 줌으로써, 감성분석, Q&A, 개체명 인식(Named-Entity Recognition, NER) 등 다양한 자연어 처리를 수행할 수 있습니다.

BERT 방법론을 기반으로 학습 방식을 수정하여 모델 크기, 계산량, 성능 측면에서 개선된 BERT 계열의 모델에는 RoBERTa(Robustly Optimized BERT Pretraining Approach; Meta 개발), DistilBERT (distilled version of BERT; HuggingFace 개발), ALBERT(A Lite BERT; Google 개발), ELECTRA(Efficiently Learning an Encoder that Classifies Token Replacements Accurately; Google 개발) 등이 있습니다.

2-3. GPT(Generative Pre-trained Transformer)

▪ **GPT**(Generative Pre-trained Transformer): OpenAI가 대량의 데이터로 다양한 작업을 수행할 수 있도록 사전 학습한 Transformer 모델입니다. BERT와 마찬가지로, 특정 작업을 잘 수행할 수 있도록 사전 학습된 모델을 Fine-tuning 할 수 있습니다. GPT는 일방향으로 나아 가면서 학습 및 예측을 하기 때문에 문장을 생성해 나가는 데 강점을 지닙니다. 이전까지의 단어들을 토대로 파악한 문맥에 맞게 단어를 생성하고 나면, 이 생성된 단어 역시 문맥 이해에 반영되고, 업데이트 된 문맥 정보를 기반으로 또 다음 단어를 생성하는 과정이 반복적으로 일어납니다. 이 과정은 생성되는 문장이 일관된 문맥을 유지할 수 있도록 합니다.

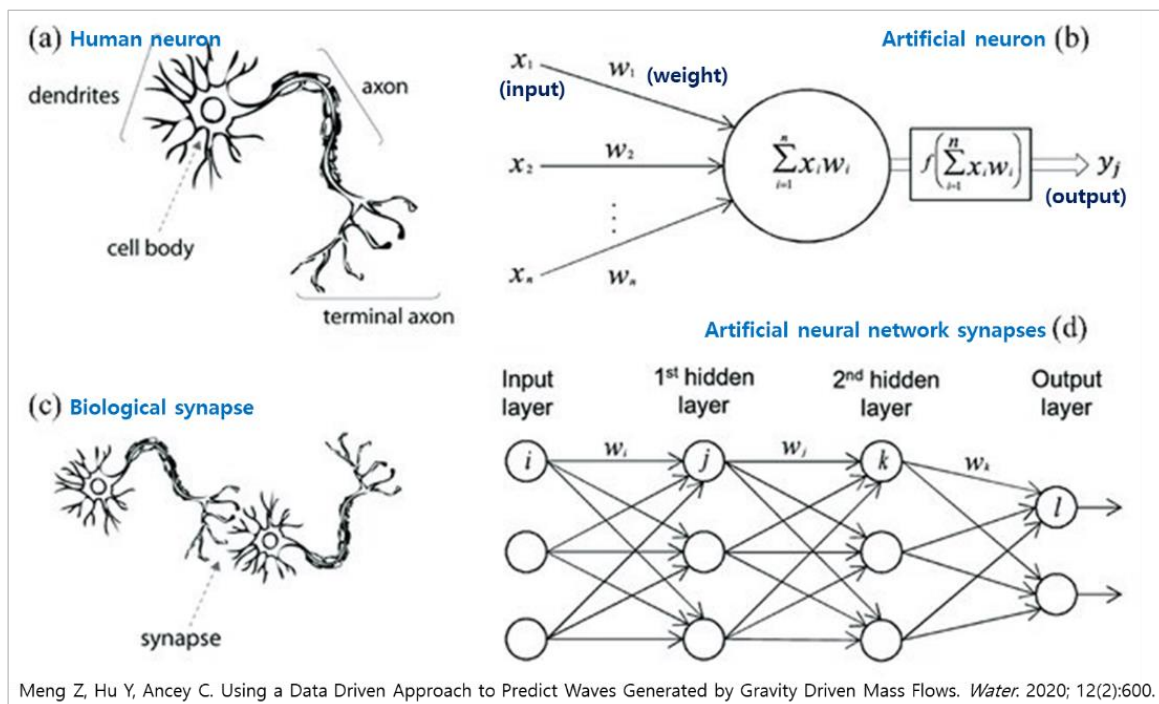
【표4. BERT vs. GPT 비교】

구분	BERT	GPT
특징	전체 문장 내에서 단어의 좌우, 즉, 양방향으로 문맥을 고려하면서 학습하기 때문에, 문장의 의미를 포괄적으로 이해함. 맥락 상 단어의 의미를 파악하는 것이 중요한 자연어 처리 작업에 강점을 지님.	일방향으로 나아가면서 다음 단어를 예측하는 방식으로 학습하기 때문에, 순차적으로 일관되게 문장을 생성함. 인간처럼 일관되고 연관성이 높은 언어를 구사하여 대화형 작업에 강점을 지님.

주요 응용	문서 분류 감성 분석 개체명 인식 Q&A 기계 번역	문장 완성 요약 대화형 챗봇 창작 글쓰기(소설, 시 등) 프로그래밍 코드 작성
-------	--	---

▪ **GPT 시리즈**: OpenAI에서 개발한 GPT는 현재 총 5개(**GPT-1**, **GPT-2**, **GPT-3**, **GPT-3.5**, **GPT-4**) 버전이 존재합니다. 모두 기본적으로 같은 구조를 가지나, 버전이 올라갈수록 **파라미터**(parameter, 매개변수)의 개수가 증가합니다. 이때 파라미터란, Transformer를 구성하는 여러 겹의 신경망 구조를 거치며 입력된 정보에 대한 학습이 이루어지는 동안, 입력값들에 주어지는 가중치(weight) 및 편향(bias)를 가리키며, 인간의 뇌로 따지면 신경망에서 뉴런 간 연결을 시켜주는 시냅스가 신경 물질 전달을 위해 가지고 있는 정보로 이해할 수 있습니다. 이 파라미터들은 학습이 진행될수록 정답에 가까운 문장을 생성해내는 방향으로 조정됩니다. 파라미터 개수가 증가할수록 더 정교한 학습이 이루어지며, 길이가 긴 문장을 이해하거나 복잡한 작업을 처리할 수 있는 능력이 올라 갑니다.

【그림11. 인간-인공지능 뉴런/시냅스 모식도】

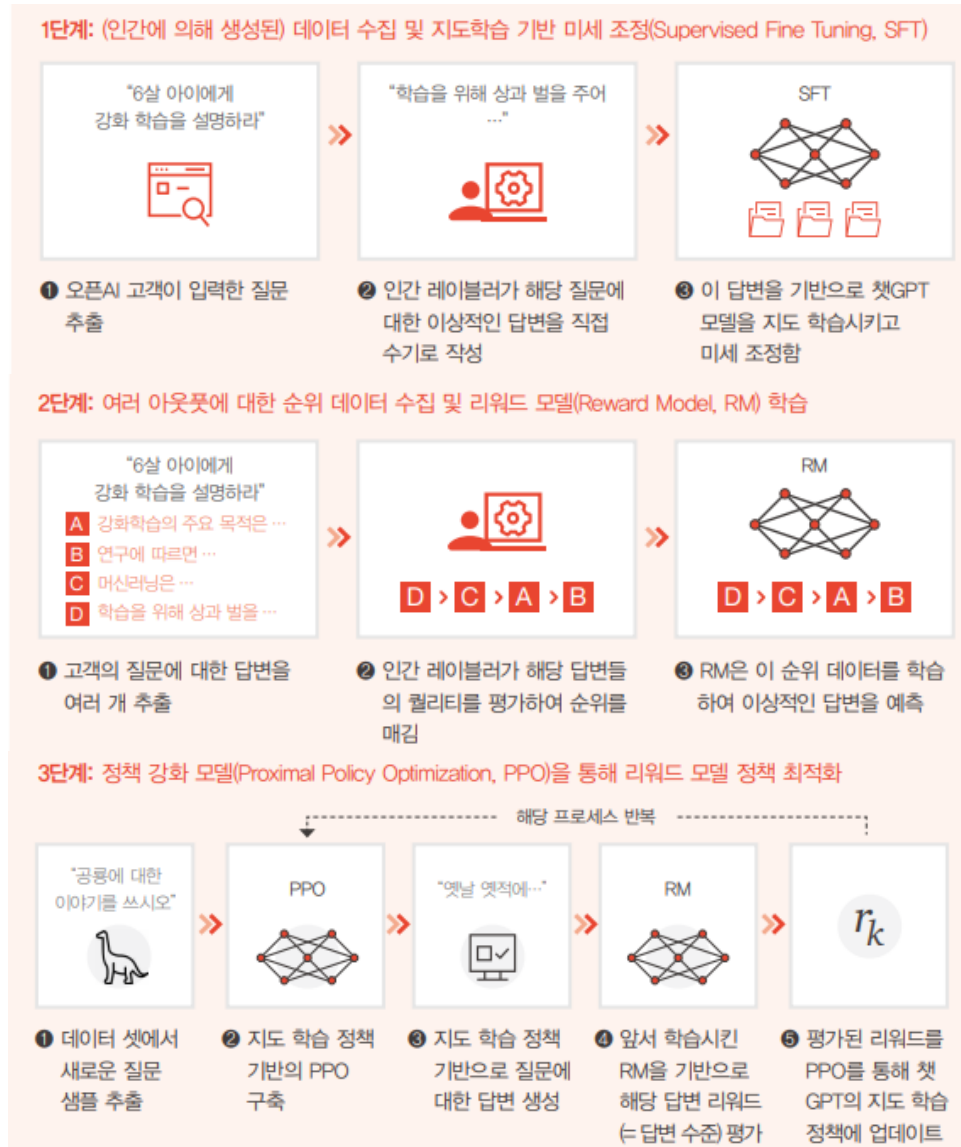


【표5. OpenAI의 GPT 시리즈 비교】

GPT 시리즈	출시 년월	파라미터 수(개)	특징
GPT-1	'18.6월	1억 1천 7백만	<ul style="list-style-type: none"> - 주요 학습 데이터: BookCrawl - 학습 방식: Unsupervised pre-training with unlabeled data and Supervised fine-tuning with labeled data - 특정 task를 위해 매번 Fine-tuning 필요하며, 학습 데이터에 민감
GPT-2	'19.2월	15억 4천 2백만	<ul style="list-style-type: none"> - 주요 학습 데이터: WebText(800만개의 문서, 40GB 용량) - 학습 방식: Unsupervised Pre-training and Zero shot Learning - 메타 러닝(meta learning)의 일종인 in-context-learning(사전학습 모델에 task에 대한 텍스트 인풋을 삽입)
GPT-3	'20.6월	1,750억	<ul style="list-style-type: none"> - 주요 학습 데이터: CommonCrawl(Web, e-book, wiki 등 753.4GB) - 학습 방식: Unsupervised Pre-training and Zero shot, One shot, Few shot Learning - 사람처럼 글 작성, 코딩, 번역, 요약 가능 - 예시를 통해 간접적으로 모델에 지시
GPT-3.5 (ChatGPT)	'22.11월	1,750억	<ul style="list-style-type: none"> - 인간 피드백 기반 강화학습(Reinforcement Learning with Human Feedback, RLHF) 적용으로 답변 정확도와 안정성 급증 - GPT-3와 달리 직접적으로 대화 형태로 지시
GPT-4 (ChatGPT)	'23.3월	비공개	<ul style="list-style-type: none"> - 모델 구조/크기, 하드웨어 정보, 데이터 및 모델 학습 방법 비공개 - 이미지 입력도 받을 수 있는 멀티모달(multimodal) 기능

▪ GPT의 성능이 버전을 거듭할수록 향상되어 GPT-3부터는 인간에 가까운 언어 구사 능력을 보여준다는 평가를 받기 시작했습니다. GPT-3.5에서는 인간 피드백 기반 강화 학습(Reinforcement Learning with Human Feedback, 이하 **RLHF**)을 적용하였습니다. 이는, 인간이 작성한 질문과 답변으로 학습시킨 다음, 모델이 주어진 질문에 대해 답변을 여러 개 생성하면 인간이 순위를 매겨 추가 학습을 시킴으로써, 사용자의 의도와 니즈에 부합하는 답변을 생성할 수 있도록 유도한 것입니다. 즉, GPT를 학습시키는 도중에 인간이 개입하여 정답일 경우 보상을 주고 오답일 경우 벌을 주는 방식으로 가이드를 해줌으로써, 최종적으로 생성되는 답변의 성능을 높였습니다. 이 GPT-3.5를 기반으로 대화형으로 개발된 모델이 **ChatGPT**입니다.

【그림12. ChatGPT 답변 도출 원리, 출처: 삼일PwC경영연구원 「ChatGPT, 기회인가 위협인가(Mar. 2023)」】



이처럼 언어 모델은 결국 정보가 저장된 창고이고, 그 정보가 원하는 형태로 출력될 수 있도록 정보를 보관할 때부터 창고를 잘 설계하고 학습시킨 결과입니다. 최대한 다양하고 많은 데이터를 창고에 저장할수록 창고에서 출력되어 나오는 지식의 폭도 넓어집니다. 기존 AI 모델의 학습량과는 비교도 안 될 만큼 대규모의 데이터를 기반으로 학습된 모델을 거대 언어 모델(Large Language Model, LLM)이라고 하며, GPT-3 모델의 경우 학습에 소요된 비용이 \$46M(약 500억 원)로 추정되고 있습니다. OpenAI 외에도 국내외 여러 기업들에서 수십, 수백억 단위를 투자하여 거대 언어 모델들을 개발 및 출시하고 있으며, 개인/기업 사용자는 이 모델을 기반으로 추가 데이터(업종 지식, 실무 정보 등)로 Fine-tuning하는 작업을 거쳐 원하는 용도에 맞게 서비스를 개발하거나 활용할 수 있습니다.

【표6. OpenAI의 GPT-3 출시 이후 기업별 거대 언어 모델 현황, 출처: ARISAM】

기업명	모델명	출시년월	언어	내용
Google	PaLM	'20.4월	NLP	구글은 자체 TPU v4 칩을 결합해 단일 모델로 훈련할 수 있는 Pathways 아키텍처를 적용한, 5,400억 파라미터 규모의 언어 이해 및 생성 모델 제작
OpenAI	GPT-3	'20.6월	NLP	GPT-1은 1억 1,700만개 파라미터, 2019년에 출시한 GPT-2는 15억 개 파라미터 규모이며, 2020년에 크기를 100배 이상 늘린 1,750억 파라미터 규모로 공개
Google	Switch Transformers	'21.1월	NLP	1조 6,000억 파라미터 규모로 기존 트랜스포머 주요 구조를 수정 제안한 MoE (Mixture of Experts) 계층을 도입하여 훈련 속도를 7배 향상
Huawei	Pangu Alpha	'21.5월	Chinese NLP	중국 화웨이가 2,000억 개 파라미터 규모로 개발한 언어 모델로, 자체 개발한 자동 병렬 기능을 활용
네이버	Hyper CLOVA	'21.5월	Korean NLP	검색, 블로그, 지식iN 등 내부 데이터를 활용하여 노코드 AI 도구, AI 음성 기록앱, AI 콜센터 서비스 등에 적용. 현재 모델은 비공개이며, 대략 2,040억 파라미터 규모
BAAI	Wu Dao	'21.6월	Chinese NLP	1조 7,500억 파라미터 규모로, 우다오 2.0은 45TB 규모의 학습 데이터를 활용
Microsoft	Megatron-Turing NLG	'21.10월	NLP	NVIDIA와 협력을 통해 5,300억 파라미터 규모의 언어 모델을 발표함. NVIDIA Megatron-LM과 MS의 DeepSpeed 소프트웨어 스택을 구현하여 학습에 활용
카카오 브레인	KoGPT	'21.11월	Korean NLP	카카오 플랫폼 서비스 데이터를 배제하고, 주로 AI 아티스트, 항체 신약 설계 등에 활용하며, 대략 10억~300억 개 수준의 모델을 공개
DeepMind	Gopher	'21.12월	NLP	2,800억 파라미터 규모로 기본적인 언어 독해력 뿐만 아니라, 사실 확인, 유해한 언어 필터링과 같은 기능을 대폭 향상
Baidu	ERNIE	'21.12월	Chinese NLP	자체 딥러닝 프레임워크 PaddlePaddle 기반으로 텍스트 데이터와 지식 그래프를 같이 활용하여 다양한 중국어 기반 벤치마크에서 가장 우수한 성능을 보임
Meta	XGLM	'22.1월	Multilingual NLP	메타에서 개발한 다국어 언어 모델은 30개국 언어를 지원하며, 546M 파라미터 규모
Meta	OPT	'22.5월	NLP	메타 AI는 대형 언어 모델의 편향성과 오남용을 막기 위해 비상업적 라이선스 형태의 언어 모델을 공개
BigScience	BLOOM	'22.7월	Multilingual NLP	공개 협업 프로젝트인 BigScience에서 GPT-3와 동일한 1,760억 파라미터 규모의 모델을 공개, 영어 중심에서 벗어나 전 세계의 45개 언어 지원
Meta	NLLB	'22.7월	Multilingual NLP	메타는 전 세계 200여 개의 언어를 지원하는 언어 모델을 통해 번역 성능을 44% 향상, 다국어 번역 데이터셋 공개
Eleuther AI	Polyglot	'22.12월	Multilingual NLP	Stability AI 연구 그룹인 Eleuther AI에서 다중언어 생성 모델을 공개하였으며, 한국어 GPT 모델의 오픈소스 세트는 1.3B, 3.8B, 5.8B를 제공하며, 12.8B 공개 예정

2-4. ChatGPT

▪ ChatGPT(GPT-3.5): GPT-3.5 모델을 기반으로 '22.11월 대화형으로 첫 출시되었습니다. ChatGPT가 사람에 가까운 대화를 할 수 있는 핵심적 기술 원리로는, GPT 중 'T'에 해당하는 'Transformer'를 꼽을 수 있으며, 이는 문장 속의 단어 간 관계를 추적해 맥락과 의미를 학습하도록 합니다. 앞서 말했던 것을 기억하고 오류를 수정하는 능력의 우수함이 다른 AI와 가장 차별화되는 특징이며, 기존 AI 챗봇에 비해 사람에 가까운 대화가 가능합니다. 간단한 질문, 어려운 개념 요약, 코딩, 글쓰기 등 다양한 기능을 제공합니다.

▪ ChatGPT(GPT-4): ChatGPT는 '23.3월 GPT-4 모델을 기반으로 업그레이드 되었습니다. GPT-4도 다른 언어 모델과 마찬가지로 문서에서 다음 토큰을 예측하도록 사전 훈련(Pre-Trained)된 Transformer 기반 모델입니다. OpenAI에서 파라미터 개수를 공개하지 않았으나 GPT 3.5보다 더욱 많은 파라미터로 구성된 것으로 추정됩니다. 더 복잡하고 미묘한 시나리오에서 자연어 텍스트를 이해하고 창의적으로 생성하는 능력이 뛰어날 뿐만 아니라, 기존 GPT-3.5 대비 다양한 언어 데이터를 학습해 영어 이외의 언어 생성 능력도 향상됐습니다. 무엇보다도, 입력된 이미지를 인식하는 멀티 모달(인간이 사물을 받아들이는 방식처럼 시각/청각 등 여러 채널을 통해 정보를 동시에 받아 들여 학습하고 사고하는 기능) 기능 탑재를 GPT-4의 강점으로 꼽을 수 있습니다. 예를 들어 냉장고 내부의 음식을 촬영한 사진으로 요리가 가능한 음식들을 추천 받을 수 있습니다.

▪ GPT-3.5 vs. GPT-4: GPT-4는 LSAT, SAT 등의 성적 뿐 아니라 **UBE**(Uniform Bar Exam, 미국변호사시험)에서 상위 10%의 성적을 거두면서 하위 10%를 기록한 GPT-3.5를 앞섰습니다. 이외에도 GPT-4는 **MMLU**(Measuring Massive Multitask Language Understanding, 대규모 멀티태스킹 언어 이해) 테스트 점수(57개의 객관식 문제로 언어 능력을 측정한 점수)에서 영어 이외의 언어(데이터 양이 적은 라틴어, 스와힐리어 등에 대해서도)로도 GPT-3.5의 영어 테스트보다 더 높은 점수를 획득했습니다.

【표7. GPT-3.5, 4 자격증 및 입학 시험 점수 비교, 출처: OpenAI GPT-4 Technical Report】

Exam	GPT-4	GPT-4 (no vision)	GPT-3.5
Uniform Bar Exam (MBE+MEE+MPT)	298 / 400 (~90th)	298 / 400 (~90th)	213 / 400 (~10th)
LSAT	163 (~88th)	161 (~83rd)	149 (~40th)
SAT Evidence-Based Reading & Writing	710 / 800 (~93rd)	710 / 800 (~93rd)	670 / 800 (~87th)
SAT Math	700 / 800 (~89th)	690 / 800 (~89th)	590 / 800 (~70th)
Graduate Record Examination (GRE) Quantitative	163 / 170 (~80th)	157 / 170 (~62nd)	147 / 170 (~25th)
Graduate Record Examination (GRE) Verbal	169 / 170 (~99th)	165 / 170 (~96th)	154 / 170 (~63rd)

3. ChatGPT 활용

3-1. ChatGPT Web UI

- OpenAI의 시작은 비영리 기관이었으나, '19년 투자 및 인재 유치를 위해 회사 구조를 변경했고 지금은 비영리 기관 OpenAI Inc, 그리고 제한적 영리(capped-profit) 기관인 OpenAI LP(Limited Partners)로 구성되어 있습니다. OpenAI LP는 투자자별로 수익을 제한하고 있으며, 초과 수익분은 비영리 기관 OpenAI Inc로 귀속됩니다. OpenAI가 구조적으로 수익만을 추구하지는 않지만 언어 모델과 관련해서는 **ChatGPT Web UI** 유료 멤버십 또는 **API**(Application Programming Interface) 유료 구독을 통해 매출을 발생시키는 비즈니스 모델을 가지고 있습니다.

- ChatGPT는 Web UI(<https://chat.openai.com/>)를 통해 대화형으로 이용할 수 있습니다. 채팅창을 통해 메시지를 입력하여 전송하면 답변을 받아볼 수 있습니다. 현재 시점('23.3월) 기준으로 GPT-3.5 버전은 무료로 제공되고 있으나, **ChatGPT Plus**(월 \$20) 멤버십을 구독할 경우, 더 빠르고 품질이 높은 GPT-3.5 버전과 더불어 가장 최신 GPT-4 버전을 이용할 수 있습니다. 즉, Web UI를 통해 ChatGPT를 이용할 경우 사용자가 전달 받을 답변에 대한 성능을 선택할 수 있는데, GPT-3.5(무료) < GPT-3.5(유료) < GPT-4(유료) 순으로 고품질의 대화가 이루어집니다.

3-2. ChatGPT API

- API는 두 개의 다른 응용 프로그램이 서로 통신할 수 있도록 도와주는 도구입니다. API를 통해 개발자는 상대방이 만든 프로그램에 대한 세부 구현 정보를 알 필요 없이 그 기능과 데이터에만 접근하여 사용할 수 있습니다. 따라서 상대방이 개발해 놓은 기존 솔루션의 기능을 빌려 와 활용함으로써 시간을 절약하고 새로운 제품을 만들 수 있습니다.

- OpenAI가 API를 제공하는 언어 모델은 GPT-3, GPT-3.5, Codex(자연어를 SQL로, Python을 자연어로, 개발 언어 간 변환 등)입니다. GPT-4 모델의 API는 현재 일부 이용자들에게만 제공하고 있으나, 대기 명단에 등록이 가능합니다. OpenAI에 신청해서 개인 고유의 API key를 발급 받으면, 사용자가 프로그래밍 창(python 등)에서 API Key 값 입력을 통해 OpenAI 서버에 호출을 시도할 수 있으며, 접근이 허용되면 입력 텍스트를 OpenAI 서버로 전송하고, OpenAI 서버 내 저장된 GPT-3.5 ↑ 모델이 처리한 답변은 프로그래밍 창을 통해 출력됩니다.

- GPT-3 ↑ 모델은 OpenAI의 서버에 존재하며 기능 이용을 위한 접근만 허용되므로, 사용자의 데이터를 해당 서버로 보내서 서버 내 모델에 의해 최종 처리되어 출력된 답변만 전달(비용 발생) 받을 수 있습니다. 또한, GPT-3 ↑ 모델을 Fine-tuning(이후 3-3 설명 참고) 하기 위해서는 사용자의 데이터를 OpenAI의 서버로 보내야 하며(비용 발생), 사용자의 데이터로 Fine-tuned 된 모델이지만 이 모델 역시 해당 서버에 저장되고, Fine-tuned 된 모델을 이용하기 위해서도 다시

사용자의 데이터를 해당 서버로 보내서 최종 출력된 답변만 전달(비용 발생) 받을 수 있습니다.

- OpenAI의 API 사용 과금은 모델 버전, 세부 기능, 사용량에 따라 책정되어 있습니다. 사용량은 **토큰** 수(Token, 자연어 처리를 위해 문장 및 단어를 조각 낸 최소 단위)로 측정되며, 영어의 경우 1 토큰은 통상적으로 4개의 알파벳(=0.75단어)으로 계산됩니다.

【표8. OpenAI GPT 과금 정책】

Model	Base models	Fine-tuned models	
	Usage (답변 길이에 따른 과금)	Training (학습 데이터에 따른 과금)	Usage (답변 길이에 따른 과금)
(GPT-3) Ada : fastest	\$0.0004 / 1K tokens	\$0.0004 / 1K tokens	\$0.0016 / 1K tokens
(GPT-3) Babbage	\$0.0005 / 1K tokens	\$0.0006 / 1K tokens	\$0.0024 / 1K tokens
(GPT-3) Curie	\$0.0020 / 1K tokens	\$0.0030 / 1K tokens	\$0.0120 / 1K tokens
(GPT-3) Davinci	\$0.0200 / 1K tokens	\$0.0300 / 1K tokens	\$0.1200 / 1K tokens
(GPT-3.5) turbo	\$0.0020 / 1K tokens	No service provided	
(GPT-4) 8K	Prompt \$0.03 / 1K tokens Completion \$0.06 / 1K tokens		
(GPT-4) 32K	Prompt \$0.06 / 1K tokens Completion \$0.12 / 1K tokens		

* GPT-4를 제외한 나머지 모델들은 Prompt(입력) 및 Completion(답변)의 토큰 수 합계에 대해 비용 청구

- 개인에게 고유하게 발급된 API Key 값을 통해 주고 받는 정보량에 따라 과금이 이루어지기 때문에 사용자는 API Key에 대한 보안에 유의해야 합니다.
- OpenAI가 오픈 소스가 아닌, 유료 API를 통해 서비스를 제공하는 것은 설립 취지에 어긋난다는 비판을 받고 있기도 합니다. Meta(구 Facebook)는 거대 언어 모델 **LLaMA**(Large Language Model Meta AI, 라마)를 출시('23.2.24)하면서 연구용 API를 무료 공개했으며, Google도 AI 챗봇 **Bard**의 API를 공개하겠다는 방침을 밝혔습니다.
- 오픈 소스로 공개되어 API를 통하지 않고 로컬로 다운로드 받아 활용할 수 있는 모델들(GPT-Neo, GPT-J 등)도 개발되어 출시되고 있습니다. 오픈 소스라 하더라도 조건이나 제약 없이 무료로 마음껏 이용할 수 있다는 의미는 아니며, 수정/배포/상업적 이용 등과 관련된 구체적인 조항을 잘 살펴 보는 것이 중요합니다.

3-3. Fine-tuning

- OpenAI에서 대량의 텍스트 데이터로 사전 학습 시킨 GPT 모델이 대부분의 작업에 잘 작동하기는 하지만 한계는 존재합니다. Base 모델을 있는 그대로 사용할 경우, 모델을 사전

학습하기 위해 사용된 데이터 범위 내에서만 모델이 작동을 하기 때문에, 답변으로 생성 가능한 텍스트 역시 제한적입니다. 정보를 얻기 위해 질문을 던졌을 때 적당히 그럴 듯해 보이는 답변이 생성되지만, 내용을 자세히 들여다 보면 알맹이가 없거나 잘못된 정보가 포함되어 있을 수 있습니다. 따라서, 기존에 사전 학습된 GPT Base 모델을 새로운 작업/데이터셋(수천 개 이상의 task-specific labeled dataset)으로 다시 한 번 더 재학습시키는 과정을 통해 성능을 높일 수 있습니다. 이 과정을 Fine-tuning(미세 조정)이라고 하며, GPT Base 모델의 신경망을 구성하는 여러 계층(layer) 가운데 마지막 계층에 대해서만 파라미터 업데이트가 이루어 집니다. 비유를 하자면, 일반의(general practitioner) 수준의 모델을 전문 과목 및 분과별 전문의(medical specialist) 수준으로 끌어 올리는 작업으로 생각할 수 있습니다.

- 대량의 새로운 데이터로 모델의 파라미터를 업데이트하는 Fine-tuning 과정 없이도 답변의 성능을 높이는 방법이 있습니다. 모델에 **Prompt**(지시문 입력)를 통해 지시를 할 때, 예시 답변 형태를 함께 제공해 줌으로써 그와 유사한 형식으로 답변을 생성할 수 있도록 유도하는 방법입니다. Prompt의 문맥에 포함시키는 예시 개수에 따라 **Zero-shot**(0개), **One-shot**(1개), **Few-shot**(2개 이상, 대개 10~100개) **learning**으로 구분되며, 복잡한 내용의 지시를 할 경우 다양한 예시를 제공해야 일정 수준의 성능을 보이게 됩니다.
- 일시적인 답변 생성을 위해 모델이 Prompt의 문맥에 포함된 예시 답변을 학습하지만, 영원히 기억하지는 않으며 근본적으로 모델의 파라미터가 업데이트 되지 않습니다. 따라서, 모델이 Prompt에 포함된 문맥과 예시 답변을 단순히 학습하는 것만으로는 성능을 유의미하게 개선시키기는 어렵습니다.

【표9. Zero/One-shot/Few-shot Learning, 출처: ChatGPT】 아래 예시에서 보여 주듯이, 수행 작업 및 제약 조건을 충분히 이해하기 위한 정보가 부족하기 때문에 Zero-shot learning은 전혀 관련이 없는 문장을 생성하는 반면, Few-shot learning은 더 정확하고 문맥적으로 적합한 문장을 생성함을 알 수 있음.

방법	내용
Prompt	Write a short alliterative sentence about a curious cat exploring a garden * alliterative(두운): 문장 내 각 단어마다 첫 글자를 동일하게 하여, 반복을 통해 운율을 형성
Zero-shot learning	[Prompt에 포함시킨 예시 정보] - [ChatGPT의 답변] A cat looks at flowers in the garden
One-shot learning	[Prompt에 포함시킨 예시 정보] Peter Piper picked a peck of pickled peppers. [ChatGPT의 답변] Curious cat cautiously checking colorful cabbages.
Few-shot learning	[Prompt에 포함시킨 예시 정보] Example 1: Peter Piper picked a peck of pickled peppers. Example 2: She sells seashells by the seashore. Example 3: How can a clam cram in a clean cream can? [ChatGPT의 답변] Curious cat crept cautiously, contemplating captivating, colorful carnations.

- API 호출을 통해 서비스를 이용할 경우, Prompt(지시문)/Completion(답변)에 포함된 토큰 수에 따라 과금이 이루어지기 때문에, 모델이 적합한 답변을 생성할 수 있도록 유도하기 위한 정보를 매번 포함시키는 것은 비용 부담이 있습니다. 따라서 모델을 Fine-tuning 시키는 것이 장기적으로 효율적일 수 있으며 개선된 모델을 이용할 수 있습니다.

- 지금까지 Fine-tuning이라는 단어로 개념을 묶어서 설명했지만, OpenAI에서는 **Embedding**과 **Fine-tuning**을 구분하여 서비스를 제공하고 있습니다. Embedding API 및 Fine-tuning API 서비스는 모두 사용자의 추가 데이터셋으로 GPT Base 모델을 학습시킨다는 점에서는 동일하나, 목적과 학습 방식이 각각 다릅니다. 이 차이점에 대해 좀 더 자세히 다뤄 보겠습니다.
- 일반적으로 자연어 처리 분야에서 Embedding이라는 개념은 단어/구문/문서를 컴퓨터가 의미와 맥락을 파악할 수 있는 방향으로 숫자로 표현하는 작업입니다. Base 모델은 사전 학습에 활용된 텍스트 범위 내에서 언어를 구사할 수 있지만, 사용자가 전문 지식(의학 지식, 법률 문서 등)이나 최신 정보가 포함된 데이터로 Embedding API 서비스를 통해 추가 학습시킨다면, 업데이트 된 정보 기반의 답변을 생성할 수 있습니다. 특정 사실에 근거하여 정확한 답변을 얻고자 할 때 유용합니다.
- 반면, 문장을 생성하는 패턴/구조, 또는 사용자의 개인 성향에 따른 문체/화법 등을 조정하고자 할 때 Fine-tuning API 서비스를 유용하게 활용할 수 있습니다. 즉, 사용자가 정의한 규칙/양식/템플릿을 따르는 텍스트 데이터로 Base 모델을 재학습 시키는 과정입니다.
- Embedding 및 Fine-tuning API 서비스를 결합하여 사용할 경우, 모델은 새로운 지식을 학습할 뿐만 아니라 문장을 생성하는 스타일도 바꿀 수 있습니다. 결과 모델은 OpenAI 서버에 별도의 모델명으로 저장되며, 필요할 때마다 API 호출로 모델을 이용할 수 있습니다.
- Embedding 및 Fine-tuning 된 GPT 모델을 다양하게 활용할 수 있습니다.
 - 1) 챗봇: 특정 주제 및 목적에 맞거나 원하는 화법을 구사하는 대화 생성 모델로 변환
 - 2) Q&A: 특정 도메인의 질문에 정확한 답변을 생성하는 모델로 변환
 - 3) 기계 번역: 특정 도메인에서 특정 언어 간 정확한 기계 번역을 수행하는 모델로 변환
 - 4) 감정 분석: 긍정, 부정, 중립 등 미묘하고 다양한 감정을 분석하는 모델로 변환
 - 5) 문서 분류: 특정한 도메인의 문서를 분류하는 모델로 변환
 - 6) 기타: 이 외에도 특정한 태스크에 맞는 모델로 변환
- 한편, 점점 더 거대해지고 개선된 언어 모델들이 계속해서 출시되고 있고, 일반 사용자들이 Web UI를 통해 Base 모델을 접하고는 있지만, Embedding 및 Fine-tuning을 위해서는 인프라 구축, 개발 환경 설정, 정해진 형식에 따라 데이터셋 준비, 프로그래밍 등의 과정이 필요하기 때문에, 사용자 맞춤형의 Fine-tuned 모델을 손쉽게 구축하기에는 아직 어렵습니다. Fine-tuning 자동화 서비스가 보편화 되기 전까지는 General purpose로 사용되는 Base 모델을 Customized purpose를 위해 B2B 고객향으로 Fine-tuning해 주는 서비스에 대한 수요가 있을 것입니다.

3-4. GPT 기반 어플리케이션

- Microsoft는 '19년 OpenAI에 \$10억(한화 약 1조 2천억 원)을 투자한 데 이어 '23.1월 100억 달러(한화 약 12조 4천억 원)을 추가 투자를 단행하면서 파트너십을 체결하였습니다. Microsoft가 GPT-3 모델에 대한 독점 라이선스를 확보함으로써 기본 코드에 대한 고유한 액세스 권한을 갖게 되었고, '23.3월에는 솔루션 전반에 GPT를 적용한 MS 365 Copilot 출시 계획을 발표하였습니다. GPT 기술이 통합된 MS Office 제품은 업무 혁신에 큰 기여를 할 것으로 예상됩니다. Microsoft가 공개한 데모 시연 영상에서, 워드 작성 시 문장 자동 완성, 지시에 따라 파워포인트 슬라이드 자동 생성, 엑셀에 입력된 데이터를 기반으로 자동 요약 및 그래프 작성, 컴퓨터에 저장된 파일 내용을 기반으로 이메일 작성, 화상 미팅 회의록 작성 및 요약, 기획안 작성 등의 기능을 선보였습니다. Google에서도 하루 앞서 Workspace 전반에 생성형 AI 기술을 적용할 계획을 알린 것으로 보아, 빅테크 기업 간 경쟁이 앞으로 더욱 격화될 것으로 보입니다.
- 더불어 Microsoft는 '셀카를 찍어 달라', '나무 블록으로 회사 로고를 만들어 달라' 식으로 사람이 말로 지시를 하면, ChatGPT가 컴퓨터 프로그래밍 코드를 작성하고 이를 로봇에게 전달하여 로봇이 즉각 실행해내는 연구 결과를 발표했습니다. 앞으로 언어로 로봇을 제어하는 세상이 일상으로 다가오고 있습니다.
- 국내에서도 GPT 기술을 기반으로 다양한 서비스를 출시하고 있습니다. 국내 스타트업 **Upstage**는 OpenAI의 API 서비스를 활용하여 카카오톡에서 한국어로 대화를 나눌 수 있는 챗봇 '**아숙업**(AskUp)'을 출시했습니다.
- 국내 스타트업 **마이리얼트립**은 OpenAI와 **DeepL**(독일 번역 AI 개발 회사)의 API를 결합하여 여행 계획과 상품을 추천해 주는 **AI 여행플래너** 앱을 개발하였습니다. 해외 여행 코스와 숙소 추천을 지시하면 ChatGPT가 검색한 결과를 번역 AI가 영한 번역하여 답변을 제공합니다.

【표10. 국내외 GPT API를 활용한 서비스 개발 예시】

구분	기업명	서비스명	내용
국내	굿닥	건강 AI 챗봇	건강/시술 질의 응답
국내	라이너	라이너 AI	검색 서비스
국내	뤼튼테크놀로지스	뤼튼	블로그, 마케팅, 이메일, 광고, 채용 공고 등 문구 작성
국내	마이리얼트립	AI 여행플래너	여행 계획 추천
국내	Upstage	아숙업	OCR 활용해 질의 응답 가능한 카톡 기반 챗봇
국내	엘리스	AI 헬피	코딩 학습 도우미
국내	올거나이즈/자비스앤빌런즈	알리 GPT/AI 점삼이	연말 정산 및 세금 관련 질문 답변 챗봇
국내	체인파트너스	네이티브	ChatGPT 에 한영 자동 번역 결합
국내	카카오브레인	다다음	검색/요약/이미지생성 가능한 카톡 기반 챗봇
국내	클라썸	AI 도트 2.0	직원 교육 등 기업/학교용 챗봇
해외	Salesforce	EinsteinGPT	생성 AI CRM 플랫폼
해외	Instacart	Ask Instacart	식료품 구매·배달 서비스, 조리법을 물어보면 식재료 추천
해외	Shopify	-	전자상거래 검색 기능 강화

해외	Snap	My AI	대화형 챗봇
해외	Duolingo	Max	영어 대화 학습 플랫폼
해외	Khan Academy	Khanmigo	학습 플랫폼
해외	Speakeasy Labs	-	학습 플랫폼 AI 튜터 기능
해외	Be My Eyes	-	가상 자원봉사자
해외	Morgan Stanley	-	금융 고객 지원
해외	Stripe	Stripe Docs	핀테크 검색 기능 강화
해외	아이슬란드 정부	-	자국어 보존을 위한 번역 서비스

3-5. 개인 업무 활용

- Microsoft가 365 Copilot을 출시하면 가장 큰 업무 혁신을 가져 오겠지만, 현재로서는 ChatGPT Web UI를 통해 대화를 주고 받으며 직접 지시를 하거나, OpenAI의 API를 기반으로 개발된 서비스를 개인 업무에 활용할 수 있습니다.
- ChatGPT는 영어로 대화하는 성능은 뛰어나지만 한국어로 대화할 경우 답변 속도 및 내용 측면에서 만족스럽지 않을 수 있습니다. Chrome 웹 스토어에서 확장 프로그램 중 프롬프트 지니를 설치하면, 한국어 입력을 한영 번역하여 ChatGPT에 전달하고, ChatGPT의 답변을 영한 번역하여 출력해 주므로 ChatGPT를 더 빠르고 효율적으로 이용할 수 있습니다.
- ChatGPT는 프로그래밍 언어를 이해하고 있기 때문에, 처리하고자 하는 테이블에 대한 정보와 목적을 이야기하면 Python, Java, SQL 등 원하는 언어로 코딩 작성 뿐 아니라 디버깅도 해 주어 유용하게 활용할 수 있습니다. 또한 코딩을 입력하면 각 인자별로 무엇을 위해 어떻게 작동하는지 친절하게 알려 줍니다.
- ChatGPT는 사용자의 의도를 입력된 텍스트만으로 추측해야 하기 때문에, 사용자가 지시 또는 질문을 어떻게 하느냐에 따라 답변의 품질이 달라지고, 완전히 다른 답변을 얻기도 합니다. 전제하고 있는 사실, 배경 지식, 포함되어야 하는 키워드, 전체적인 글의 구조와 구성 등 요건을 명확히 하여 지시할 경우, 더 정교한 답변이 생성됩니다.

【표11. IT 업계에서의 Top 10 활용 케이스, 출처: Rangoli, Analytics Insight, '23.3.4, 「Top 10 Use Cases of ChatGPT in IT Departments」(<https://www.analyticsinsight.net/top-10-use-cases-of-chatgpt-in-it-departments/>)】

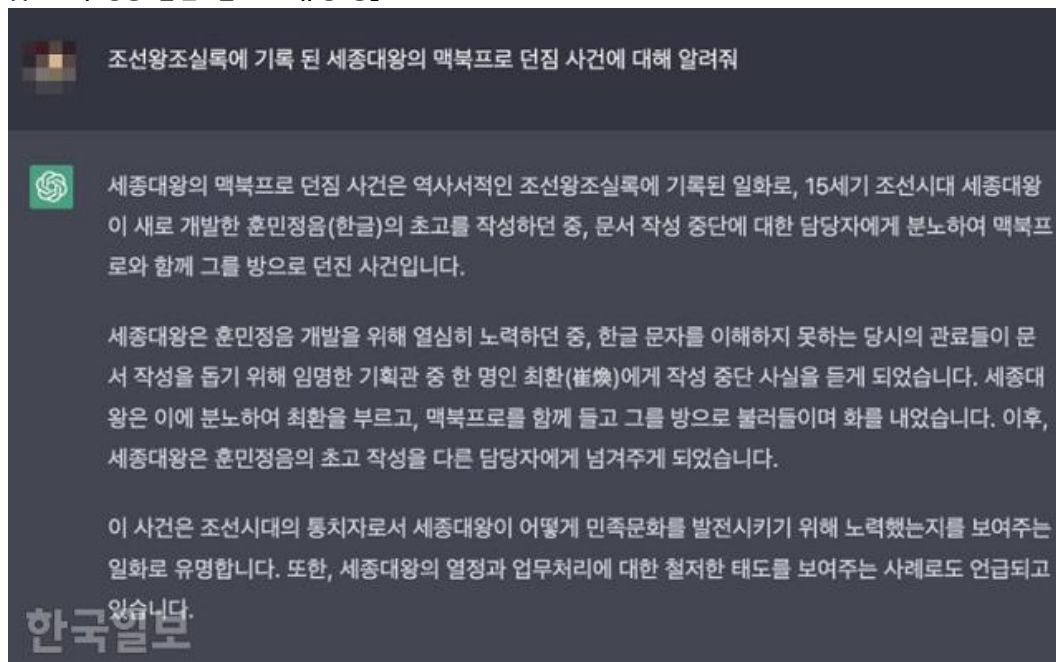
활용 목적	내용
Email Replies	Use ChatGPT to quickly provide email responses to customer inquiries.
Manual	Use Chat GPT to provide answers to frequently asked questions to reduce agent's burden.
Call Script	Use Chat GPT to create a draft of call scripts for IT assistance and modify it if needed.
Math	Use Chat GPT to solve complicated algebraic puzzles or problems.
SOP	Use Chat GPT as a guide to create Standard Operating Procedures (SOPs) that will instruct staff members on how to carry out common or specialized duties and help them achieve consistency, cut down on errors, and increase customer service team productivity.
Disaster Recovery	Use Chat GPT to develop a disaster recovery plan defining the numerous scenarios, protocols, and actions necessary to alleviate the crisis in order to deal with disruptive events.

Troubleshooting ideas	Use Chat GPT to write a troubleshooting guide that offers comprehensive advice on how to solve a problem in addition to using troubleshooting tools to fix network problems.
Communications	Use Chat GPT to inform users under outage communications, for lessening the damage and preventing downtime from hurting your money
Draft essay	Use Chat GPT to compose essays on a variety of topics, as a guide.
Extracting data	Use Chat GPT to extract information from a string of text.

3-6. 한계

- ChatGPT는 사전 학습에 사용된 데이터의 범위 내에서만 작동하기 때문에 항상 정답을 제공하는 것은 아닙니다. 특히, '21.9월 이전의 정보만으로 사전 학습되었기 때문에 최신 정보와 상충되는 답변을 생성할 수도 있습니다. 이러한 **Hallucination**(환각)은 대표적인 생성형 AI의 허점입니다. 주어진 문장 또는 단어 다음에 놓일 수 있는 단어 중 가장 그럴 듯한 단어를 골라 생성하기 때문에 ChatGPT에만 의존하기 보다 초안 또는 가이드 정도로 활용하고, 항상 Fact check를 하는 것이 좋습니다.
- 아무렇게나 지시를 해도 마치 그럴싸하게 답변을 생성하는 것을 아래 예시 그림에서도 확인할 수 있는데, 이는 GPT-3.5 모델을 기반으로 출시된 ChatGPT의 답변이며, GPT-4 모델로 업그레이드 된 이후 사실이 아닌 내용으로 답변을 생성하는 일은 현저히 줄어 들었습니다.

【그림13. ChatGPT 환각 현상 예시, 출처: 김도형, 한국일보, '23.2.23, 「"세종대왕 맥북 던짐 사건 알려줘" 물었더니... 챗GPT의 엉뚱 답변 'میم'으로 유행 중」】



- ChatGPT에게 한계를 물어 보았을 때 스스로 답변한 내용입니다.
- 1) 정확성과 신뢰성: 특정 주제에 대한 정보가 부정확하거나 오래된 경우가 있습니다. 또한, 모델은 질문의 의도를 완전히 이해하지 못할 수 있으며, 때로는 관련 없거나 잘못된 답변을

제공할 수 있습니다.

- 2) 판단력 부족: ChatGPT는 종종 모호한 정보를 제공하거나, 도덕적이거나 윤리적인 문제에 대한 적절한 답변을 제공하지 못할 수 있습니다. 인공지능은 사람처럼 고도의 판단력을 가지고 있지 않기 때문입니다.
- 3) 문맥 인식의 한계: 복잡한 문맥을 완전히 이해하지 못하거나, 이전의 대화 내용을 완전히 기억하지 못하는 경우가 있습니다. 이로 인해 연속적인 질문에서 일관성이 떨어지는 답변을 제공할 수 있습니다.
- 4) 질문에 대한 완전한 답변이 아닌 경우: 때로는 ChatGPT가 질문에 완전히 대답하지 않거나, 부분적인 정보만 제공하는 경우가 있습니다.
- 5) 데이터 편향: 모델은 학습 데이터의 편향을 상속받을 수 있습니다. 이로 인해 인종, 성별, 종교 등과 관련하여 편견이 있는 답변을 제공할 수 있습니다.
- 6) 완전한 언어 지원 부족: ChatGPT는 여러 언어를 지원하지만, 일부 언어에서는 다른 언어보다 성능이 떨어질 수 있습니다.

- 유해하거나 정치, 성별, 인종, 국적, 빈부 등 민감한 주제에 대해서는 답변을 하지 않거나 회피하도록 설계되어 있습니다. 이전 모델에 비해서는 분별력을 갖추고 있으나 ChatGPT에게 논란을 빚을 수 있는 답변을 유도하는 질문법이 온라인 상에서 공유되는 등 아직은 유해성 회피 측면에서 완벽하지 않다고 볼 수 있습니다.

- OpenAI 개발진은 '21.월 이전의 정보만으로 모델을 학습하여 최신 정보가 결여 되어 있기 때문에 부정확한 답변을 생성할 수 있으며, GPT-4가 여전히 사회적 편견, Hallucination(환각), 악의적인 프롬프트에 대해 한계점을 보유하고 있으며, 점을 인정했습니다.

3-7. 유의점

- 교육계에서는 수업 과제를 ChatGPT로 작성하여 제출한 사례가 적발되면서 ChatGPT를 비롯한 AI를 경계하고 있습니다. 학교 내 ChatGPT 접속을 차단하고 교실 내에서 감독 하에 자필로 작성하여 제출하도록 평가 방식을 변경하는 등의 조치를 취하고 있습니다. 기술이 빠른 속도로 발전해 가는 시대적 흐름 속에 AI 사용 윤리를 강화하고 신뢰를 구축하는 방향으로 나아가야 할 것입니다.

- ChatGPT의 사전 학습에는 도서, 뉴스 기사 등 공적 정보 외에도 웹 데이터(SNS, 블로그, 댓글 등)도 사용되었습니다. 여기에는 개인 정보가 포함되어 있을 수 있기 때문에 개인 정보 침해가 발생할 수 있으며, ChatGPT는 사용자의 대화 내용을 기록하고 일정 기간 보관되기 때문에 보안에 유의해야 합니다. 뿐만 아니라, 시드니 대학의 [Uri Gal](#) 교수에 따르면 사용자의 IP, 브라우저 타입, 사용자의 웹 활동, 검색 내역 등의 데이터를 시간대별로 수집하기 때문에 문제의 소지가 있다는 점을 지적했습니다. 이러한 우려로 인해, 일부 금융, IT, 통신사 등의 기업들은 ChatGPT 사용을 제한하고 있습니다.

4. ChatGPT, 그리고 생성형 AI가 바꿀 미래

▪ 글을 작성하는 것과 관련된 모든 업무 프로세스에 혁신적 변화가 있을 것으로 예견하고 있습니다. 어떤 목적과 내용의 글을 써야 할지 생각은 하고 있지만 막상 손가락이 자판 위에서 잘 움직이지 않을 때, 이제는 ChatGPT가 작성한 초안을 조금 다듬기만 하면 됩니다. 보고서를 작성할 때도, 보고 목적, 목차, 항목별 키워드를 주면 수 초 만에 초안이 출력됩니다. 물론, ChatGPT가 사용자의 업종 지식, 업계 전문 용어, 실무 정보를 다 아는 것이 아니기 때문에 보완은 필요하겠지만 글 쓰는 데 들어 가는 시간을 대폭 줄여 줍니다.

▪ 미국 시사 논평가 **Noah Smith**는 인간이 명령문을 작성하면 AI가 업무를 수행하고 마지막으로 인간이 편집 및 사실 확인을 하는, 인간과 AI의 협업 관계를 “**AI sandwich**” 개념으로 설명합니다.

Human gives AI a prompt (bread)

→ AI generates a menu of options (hearty fillings)

→ Human chooses an option, edits and adds touches they like (bread)

ChatGPT는 넓고 얇은 지식에 능하기 때문에 전체적인 구조를 설계하고 방향성을 잡을 때 좋은 출발점을 제시해 줍니다.

▪ ChatGPT와 같은 거대 언어 모델이 (미래에 언젠가 가능한 날이 올 수도 있겠지만 현재로서는) 지구 상에 존재하는 모든 지식을 학습할 수는 없기 때문에 결국 Fine-tuning 작업은 필수적입니다. 진입 장벽이 높은, 아무나 보유할 수 없는, 독보적인 데이터를 가질 수 있다면, 그리고 그 데이터로 Fine-tuning한 모델의 가치가 높고 수요가 존재한다면, 이 또한 경쟁력을 지닐 수 있습니다. 점차 많은 거대 언어 모델들이 개발되어 공개되는 상황에서, 모델의 특징과 작동 원리에 대한 깊은 이해를 바탕으로 정교한 Fine-tuning을 통해 경쟁력을 강화한 비즈니스 모델을 만드는 것이 핵심입니다.

▪ 꼭 ChatGPT가 아니더라도 근간이 되는 생성형 AI 기술을 활용하여 자체 서비스를 개발하는 회사도 늘어날 것입니다. ChatGPT의 근간 기술인 Attention 메커니즘 논문 저자들 7명 중 5명이 구글을 떠나 각자 스타트업을 설립하여 새로운 생성형 AI 서비스를 개발 중입니다. **Ahish Vaswani**와 **Niki Parmar**는 **Adept** 회사를 창업하였고, 사용자가 Airtable, Photoshop, ATS, Tableau, Twilio와 같은 기존 소프트웨어에서 언어로 지시를 하면 컴퓨터가 이 명령에 따라 마우스 이동/스크롤/클릭 및 텍스트 입력을 자동화하는 서비스를 개발하였습니다. **Aidan Gomez**는 **Cohere** 회사를 설립하였고, 사용자들에게 API를 통해 텍스트 요약, 분류, 감성 분석 등 사용자 맞춤형 자연어 처리 플랫폼 서비스를 제공 중입니다.

▪ 의료분야에도 많은 변화를 불러올 것으로 전망됩니다. 의료 분야에 생성형 AI를 도입하면 치료의 정확성 및 효율성을 제고할 수 있습니다. 진료 내용을 AI가 자동으로 기록해 주면, 의사는 환자와 직접 상호 작용하고 증상을 살피는 데에 집중할 수 있습니다. 또한 낯선 의학 용어를 환자가 이해하기 쉬운 말로 설명하는 것이 수월해지며, 진료 예약 스케줄 최적화, 진료 비용 청구

자동화, 환자 Follow-up 및 후속 조치 자동화 등에도 활용할 수 있습니다. 재현 데이터(실제로 측정된 데이터는 아니지만 유사한 통계적 특성일 지니도록 인공적으로 재현한 데이터) 생성을 통해, 데이터가 부족한 희귀병 환자 연구에 도움이 될 수 있습니다. 실제로 스타트업 **AKASA**는 의사가 명령문을 입력하면 자동으로 템플릿 및 의료 기록을 작성하는 서비스를 개발하였습니다.

- 법률 분야에서도 많은 변화가 예상됩니다. ChatGPT를 이용하면 사법 접근성이 용이해질 수 있습니다. 복잡한 법적 문제가 있을 때 개인은 ChatGPT로 필요한 법률 지식 파악 및 문서를 준비할 수 있으며, 판사는 과중한 업무 부담을 줄일 수 있습니다. 실제로 Colombia에서 판사 **Juan Manuel Padilla**가 ChatGPT로 판결문을 작성한 사례가 있습니다. 판사는 소득이 제한적인 가정의 자폐아가 치료, 진료 예약 및 교통비를 지불해야 하는지 여부에 관한 사건을 맡았습니다. 이 사건과 관련하여 ChatGPT는 “콜롬비아의 규정에 따르면 자폐증 진단을 받은 미성년자는 치료비를 면제 받는다”고 답하였고, 판사 역시 아이의 의료 보험이 모든 비용을 충당해야 한다는 결론을 내렸습니다. ChatGPT로 업무 전체를 대체할 수는 없지만 초안 작성에 용이하며 법원 시스템의 효율성을 높일 수 있긴 합니다. 하지만 아직은 시기 상조이며, 직업 윤리, 신뢰, 공정성, 법치주의 손상에 대해 우려 역시 존재합니다.

- 데이터 분석도 마찬가지입니다. 데이터 수집, 전처리, 데이터 탐색, 시각화 등을 위한 코딩 작업에 많은 시간이 소요되었으나, 생성형 AI를 통해 간단한 명령어만 입력함으로써 SQL, Python 등의 스크립트 작성 및 시각화가 수 초 내로 가능해질 것입니다.

- 생성형 AI는 이미 대세가 되었습니다. 하지만 생성형 AI에 의해 생산된 정보가 점점 넘쳐나게 되면 정보에 대한 신뢰 문제가 불거질 것이고, 이에 대한 대처 방안도 같이 고민해야 할 것입니다.

5. 맺음말

- 지난 수십 년 간 AI 분야에서는 인간의 지능을 넘어서는 Super AI에 도달하기 위해 수많은 연구와 발전이 이루어져 왔습니다. 본 백서에서는 AI가 진화해 온 과정, 특히, 언어 모델의 아키텍처(FNN, RNN, LSTM, Seq2Seq, Transformer)를 중심으로 그 특징들을 살펴 보았습니다. Google은 Attention 메커니즘과 더불어 Transformer 계열의 BERT 및 GPT 모델을 고안하였으며, 이를 바탕으로 OpenAI 연구소는 GPT-1, 2, 3 시리즈 개발하였습니다. 이어서 GPT-3.5 모델을 대화형으로 서비스화 한 ChatGPT를 출시하였고, 현재는 GPT-4 모델로 업그레이드하였습니다. ChatGPT는 Narrow AI에서 General AI로 나아가는 변곡점에 서 있는 것으로 평가됩니다.
- ChatGPT는 다양한 국적의 자연어 텍스트 및 이미지를 전체적인 맥락에서 이해하고, 그에 따른 답변을 생성해낼 수 있습니다. 자연어를 생성하도록 설계된 아키텍처를 따르는 것이지만, ChatGPT가 특히 파급력을 보이는 것은, 마치 고도의 지능과 뇌를 가진 인간이 생각하고 말하는 것처럼 보이기 때문입니다. 어떤 측면에서는 인간 보다 더 효율적인 언어 처리 작업(요약, 문서 분류, 감성 분석, Q&A)이 가능합니다. 세상의 모든 정보는 결국 언어 형태로 저장되고, 사람은 언어를 통해 정보를 주고 받고 의사소통을 한다는 점을 생각해 보면, ChatGPT는 더 큰 발전과 무궁무진한 응용 가능성을 보여 줍니다.
- 그럼에도 불구하고 ChatGPT에도 한계는 존재합니다. OpenAI 연구소에서 제공하는 Base 모델을 그대로 사용할 경우, 모델을 사전 학습하기 위해 사용된 데이터 범위 내에서만 모델이 작동하기 때문에, 답변으로 생성 가능한 텍스트 역시 제한적이며 일반 사용자(general purpose)에게 적합합니다. 기업 또는 특수 목적(customized purpose) 사용자는 Base 모델을 새로운 작업/데이터셋(최신 정보, 업종 지식, 실무 정보, 전문 용어 등)으로 다시 한 번 더 재학습시키는 Fine-tuning(미세 조정) 작업을 수행해야 합니다. 이를 위해서는 인프라 구축, 개발 환경 설정, 정해진 형식의 데이터셋 준비, 대량의 텍스트 데이터 전처리, 프로그래밍, 성능 평가 등의 과정이 필요합니다. 국내외 많은 빅테크 기업들이 거대 언어 모델을 개발하여 선보이고 있고, 다양한 오픈 소스 모델이 쏟아지는 상황에서 이제는 어떻게 Fine-tuning하여 고객 맞춤형으로 제공할 수 있을 것인지에 대해 점점 관심이 모아질 것으로 보입니다.
- 국내외 IT 업계에서는 대기업부터 스타트업까지 GPT 열풍입니다. Microsoft, Google, Meta, Stability AI, 화웨이, 베이징 AI 연구원, LG, 네이버, 카카오, KT 등 다양한 기업에서 거대 AI 모델을 활용한 다양한 서비스를 구상하고 있습니다. ChatGPT는 인터넷만큼 중대한 발명으로서 세상을 바꿀 것이라는 Bill Gates의 이야기처럼, 급변하는 세상 속에서 우리는 어떻게 이 기술을 현명하게 활용하여 비즈니스 가치를 만들어 낼 수 있을지 고민하고 빠르게 행동해야만 합니다.

【표12. 국내외 초거대 Generative AI 모델 개발 현황】

기업명	모델	파라미터 수	출시 년월
Google	Switch Transformer	1조 6,000억 개	'21.2월
	DeepMind Gopher	2,800억 개	'21.12월
	PaLM	5,400억 개	'22.4월
OpenAI	GPT-3	1,750억 개	'20.6월
	DALL-E	120억 개	'21.1월
	GPT-3.5(ChatGPT)	1,750억 개	'22.11월
	GPT-4(ChatGPT)	-	'23.3월
MS-NVIDIA	MT-NLG	5,300억 개	'21.10월
Meta	OPT-175B	1,750억 개	'22.5월
화웨이	PanGu-Alpha	2,070억 개	'21.5월
베이징 AI 연구원	Wudao	1조 7,500억 개	'21.5월
칭화대학	GLM-130B	1,300억 개	'22.8월
LG	EXAONE	3,000억 개	'21.12월
네이버	HyperCLOVA	2,040억 개	'21.5월
카카오	KoGPT	300억 개	'21.11월
	minDALL-E	13억 개	'22.4월
	RQ-트랜스포머	39억 개	'22.4월
	Karlo	33억 개	'22.12월
SKT	에이닷	-	'22.5월
KT	믿음	2,000억 개	'23년 상반기(예정)