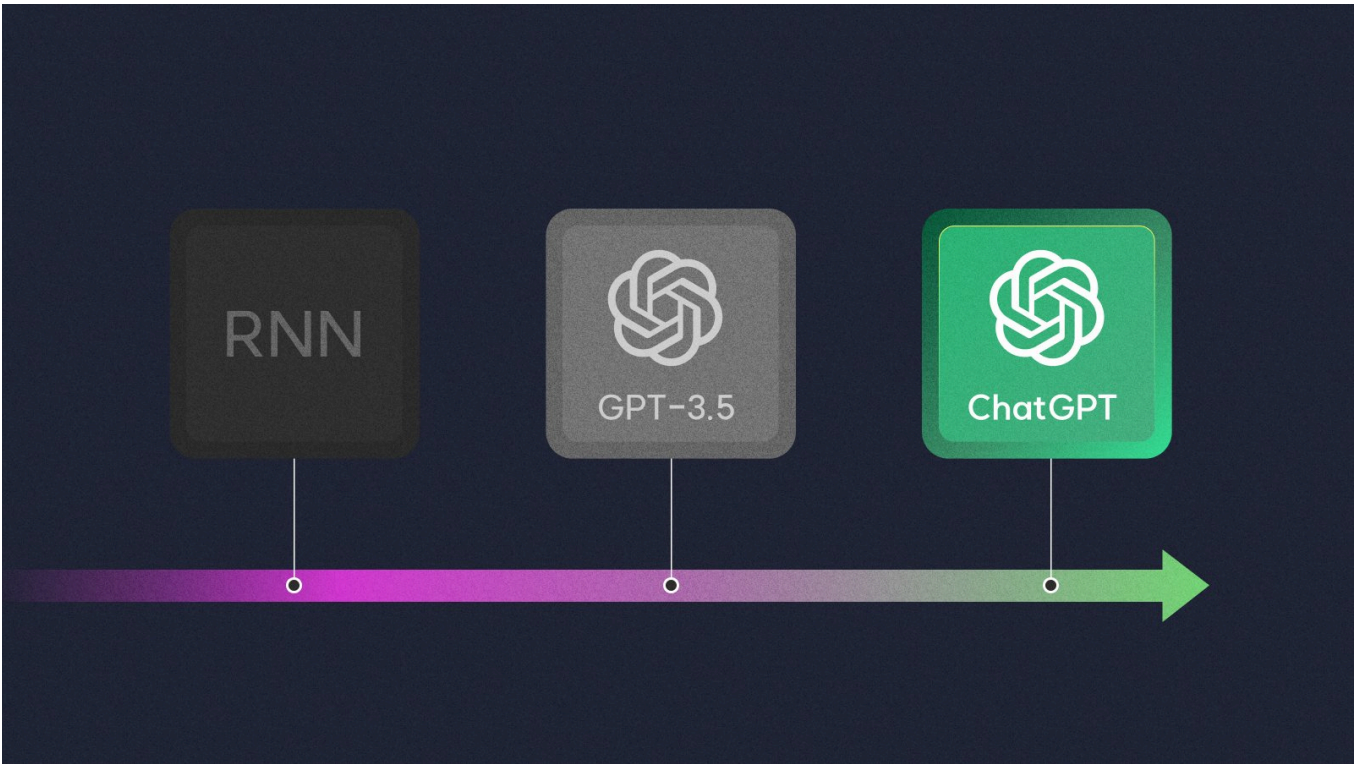


# GPT 시리즈와 발전 과정

2023/08/24 | 4 mins



Writer	✓	(이 콘텐츠는 업스테이지가 제작한 '모두를 위한 ChatGPT UP!' 강의 중 "GPT 시리즈와 발전 과정"의 내용을 바탕으로 작성되었습니다.)
이런 분이 읽으면 좋아요!	✓	
이 글로 확인할 수 있는 내용	✓	하루가 다르게 발전하는 AI, 그중에서도 가장 대중적으로 널리 알려진 ChatGPT가 탄생하기까지 GPT 시리즈는 어떻게 발전해 왔을까요? 기본적인 언어 모델의 개념부터 RNN(순환 신경망)에서 ChatGPT 시대를 마주하기까지 약 5년간의 여정을 살펴봅니다.
목차	✓	

## 언어 모델 (Language model)

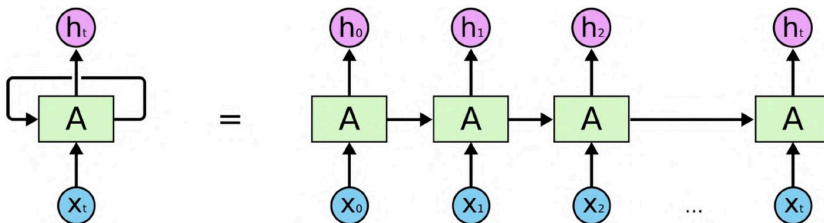
GPT(Generative Pre-trained Transformer)는 OpenAI에서 개발한 대규모 언어 모델로 다양한 자연어 처리 작업에 사용되고 있습니다. 그 때문에 먼저 언어 모델에 대해 이해해야 GPT의 발전 과정을 살펴보는 데 도움이 될 수 있는데요. 언어 모델이 답변을 생성할 때에는 보통 **다음 단어를 맞추는** 방식으로 풀어내곤 합니다. 아래의 문제를 예시로 들어보겠습니다.

Q. 빈칸에 들어갈 단어로 알맞은 것은 무엇일까요?

“오늘 참가한 [ ] 은/는 힘들었지만 정말 보람찬 일이었다.”

- (1) 달리기
- (2) 낮잠
- (3) 축제

빈칸에 무엇이 들어갈지 알아 맞추는 문제인데요. 이러한 방식을 언어 모델에도 그대로 적용하게 됩니다. 이 경우 사람이 일일이 정답 데이터를 생성해주지 않아도 모델 스스로 단어나 문장의 구조를 이용해 정답 데이터를 무수히 생성할 수 있다는 것이 장점입니다. 따라서 언어 모델링은 이러한 **Self-supervised learning**의 특징을 갖고 있어 pre-trained 모델을 만드는데 유리하다고 볼 수 있습니다.



Recurrent Neural Networks (출처: <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>)

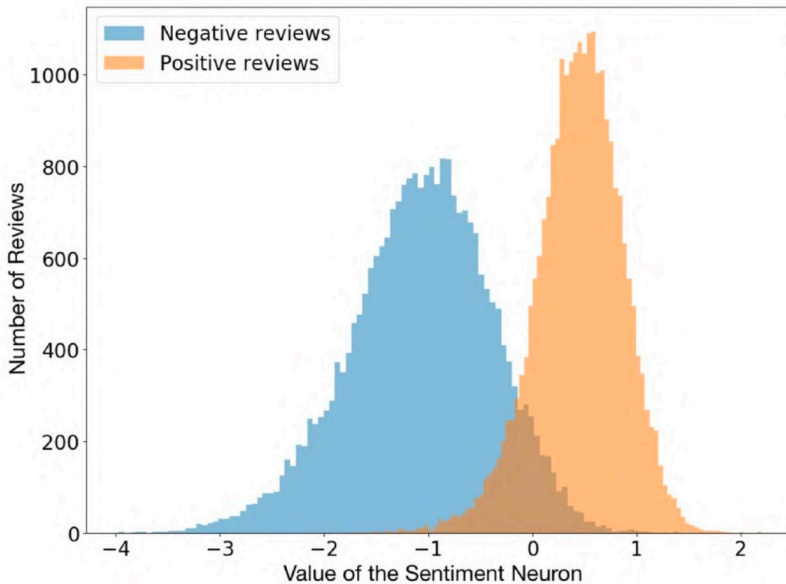
딥러닝 초창기에는 언어 처리 모델을 “RNN” 아키텍처(모델의 구조, 연산의 틀)로 만들었습니다. RNN은 노드 사이의 연결고리가 cycle을 이룬다고 해서 붙여진 이름으로, Recurrent Neural Network를 의미합니다. 이러한 특성으로 인해 **자연어와 같은 sequence 형태의 데이터를 다루는데 특화되어** 있습니다.

그렇다면 앞서 보여드리 예시와 같이 다수화 다른 단어 만츠기가 ChatGPT로

업스테이지 홈페이지에 오신 것을 환영합니다. 웹사이트를 원활하게 표시하기 위해 쿠키를 사용합니다. [www.upstage.ai](http://www.upstage.ai)를 계속 이용하려면 쿠키 사용에 동의해야 합니다.

## GPT 시리즈와 발전 과정

### Emergence (2017년 4월)



Sentiment neuron (출처: <https://openai.com/research/unsupervised-sentiment-neuron>)

2017년, OpenAI에서는 언어 모델을 RNN(Recurrent Neural Network)으로 만들고 있었습니다. 이 과정에서 특정 뉴런이 감성 분석을 하고 있음을 발견하게 되는데요. 이로 인해 의도하지 않았던 능력이 언어 모델링 과정에서 생기게 되는것인가 라는 가설이 등장합니다.

#### <감성 분석이란?>

- 인공지능 기술을 활용하여 텍스트의 내용을 분석하고, 이로부터 추출된 감정이나 의견을 판단하는 과정
- 주로 영화 리뷰, 온라인 게시글 등 텍스트 데이터를 대상으로 이루어지며, AI가 사람처럼 문장을 이해하고 어떤 감정이 담겼는지 파악하여 긍정, 부정, 중립 등을 구분해내는 것

업스테이지 홈페이지에 오신 것을 환영합니다. 웹사이트를 원활하게 표시하기 위해 쿠키를 사용합니다. [www.upstage.ai](https://www.upstage.ai)를 계속 이용하려면 쿠키 사용에 동의해야 합니다.

## Transformer (2017년 6월)

2017년에는 RNN(Recurrent Neural Network), CNN(Convolutional Neural Network)과 유사한 아키텍처의 일종인 Transformer가 등장합니다. 이것의 핵심은 **항목과 항목 사이의 연관성을 나타내는 'Attention'**이라는 연산인데, 이에 따라 "Attention is all you need"라는 구글 브레인팀의 논문이 등장할 정도로 중요성이 대두되었습니다. Transformer는 기존의 RNN 등에 비해 계산 효율과 결과의 품질이 좋았기 때문에 이후 비전, 추천, 바이오 등 다른 모든 분야에서 쓰는 기술이 될 정도로 큰 영향력을 미치게 되었습니다.

## GPT (2018년 6월)

1년 후, Generative Pre-training Transformer(GPT)가 처음으로 등장했습니다. 이는 앞서 설명했던 Self-supervised learning의 방식으로 언어 모델을 만든 것이라고 이해할 수 있습니다. GPT는 **Pretraining-finetuning 패러다임**의 대표적인 논문으로 꼽히기도 하는데, 이는 GPT의 등장으로부터 큰 규모의 언어 모델링을 통해 사전학습된 모델을 만들고, 이 모델을 각 task에 맞는 작은 데이터셋으로 학습하는 finetuning의 과정을 거치면 다양한 NLP 태스크에서 우수한 성능을 보인다는 것을 보여줬기 때문입니다.

### <파인튜닝이란?>

- 사전 학습된 모델(pre-trained model)을 기반으로 특정 도메인이나 작업에 적합하게 성능을 개선하는 작업
- 대용량 데이터셋을 사용해서 사전에 학습된 모델을 재사용하면, 새로운 작업이나 도메인에서 모델의 학습 시간을 줄이고, 데이터가 제한된 경우에도 성능을 향상시킬 수 있다는 것이 핵심 아이디어

## GPT-2 (2019년 2월)

업스테이지 홈페이지에 오신 것을 환영합니다. 웹사이트를 원활하게 표시하기 위해 쿠키를 사용합니다. [www.upstage.ai](https://www.upstage.ai)를 계속 이용하려면 쿠키 사용에 동의해야 합니다.

## Language Models are Unsupervised Multitask Learners

Alec Radford<sup>\*1</sup> Jeffrey Wu<sup>\*1</sup> Rewon Child<sup>1</sup> David Luan<sup>1</sup> Dario Amodei<sup>\*\*1</sup> Ilya Sutskever<sup>\*\*1</sup>

### Abstract

Natural language processing tasks, such as question answering, machine translation, reading comprehension, and summarization, are typically approached with supervised learning on task-specific datasets. We demonstrate that language models begin to learn these tasks without any explicit supervision when trained on a new dataset of millions of webpages called WebText. When conditioned on a document plus questions, the answers generated by the language model reach 55 F1 on the CoQA dataset - matching or exceeding the performance of 3 out of 4 baseline systems without using the 127,000+ training examples. The capacity of the language model is essential to the success of zero-shot task transfer and increasing it improves performance in a log-linear fashion across tasks. Our largest model, GPT-2, is a 1.5B parameter Transformer that achieves state of the art results on 7 out of 8 tested language modeling datasets in a zero-shot setting but still underfits WebText. Samples from the model reflect these improvements and contain coherent paragraphs of text. These findings suggest a promising path towards building language processing systems which learn to perform tasks from their naturally occurring demonstrations.

competent generalists. We would like to move towards more general systems which can perform many tasks – eventually without the need to manually create and label a training dataset for each one.

The dominant approach to creating ML systems is to collect a dataset of training examples demonstrating correct behavior for a desired task, train a system to imitate these behaviors, and then test its performance on independent and identically distributed (IID) held-out examples. This has served well to make progress on narrow experts. But the often erratic behavior of captioning models (Lake et al., 2017), reading comprehension systems (Jia & Liang, 2017), and image classifiers (Alcorn et al., 2018) on the diversity and variety of possible inputs highlights some of the shortcomings of this approach.

Our suspicion is that the prevalence of single task training on single domain datasets is a major contributor to the lack of generalization observed in current systems. Progress towards robust systems with current architectures is likely to require training and measuring performance on a wide range of domains and tasks. Recently, several benchmarks have been proposed such as GLUE (Wang et al., 2018) and decaNLP (McCann et al., 2018) to begin studying this.

Multitask learning (Caruana, 1997) is a promising framework for improving general performance. However, multitask training in NLP is still nascent. Recent work reports modest performance improvements (Yogatama et al., 2019) and the two most ambitious efforts to date have trained on a total of 10 and 17 (dataset, objective)

출처: [Language Models are Unsupervised Multitask Learners](#)

GPT-2는 기존 모델의 크기를 키우고 (117M → 1.5B) 학습 데이터의 양을 늘려서 (4GB → 40GB) 탄생한 버전입니다. 하지만 OpenAI는 생성에 탁월한 능력을 가진 GPT-2가 가짜 정보를 다량 생성할 위험성이 크다고 판단하여 외부에 공개하지 않기도 했는데요. GPT-2는 언어 생성 능력뿐만 아니라 또 다른 영향력을 시사하는 emergence를 보였습니다.

## “Emergence”: Zero-shot learning

GPT-2의 등장은 어떤 새로운 가능성을 보여주었을까요? 바로 모델이 예시를 전혀 보지 않고도 새로운 태스크를 수행하는 “Zero-shot learning”의 개념입니다. 이를 Unsupervised multitask learners라고도 부릅니다. 초기에는 언어 모델로 출발했지만 독해, 번역, 요약, Q&A 등 다른 다양한 태스크를 수행할 수 있는지에 대한 의문을 풀어가고자 여러 실험이 진행되었습니다.



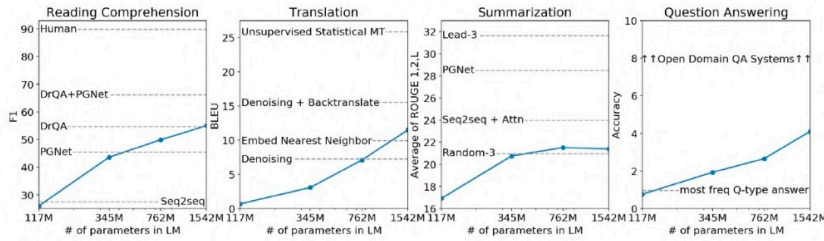


Figure 1. Zero-shot task performance of WebText LMs as a function of model size on many NLP tasks. Reading Comprehension results are on CoQA (Reddy et al., 2018), translation on WMT-14 Fr-En (Artetxe et al., 2017), summarization on CNN and Daily Mail (See et al., 2017), and Question Answering on Natural Questions (Kwiatkowski et al., 2019). Section 3 contains detailed descriptions of each result.

	LAMBADA (PPL)	LAMBADA (ACC)	CBT-CN (ACC)	CBT-NE (ACC)	WikiText2 (PPL)	PTB (PPL)	enwik8 (BPF)	text8 (BPC)	WikiText103 (PPL)	1BW (PPL)
SOTA	99.8	59.23	85.7	82.3	39.14	46.54	0.99	1.08	18.3	<b>21.8</b>
117M	<b>35.13</b>	45.99	<b>87.65</b>	<b>83.4</b>	<b>29.41</b>	65.85	1.16	1.17	37.50	75.20
345M	<b>15.60</b>	55.48	<b>92.35</b>	<b>87.1</b>	<b>22.76</b>	47.33	1.01	<b>1.06</b>	26.37	55.72
762M	<b>10.87</b>	<b>60.12</b>	<b>93.45</b>	<b>88.0</b>	<b>19.93</b>	<b>40.31</b>	<b>0.97</b>	<b>1.02</b>	22.05	44.575
1542M	<b>8.63</b>	<b>63.24</b>	<b>93.30</b>	<b>89.05</b>	<b>18.34</b>	<b>35.76</b>	<b>0.93</b>	<b>0.98</b>	<b>17.48</b>	42.16

Table 3. Zero-shot results on many datasets. No training or fine-tuning was performed for any of these results. PTB and WikiText-2 results are from (Gong et al., 2018). CBT results are from (Bajgar et al., 2016). LAMBADA accuracy result is from (Hoang et al., 2018) and LAMBADA perplexity result is from (Grave et al., 2016). Other results are from (Dai et al., 2019).

SOTA란 state-of-the-art, 즉 현존하는 제일 좋은 모델.  
볼드체가 더 좋은 점수. ACC제외하고는 낮을수록 좋음.

출처: [Language Models are Unsupervised Multitask Learners](#)

위의 논문에도 언급되었듯이 파라미터(매개변수) 수를 늘릴수록 Zero-shot의 성능이 올라가고, 특정 태스크에서는 기존의 SOTA(state-of-the-art, 현존하는 제일 좋은 모델) 모델을 능가하는 것이 실제로 확인할 수 있었다는 것이 특이점이었습니다.

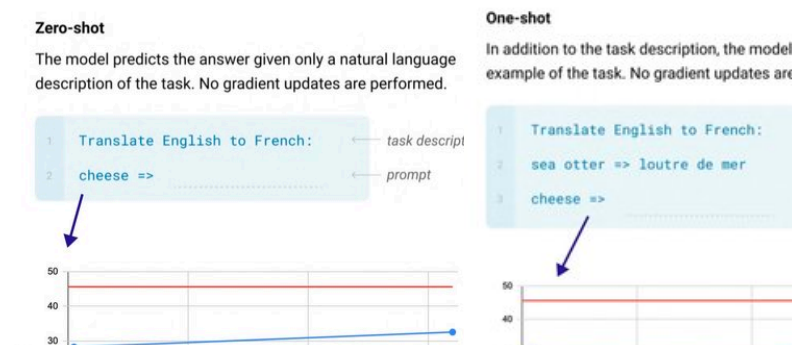
## GPT-3 (2020년 6월)

이처럼 실험을 통해 GPT의 여러 능력을 확인한 뒤로 한 번 더 크기를 키운 것이 2020년에 등장한 GPT-3입니다. 모델은 1.5B에서 175B로 커졌으며, 데이터 또한 600GB 이상이 들어갔습니다. 이렇게 앞선 시리즈보다 많은 데이터로 pretraining을 했기에 더욱 놀라운 생성 능력을 갖추게 되었는데요. GPT-3 역시 지식과 학습 없이도 태스크를 “배우는” 능력(few-shot learners)을 포함하여 여러 측면으로 emergence를 확인할 수 있었습니다. 이전 버전에서는 태스크를 수행하는 것에 그쳤지만, 스스로 태스크를 배우는 능력까지 보이게 된 것입니다.



(In-context learning 이전에 태스크별 예시를 모델에 입력해 파인튜닝한 과정을 도식화한 것 / 출처: [Language Models are Few-Shot Learners](#))

GPT-3에서 나타난 emergence를 In-context learning이라고도 부르는데요. In-context learning 이전에는 태스크의 예시를 모델에 입력하여 파인튜닝을 하는 것이 필요했습니다. 이 경우 태스크별로 모델과 데이터가 필요하다는 점에서 한계가 있었는데, GPT-2부터 Zero-shot learning이 가능해지며 프롬프트에 예시 몇 개(few-shot)을 넣어주면 모델 업데이트 없이도 새로운 태스크를 수행하게 되었습니다.



(출처: [Language Models are Few-Shot Learners](#))

## GPT-4 출시 전, 2021~2022년

GPT-3 이후, GPT-4 출시에 대한 사람들의 기대감은 점점 커졌습니다.

GPT-4의 정식 출시 이전에 업계에서는 주목할 만한 것들이 크게 네 가지

업스테이지 홈페이지에 오신 것을 환영합니다. 웹사이트를 원활하게 표시하기 위해 쿠키를 사용합니다. [www.upstage.ai](http://www.upstage.ai)를 계속 이용하려면 쿠키 사용에 동의해야 합니다.

모델인 InstructGPT입니다. InstructGPT는 기존의 GPT와 달리 모델에게 직접적으로 지시를 할 수 있으며, 일련의 지시사항에 따라 사용자의 의도에 맞는 답을 할 수 있도록 설계된 언어 모델로써 큰 주목을 받았습니다.

- CLIP (2021년 1월): “zero-shot” 이미지 분류
- DALL-E(2021년 1월): 주어진 텍스트로부터 이미지 생성
- Codex(2021년 8월): 코드 생성을 위한 모델
- InstructGPT (2022년 1월)  
: 명령에 대한 파인튜닝과 강화 학습이 이루어진 모델. 기존의 GPT는 프롬프트 엔지니어링을 통해 모델이 특정 태스크를 잘 수행할 수 있는 조건이나 예시를 들어줘야 했다면, InstructGPT는 간단한 자연어 지시문 만으로도 사용자가 요청한 대로 결과물을 생성함

○ <일반 언어 모델이 지시사항에 대해 생성한 내용>

💬 “ChatGPT에 대해 설명해줘”

→ BERT에 대해 설명해줘 / GPT에 대해 설명해줘

<Instruction fine-tuning을 적용한 언어 모델이 지시사항에 대해 생성한 내용>

💬 “ChatGPT에 대해 설명해줘”

→ ChatGPT는 OpenAI에서 개발한 자연어 처리 모델 중 하나입니다. 이 모델은 GPT (Generative Pre-trained Transformer) 아키텍처를 기반으로 하며, 대화 기반의 인공지능 모델로 사용됩니다. ChatGPT는 사전에 대규모의 데이터로 사전 훈련된 후, 다양한 대화 데이터를 기반으로 Fine-tuning 과정을 거쳐 최적화됩니다. 이를 통해 ChatGPT는 사용자와의 대화에서 자연스러운 응답을 생성하고 다양한 주제에 대해 대화를 나눌 수 있습니다.

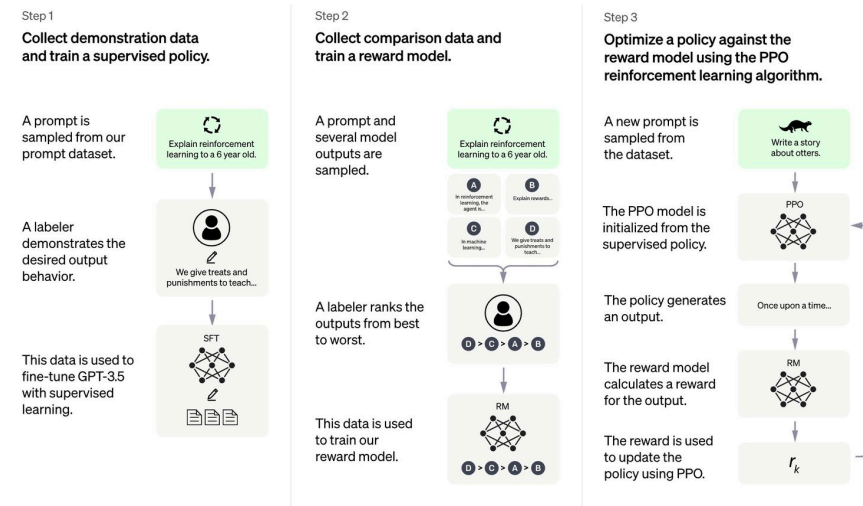


**GPT-3.5는 GPT-3에 코드 데이터와 Instruction fine-tuning이 추가된 버전입니다.** 많은 연구자들의 추측에 따르면 이러한 방식이 모델에 직접적으로 영향을 미치는지는 알 수 없지만, Code 데이터를 추가함으로 인해 GPT의 추론 능력과 긴 입력에 대한 이해가 올라가는 것을 관찰했다고 합니다.

이외에도 GPT-3.5에는 Instruction fine-tuning이 적용되어 있는데, 명령에 대한 파인튜닝과 강화학습을 하면 사용자의 의도를 더 잘 파악하고 답변한다는 것에 착안한 InstructGPT(2022년 1월)의 실험 방식이 가미되어 있습니다.

## ChatGPT (2022년 11월)

2022년에 등장한 ChatGPT는 AI의 대중화를 이끈 모델 중 하나입니다. 이는 GPT-3.5를 파인튜닝한 것으로, OpenAI에서는 InstructGPT와 학습 방식이 유사하기 때문에 “sibling model”이라고도 부릅니다.



출처: [OpenAI Blog](#)

ChatGPT 모델의 학습 방식을 살펴보면, 첫 번째 단계에서는 지시 프롬프트와 데이터셋으로 이루어진 Demonstration data를 넣어줍니다. 여기서 라벨러는 지시 프롬프트에 적합하다고 여겨지는 행동을 라벨링 하는데, 이렇게 모아진 데이터셋은 SFT(Supervised Fine Tuning) 모델

업스테이지 홈페이지에 오신 것을 환영합니다. 웹사이트를 원활하게 표시하기 위해 쿠키를 사용합니다. [www.upstage.ai](http://www.upstage.ai)를 계속 이용하려면 쿠키 사용에 동의해야 합니다.

다음 단계에서는 유저의 선호도에 대한 보상 모델(Reward model, RM)을 활용해 ChatGPT를 강화학습(Reinforcement learning, RL)으로 업데이트합니다. 이러한 방식을 거쳐 ChatGPT는 보다 다양하고 유연한 대화를 제공할 수 있게 되었습니다.



RNN부터 ChatGPT까지의 여정

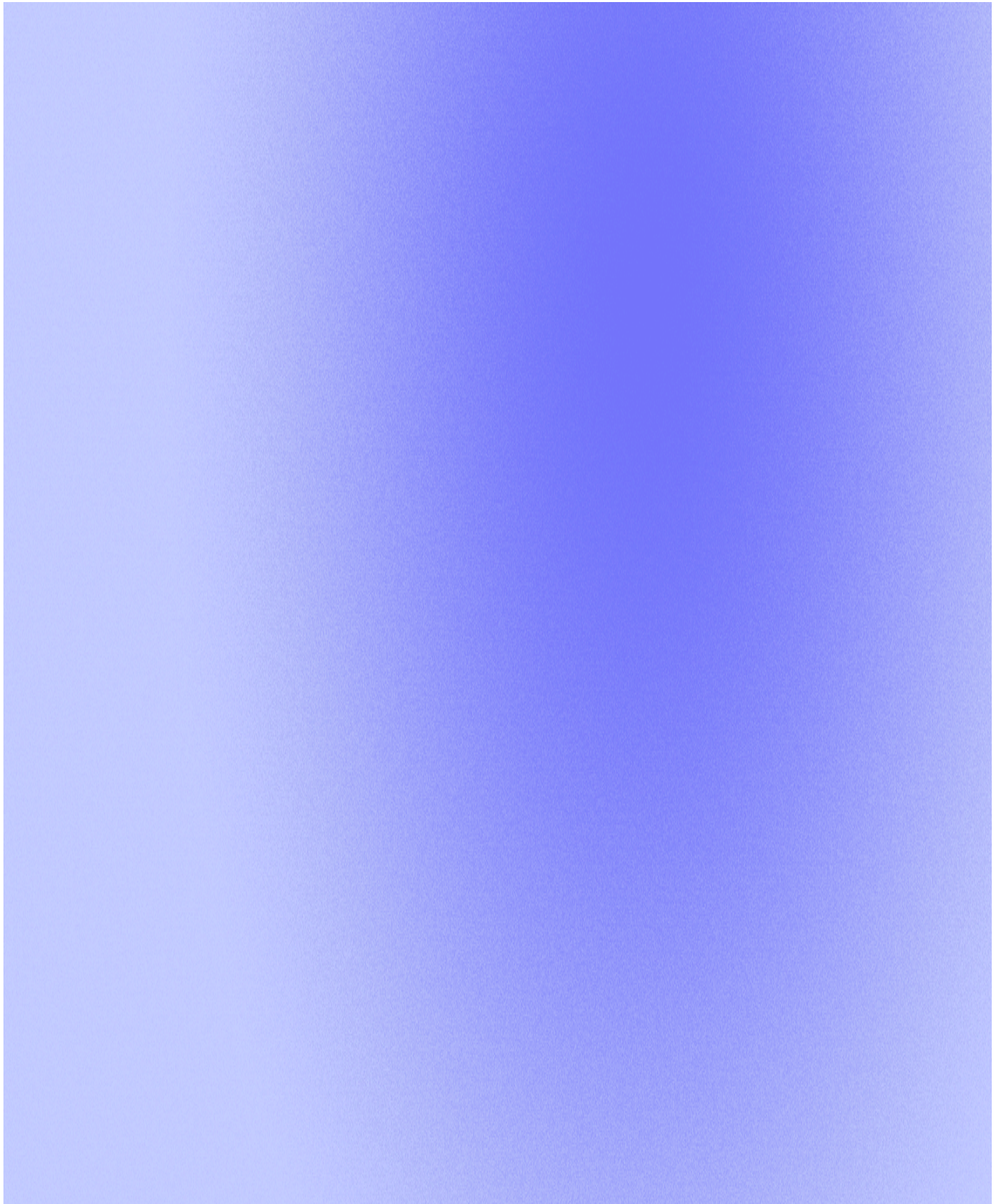
이렇게 RNN부터 ChatGPT까지, GPT 시리즈의 긴 여정을 함께 살펴보았습니다. ChatGPT의 뒤를 이을 Next GPT를 향한 움직임에는 어떤 것들이 있을까요? 미래에는 ChatGPT가 어떻게 활용될지, 또 발전을 위해 필요한 다양한 측면은 무엇이 있을지에 대해 보다 자세한 내용은 [웨비나 다시보기](#) 페이지에서 확인하실 수 있습니다.

Next GPT를 향한 움직임들

 [웨비나 시청하기](#)



업스테이지 홈페이지에 오신 것을 환영합니다. 웹사이트를 원활하게 표시하기 위해 쿠키를 사용합니다. [www.upstage.ai](https://www.upstage.ai)를 계속 이용하려면 쿠키 사용에 동의해야 합니다.



업스테이지 홈페이지에 오신 것을 환영합니다. 웹사이트를 원활하게 표시하기 위해 쿠키를 사용합니다. [www.upstage.ai](https://www.upstage.ai)를 계속 이용하려면 쿠키 사용에 동의해야 합니다.