

The background of the image is a complex, abstract network diagram. It consists of numerous nodes of varying sizes and colors (dark blue, light blue, and grey) connected by a web of thin, light grey lines. Some nodes are highlighted with larger, concentric circles. The overall aesthetic is clean and modern, suggesting a digital or technological theme.

FINE-TUNING

FINE-TUNING

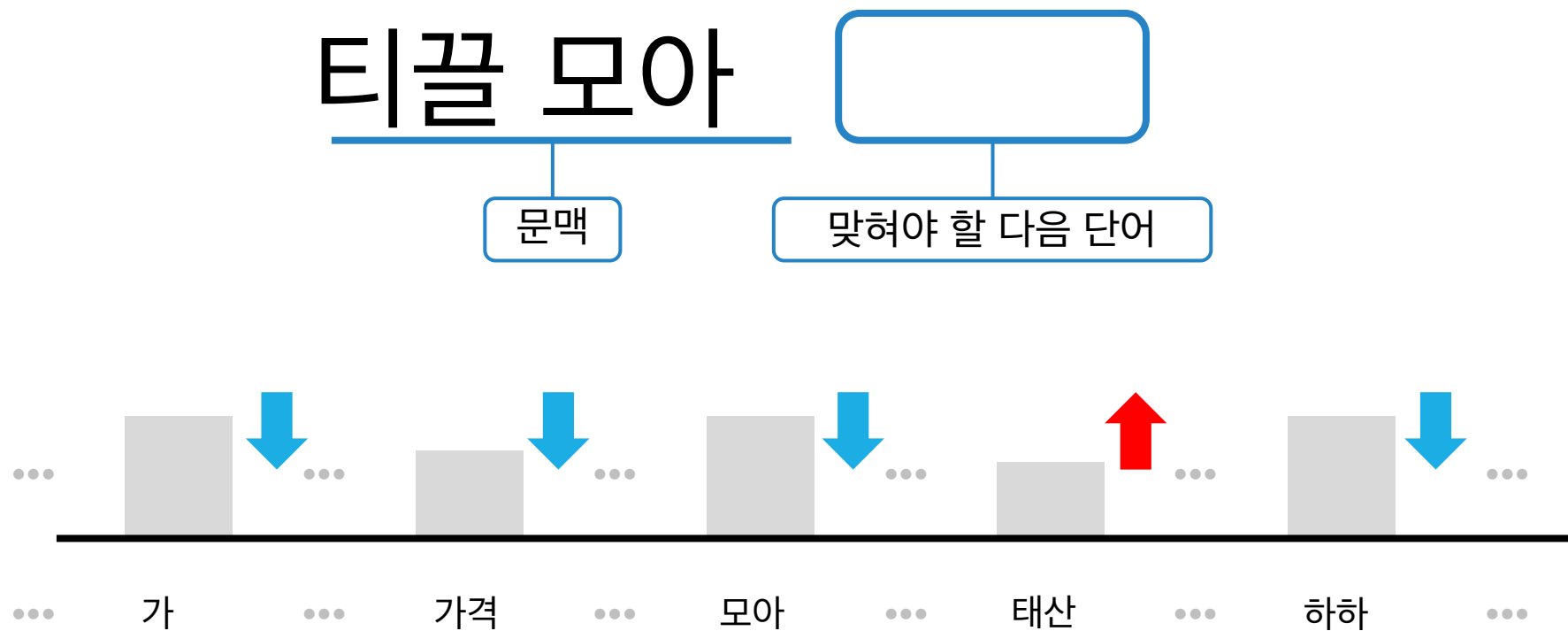
- 프리트레이닝을 마친 언어 모델 위에 작은 모듈을 더 쌓아 전체 문서 분류, 개체명 인식 등 다운스트림 데이터로 업데이트 하는 과정

UPSTREAM TASK

- 트랜스퍼 러닝의 대표적인 방법으로 다음 다음 단어 맞추기, 빈칸 채우기 등 대규모 말뭉치의 문맥을 이해하는 과정

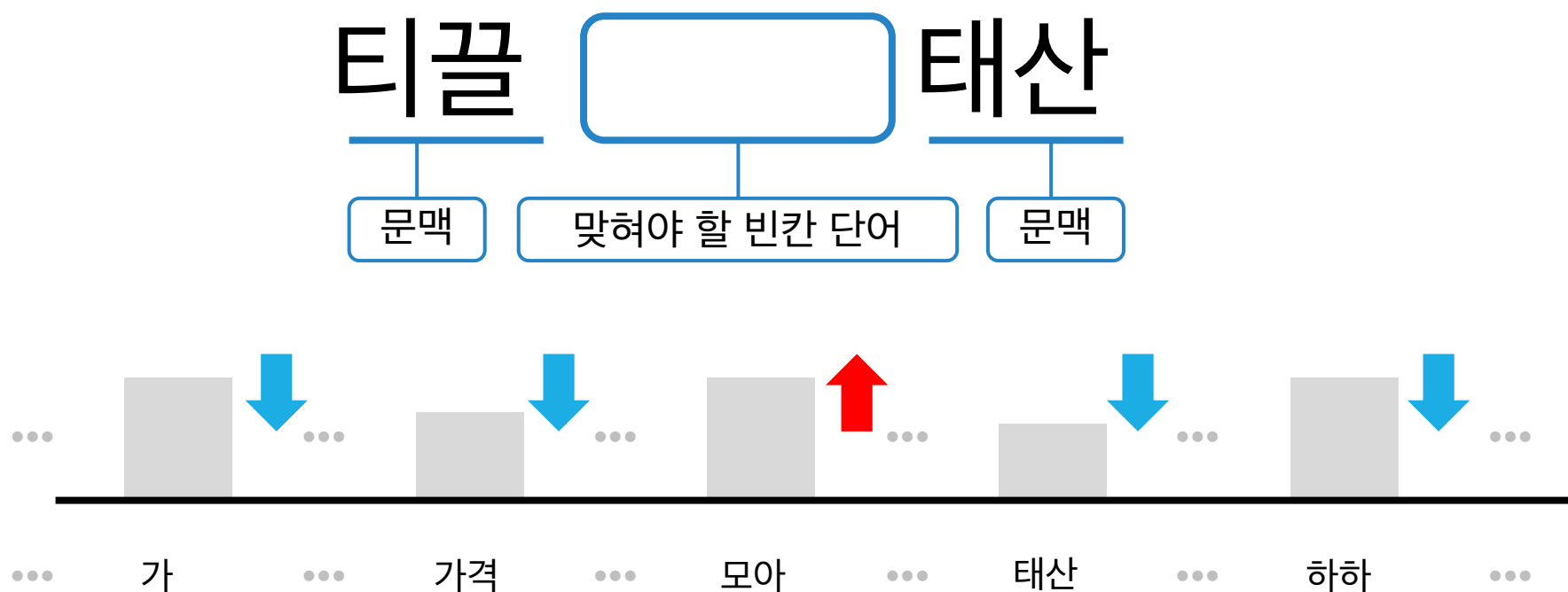
UPSTREAM TASK

- 다음 단어 맞추기



UPSTREAM TASK

- 빈칸 채우기



DOWNSTREAM TASK

- 트랜스퍼 러닝의 대표적인 방법으로 문서 분류, 개체명 인식, 문장 생성 등 우리가 풀고자 하는 자연어 처리의 구체적인 문제들을 풀어가는 과정

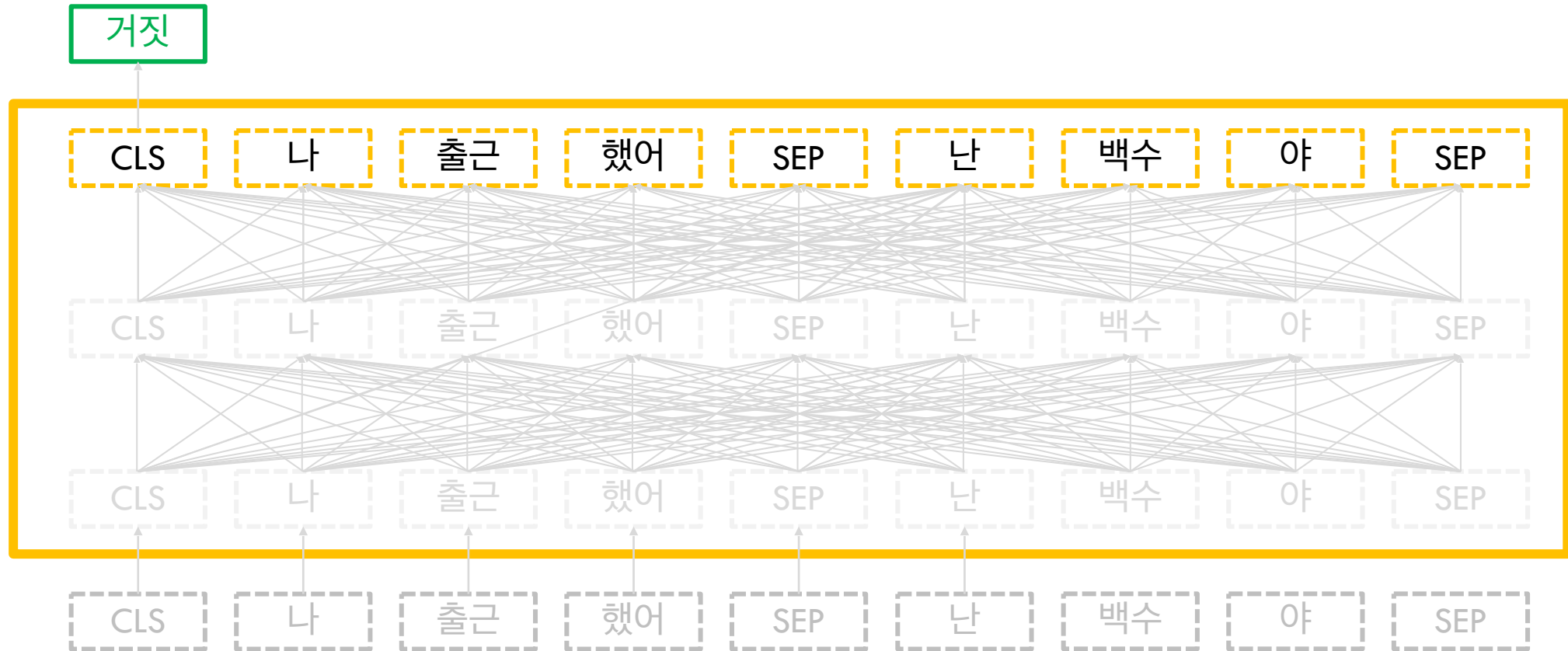
DOWNSTREAM TASK

- 문서 분류



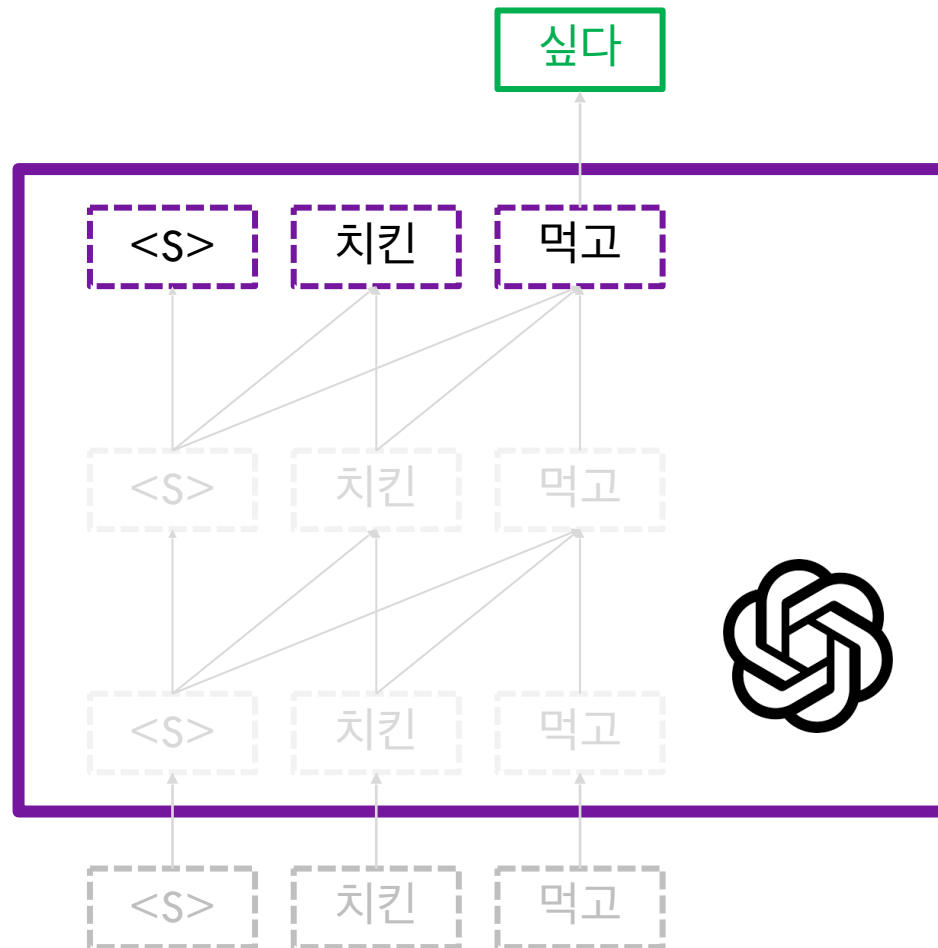
DOWNSTREAM TASK

- 자연어 추론

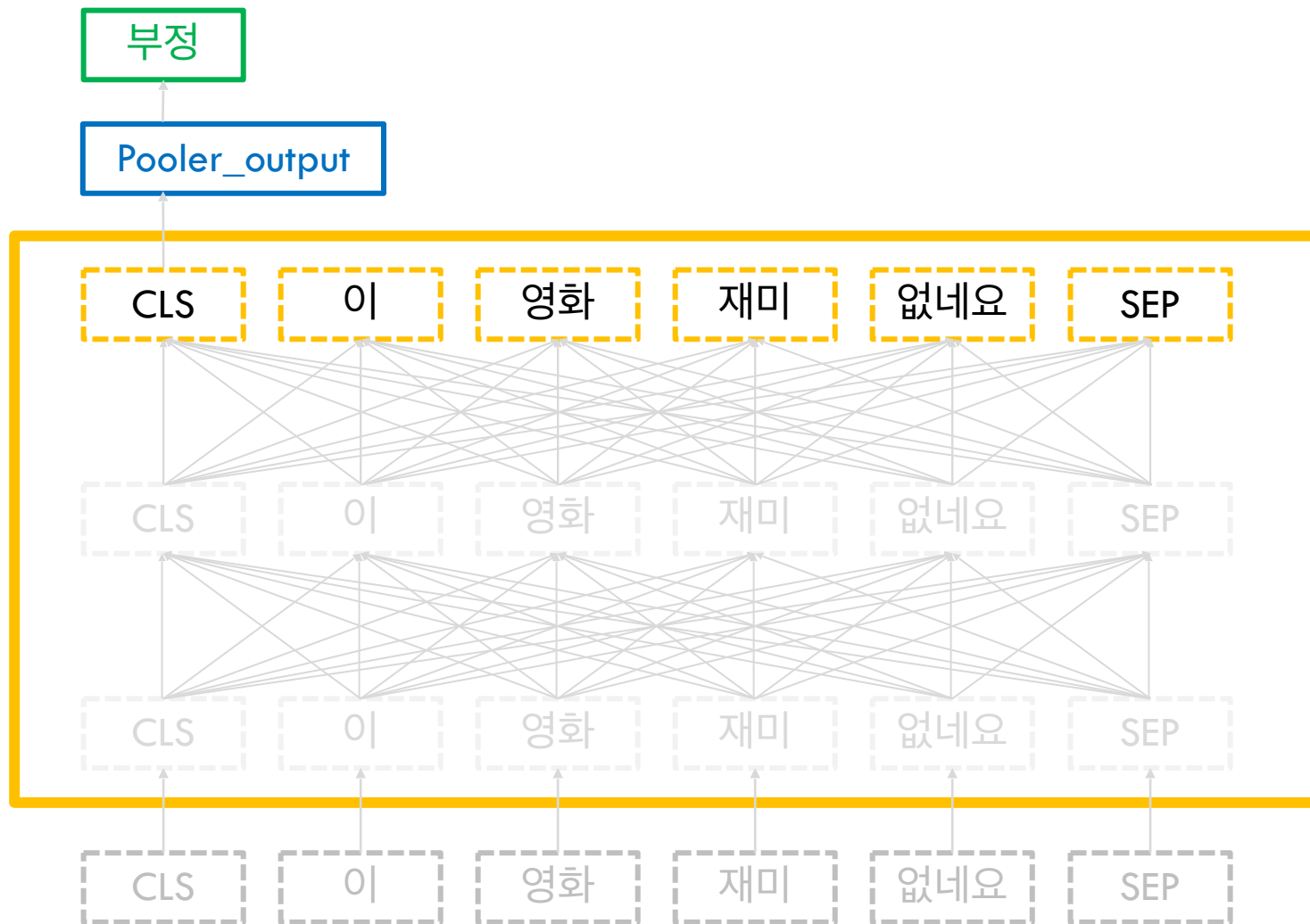


DOWNSTREAM TASK

- 문장 생성



FINE-TUNING



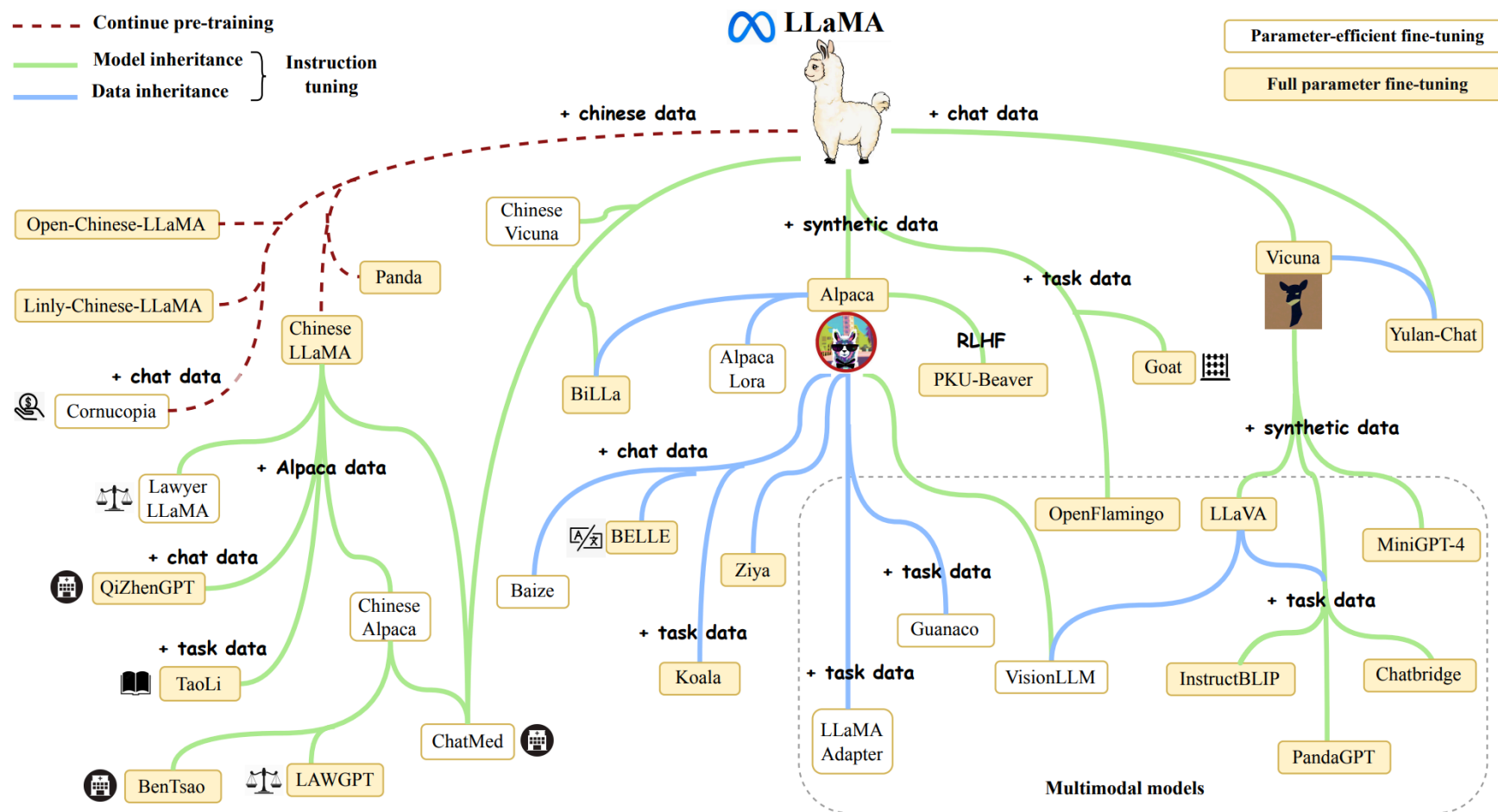
FINE-TUNING



LLAMA

- Meta에서 공개한 오픈소스 LLM
- 70억에서 650억 파라미터까지 다양한 모델을 공개하였으며, Meta가 공개한 내용에 따르면 대부분의 NLP 벤치마크에서 130억 파라미터 모델의 성능이 훨씬 더 큰 GPT-3(1,750억 파라미터 포함)의 성능을 초과
- Version 1 에서는 연구용으로 사용 할 수 있었지만, Version 2 부터는 상용으로 사용 가능
- 거의 모든 오픈소스 LLM이 이 모델을 기반으로 하고 있음

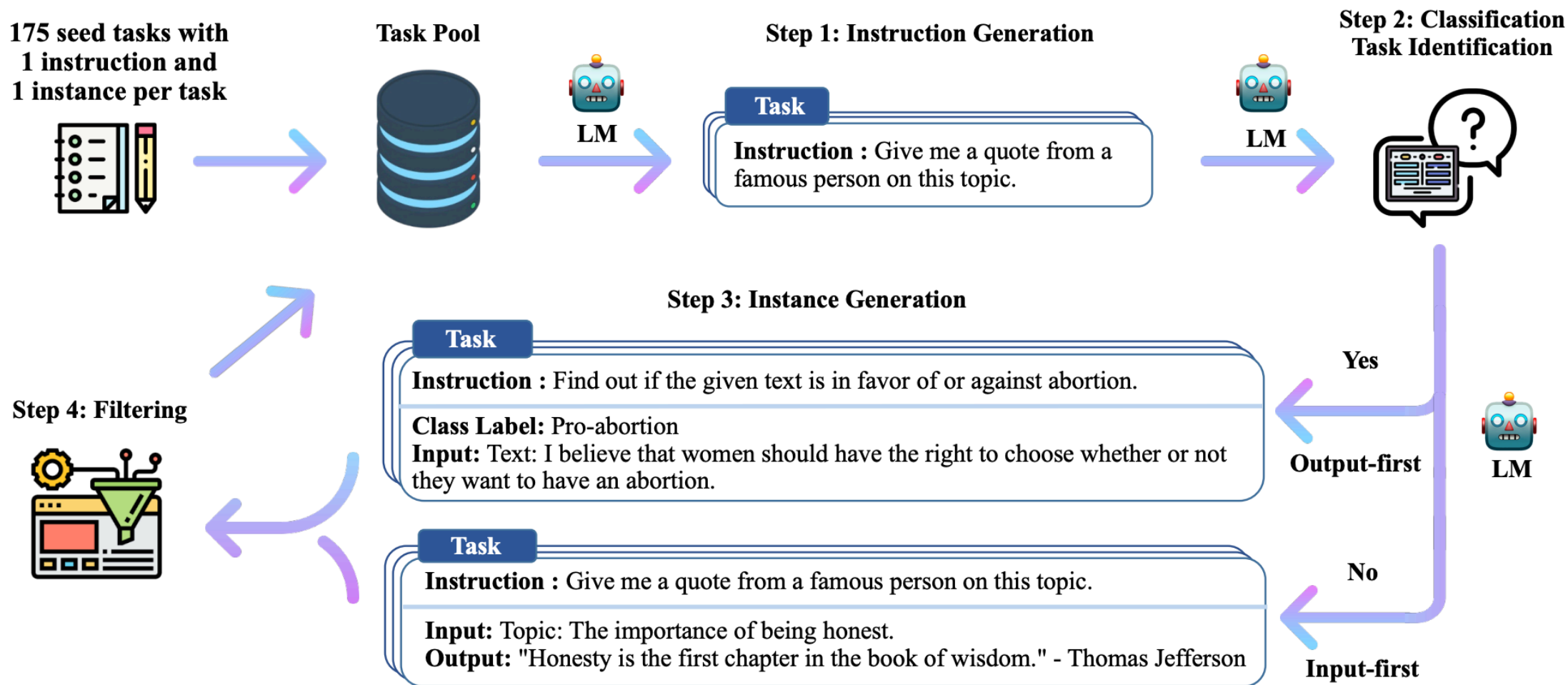
LLaMA



ALPACA

- LLaMA 모델을 기반으로 Stanford 대학에서 Instruct tuning 을 통해 성능을 크게 향상시킨 모델
- Self-instruct 라는 방법을 이용하였으며, 이는 사람이 작성한 instruction seed 세트를 기반으로 LLM 에게 학습에 필요한 instruction 세트를 생성해서 학습시키는 방법
- 52,000개의 instruct 데이터 세트만으로 fine-tuning 한 7B의 Alpaca 모델이 175B의 GPT-3.5의 성능을 상회했다고 함 (약 \$600 정도의 비용을 사용)
- 다만, 환각이 높고 학습된 데이터에 편향된 결과가 나타난다고 함

SELF-INSTRUCT



KOALPACA

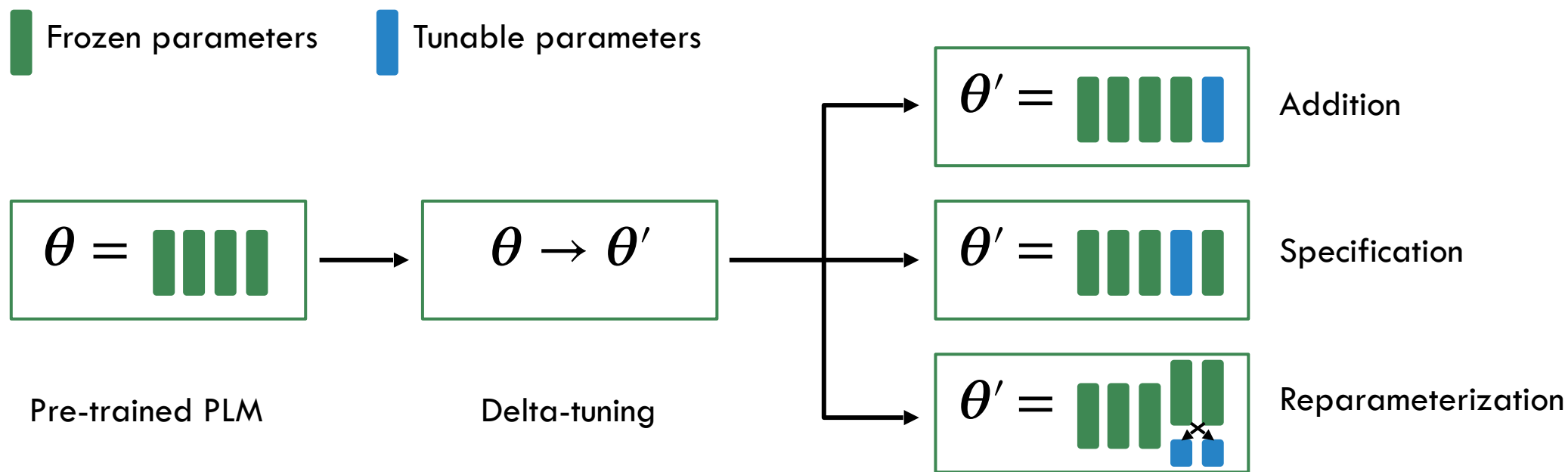
- Polyglot-ko 모델을 기반으로 Stanford Alpaca 모델을 학습한 방식과 동일한 방식으로 학습을 진행한 한국어 LLM
- Polyglot-ko 모델은 EleutherAI 의 컴퓨팅 자원과 TuNiB AI에서 수집한 1.2TB 규모의 한국어 데이터를 바탕으로 학습한 한국어 LLM

PEFT(PARAMETER EFFICIENT FINE-TUNING)

- 모델의 전체 파라미터를 다 학습 시키지 않고 매우 적은 양의 파라미터만 학습하여 빠른 시간 내에 새로운 문제를 거의 비슷한 성능으로 학습시키는 방법

.

PEFT (PARAMETER EFFICIENT FINE-TUNING)



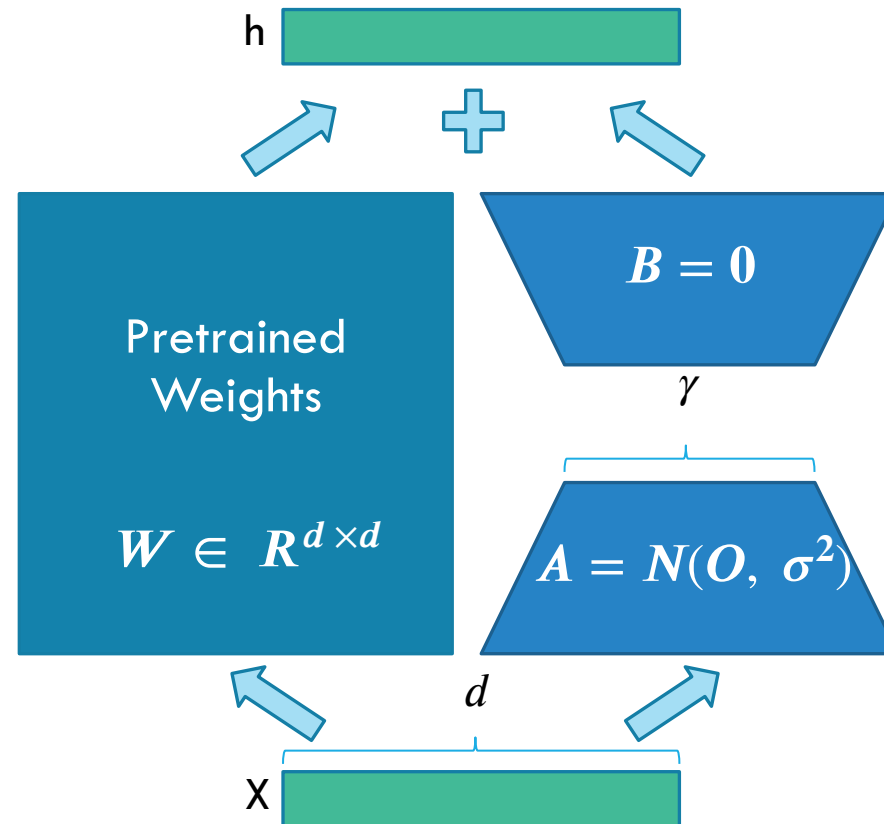
LORA

(LOW-RANK ADAPTATION OF LARGE LANGUAGE MODELS)

- Parameter Efficient Fine-Tuning(PEFT) 방법의 하나로, 매우 적은 양의 파라미터만 학습하여 빠른 시간 내에 새로운 문제를 거의 비슷한 성능으로 학습시키는 방법
- 고정된 weight를 갖는 pretrained model에 학습이 가능한 rank decomposition 행렬을 삽입한것으로, 행렬의 차원을 r 만큼 줄이는 행렬과 다시 원래 크기로 키워주는 행렬의 곱으로 이루어짐

LORA

(LOW-RANK ADAPTATION OF LARGE LANGUAGE MODELS)



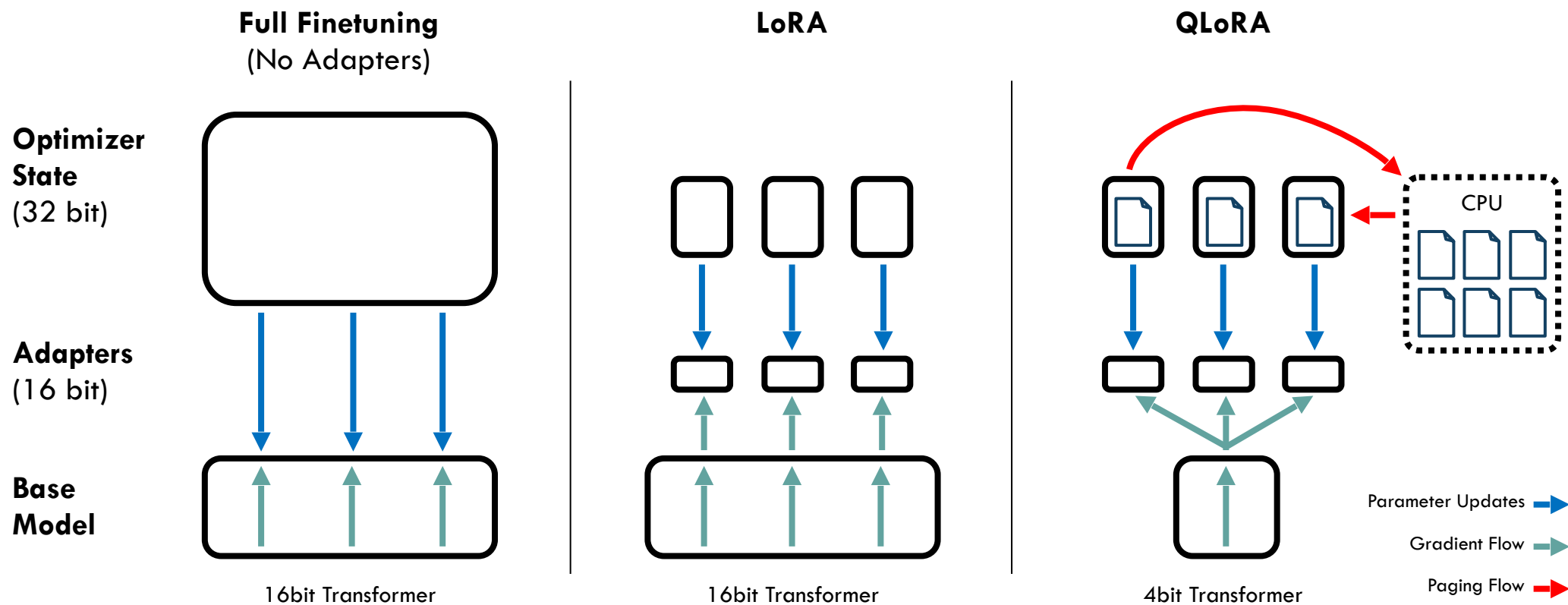
QLORA

(EFFICIENT FINETUNING OF QUANTIZED LLMS)

- Parameter Efficient Fine-Tuning(PEFT) 방법의 하나로, 매우 적은 양의 파라미터만 학습하여 빠른 시간 내에 새로운 문제를 거의 비슷한 성능으로 학습시키는 방법
- 고정된 weight을 갖는 pretrained model에 학습이 가능한 rank decomposition 행렬을 삽입한것으로, 행렬의 차원을 r 만큼 줄이는 행렬과 다시 원래 크기로 키워주는 행렬의 곱으로 이루어짐

QLORA

(EFFICIENT FINETUNING OF QUANTIZED LLMS)



HUGGING FACE

- 다양한 트랜스포머 모델을 올리고 쉽게 학습 시키거나 사용해 볼 수 있는 API를 제공하는 서비스로, 학습 스크립트나 추론 스크립트를 만드는 수고를 크게 덜 수 있음
- 인공지능의 모델을 위한 깃허브