# GEN AI 인텐시브 과정

강사 장철원

DAY1     DAY2     DAY3     DAY4     DAY5     DAY6     DAY7     DAY8

| LLM Basic Concept | Transformers paper review | Transformers LangChain LangGraph | LLM service develop | Final Project |

❑Full Fine Turning

❑Adapter

❑Adapter Fusion

❑LoRA

# GEN AI 인텐시브 과정

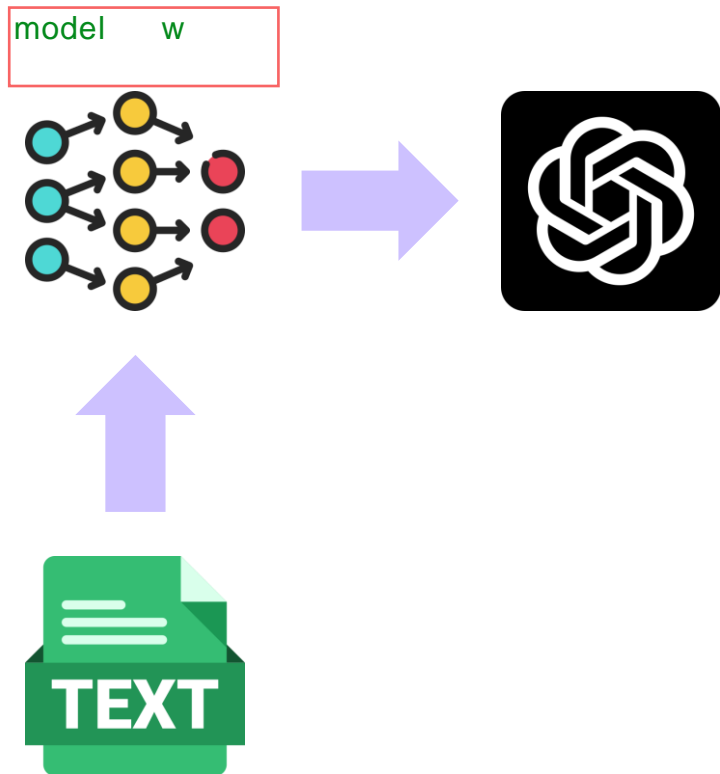Section 1. Full Fine Tuning

## Section 1-1.  Full Fine Tuning의 개념

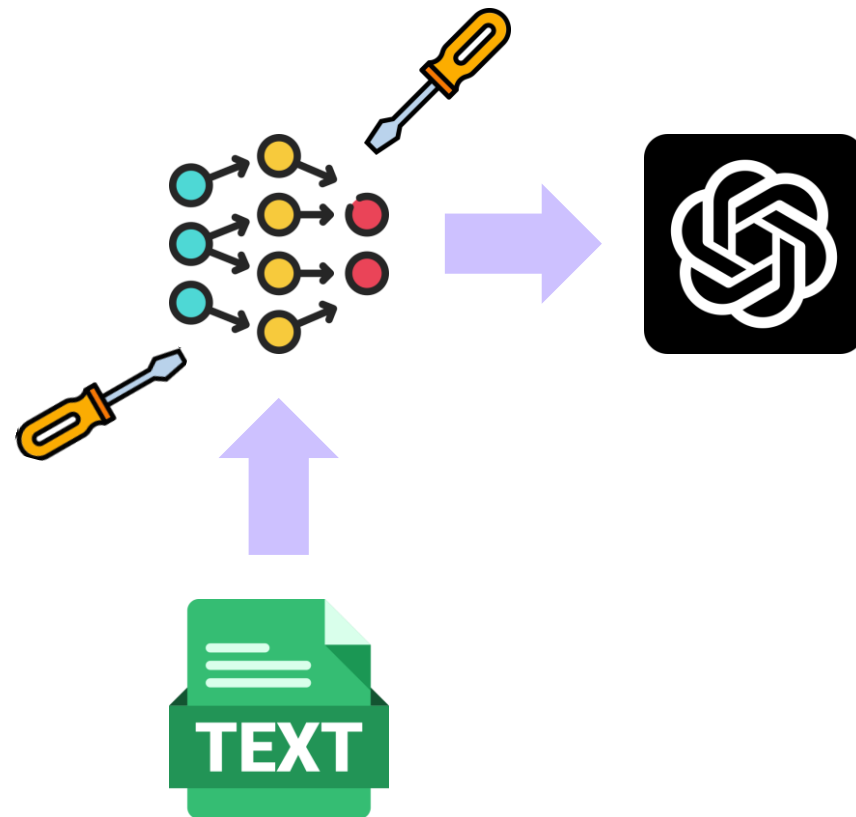# Full Fine Tuning의 개념

- 파인 튜닝(Fine Tuning)

  - 사전 학습이 끝난 모델에 대해 또 다른 Task에 대해 데이터를 추가시켜 재학습 하는 과정

  - LLM을 구성하는 파라미터의 개수: 수억 ~ 수십억개


- 풀 파인 튜닝(Full Fine Tuning)

  - 모델을 구성하는 전체 파라미터를 모두 재학습 시키는 방법

# Full Fine Tuning의 개념



## Pre-Trained Model

## Fine-Tuning Model

# Full Fine Tuning의 한계

- 수 십억개의 파라미터를 모두 재학습 시켜야 하므로...

  - 학습 시간이 오래걸림

  - Task 개수만큼의 모델이 필요하므로 요구되는 저장 공간이 넓어짐

  - 효율성과 확장성이 떨어짐.

# PEFT(Parameter Efficient Fine Tuning)

- 모델 전체를 다시 학습하지 않고, 일부만 수정하여 효율적인 학습 추구

    - 모델을 구성하는 파라미터의 대부분을 동결(freeze) 시키고,

    - 작은 학습 가능한 모듈을 추가함으로써 파인튜닝 하는 방법

# GEN AI 인텐시브 과정

Section 2. Adapter

**Section 2-1. Adapter의 개념**

# Adapter의 개념 - 기존 트랜스포머 구조에 탈부착 가능한 모듈



**AdapterHub: A Framework for Adapting Transformers**

Jonas Pfeiffer et al, 2020

트랜스포머 인코더 모델인 BERT 기준으로 설명
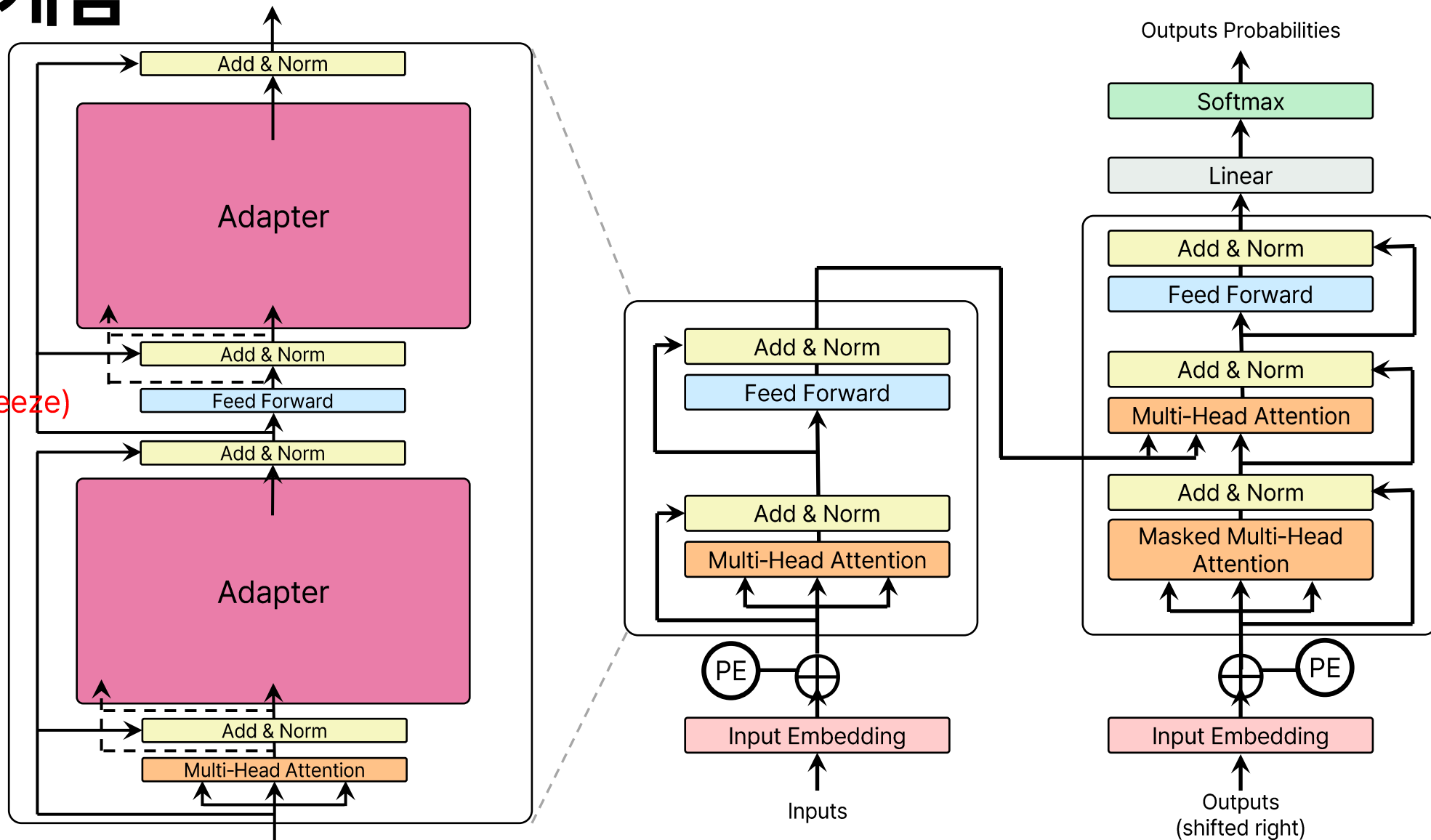
# Adapter의 개념

- LLM에 포함된 전체 파라미터를 튜닝하지 않아도, 사전 학습된 모델을 효율적으로 수정하거나 확장할 수 있는 방법

- 기존 트랜스포머 구조에 외부 모듈을 추가하는 방법

    - 어댑터만 학습하면 되므로 확장성, 모듈성이 뛰어남

    - 여러 어댑터를 조합해서 사용 가능(이후 AdapterFusion에서 설명)

- 기존 풀 파인튜닝의 단점을 보완하기 위해 등장한 방법

# Adapter의 개념



전체 파라미터를
학습하는게 아니라
Adapter 내부의
파라미터만 학습

기존 파라미터들은 동결(freeze)

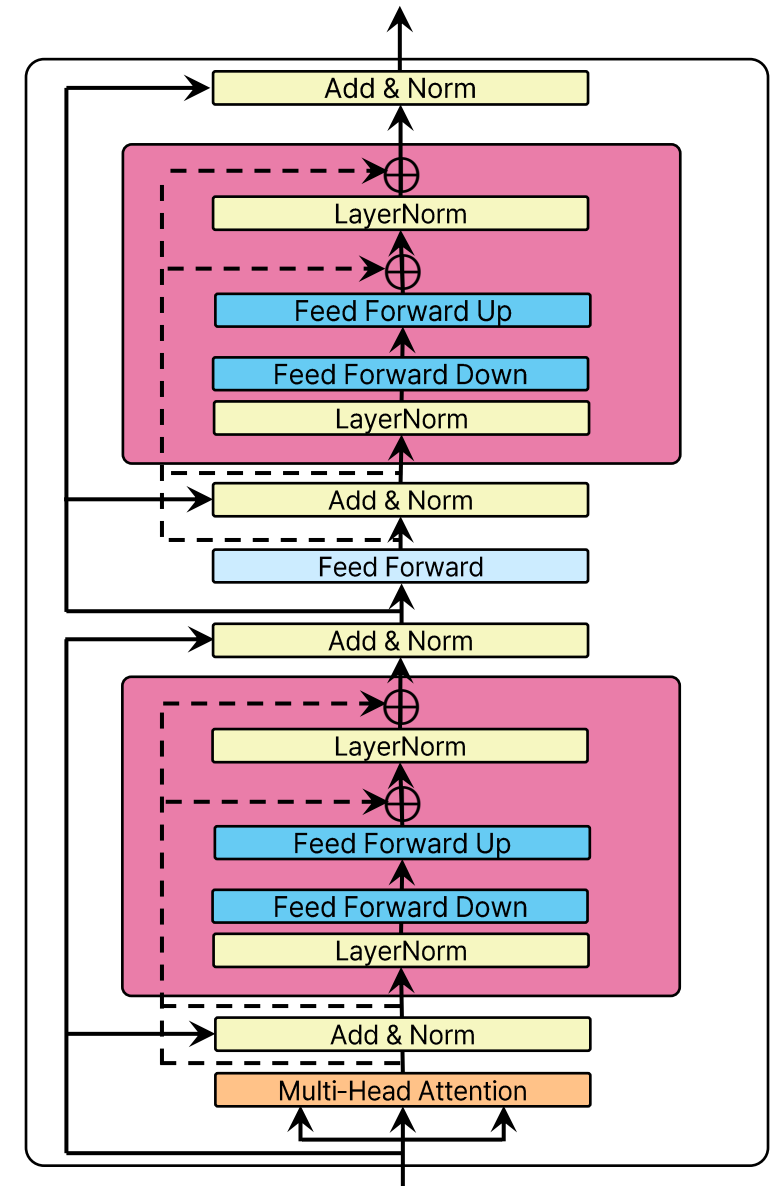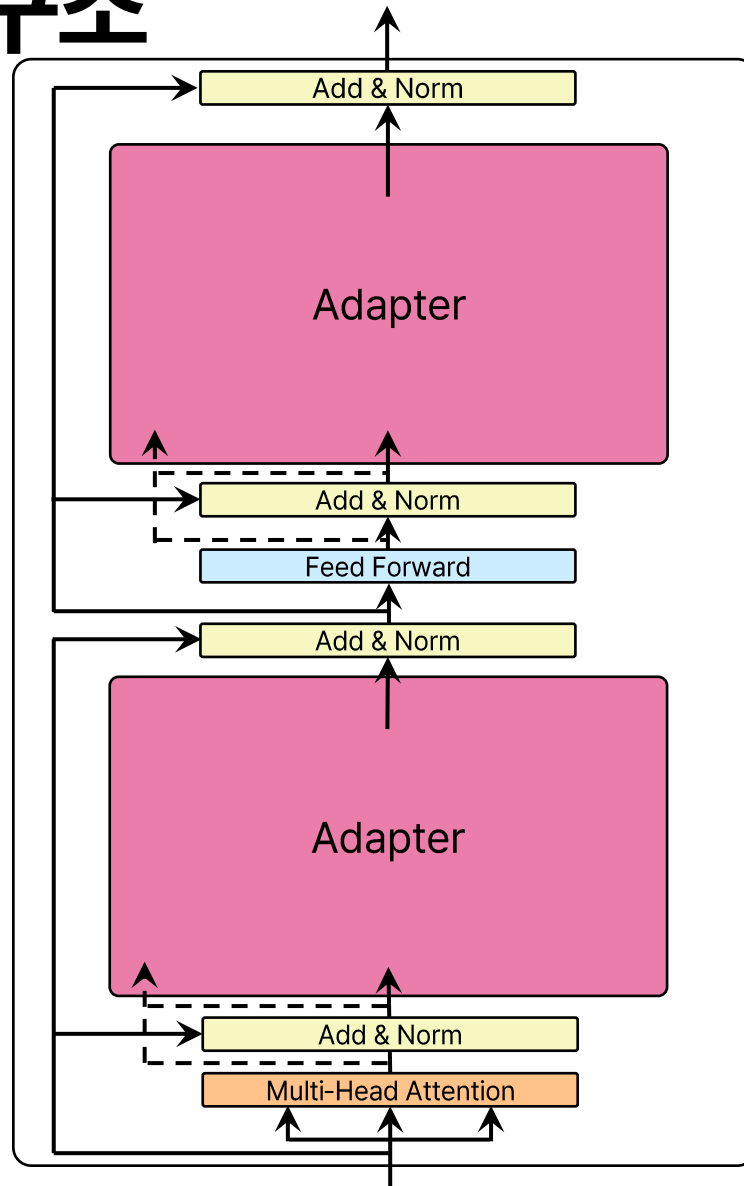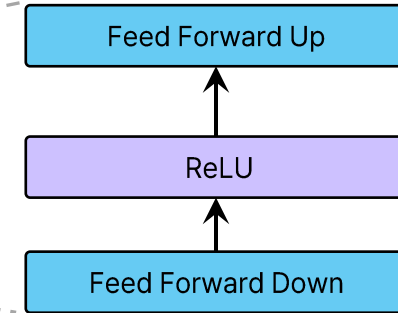전체 파라미터를
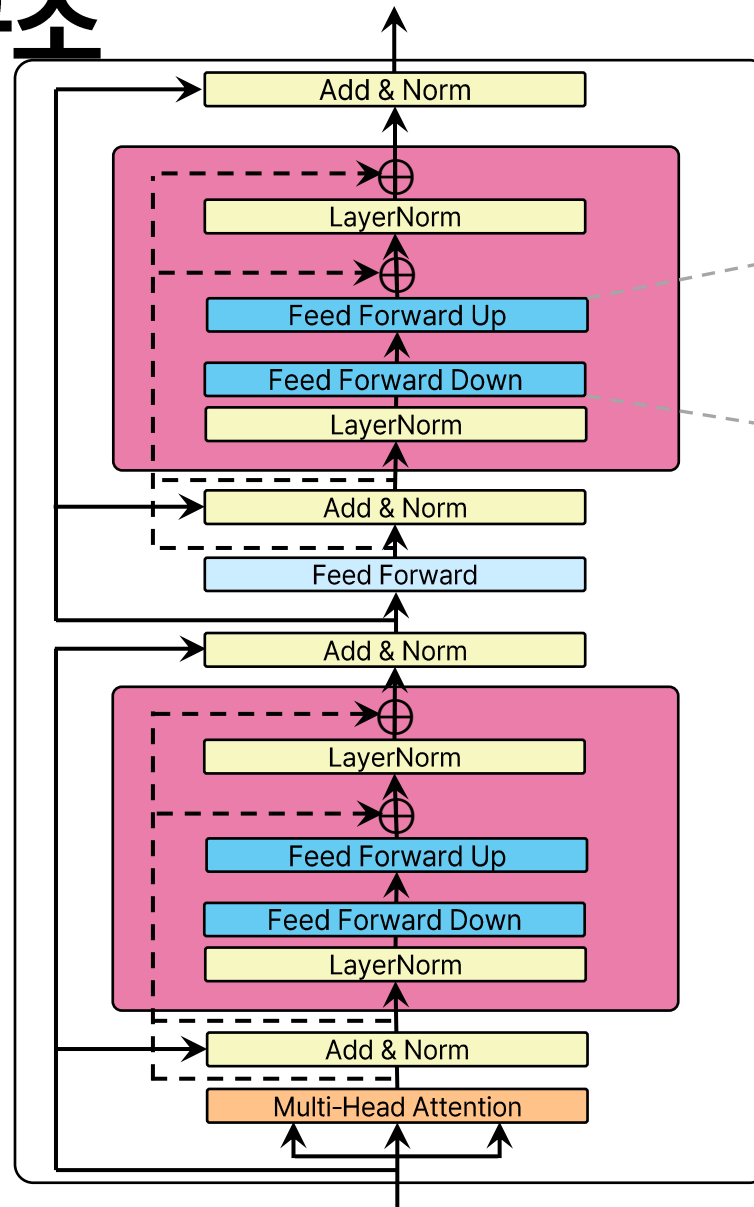학습하는게 아니라
Adapter 내부의
파라미터만 학습

# GEN AI 인텐시브 과정

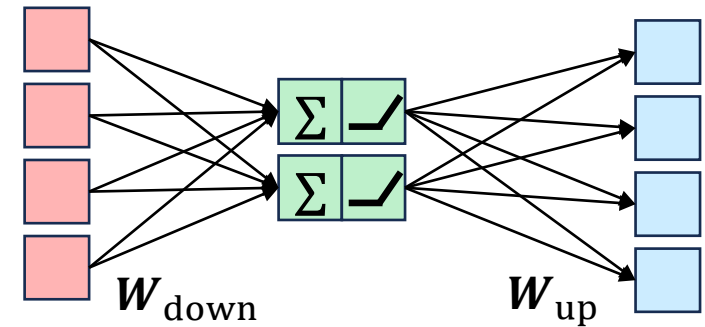Section 2. Adapter

**Section 2-2. Adapter의 구조**

# Adapter의 구조

# Adapter의 구조



$$W_{\text{up}}[\max(\mathbf{0}, W_{\text{down}}\mathbf{x} + \mathbf{b}_{\text{down}})] + \mathbf{b}_{\text{up}}$$

auto encoder

$W_{\text{down}}$

$W_{\text{up}}$

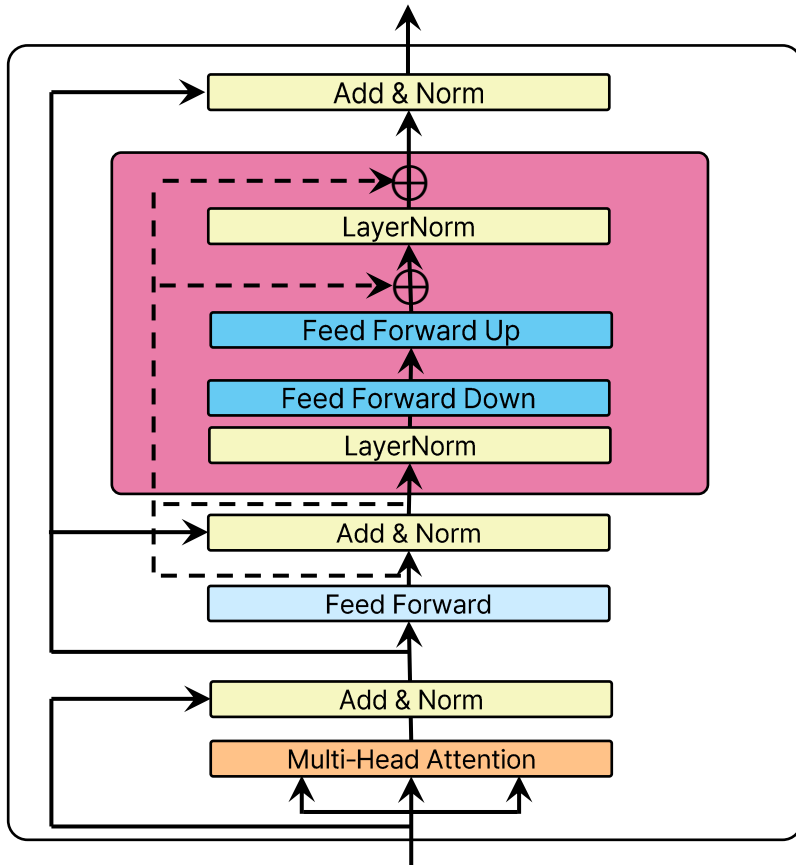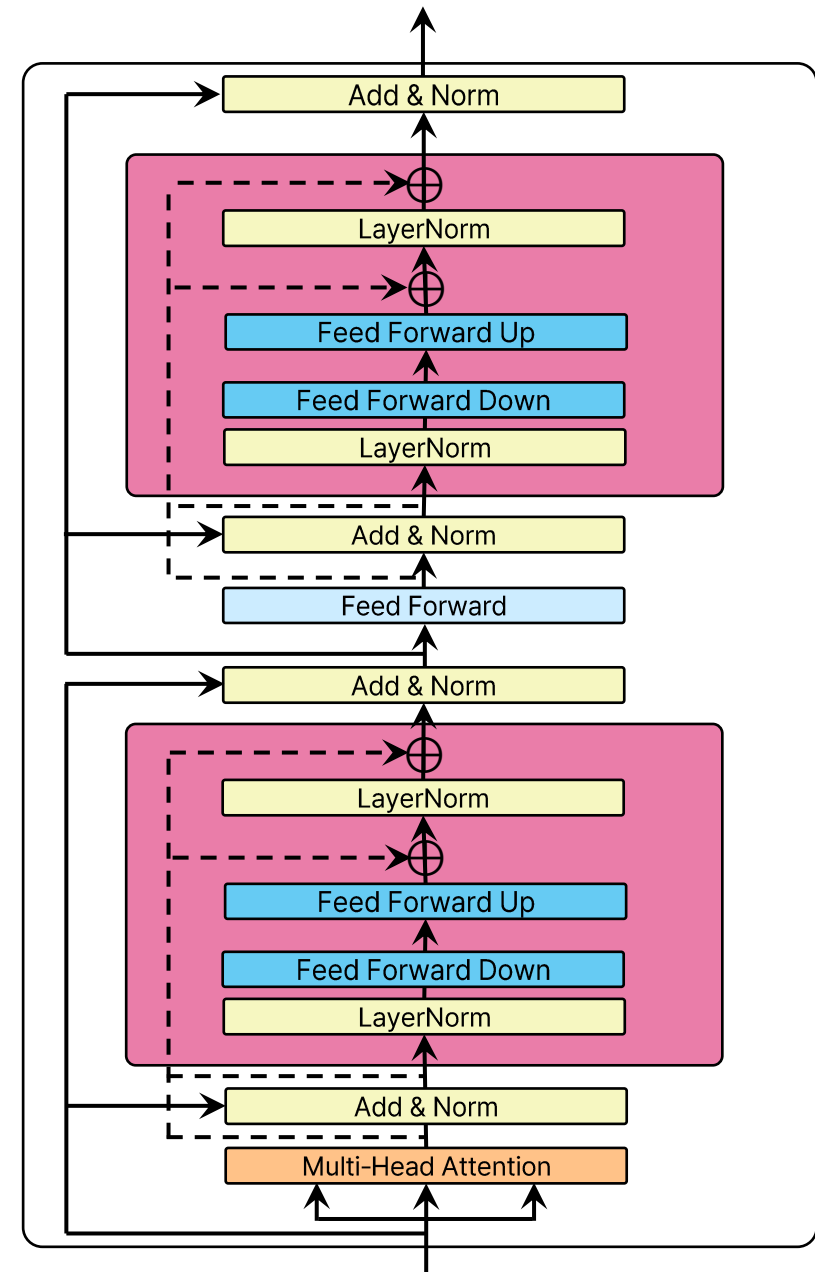Feed Forward down | ReLU | Feed Forward Up

# GEN AI 인텐시브 과정

Section 2. Adapter

**Section 2-3. Adapter의 종류**

# Adapter의 종류



**Pfeiffer Architecture**

**Houlsby Architecture**

# Pfeiffer Architecture

- Adapter 삽입 위치에 따라 Pfeiffer Architecture와 Houlsby Architecture로 나뉨

- Pfeiffer Architecture

  - 하나의 Adapter만 삽입(Feed Forward 단계 이후)

  - 하나의 Adapter만 다루므로 구조가 단순하고 효율적

  - 학습해야할 파라미터가 적음

  - 작업이 단순하거나 리소스가 제한적일 때 효과적으로 사용 가능

# Houlsby Architecture

- Houlsby Architecture

  - 두 개의 Adapter만 삽입(Multi-head attention, Feed Forward 단계 이후)

  - Pfeiffer보다 구조가 복잡하고 학습해야할 파라미터가 많음

  - Adapter를 두 개 사용하므로 다양한 표현을 학습할 수 있어 성능이 더 좋은 경우가 많음

  - 복잡한 Task에 효과적

# GEN AI 인텐시브 과정

Section 3. AdapterFusion

## Section 3-1. AdapterFusion의 개념

# AdapterFusion의 개념



**AdapterFusion:**
**Non-Destructive Task Composition for Transfer Learning**
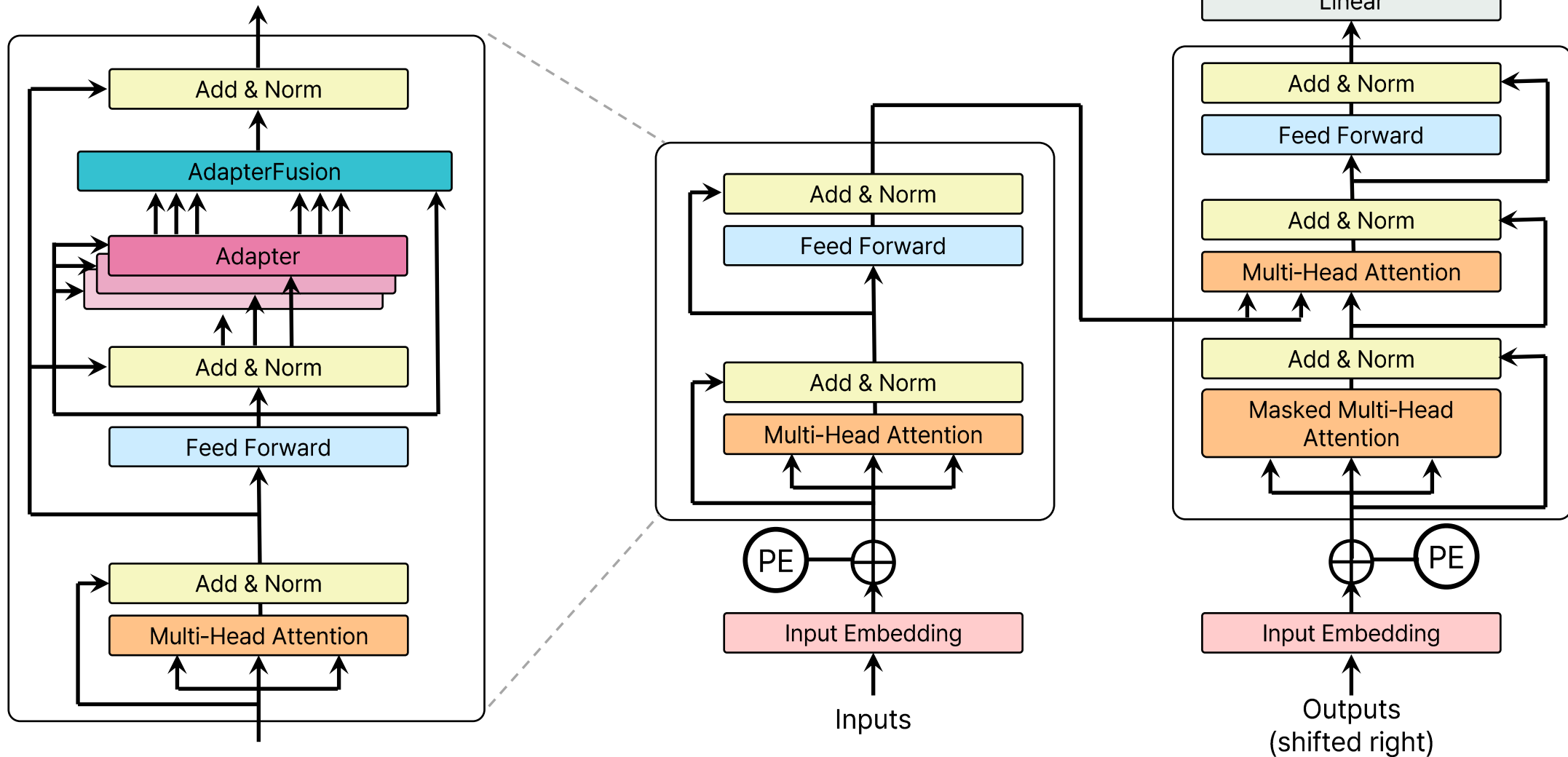
Jonas Pfeiffer et al, 2020

<span style="color:red">트랜스포머 인코더 모델인 BERT 기준으로 설명</span>

# AdapterFusion의 개념

- <u>여러 개의 Adapter들을 동시에 활용하는 방법</u>

- 지도학습의 앙상블 학습(ensemble learning)과 비슷한 원리

- 각각의 Adapter들은 이미 사전 학습되어 있는 상황이고, 더이상 학습하지 않고 Freeze 상태 유지

- 이후 AdapterFusion 레이어만 새롭게 추가하고 AdapterFusion 레이어만 학습함

  - AdapterFusion 레이어에는 이전에 각 Adapter들의 출력값을 조합하는 파라미터로 구성되어 있음
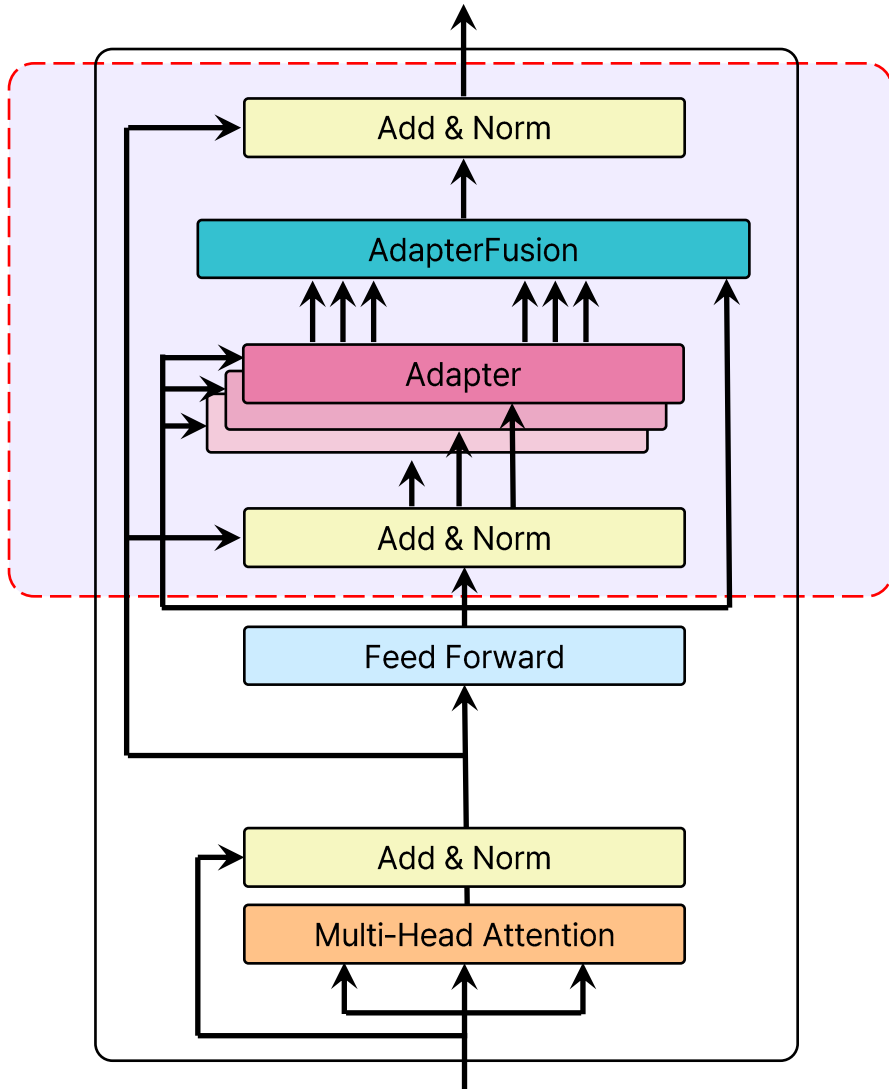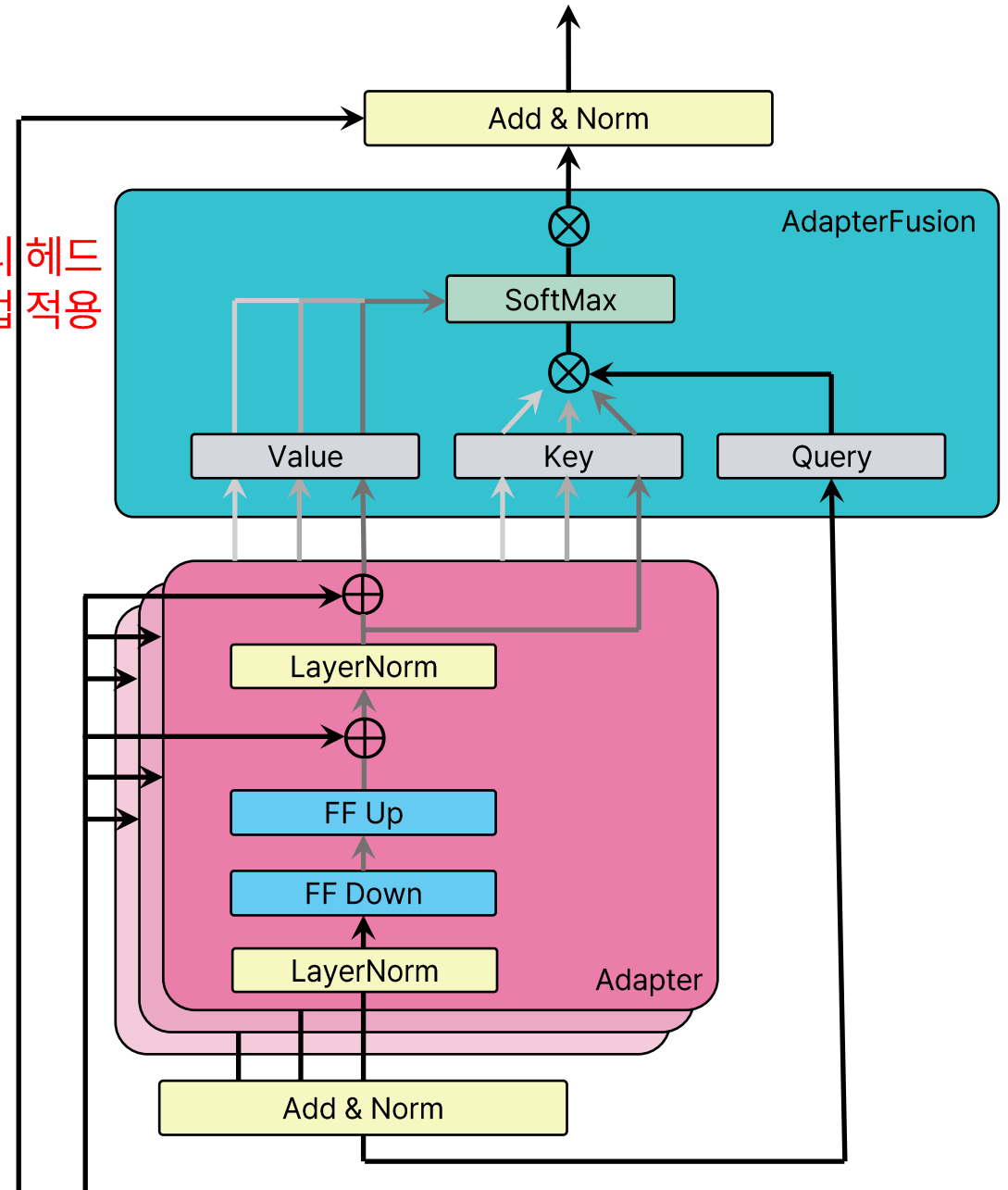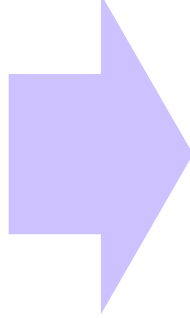
# AdapterFusion의 개념

# GEN AI 인텐시브 과정

Section 3. AdapterFusion

## Section 3-2. AdapterFusion의 구조

# AdapterFusion의 구조



기존의 멀티 헤드
어텐션 방법 적용

# GEN AI 인텐시브 과정

Section 4. LoRA

**Section 4-1. LoRA의 개념**

# LoRA의 개념



LoRA:
Low-Rank Adaptation of Large Language Models

Edward Hu et al, 2021

Rank :               feature

존재감          성능          가격
6             3            5

4             2

Rank : 2
(              =        *2)

Rank                - >

# LoRA의 개념

- 등장 배경

    - 풀 파인 튜닝 방식처럼 전체 파라미터를 재학습 시키기도 싫고...

    - 앞서 배운 Adapter처럼 외부 모듈을 트랜스포머 내부에 삽입하기도 싫고...


- LoRA는 기존 모델의 내부 파라미터를 수정하는 방식

    - 풀 파인 튜닝 방식처럼 모든 파라미터를 학습시키는 것이 아니라,

    - 기존 파라미터에 추가되는 보정 행렬만 학습

    - 보정 행렬의 Rank가 작아서 학습해야할 파라미터 개수 감소

# LoRA의 개념



Multi-Head Attention

LoRA는 트랜스포머 내부 구조에서 Attention이나 Feed Forward 과정에 적용할 수 있음
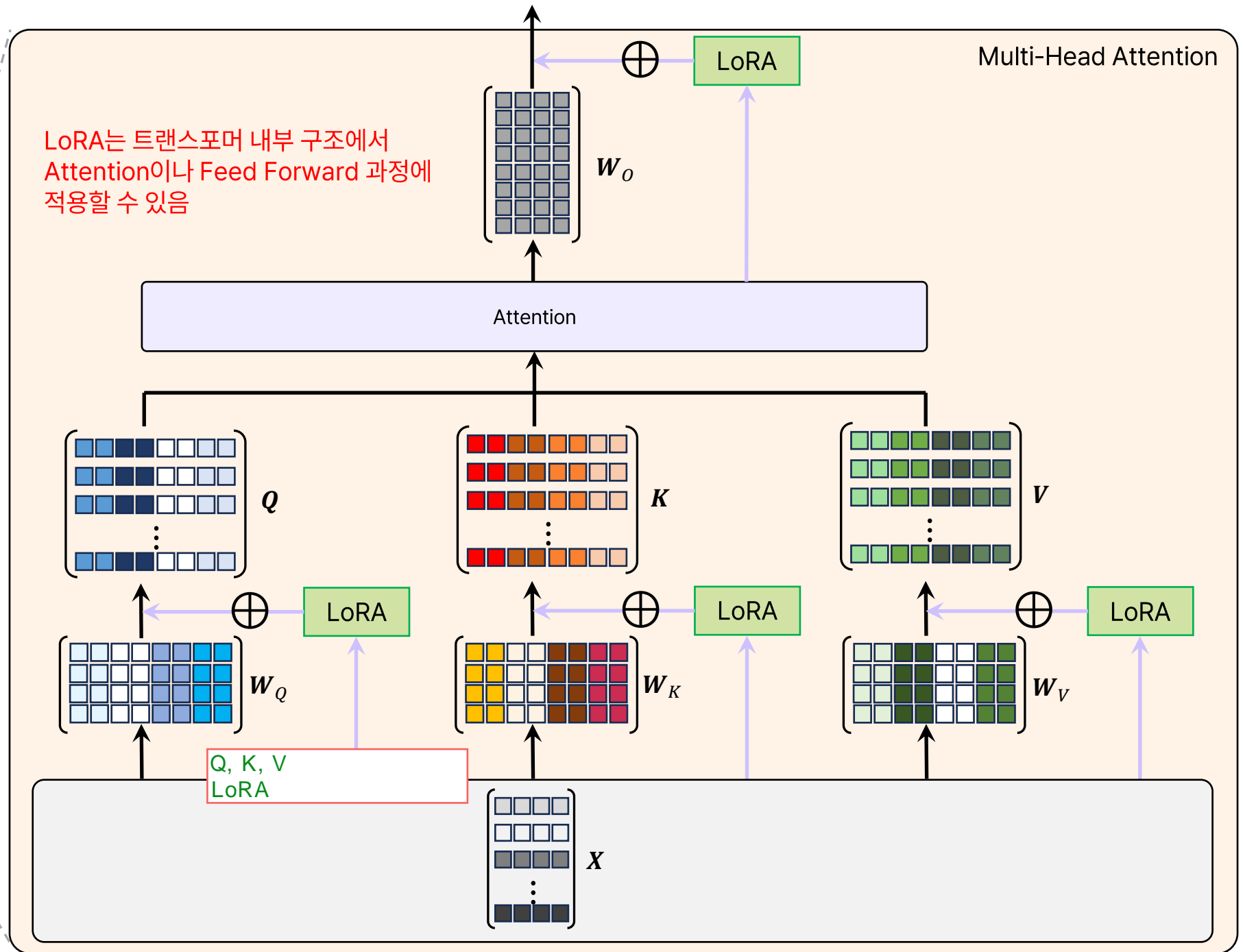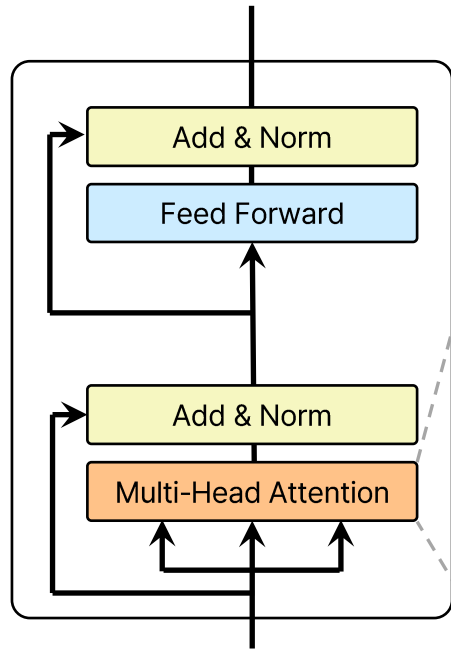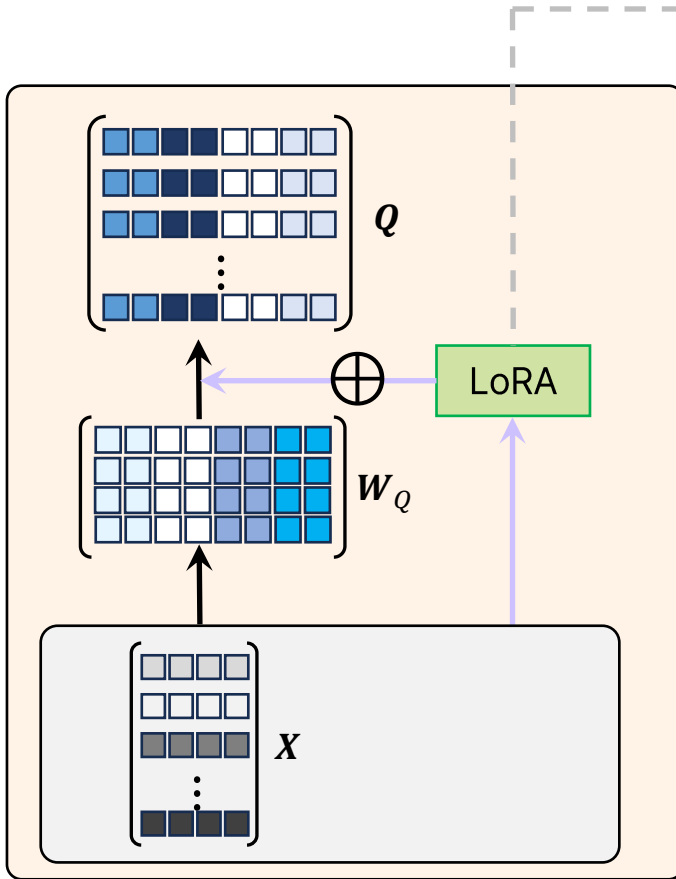
# GEN AI 인텐시브 과정

Section 4. LoRA

**Section 4-2. LoRA의 구조**

# LoRA의 구조

rank  - > full rank

$$X(W_Q + \Delta W) = Q$$

$n \times d \quad d \times (h \cdot d_h) \quad d \times (h \cdot d_h)$

$$W_Q + BA$$

$, r \ll \min(d, h \cdot d_h)$

rank



$W_Q$    $B$

$$+ \quad A$$

$d \times (h \cdot d_h) \qquad d \times r \qquad r \times (h \cdot d_h)$

행렬 BA의 Rank는 r

# 감사합니다.

# Q & A