

AI 모델 개발 Part 1 (심화 강의)

목차

- Part 1: AI 모델 아키텍처 설계
 - 1.1.1 CNN의 생물학적 영감
 - 1.1.2 컨볼루션 연산의 수학적 원리
 - 1.1.3 CNN의 핵심 하이퍼파라미터
 - 1.1.4 CNN 아키텍처의 진화
 - 1.1.5 CNN의 고급 기법과 최적화
 - 1.2.1 순환 신경망의 한계와 LSTM의 등장
 - 1.2.2 LSTM의 게이트 메커니즘: 정보 흐름 제어
 - 1.2.3 RNN의 고급 기법과 최적화
 - 1.3.1 Diffusion Model Forward Process
 - 1.3.2 Diffusion Model Reverse Process
 - 1.3.3 Diffusion Model의 고급 기법
 - 1.3.4 Diffusion Model의 장점과 응용
 - 1.4.1 Attention is All You Need
 - Attention 메커니즘의 철학적 배경
 - 1.4.2 Multi-Head Attention의 다관점 학습
 - 1.4.3 Transformer의 고급 기법과 최적화
 - 1.4.4 위치 인코딩과 순서 정보
 - 1.5.1 Transformer 아키텍처 패밀리 : BERT
 - 1.5.2 Transformer 아키텍처 패밀리 : GPT
 - 1.5.3 Transformer 아키텍처 패밀리 : T5
 - Encoder-Decoder: T5 (Text-to-Text Transfer Transformer)
 - 1.5.4 모델 아키텍처 선택 전략
- Part 2: Explainable AI (XAI)
 - XAI의 필요성과 사회적 요구
 - 2.1.1 AI의 블랙박스 문제와 해석 가능성
 - 2.1.2 XAI의 사회적 요구
 - 2.2.1 XAI 기법 심층 분석: PDP
 - Partial Dependence Plot (PDP): 전역적 특성 이해
 - 2.2.2 XAI 기법 심층 분석: ICE
 - Individual Conditional Expectation (ICE): 개별적 특성 이해
 - 2.2.3 XAI 기법 심층 분석: Grad-CAM
 - Grad-CAM: 시각적 주의 메커니즘
 - 2.2.4 XAI 기법 심층 분석: LIME & SHAP
 - LIME (Local Interpretable Model-agnostic Explanations)
 - 2.2.5 XAI 기법 선택 가이드
 - 2.2.6 성능과 해석성의 균형

AI 모델 개발 Part 1

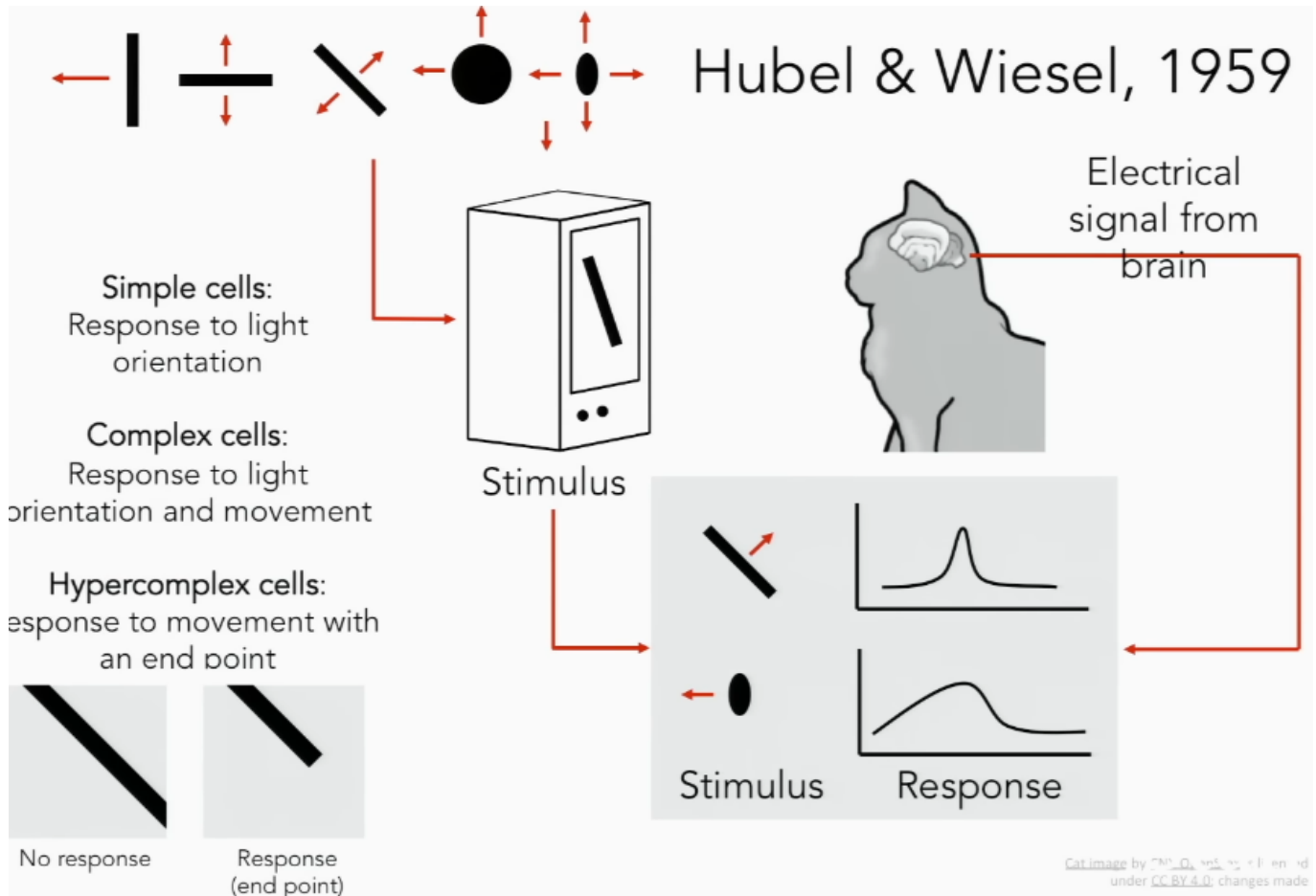
아키텍처의 심층 이해와 설명 가능한 AI

"모델을 이해하고, 설명할 수 있어야 진정한 AI 개발자"

Part 1: AI 모델 아키텍처 설계

1.1.1 CNN의 생물학적 영감

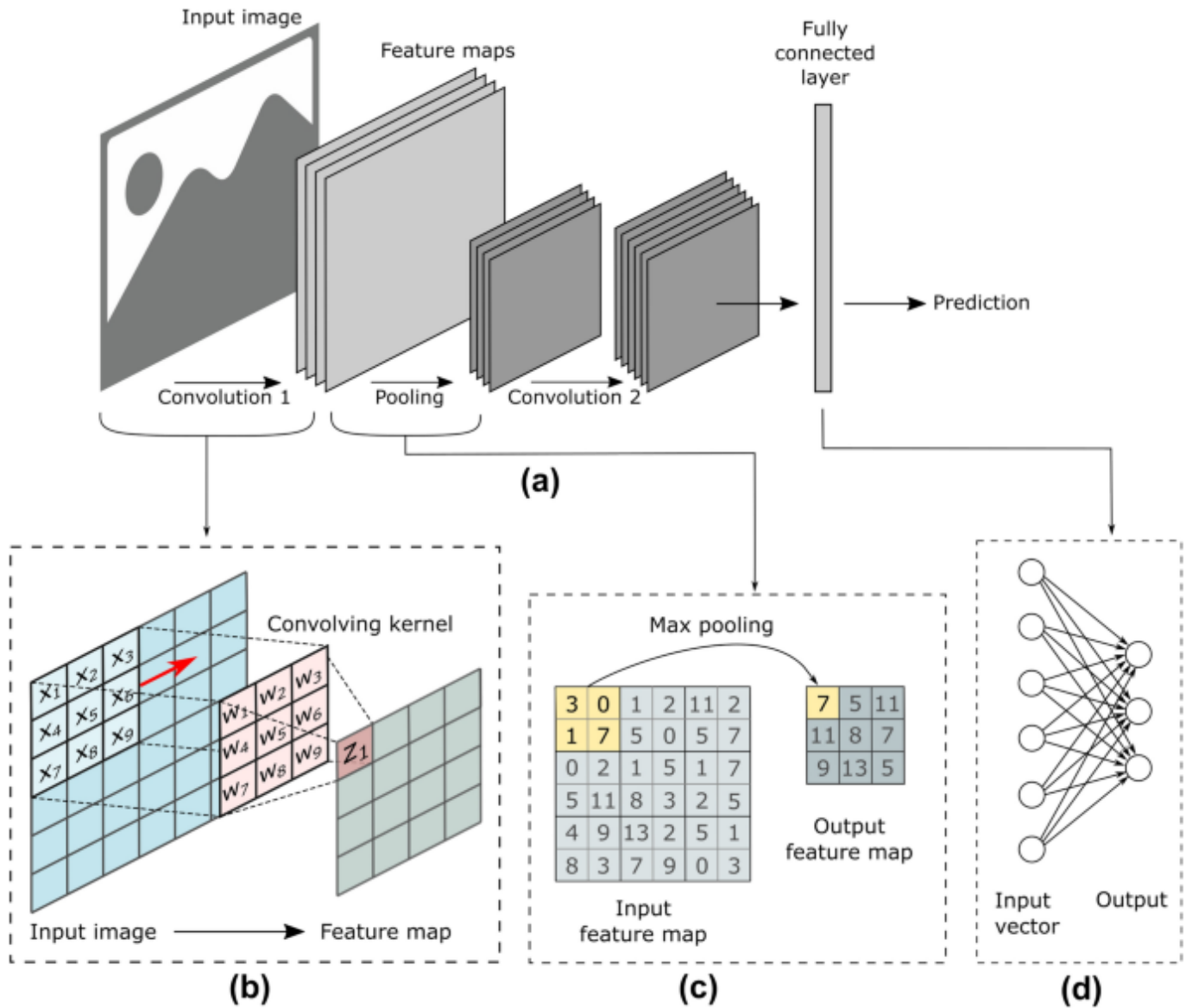
CNN(Convolutional Neural Networks)의 설계는 1960년대 Hubel과 Wiesel의 고양이 시각 피질 연구에서 영감을 받았습니다. 이들은 시각 피질의 뉴런들이 계층적으로 조직되어 있으며, 각 층이 점진적으로 복잡한 시각적 특징을 인식한다는 것을 발견했습니다. CNN은 이러한 생물학적 원리를 컴퓨터 과학적으로 구현한 것입니다.



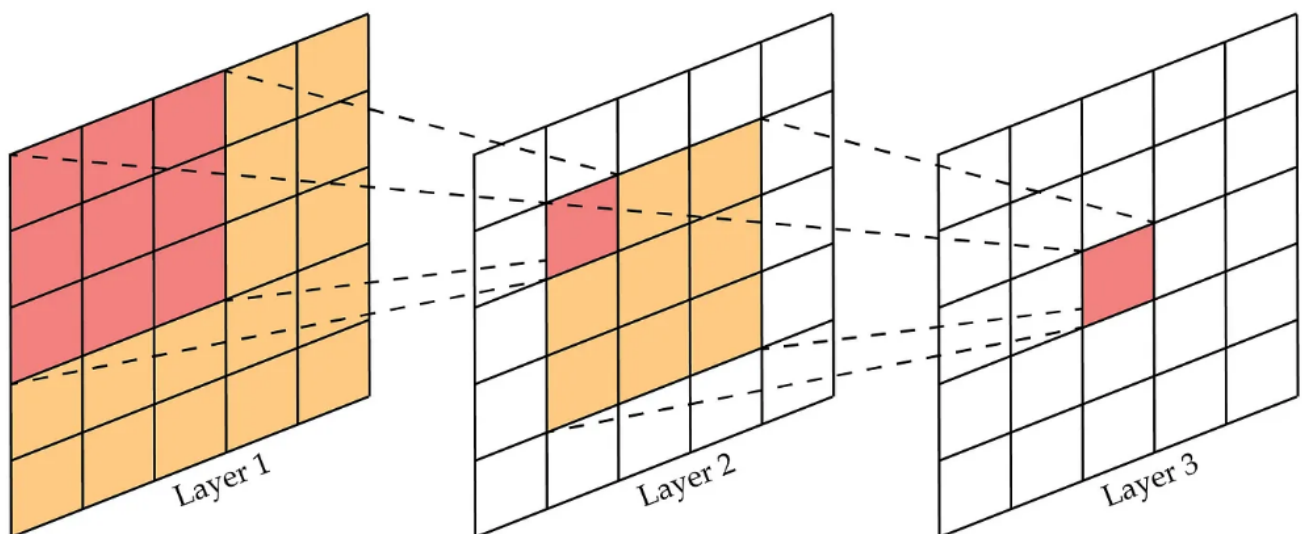
- 단순 세포(Simple Cells): 특정 방향의 가장자리(Edge)를 탐지하는 역할. CNN의 컨볼루션 필터에 해당.
- 복합 세포(Complex Cells): 위치에 무관하게 특징을 탐지하는 역할. CNN의 풀링(Pooling) 레이어에 해당.
- 계층적 처리: 저수준 특징(선, 점)에서 고수준 특징(객체의 부분, 전체)으로 점진적으로 추상화.

1.1.2 컨볼루션 연산의 수학적 원리

컨볼루션은 신호 처리에서 두 함수의 관계를 나타내는 연산으로, CNN에서는 이산 컨볼루션을 사용합니다. 입력 이미지와 필터(커널)의 연산을 통해 특징 맵(Feature Map)을 생성합니다.



Receptive Field in Convolutional Networks



- 필터(커널): 학습 가능한 가중치 행렬로, 이미지에서 특정 패턴(엣지, 질감, 코너 등)을 탐지하는 역할을 합니다.

- 특징 맵: 컨볼루션 연산의 결과물로, 입력 이미지에서 필터가 탐지한 특징의 위치와 강도를 나타냅니다.
- 수용 영역(Receptive Field): 각 뉴런이 반응하는 입력 이미지의 영역으로, 깊은 층으로 갈수록 넓어져 더 추상적인 특징을 인식하게 됩니다.

1.1.3 CNN의 핵심 하이퍼파라미터

핵심 하이퍼파라미터의 영향

- 스트라이드(Stride): 필터가 이미지를 가로지르는 간격입니다. 스트라이드가 클수록 특징 맵의 크기가 줄어들어 다운샘플링 효과와 계산량 감소를 가져옵니다. (e.g., Stride 1: 세밀한 특징 추출, Stride 2: 공간 해상도 감소)
- 패딩(Padding): 이미지 경계 부분의 정보 손실을 막기 위해 이미지 주변에 픽셀을 추가하는 기법입니다. 'Valid' 패딩은 패딩이 없고, 'Same' 패딩은 출력 크기를 입력과 같게 유지합니다. 경계 효과를 최소화하고 가장자리 정보를 보존하는 데 중요합니다.
- 필터 크기(Kernel Size): 3x3, 5x5, 7x7 등 필터의 크기로, 큰 필터는 넓은 수용 영역을 가지지만 계산량이 증가합니다.
- 채널 수(Channels): 각 레이어의 필터 개수로, 깊은 레이어일수록 더 많은 채널을 사용하여 복잡한 특징을 학습합니다.

출력 크기 계산 공식

$$\text{Output_Size} = (\text{Input_Size} - \text{Kernel_Size} + 2 * \text{Padding}) / \text{Stride} + 1$$

1.1.4 CNN 아키텍처의 진화

CNN 아키텍처의 발전사

- LeNet-5 (1998): CNN의 시초. 손글씨 숫자 인식에 사용. 단순한 구조와 Sigmoid 활성화 함수.
- AlexNet (2012): GPU 활용, ReLU 활성화 함수, 드롭아웃, 데이터 증강 기법 도입으로 딥러닝 부활의 신호탄. 대규모 이미지 분류 태스크에서 성능 입증.
- VGGNet (2014): 3x3의 작은 필터를 연속 사용하며 네트워크 깊이의 중요성 입증. 균일한 아키텍처 설계 철학.
- ResNet (2015): 잔차 연결(Residual Connection)이라는 혁신적 구조를 통해 100층 이상의 매우 깊은 네트워크 학습을 가능하게 했습니다. 입력 신호가 일부 레이어를 건너뛰도록(Skip) 허용하여, 깊은 망에서 발생하는 그래디언트 소실 문제를 획기적으로 해결했습니다.
- Inception (2014-2016): 다양한 크기의 필터를 병렬로 사용하는 Inception 모듈 도입. 1x1 컨볼루션으로 차원 축소, 계산 효율성과 표현력 균형.
- DenseNet (2017): 모든 레이어를 직접 연결하는 Dense Block 구조. 특징 재사용과 그래디언트 흐름 개선.
- EfficientNet (2019): Compound Scaling을 통한 체계적 모델 확장. 깊이, 너비, 해상도를 균형있게 조정.

```
# ResNet의 핵심: Residual Block.  $H(x) = F(x) + x$ 
# DenseNet의 핵심: Dense Block.  $x_l = H_l([x_0, x_1, \dots, x_{l-1}])$ 
```

1.1.5 CNN의 고급 기법과 최적화

현대적 CNN 기법들

- Batch Normalization (2015): 미니배치 단위로 정규화하여 내부 공변량 이동(Internal Covariate Shift) 문제 해결. 학습 속도 향상과 그래디언트 소실 방지.
- Group Convolution: 채널을 그룹으로 나누어 병렬 처리. ResNeXt, MobileNet에서 사용하여 계산 효율성 증대.
- Depthwise Separable Convolution: 깊이별 컨볼루션과 점별 컨볼루션을 분리. MobileNet의 핵심으로 계산량 대폭 감소.
- Attention Mechanism: Self-Attention을 CNN에 도입. SENet의 Squeeze-and-Excitation, CBAM의 Channel & Spatial Attention.

실무 구현 예시

```
# Depthwise Separable Convolution 구현 (TensorFlow)
from tensorflow.keras.layers import Conv2D, DepthwiseConv2D

def depthwise_separable_conv(x, filters, kernel_size):
    # Depthwise convolution
    x = DepthwiseConv2D(kernel_size, padding='same', use_bias=False)(x)
    # Pointwise convolution
    x = Conv2D(filters, 1, padding='same', use_bias=False)(x)
    return x

# Batch Normalization 효과
from tensorflow.keras.layers import BatchNormalization, ReLU

def conv_bn_relu(x, filters, kernel_size):
    x = Conv2D(filters, kernel_size, padding='same')(x)
    x = BatchNormalization()(x)
    x = ReLU()(x)
    return x

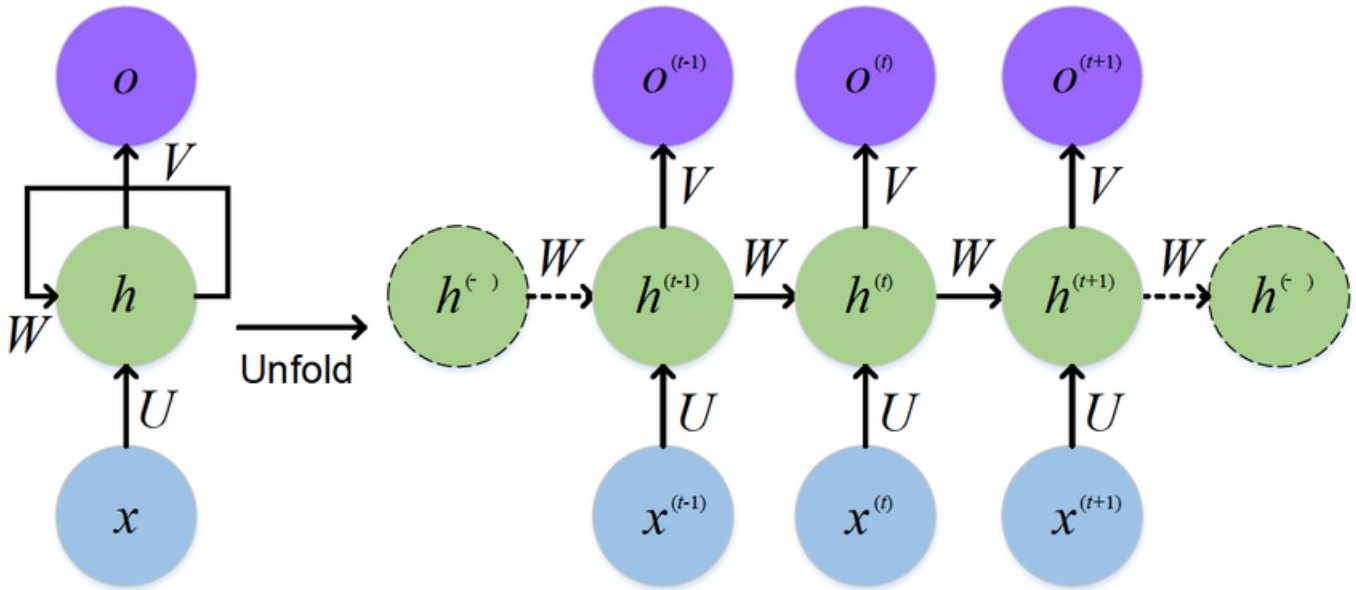
# PyTorch 버전
import torch.nn as nn

class DepthwiseSeparableConv(nn.Module):
    def __init__(self, in_channels, out_channels, kernel_size):
        super().__init__()
        self.depthwise = nn.Conv2d(in_channels, in_channels, kernel_size,
                                     padding=kernel_size//2, groups=in_channels)
        self.pointwise = nn.Conv2d(in_channels, out_channels, 1)

    def forward(self, x):
        x = self.depthwise(x)
        x = self.pointwise(x)
        return x
```

1.2.1 순환 신경망의 한계와 LSTM의 등장

기존 순환 신경망(RNN)은 시퀀스가 길어질수록 앞쪽의 정보가 뒤로 전달되지 못하는 장기 의존성 문제(Long-term Dependency Problem)가 있었습니다. 이는 역전파 시 그래디언트가 소실되거나 폭발하기 때문입니다.



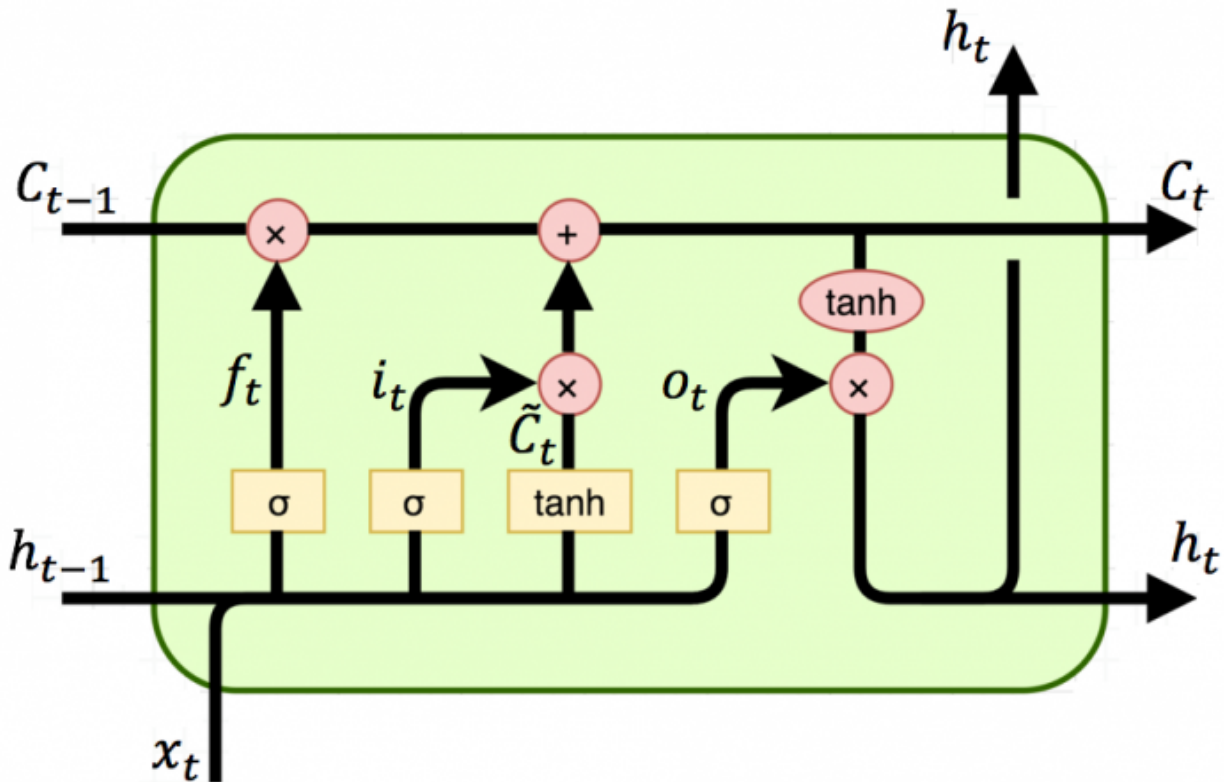
그래디언트 소실 문제의 수학적 분석:

$$\partial L / \partial h_s = \partial L / \partial h_t \cdot \prod_{(k=s+1 \text{ to } t)} \partial h_k / \partial h_{(k-1)}$$

연쇄 법칙에 의한 그래디언트의 곱셈 구조로 인해, 가중치가 1보다 작으면 지수적으로 감소(소실), 1보다 크면 지수적으로 증가(폭발)하여 장기 의존성 학습이 어려워집니다. 이런 문제에 대한 해결 방법으로 LSTM(Long Short-Term Memory)이 등장하게 됩니다.

1.2.2 LSTM의 게이트 메커니즘: 정보 흐름 제어

LSTM은 셀 상태(Cell State)라는 별도의 정보 고속도로를 두고, 세 개의 정교한 게이트를 통해 정보를 선택적으로 제어합니다. 이는 컴퓨터의 메모리 시스템과 유사한 개념입니다.



- Forget Gate (f_t): 이전 셀 상태에서 어떤 정보를 버릴지 결정합니다. $(\sigma(W_f \cdot [h_{t-1}, x_t] + b_f))$
- Input Gate (i_t): 현재 입력과 이전 은닉 상태에서부터 새로운 정보 중 어떤 것을 셀 상태에 새로 저장할지 결정합니다. $(\sigma(W_i \cdot [h_{t-1}, x_t] + b_i))$
- Output Gate (o_t): 업데이트된 셀 상태 정보 중 어떤 것을 다음 시점의 은닉 상태로 출력할지 결정합니다. $(\sigma(W_o \cdot [h_{t-1}, x_t] + b_o))$

LSTM의 핵심: Cell State 업데이트
 $C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$ (\tilde{C}_t 는 새로운 후보 값)

이러한 게이트 메커니즘 덕분에 LSTM은 시계열 데이터나 자연어처럼 장기적인 맥락이 중요한 태스크에서 뛰어난 성능을 보입니다.

1.2.3 RNN의 고급 기법과 최적화

RNN 최적화 기법들

- Bidirectional RNN: 과거와 미래 정보를 모두 활용. 시퀀스의 전체 맥락을 고려하여 성능 향상.
- Stacked RNN: 여러 RNN 레이어를 쌓아 복잡한 패턴 학습. 깊은 네트워크의 표현력 활용.
- Dropout in RNN: 순환 연결에 드롭아웃 적용. Zoneout, Recurrent Dropout 등 다양한 기법.
- Gradient Clipping: 그래디언트 폭발 문제 해결. 임계값을 넘는 그래디언트를 제한.

Attention 기반 RNN

- Bahdanau Attention: 인코더-디코더 구조에서 동적 컨텍스트 벡터 생성.
- Luong Attention: 다양한 Attention 스코어 함수 제안.

- Self-Attention: 시퀀스 내 모든 위치 간 관계 모델링.

실무 구현 예시

```
# Attention 메커니즘 구현
import torch
import torch.nn.functional as F
import math

def attention_mechanism(query, keys, values):
    # Attention score 계산
    scores = torch.matmul(query, keys.transpose(-2, -1))
    attention_weights = F.softmax(scores / math.sqrt(d_k), dim=-1)

    # Weighted sum
    context = torch.matmul(attention_weights, values)
    return context, attention_weights

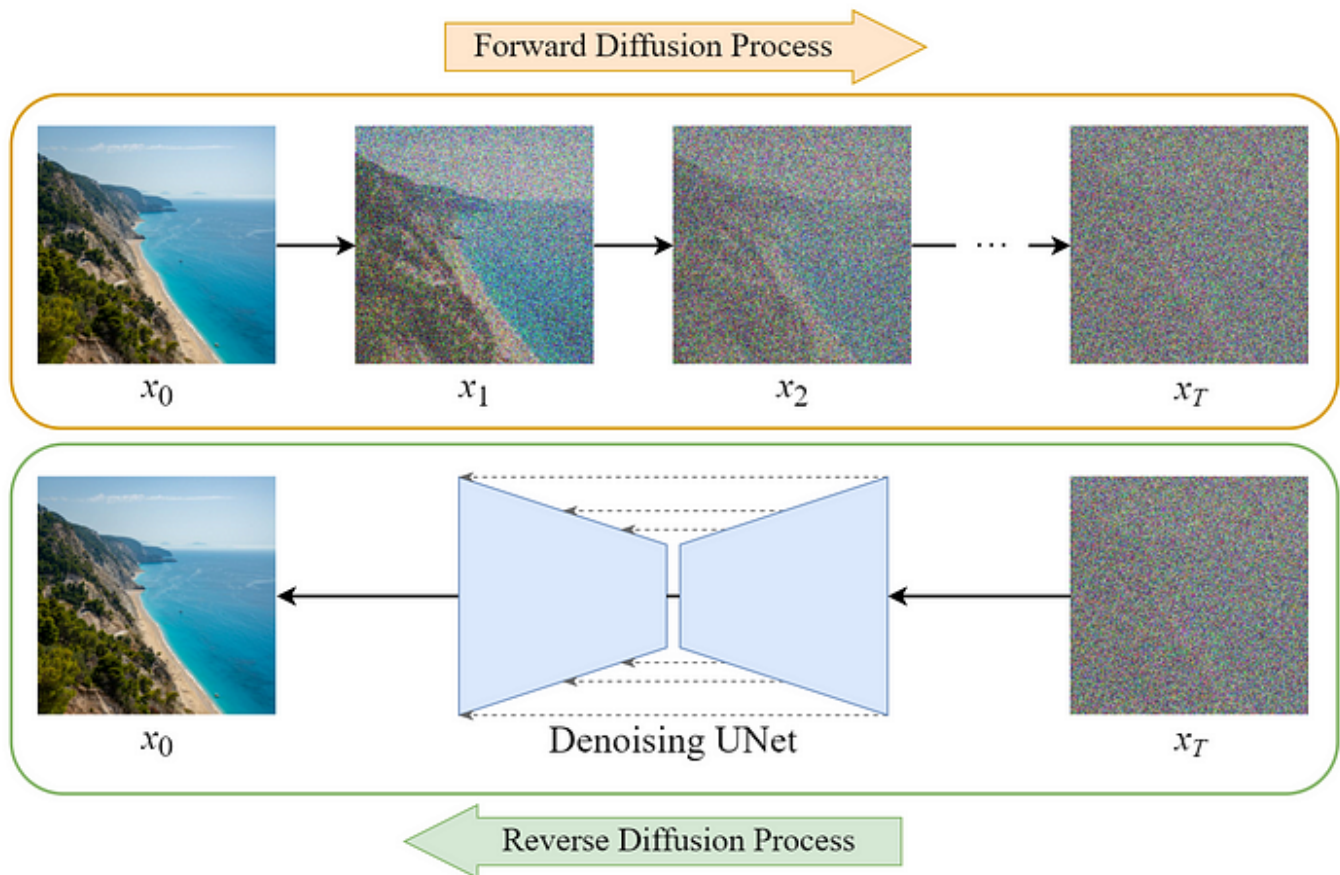
# Gradient Clipping
torch.nn.utils.clip_grad_norm_(model.parameters(), max_norm=1.0)
```

1.3.1 Diffusion Model Forward Process

물리학의 확산(Diffusion) 과정에서 영감을 받은 생성 모델로, GAN을 뛰어넘는 고품질의 다양한 샘플 생성 능력으로 현재 생성 AI의 주류 기술로 자리 잡았습니다.

Forward Process (확산 과정)

원본 이미지(x_0)에 정해진 스케줄(β_t)에 따라 점진적으로 가우시안 노이즈를 추가하여, 최종적으로는 완전한 노이즈 이미지(x_T)로 만듭니다. 이 과정은 정해져 있어 학습이 필요 없습니다.



수학적 분석

- 노이즈 스케줄: β_t 는 0에서 1로 증가하는 함수. 선형, 코사인, 스케줄 등 다양한 형태.
- 중간 분포: $q(x_t|x_0) = N(x_t; \sqrt{\alpha_t} x_0, (1-\alpha_t) I)$ where $\bar{\alpha}_t = \prod(1-\beta_i)$.
- 물리학적 해석: 엔트로피 증가 과정의 모델링, 열평형 상태로의 자발적 진행.

실무 구현 예시

```
# Forward Process 구현
import torch

def forward_diffusion(x_0, t, noise_schedule):
    # 노이즈 스케줄에서  $\bar{\alpha}_t$  계산
    alpha_bar = noise_schedule.get_alpha_bar(t)

    # 노이즈 생성
    noise = torch.randn_like(x_0)

    # 노이즈가 추가된 이미지
    x_t = torch.sqrt(alpha_bar) * x_0 + torch.sqrt(1 - alpha_bar) * noise
    return x_t, noise
```

1.3.2 Diffusion Model Reverse Process

노이즈 이미지(x_T)에서 시작하여, 각 시간 단계(timestep)에서 추가된 노이즈를 예측하고 제거하는 과정을 학습합니다. 이 역방향 과정을 통해 노이즈로부터 새로운 이미지를 창조합니다.

$$p_{\theta}(x_{t-1}|x_t) = N(x_{t-1}; \mu_{\theta}(x_t, t), \Sigma_{\theta}(x_t, t))$$

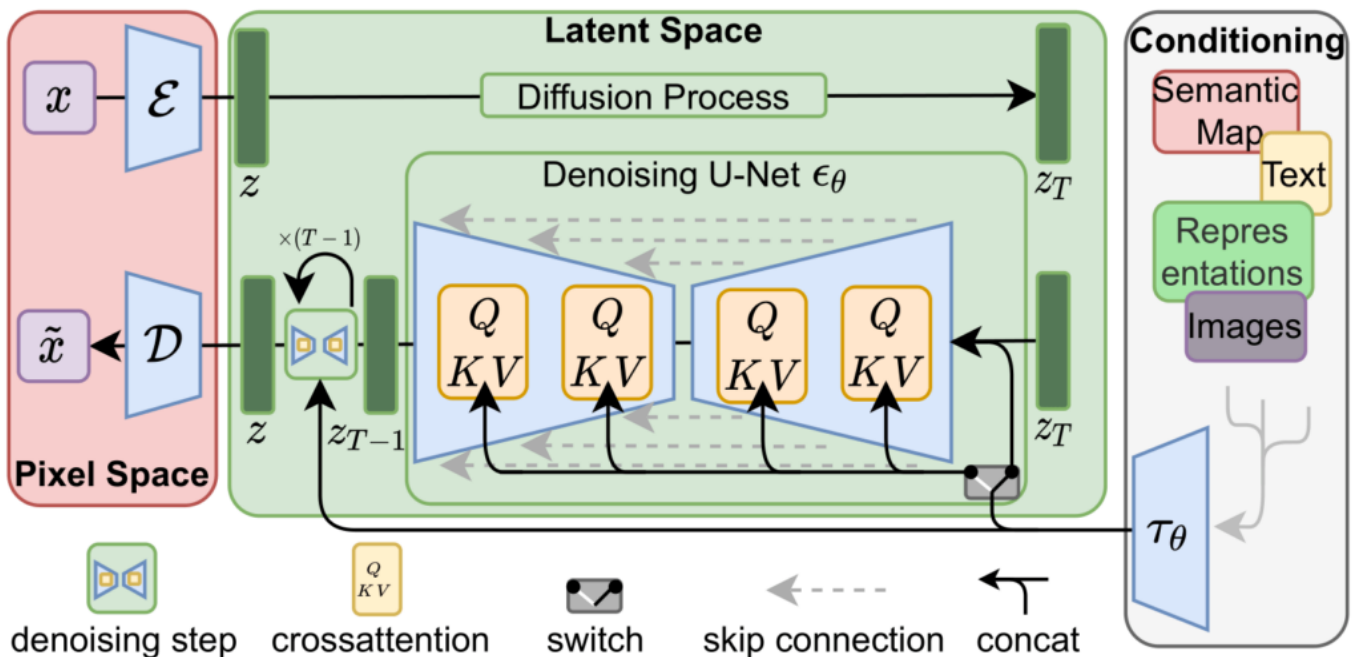
DDPM의 핵심 아이디어

- 노이즈 예측: 네트워크가 추가된 노이즈 ϵ 를 직접 예측하도록 학습.
- U-Net 구조: 인코더-디코더 구조로 다중 스케일 특징 활용.
- 시간 임베딩: Sinusoidal Positional Encoding으로 시간 정보 주입.

손실 함수

$$L = E[||\epsilon - \epsilon_{\theta}(x_t, t)||^2]$$

DDPM (Denoising Diffusion Probabilistic Models)은 이 과정에서 추가된 노이즈 자체를 예측하도록 U-Net 기반의 네트워크를 학습시켜 안정성과 성능을 크게 향상시켰습니다. GAN과 달리 학습이 안정적이고 모드 붕괴 문제가 없다는 장점이 있습니다.



실무 구현 예시

```
# Reverse Process 구현
import torch

def reverse_diffusion(x_t, t, model, noise_schedule):
    # 모델로 노이즈 예측
    predicted_noise = model(x_t, t)
```

```

    # 노이즈 제거
    alpha_bar = noise_schedule.get_alpha_bar(t)
    x_prev = (x_t - torch.sqrt(1 - alpha_bar) * predicted_noise) /
    torch.sqrt(alpha_bar)

    return x_prev

```

1.3.3 Diffusion Model의 고급 기법

최신 Diffusion 기법들

- DDIM (Denoising Diffusion Implicit Models): 확률적 과정을 결정적 과정으로 변환. 샘플링 속도 대폭 향상.
- Classifier-Free Guidance: 조건부 생성에서 분류기 없이 조건 정보 활용. CFG 스케일로 생성 품질 조절.
- Stable Diffusion: 잠재 공간에서의 Diffusion. 메모리 효율성과 생성 속도 개선.
- ControlNet: 추가 조건(스케치, 포즈 등)을 통한 정밀한 제어.

실무 구현 예시

```

# Classifier-Free Guidance 구현
import torch

def classifier_free_guidance(x_t, t, model, condition, cfg_scale=7.5):
    # 조건부 예측
    pred_cond = model(x_t, t, condition)
    # 무조건부 예측
    pred_uncond = model(x_t, t, None)

    # CFG 적용
    pred = pred_uncond + cfg_scale * (pred_cond - pred_uncond)
    return pred

# DDIM 샘플링
def ddim_sampling(model, x_T, timesteps, eta=0.0):
    for t in reversed(timesteps):
        # 노이즈 예측
        noise_pred = model(x_T, t)
        # DDIM 업데이트
        x_T = update_step_ddim(x_T, t, noise_pred, eta)
    return x_T

```

1.3.4 Diffusion Model의 장점과 응용

기존 생성 모델과의 비교

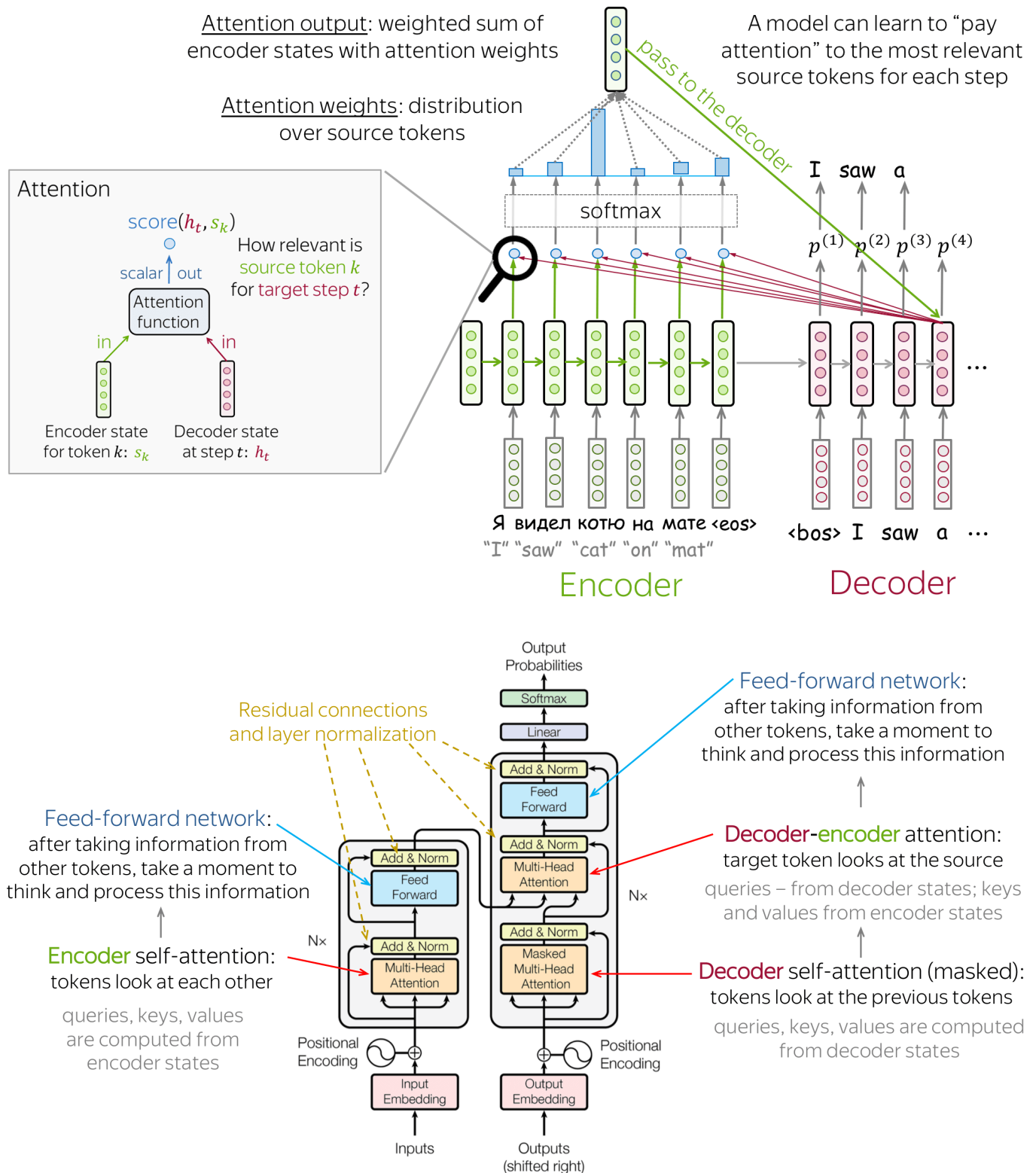
- GAN 대비 장점: 모드 붕괴 문제 해결, 안정적인 학습 과정, 고품질 다양한 샘플 생성.
- VAE 대비 장점: 더 선명하고 세밀한 이미지 생성, 후방 붕괴(Posterior Collapse) 문제 없음.

최신 응용 분야

- 이미지 생성: DALL-E 2, Midjourney, Stable Diffusion 등.
- 오디오 합성: WaveGrad.
- 3D 모델 생성: DreamFusion.
- 비디오 생성: Video Diffusion Models.

1.4.1 Attention is All You Need

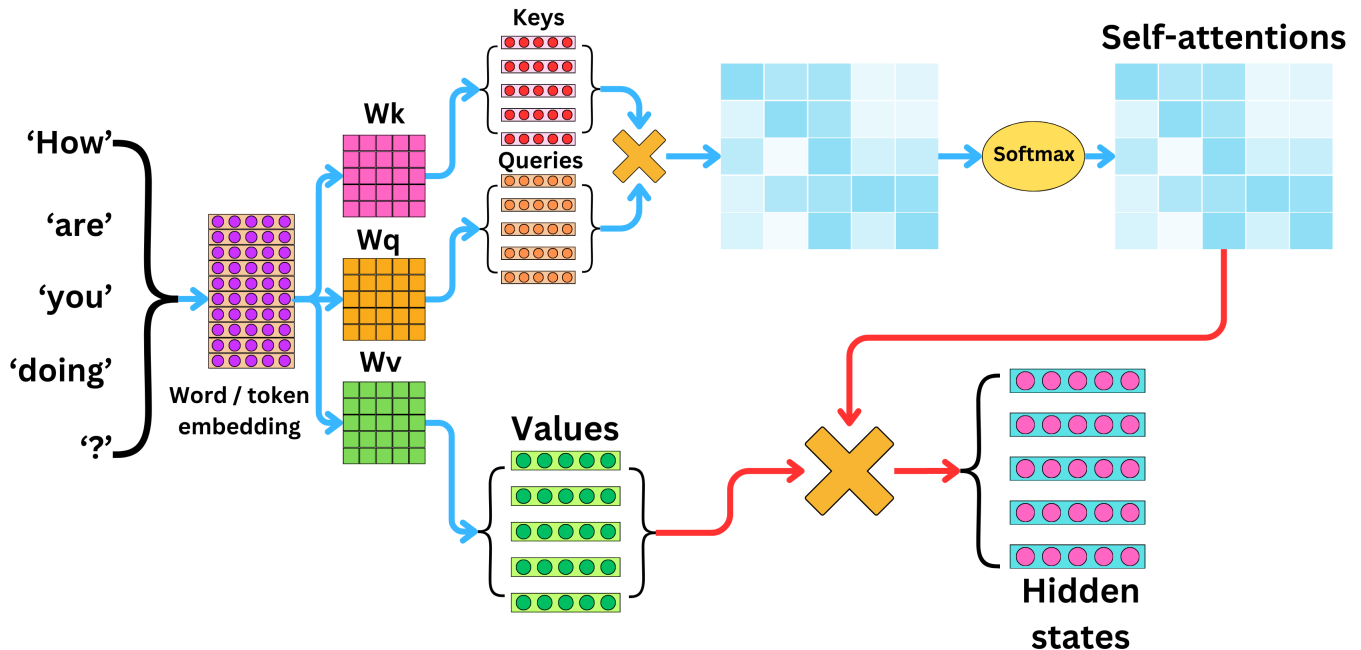
Attention 메커니즘의 철학적 배경



Transformer는 RNN의 순차적 처리 구조를 완전히 버리고, Self-Attention 메커니즘만으로 시퀀스 내 모든 요소 간의 관계를 병렬적으로 계산하는 혁신을 보여주었습니다. 이는 인간이 정보를 처리할 때 관련 있는 부분에 선택적으로 주의를 기울이는 과정을 모델링한 것입니다.

Self-Attention의 수학적 원리

모든 단어가 다른 모든 단어와 직접 상호작용하며 연관성을 계산합니다. 각 단어는 정보를 요청하는 Query(Q), 정보를 식별하는 Key(K), 실제 정보를 담은 Value(V)라는 세 가지 벡터로 표현됩니다.



$$\text{Attention}(Q, K, V) = \text{softmax}(QK^T / \sqrt{d_k})V$$

Attention의 핵심 구성요소

- Query(Q): 현재 위치에서 찾고자 하는 정보를 표현하는 벡터.
- Key(K): 각 위치가 가지고 있는 정보의 특성을 표현하는 벡터.
- Value(V): 실제로 전달될 정보를 담고 있는 벡터.
- Scale Factor ($\sqrt{d_k}$): 그래디언트 소실을 방지하기 위한 정규화 상수.

실무 구현 예시

```
# Self-Attention 구현
import torch
import torch.nn.functional as F
import math

def self_attention(query, key, value, mask=None):
    # Attention score 계산
    scores = torch.matmul(query, key.transpose(-2, -1))
    scores = scores / math.sqrt(query.size(-1))
```

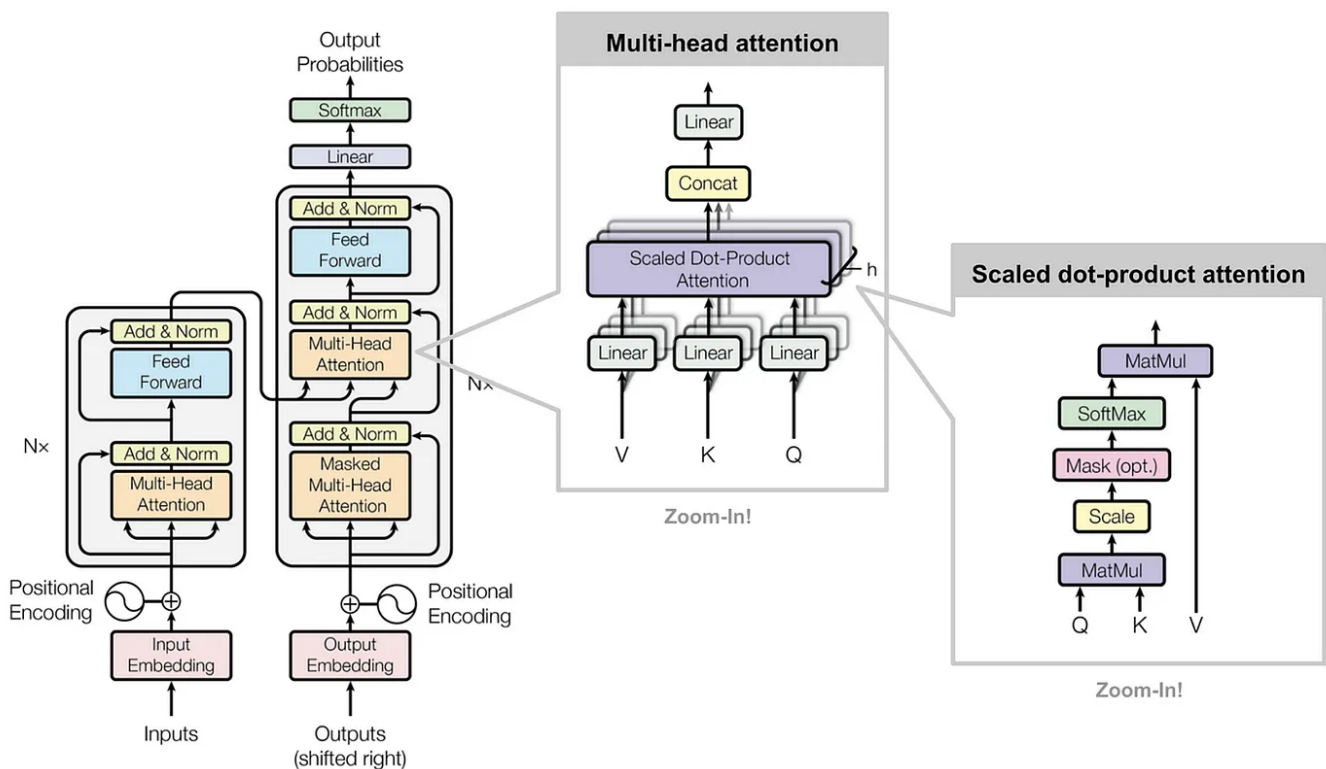
```
# Mask 적용 (선택적)
if mask is not None:
    scores = scores.masked_fill(mask == 0, -1e9)

# Softmax 적용
attention_weights = F.softmax(scores, dim=-1)

# Weighted sum
output = torch.matmul(attention_weights, value)
return output, attention_weights
```

1.4.2 Multi-Head Attention의 다관점 학습

Multi-Head Attention은 Self-Attention을 여러 개 병렬로 수행하여, 다양한 관점(구문적, 의미적 등)에서 문장을 동시에 바라보는 효과를 냅니다.



```
MultiHead(Q, K, V) = Concat(head_1, ..., head_h)W^O
head_i = Attention(QW_i^Q, KW_i^K, VW_i^V)
```

- 각 헤드의 특화된 역할: 구문적 관계 포착, 의미적 연관성 학습, 장거리 의존성 모델링, 다양한 언어적 패턴 인식.
- 이를 통해 모델은 더 풍부하고 다양한 관계를 학습할 수 있습니다.

1.4.3 Transformer의 고급 기법과 최적화

Transformer 최적화 기법들

- Relative Position Encoding: 절대적 위치 대신 상대적 위치 정보 활용. Shaw et al. (2018)의 방법.

- Layer Normalization: 각 레이어의 출력을 정규화하여 학습 안정성 향상.
- Residual Connection: 그래디언트 흐름 개선과 깊은 네트워크 학습 가능.
- Pre-Layer Norm: Layer Norm을 Attention과 FFN 앞에 배치하여 학습 안정성 향상.

효율적 Attention 기법

- Sparse Attention: 모든 토큰 쌍을 계산하지 않고 일부만 계산. Longformer, BigBird.
- Linear Attention: $O(n^2)$ 복잡도를 $O(n)$ 으로 줄이는 기법.
- Flash Attention: 메모리 효율적인 Attention 계산.

실무 구현 예시

```
# Transformer Block 구현
import torch.nn as nn

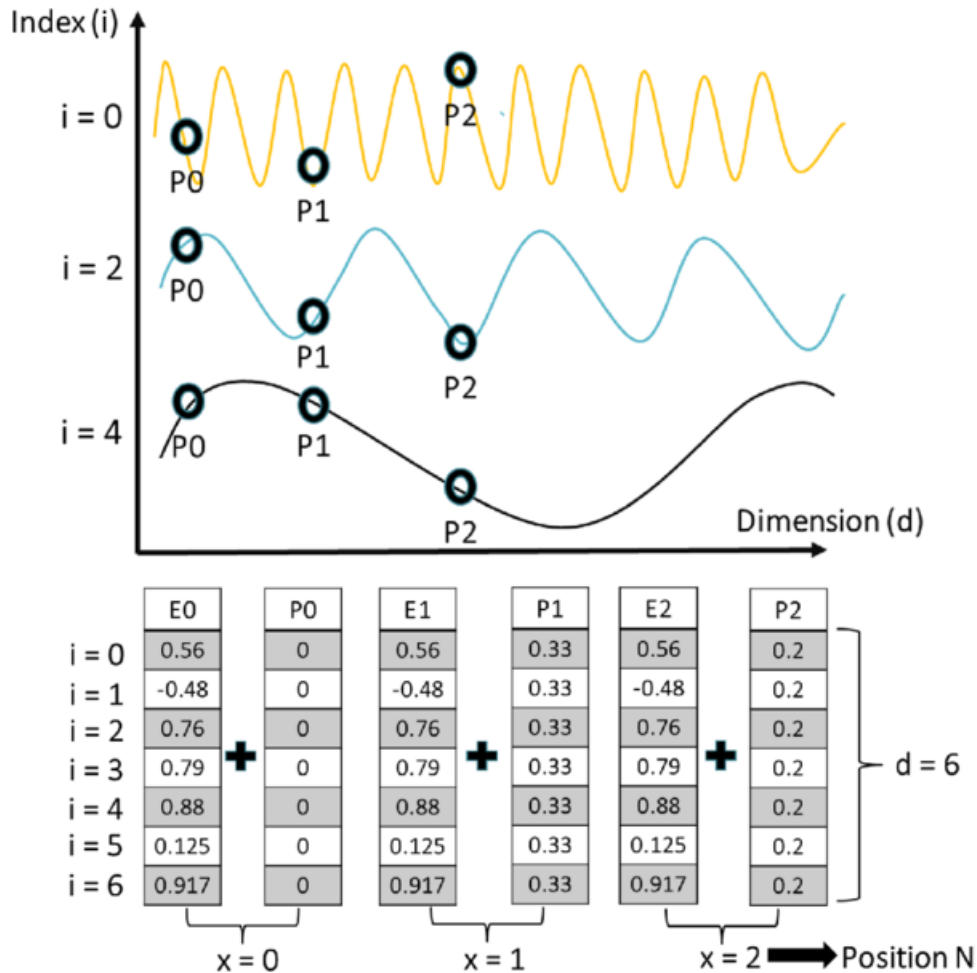
class TransformerBlock(nn.Module):
    def __init__(self, d_model, n_heads, d_ff, dropout=0.1):
        super().__init__()
        self.attention = MultiHeadAttention(d_model, n_heads)
        self.feed_forward = nn.Sequential(
            nn.Linear(d_model, d_ff),
            nn.ReLU(),
            nn.Linear(d_ff, d_model)
        )
        self.norm1 = nn.LayerNorm(d_model)
        self.norm2 = nn.LayerNorm(d_model)
        self.dropout = nn.Dropout(dropout)

    def forward(self, x, mask=None):
        # Self-Attention with residual connection
        attn_output = self.attention(x, x, x, mask)
        x = self.norm1(x + self.dropout(attn_output))

        # Feed-Forward with residual connection
        ff_output = self.feed_forward(x)
        x = self.norm2(x + self.dropout(ff_output))
        return x
```

1.4.4 위치 인코딩과 순서 정보

Transformer는 본질적으로 순서에 무관한(permutation-invariant) 모델이므로, 시퀀스 내 단어들의 순서 정보를 명시적으로 주입해야 합니다. 이를 위해 위치 인코딩(Positional Encoding)이 사용됩니다.



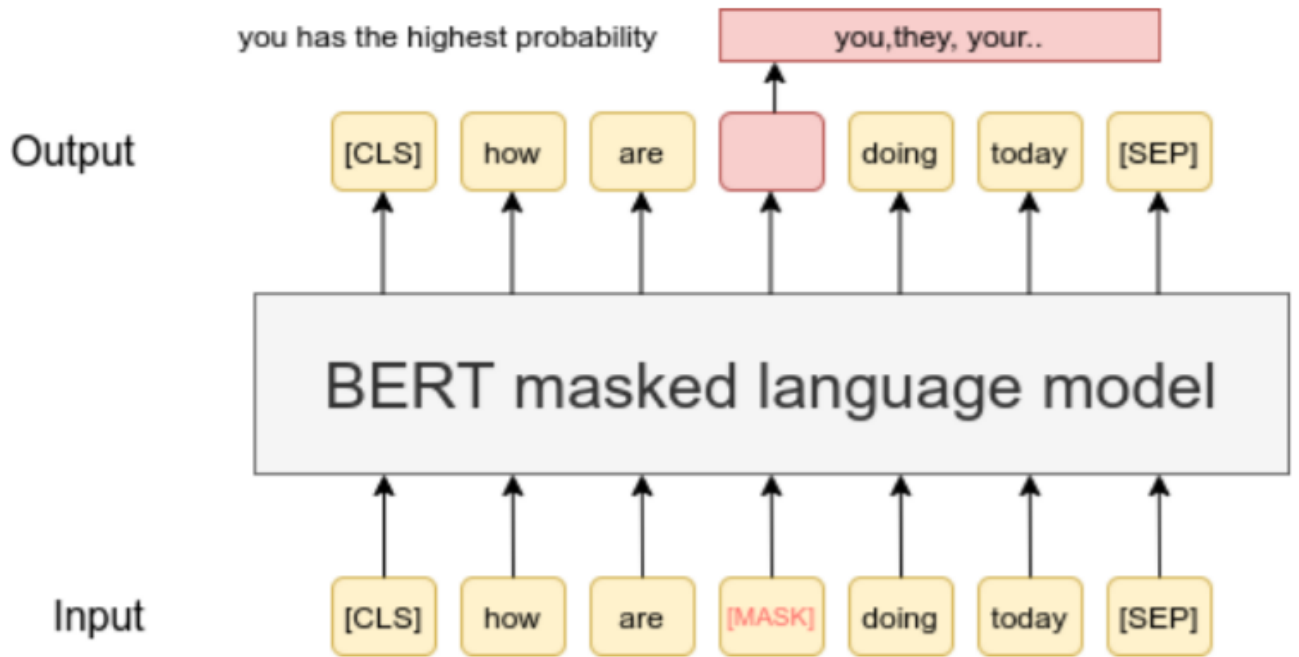
```
PE(pos, 2i) = sin(pos / 10000^(2i/d_model))
PE(pos, 2i+1) = cos(pos / 10000^(2i/d_model))
```

- 삼각함수 사용의 이유: 상대적 위치 관계 표현 가능, 임의 길이 시퀀스 처리 가능, 선형 변환을 통한 상대 위치 계산.
- 위치 인코딩은 단어 임베딩에 더해져 모델이 단어의 의미와 위치 정보를 동시에 학습하도록 돕습니다.

1.5.1 Transformer 아키텍처 패밀리 : BERT

Encoder-Only: BERT (Bidirectional Encoder Representations from Transformers)

Masked Language Model (MLM)을 통해 문장의 일부를 가리고, 문장 전체의 양방향 문맥을 모두 활용하여 가려진 단어를 예측하도록 학습합니다. 문장 전체의 의미를 깊이 이해하는 데 강점이 있어 분류, 개체명 인식 등 자연어 이해(NLU) 태스크에 주로 사용됩니다.



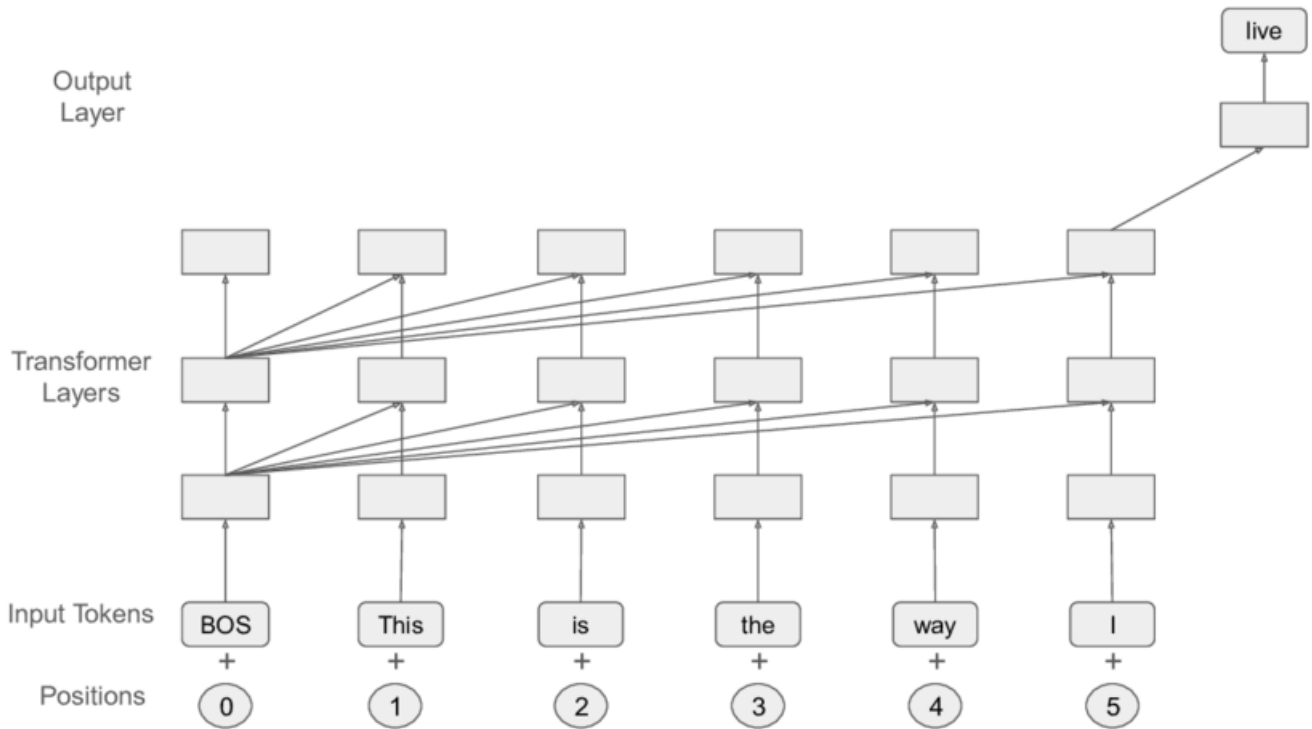
[CLS] how are [MASK] doing today. [SEP]

- Next Sentence Prediction (NSP): 두 문장이 이어지는지 예측하여 문장 간 관계 학습.
- 구조적 혁신: [CLS][CLS] 토큰 (문장 전체 표현), [SEP][SEP] 토큰 (문장 경계), Segment Embedding (문장 구분).

1.5.2 Transformer 아키텍처 패밀리 : GPT

Decoder-Only: GPT (Generative Pre-trained Transformer)

왼쪽에서 오른쪽으로, 이전 단어들을 기반으로 다음 단어를 예측하는 자기회귀적(Autoregressive) 방식으로 학습합니다. 이 구조는 자연스러운 자연어 생성(NLG)에 매우 강력하며, 별도의 파인튜닝 없이 프롬프트의 예시만으로 새로운 태스크를 수행하는 In-Context Learning이라는 놀라운 능력을 보여줍니다.



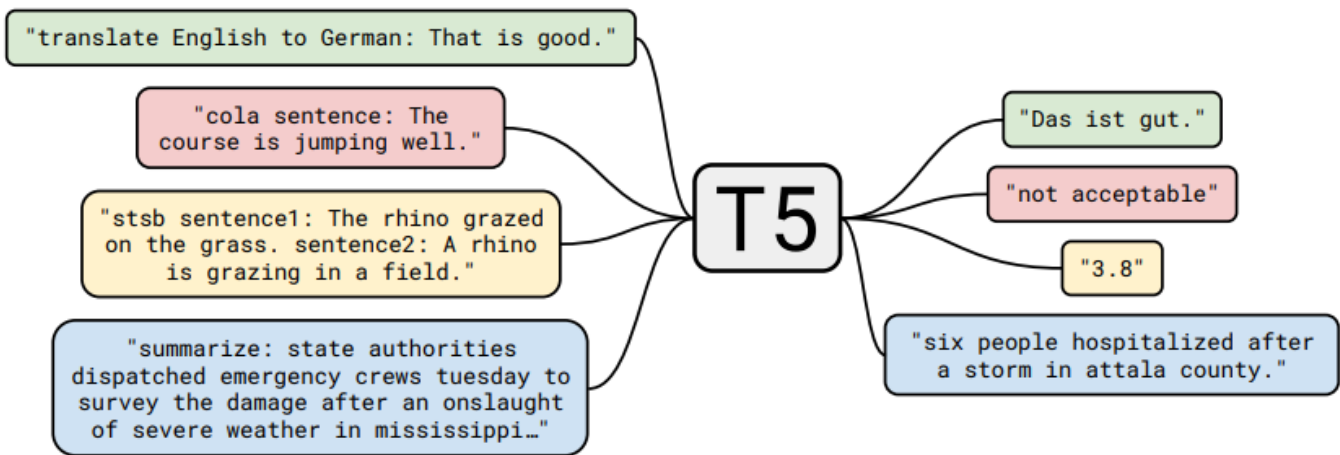
$$P(x_1, x_2, \dots, x_n) = \prod_{i=1}^n P(x_i | x_1, \dots, x_{i-1})$$

- Causal Attention: 미래 토큰 정보를 차단하여 단방향 정보 흐름 유지.
- 스케일링 법칙: 모델 크기가 커질수록 성능이 비례하여 증가하는 경향.
- Emergent Abilities: 특정 규모에서 나타나는 새로운 능력 (e.g., 추론, 코딩).

1.5.3 Transformer 아키텍처 패밀리 : T5

Encoder-Decoder: T5 (Text-to-Text Transfer Transformer)

모든 NLP 문제를 "텍스트-to-텍스트" 형식으로 통일한 프레임워크입니다. 입력 텍스트를 Encoder로 이해하고, 출력 텍스트를 Decoder로 생성하는 구조로 번역, 요약 등 입력과 출력이 모두 필요한 Seq2Seq 태스크에 효과적입니다.



- 통합 프레임워크의 장점: 일관된 사전 훈련 및 파인튜닝 절차, 태스크 간 지식 전이 효과.
- 구조적 특징: Relative Position Embedding, Pre-Layer Norm.

- 사전 훈련 목표: Span Corruption (연속된 토큰 스패น 마스크).

1.5.4 모델 아키텍처 선택 전략

문제 유형별 최적 아키텍처 선택

- 이미지: 분류(ResNet, EfficientNet), 객체 탐지(YOLO, R-CNN), 세그멘테이션(U-Net), 생성(GAN, Diffusion)
- 자연어: 이해/분류(BERT, RoBERTa), 생성(GPT 계열), 번역/요약(T5, BART)
- 시계열: 단기 예측(LSTM, GRU), 장기 예측(Transformer 기반 모델)

하이브리드 아키텍처 설계

- 멀티모달 학습: 비전-언어 (CLIP, DALL-E), 음성-텍스트 (Wav2Vec + BERT), 센서-비전 (CNN + LSTM 결합).
- 앙상블 전략: 모델 다양성, 배깅, 부스팅.

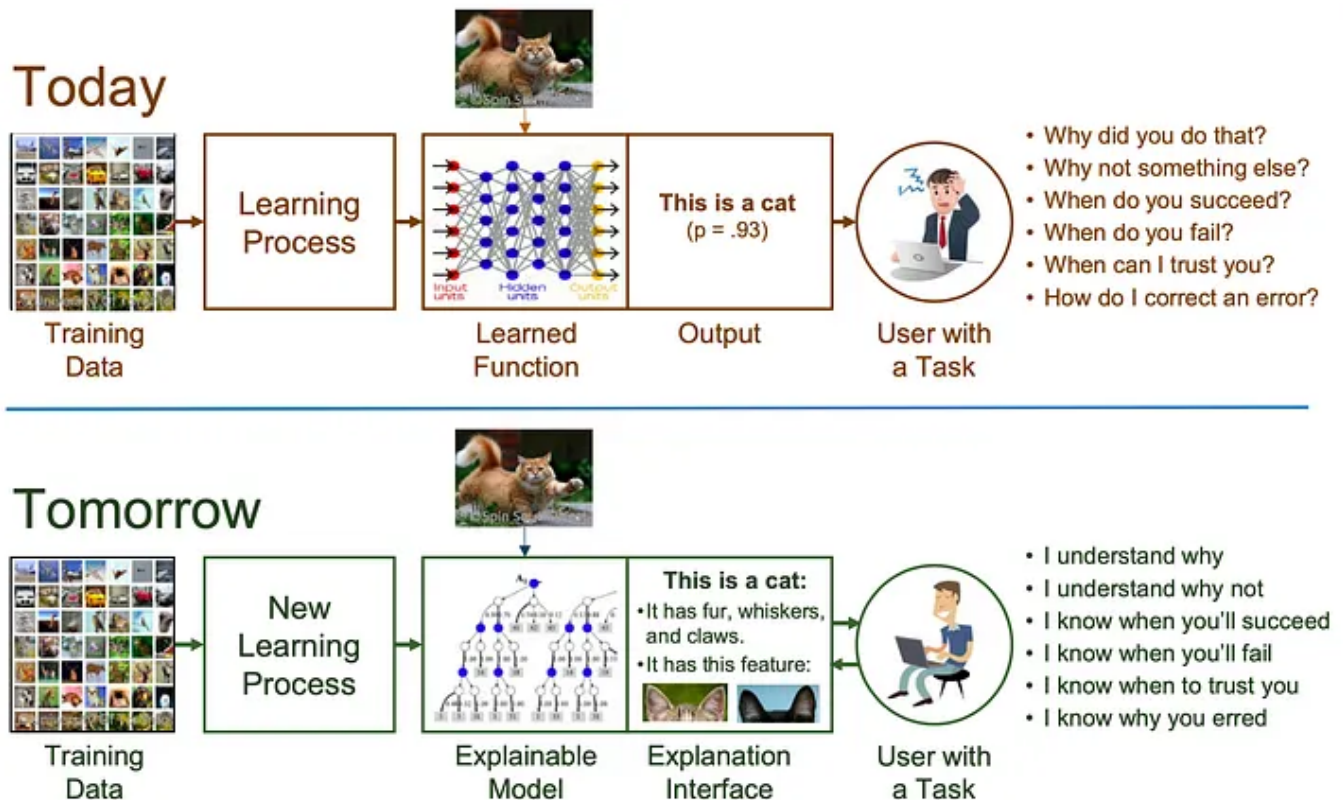
Part 2: Explainable AI (XAI)

XAI의 필요성과 사회적 요구

AI의 성능이 높아지고 사회적 영향력이 커지면서, 모델의 결정을 설명하고 신뢰할 수 있어야 한다는 요구가 커지고 있습니다. AI의 "블랙박스" 문제는 기술적 문제를 넘어 사회적, 윤리적 문제로 확장됩니다.

2.1.1 AI의 블랙박스 문제와 해석 가능성

현대 AI 모델의 성능 향상은 모델의 복잡성 증가와 함께 이루어졌습니다. 이는 성능과 해석 가능성 사이의 근본적인 긴장 관계를 만들어냅니다.



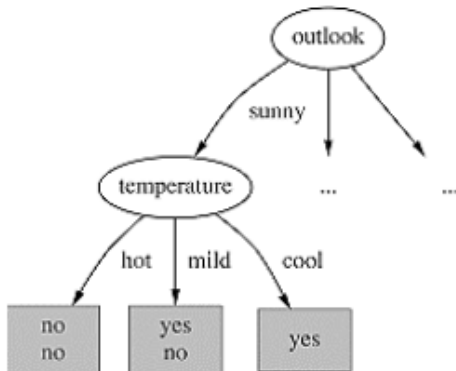
해석 가능성의 다층적 의미

- 투명성(Transparency): 모델의 내부 작동 원리 이해.
- 설명 가능성(Explainability): 특정 예측의 근거 제시.
- 신뢰성(Trustworthiness): 모델 예측에 대한 신뢰 구축.
- 공정성(Fairness): 편향과 차별 없는 의사결정.

XAI의 기술적 분류

• Intrinsic interpretability

- E.g.> Decision tree model



• Post hoc interpretability

- E.g.> Black box model



- Intrinsic vs Post-hoc: 내재적 해석 가능성(선형 모델) vs 사후 설명(복잡한 모델에 별도 설명 모듈).
- Model-specific vs Model-agnostic: 특정 모델에만 적용 vs 모든 모델에 적용 가능.
- Local vs Global: 개별 예측 설명 vs 전체 모델 동작 설명.

실무 구현 예시

```
# XAI 프레임워크 선택 가이드
def select_xai_method(model_type, data_type, explanation_scope):
    if model_type == "linear":
        return "intrinsic" # 모델 자체가 해석 가능
    elif data_type == "image":
        return "grad_cam" # 시각적 설명
    elif explanation_scope == "local":
        return "lime" # 지역적 설명
    else:
        return "shap" # 전역적 설명
```

2.1.2 XAI의 사회적 요구

- 신뢰성 및 디버깅: 개발자가 모델이 왜 특정 예측을 했는지 이해하고, 오류를 찾아 개선하기 위해 필요합니다.
- 공정성 및 윤리: 모델이 특정 집단에 대해 편향된 결정을 내리는지 확인하고 바로잡기 위해 필수적입니다.
- 규제 준수: EU의 GDPR에서 보장하는 "설명권"과 같이, 금융, 의료 등 고위험 분야에서 알고리즘 결정에 대한 설명은 법적 의무가 되고 있습니다. (e.g., 미국 알고리즘 책임법, 의료/금융 AI 규제)

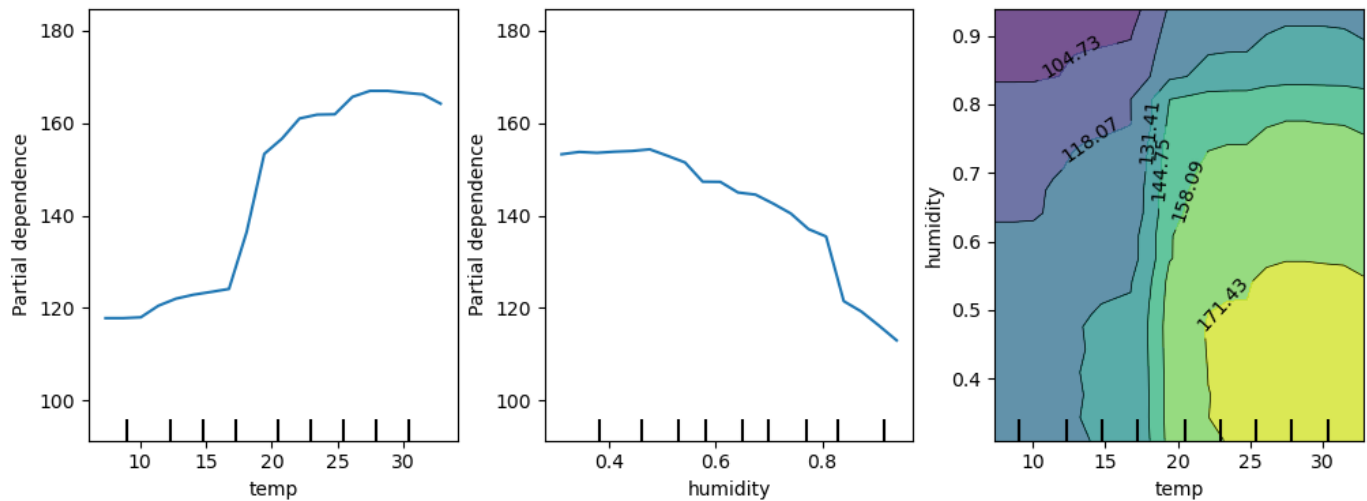
- 인간-AI 협업: 의사, 판사 등 전문직 사용자가 AI의 제안을 신뢰하고 최종 결정을 내리기 위해 그 근거를 이해해야 합니다.

2.2.1 XAI 기법 심층 분석: PDP

Partial Dependence Plot (PDP): 전역적 특성 이해

편의존성(Partial Dependence)은 관심있는 특정 입력변수를 제외한 다른 입력변수들의 값은 고정시킨 상태(상수 취급)에서 관심있는 입력변수의 값을 변화시키며(변수 취급) 예측값을 계산한 후, 그 값들의 평균을 낸 것으로 하나의 특성이 모델 예측에 미치는 평균적인(전역적) 영향을 시각화합니다. 전체적인 경향성을 파악하는 데 유용하지만, 특성 간 상호작용을 고려하지 못하고 개별 데이터의 이질적인 효과를 숨기는 한계가 있습니다.

1-way vs 2-way of numerical PDP using gradient boosting



$$PD_j(x_j) = E[f(x_j, X_C)] = (1/n) \sum f(x_j, x_C^{(i)})$$

- 해석: 특성 x_j 의 변화에 따른 모델 예측값의 평균적인 변화 추이.
- 한계: 특성 간 복잡한 상호작용을 반영하지 못함.

2.2.2 XAI 기법 심층 분석: ICE

Individual Conditional Expectation (ICE): 개별적 특성 이해

PDP의 한계를 보완하여, 각 개별 데이터 포인트가 특정 특성 변화에 어떻게 반응하는지를 하나하나의 곡선으로 보여줍니다. PDP가 이 곡선들의 평균이라면, ICE는 모든 곡선을 보여줌으로써 PDP에서는 보이지 않던 이질적인 상호작용 효과나 숨겨진 서브그룹을 발견할 수 있게 해줍니다.

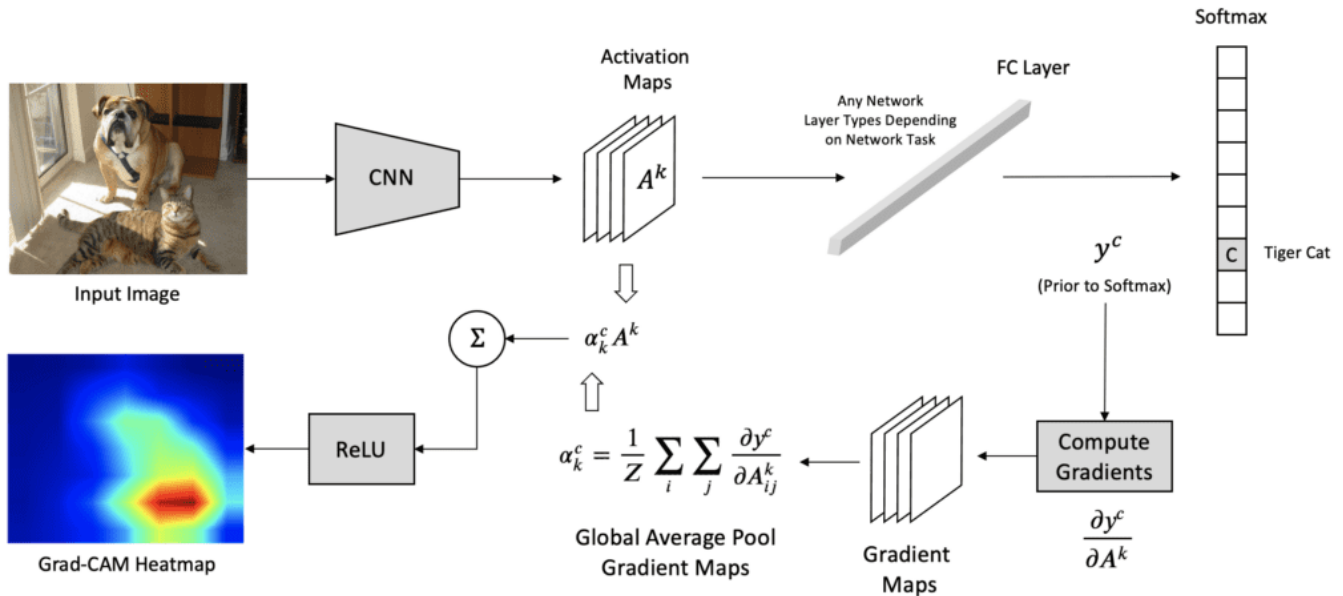
$$ICE_i(x_j) = f(x_j, x_C^{(i)})$$

- 중앙집중 ICE (Centered ICE): 기준점에서의 값을 빼서 상대적 변화를 강조하여 개별 곡선 비교 용이.
- 통찰: 데이터의 이질성 발견, 특성 간 상호작용 효과 파악, 유사한 반응 패턴을 가진 서브그룹 식별.

2.2.3 XAI 기법 심층 분석: Grad-CAM

Grad-CAM: 시각적 주의 메커니즘

GAP(Global Average Pooling) 레이어를 사용해 특정한 CAM(Class Activation Map)을 추출하는 CAM 방식을 개선한 방식으로 CNN 모델이 이미지 분류 결정을 내릴 때, 이미지의 어떤 영역에 크게 기여했는지(주목했는지)를 그래디언트(Gradient) 정보를 활용하여 히트맵으로 시각화합니다. 마지막 컨볼루션 레이어의 특성 맵에 대한 클래스 점수의 그래디언트를 계산하여 각 특성 맵의 중요도를 구하고, 이를 가중합하여 히트맵을 생성합니다. 모델이 올바른 근거로 정답을 맞혔는지, 혹은 잘못된 부분을 보고 있는지 확인할 수 있습니다.



$$\alpha_k^c = \frac{1}{Z} \sum_i \sum_j \left(\frac{\partial y^c}{\partial A_{ij}^k} \right)$$

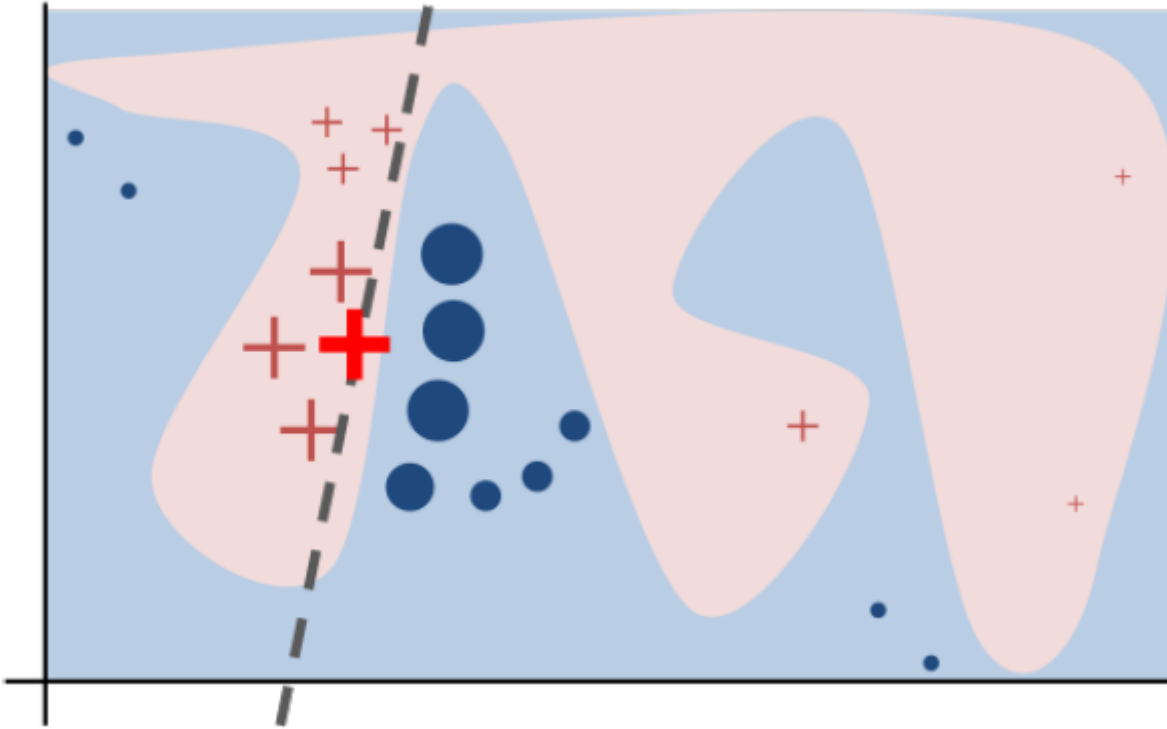
$$L_{\text{Grad-CAM}}^c = \text{ReLU} \left(\sum_k \alpha_k^c A^k \right)$$

- 해석: 히트맵의 밝은 영역이 모델 예측에 중요한 영향을 미친 부분.
- 한계: 해상도 제한, 클래스 구분성 강조.

2.2.4 XAI 기법 심층 분석: LIME & SHAP

LIME (Local Interpretable Model-agnostic Explanations)

모델에 구애받지 않는(Model-agnostic) XAI 기법의 대표주자입니다. 아무리 복잡한 블랙박스 모델이라도, 설명하고자 하는 특정 예측 지점의 주변(Local)에서는 단순한 선형 모델(Interpretable Model)로 근사할 수 있다는 아이디어에 기반합니다. 원본 데이터 주변에 교란된 샘플들을 생성하고, 이 샘플들에 대한 블랙박스 모델의 예측값을 이용해 가중 선형 회귀를 학습시켜 설명을 생성합니다. 이미지, 텍스트, 정형 데이터 등 다양한 데이터에 적용 가능합니다.



$$\xi(x) = \operatorname{argmin}_{\{g \in G\}} L(f, g, \pi_x) + \Omega(g)$$

- 작동 과정: 교란된 샘플 생성 → 모델 예측 획득 → 가중치 계산 → 선형 모델 학습.
- 장점: 범용성, 유연성, 직관성.
- 한계: 지역성 가정, 샘플링 품질, 안정성 문제.

SHAP (SHapley Additive exPlanations)

SHAP은 협력게임이론(Coalitional Game Theory) 중에서 각 게임 참여자(player)에게 수익의 배당(payout)을 공정하게 배분하는 알고리즘인 Shapley Value를 기반으로 한 XAI 기법으로, 각 특성이 예측에 기여하는 정도를 공정하게 분배합니다.

$$\phi_i = \sum_{S \subseteq N \setminus \{i\}} \frac{|S|!(|N|-|S|-1)!}{|N|!} [f(S \cup \{i\}) - f(S)]$$

- Shapley 값의 장점: 공정성, 일관성, 효율성, 가산성을 만족하는 유일한 해석 방법.
- 계산 복잡도: $O(2^n)$ 으로 특성 수에 따라 지수적 증가. TreeSHAP, KernelSHAP 등 근사 방법 사용.

실무 구현 예시

```
# LIME 구현
from lime import lime_tabular

explainer = lime_tabular.LimeTabularExplainer(
    X_train, feature_names=feature_names, class_names=class_names
)
explanation = explainer.explain_instance(X_test[0], model.predict_proba)
```

```
# SHAP 구현
import shap

explainer = shap.TreeExplainer(model)
shap_values = explainer.shap_values(X_test)
```

2.2.5 XAI 기법 선택 가이드

데이터 타입별 설명 방법

- 이미지 데이터: 전역 해석 (특성맵 시각화), 지역 해석 (Grad-CAM, LIME, SHAP).
- 텍스트 데이터: 어텐션 시각화, 토큰 중요도 (LIME, SHAP for text).
- 표 형태 데이터: 특성 중요도 (SHAP), 부분 의존성 (PDP, ICE).

이해관계자별 설명 전략

- 개발자에게는: 디버깅을 위한 상세한 기술적 설명 (e.g., 그래디언트 분석, SHAP 값, 활성화 시각화)
- 도메인 전문가에게는: 전문 지식과 비교할 수 있는 인과관계 중심의 설명 (e.g., PDP, LIME, 특성 중요도)
- 최종 사용자에게는: 직관적인 시각화와 간결한 자연어 설명, 그리고 대안 제시.

2.2.6 성능과 해석성의 균형

성능과 해석 가능성은 더 이상 트레이드오프 관계가 아니라, 신뢰할 수 있고 책임감 있는 AI 시스템을 구축하기 위한 상호 보완적인 필수 요소입니다.

해석 가능 모델 vs 고성능 모델

- 해석성 우선 상황: 고위험 의사결정 (의료, 금융, 법률), 규제 요구사항, 신뢰 구축.
- 성능 우선 상황: 추천 시스템, 이미지/음성 인식, 게임 AI.

하이브리드 접근법

- 포스트 혹 설명 (Post-hoc Explanation): 고성능 모델 + 별도 설명 모듈. 성능 저하 없이 설명성 확보.
- 해석 가능성 제약 학습: 정규화를 통한 해석성 유도, 어텐션 패턴 정규화.