

# AI 시대 필수 용어집

## Essential AI Glossary

한글 용어 101개 | 영문 용어 299개

기초 개념부터 최신 트렌드까지 완벽 정리

2026년 1월

# 한글 용어집 (가나다순)

총 101개 용어 수록

용어	정의
<b>가중치 (Weight)</b>	신경망에서 입력 신호의 중요도를 결정하는 학습 가능한 매개변수. 역전파를 통해 최적화됨
<b>강화학습 (Reinforcement Learning)</b>	에이전트가 환경과 상호작용하며 보상을 최대화하는 행동을 학습하는 ML 패러다임
<b>경사 하강법 (Gradient Descent)</b>	손실 함수를 최소화하기 위해 매개변수를 반복적으로 업데이트하는 최적화 알고리즘
<b>과적합 (Overfitting)</b>	모델이 훈련 데이터에 과도하게 맞춰져 새로운 데이터에 대한 일반화 성능이 저하되는 현상
<b>교차 검증 (Cross-Validation)</b>	데이터를 여러 부분으로 나눠 모델 성능을 평가하는 기법. K-Fold가 대표적
<b>그래디언트 소실/폭발 (Vanishing/Exploding Gradient)</b>	심층 신경망에서 역전파 시 그래디언트가 0에 수렴하거나 무한대로 발산하는 문제
<b>기계 학습 (Machine Learning)</b>	명시적 프로그래밍 없이 데이터로부터 패턴을 학습하여 예측/분류를 수행하는 AI 분야
<b>기울기 클리핑 (Gradient Clipping)</b>	그래디언트 폭발 방지를 위해 그래디언트 값을 특정 임계값으로 제한하는 기법
<b>나이브 베이즈 (Naive Bayes)</b>	베이즈 정리에 기반한 확률적 분류기. 특성 간 조건부 독립 가정
<b>노이즈 제거 오토인코더 (Denoising Autoencoder)</b>	입력에 노이즈를 추가한 후 원본을 복원하도록 학습하는 오토인코더 변형
<b>다중 레이블 분류 (Multi-label Classification)</b>	하나의 입력이 여러 클래스에 동시에 속할 수 있는 분류 문제
<b>다층 퍼셉트론 (MLP, Multi-Layer Perceptron)</b>	입력층, 은닉층, 출력층으로 구성된 순방향 신경망의 기본 구조
<b>단어 임베딩 (Word Embedding)</b>	단어를 밀집 벡터로 표현하는 기법. Word2Vec, GloVe, FastText 등
<b>드롭아웃 (Dropout)</b>	학습 시 무작위로 뉴런을 비활성화하여 과적합을 방지하는 정규화 기법
<b>딥러닝 (Deep Learning)</b>	다층 신경망을 사용한 기계학습. 비정형 데이터(이미지, 텍스트, 음성)에 강점
<b>랜덤 포레스트 (Random Forest)</b>	다수의 결정 트리를 앙상블하여 예측하는 알고리즘. 배깅 기법 활용
<b>레이블 (Label)</b>	지도학습에서 입력 데이터에 대한 정답 또는 목표값
<b>로짓 (Logit)</b>	확률의 로그 오즈. 분류 모델의 소프트맥스 이전 출력값
<b>리커런트 신경망 (RNN)</b>	순차적 데이터 처리를 위한 신경망. 이전 상태를 현재 계산에 활용
<b>마스크 언어 모델 (MLM)</b>	BERT의 사전학습 방식. 입력 토큰 일부를 마스킹하고 예측하도록 학습

<b>메타 학습 (Meta-Learning)</b>	학습하는 방법을 학습하는 접근법. Few-shot learning의 기반 기술
<b>멀티모달 (Multimodal)</b>	텍스트, 이미지, 오디오 등 여러 형태의 데이터를 동시에 처리하는 AI 시스템
<b>멀티헤드 어텐션 (Multi-Head Attention)</b>	Transformer의 핵심 구성요소. 여러 어텐션을 병렬로 수행하여 다양한 표현 학습
<b>모델 경량화 (Model Compression)</b>	프루닝, 양자화, 지식 증류 등을 통해 모델 크기/연산량을 줄이는 기법
<b>모멘텀 (Momentum)</b>	이전 그라디언트의 방향을 현재 업데이트에 반영하는 최적화 기법
<b>미세조정 (Fine-tuning)</b>	사전학습된 모델을 특정 태스크에 맞게 추가 학습하는 전이학습 방식
<b>반지도 학습 (Semi-supervised Learning)</b>	소량의 레이블 데이터와 대량의 비레이블 데이터를 함께 사용하는 학습 방식
<b>배치 정규화 (Batch Normalization)</b>	미니배치 단위로 활성화값을 정규화하여 학습을 안정화하는 기법
<b>백본 (Backbone)</b>	특징 추출을 담당하는 신경망의 기본 구조. ResNet, ViT 등
<b>벡터 데이터베이스 (Vector Database)</b>	고차원 임베딩 벡터의 유사도 검색에 최적화된 데이터베이스. RAG 시스템의 핵심
<b>변분 오토인코더 (VAE)</b>	잠재 공간을 확률 분포로 모델링하는 생성 모델. 데이터 생성 및 표현 학습에 활용
<b>분류 (Classification)</b>	입력 데이터를 사전 정의된 범주 중 하나로 할당하는 지도학습 태스크
<b>분산 학습 (Distributed Training)</b>	여러 GPU/노드에서 병렬로 모델을 학습하는 방식. 데이터/모델 병렬화
<b>비용 함수 (Cost Function)</b>	모델 예측과 실제값 간의 차이를 측정하는 함수. 손실 함수와 동의어
<b>비지도 학습 (Unsupervised Learning)</b>	레이블 없이 데이터의 내재적 구조나 패턴을 발견하는 학습 방식
<b>사전학습 (Pre-training)</b>	대규모 데이터로 일반적인 표현을 학습하는 초기 학습 단계. 전이학습의 기반
<b>서포트 벡터 머신 (SVM)</b>	최대 마진 초평면을 찾는 분류 알고리즘. 커널 트릭으로 비선형 분류 가능
<b>생성적 적대 신경망 (GAN)</b>	생성자와 판별자가 경쟁하며 학습하는 생성 모델. 고품질 이미지 생성에 활용
<b>셀프 어텐션 (Self-Attention)</b>	시퀀스 내 각 위치가 다른 모든 위치와의 관계를 계산하는 메커니즘
<b>소프트맥스 (Softmax)</b>	로짓을 확률 분포로 변환하는 활성화 함수. 다중 클래스 분류의 출력층에 사용
<b>손실 함수 (Loss Function)</b>	모델의 예측 오차를 정량화하는 함수. MSE, Cross-Entropy 등
<b>순환 신경망 (Recurrent Neural Network)</b>	시계열/순차 데이터 처리를 위한 신경망 구조. LSTM, GRU가 대표적
<b>스킵 연결 (Skip Connection)</b>	입력을 출력에 직접 더하는 연결. ResNet에서 도입, 깊은 네트워크 학습 가능하게 함
<b>슬라이딩 윈도우 (Sliding Window)</b>	고정 크기 윈도우를 이동하며 데이터를 처리하는 방식. CNN, 시계열 분석에 활용
<b>시그모이드 (Sigmoid)</b>	출력을 0~1 범위로 압축하는 활성화 함수. 이진 분류, 게이트 메커니즘에 사용
<b>시퀀스-투-시퀀스 (Seq2Seq)</b>	입력 시퀀스를 출력 시퀀스로 변환하는 모델 구조. 번역, 요약에 활용
<b>신경망 (Neural Network)</b>	뉴런을 모방한 연결 구조로 정보를 처리하는 계산 모델

심층 신경망 (DNN)	은닉층이 2개 이상인 신경망. 딥러닝의 기본 구조
아키텍처 검색 (Neural Architecture Search)	최적의 신경망 구조를 자동으로 탐색하는 AutoML 기법
앙상블 (Ensemble)	여러 모델의 예측을 결합하여 성능을 향상시키는 기법. 배깅, 부스팅, 스태킹
어텐션 메커니즘 (Attention Mechanism)	입력의 중요한 부분에 가중치를 부여하여 집중하는 메커니즘
에이전트 (Agent)	환경을 인식하고 목표 달성을 위해 자율적으로 행동하는 AI 시스템
에포크 (Epoch)	전체 훈련 데이터셋을 한 번 완전히 학습하는 단위
역전파 (Backpropagation)	출력 오차를 역방향으로 전파하여 가중치를 업데이트하는 학습 알고리즘
연합 학습 (Federated Learning)	중앙 서버에 데이터를 전송하지 않고 분산된 장치에서 모델을 학습하는 방식
오토인코더 (Autoencoder)	입력을 압축(인코딩)했다가 복원(디코딩)하도록 학습하는 비지도 신경망
온톨로지 (Ontology)	특정 도메인의 개념과 관계를 정형화한 지식 표현 체계
옵티マイ저 (Optimizer)	손실 함수를 최소화하도록 가중치를 업데이트하는 알고리즘. Adam, SGD 등
원-핫 인코딩 (One-Hot Encoding)	범주형 변수를 이진 벡터로 변환하는 인코딩 방식
위치 인코딩 (Positional Encoding)	Transformer에서 토큰의 위치 정보를 주입하는 기법. 사인/코사인 함수 사용
은닉층 (Hidden Layer)	입력층과 출력층 사이에 위치한 신경망 계층
이미지 분할 (Image Segmentation)	이미지의 각 픽셀을 클래스로 분류하는 컴퓨터 비전 태스크
인공 신경망 (Artificial Neural Network)	생물학적 뉴런의 연결 구조를 모방한 기계학습 모델
임베딩 (Embedding)	고차원 이산 데이터를 저차원 연속 벡터로 표현하는 기법
자기 지도 학습 (Self-supervised Learning)	레이블 없이 데이터 자체에서 학습 신호를 생성하는 방식. 대조 학습, 마스킹 등
자연어 처리 (NLP)	인간의 언어를 컴퓨터가 이해하고 생성할 수 있게 하는 AI 분야
자연어 이해 (NLU)	텍스트의 의미와 의도를 파악하는 NLP 하위 분야
자연어 생성 (NLG)	구조화된 데이터나 의미로부터 자연스러운 텍스트를 생성하는 기술
잔차 네트워크 (ResNet)	스킵 연결을 도입하여 매우 깊은 네트워크 학습을 가능하게 한 CNN 아키텍처
잠재 공간 (Latent Space)	데이터의 압축된 표현이 존재하는 저차원 공간. 생성 모델의 핵심 개념
적대적 공격 (Adversarial Attack)	입력에 미세한 변형을 가해 모델을 오작동시키는 공격 기법
전이 학습 (Transfer Learning)	한 태스크에서 학습한 지식을 다른 태스크에 적용하는 기법
정규화 (Regularization)	과적합 방지를 위해 모델 복잡도에 페널티를 부여하는 기법. L1, L2 정규화
정밀도/재현율 (Precision/Recall)	분류 모델의 성능 지표. 정밀도는 예측의 정확성, 재현율은 실제 양성 탐지율

<b>제로샷 학습 (Zero-shot Learning)</b>	학습 시 보지 못한 클래스에 대해 추론하는 능력
<b>지도 학습 (Supervised Learning)</b>	레이블이 있는 데이터로 입력-출력 관계를 학습하는 방식
<b>지식 그래프 (Knowledge Graph)</b>	엔티티와 관계를 그래프 구조로 표현한 지식 베이스
<b>지식 증류 (Knowledge Distillation)</b>	대형 모델(교사)의 지식을 소형 모델(학생)로 전달하는 모델 압축 기법
<b>차원 축소 (Dimensionality Reduction)</b>	고차원 데이터를 저차원으로 변환하는 기법. PCA, t-SNE, UMAP
<b>청킹 (Chunking)</b>	긴 문서를 작은 단위로 분할하는 RAG 전처리 기법
<b>초매개변수 (Hyperparameter)</b>	학습 전에 설정하는 매개변수. 학습률, 배치 크기, 레이어 수 등
<b>추론 (Inference)</b>	학습된 모델로 새로운 입력에 대해 예측을 수행하는 과정
<b>컨볼루션 신경망 (CNN)</b>	합성곱 연산을 사용하여 공간적 패턴을 학습하는 신경망. 이미지 처리에 강점
<b>컨텍스트 윈도우 (Context Window)</b>	LLM이 한 번에 처리할 수 있는 토큰의 최대 길이
<b>클러스터링 (Clustering)</b>	데이터를 유사성 기반으로 그룹화하는 비지도 학습 태스크. K-means, DBSCAN
<b>토큰 (Token)</b>	텍스트를 처리하는 기본 단위. 단어, 서브워드, 문자 등
<b>토크나이저 (Tokenizer)</b>	텍스트를 토큰으로 분할하는 도구. BPE, WordPiece, SentencePiece
<b>특성 공학 (Feature Engineering)</b>	원시 데이터에서 모델 성능을 향상시키는 특성을 추출/생성하는 과정
<b>파라미터 효율적 미세조정 (PEFT)</b>	전체 모델이 아닌 일부 파라미터만 학습하는 효율적인 미세조정 기법. LoRA, QLoRA
<b>퍼플렉시티 (Perplexity)</b>	언어 모델의 예측 불확실성을 측정하는 지표. 낮을수록 좋음
<b>편향 (Bias)</b>	뉴런의 활성화 임계값을 조절하는 학습 가능한 매개변수
<b>풀링 (Pooling)</b>	CNN에서 특징 맵의 크기를 줄이는 다운샘플링 연산. Max/Average Pooling
<b>프롬프트 엔지니어링 (Prompt Engineering)</b>	LLM에서 원하는 출력을 얻기 위해 입력 프롬프트를 설계하는 기술
<b>피처 (Feature)</b>	모델 입력으로 사용되는 데이터의 개별 속성이나 특성
<b>하이퍼파라미터 튜닝 (Hyperparameter Tuning)</b>	최적의 초매개변수 조합을 찾는 과정. Grid Search, Random Search, Bayesian Optimization
<b>학습률 (Learning Rate)</b>	가중치 업데이트의 크기를 결정하는 하이퍼파라미터
<b>할루시네이션 (Hallucination)</b>	LLM이 사실이 아닌 정보를 생성하는 현상
<b>합성곱 (Convolution)</b>	필터를 입력에 슬라이딩하며 특징을 추출하는 연산
<b>활성화 함수 (Activation Function)</b>	뉴런의 출력에 비선형성을 부여하는 함수. ReLU, Sigmoid, Tanh
<b>회귀 (Regression)</b>	연속적인 값을 예측하는 지도학습 태스크

**효율적 어텐션 (Efficient Attention)**

어텐션의  $O(n^2)$  복잡도를 줄이는 기법. Flash Attention, Linear Attention

# English Glossary (A-Z)

Total 299 terms included

Term	Definition
<b>A/B Testing</b>	Comparing two model versions on live traffic to determine which performs better
<b>Ablation Study</b>	Systematic removal of model components to understand their contribution to performance
<b>Activation Function</b>	Non-linear function applied to neuron output. Common types: ReLU, GELU, SiLU, Swish
<b>Activation Maximization</b>	Technique to visualize what patterns a neuron responds to by optimizing input
<b>Adam Optimizer</b>	Adaptive learning rate optimizer combining momentum and RMSprop. Default choice for deep learning
<b>Adversarial Training</b>	Training on adversarial examples to improve model robustness against attacks
<b>Agent</b>	AI system that perceives environment and takes autonomous actions to achieve goals
<b>Agentic AI</b>	AI systems capable of autonomous planning, tool use, and multi-step reasoning to complete tasks
<b>AI Alignment</b>	Research field ensuring AI systems behave according to human values and intentions
<b>AI Safety</b>	Research area focused on preventing harmful AI behaviors and ensuring beneficial outcomes
<b>Anchor Loss</b>	Contrastive learning loss that pulls similar samples together and pushes dissimilar apart
<b>Annotation</b>	Process of labeling data for supervised learning. Can be manual or semi-automated
<b>Anomaly Detection</b>	Identifying rare patterns or outliers that deviate from expected behavior
<b>Attention Mechanism</b>	Mechanism allowing models to focus on relevant parts of input when producing output
<b>Autoencoder (AE)</b>	Neural network that learns compressed representation by encoding then decoding input
<b>AutoML</b>	Automated machine learning - automating model selection, hyperparameter tuning, and architecture design
<b>Autoregressive Model</b>	Model that generates output sequentially, conditioning each token on previous ones. GPT-style
<b>Backpropagation</b>	Algorithm for computing gradients by propagating error backward through network

<b>Bag of Words (BoW)</b>	Text representation counting word occurrences, ignoring order and grammar
<b>Batch Normalization</b>	Normalizing layer inputs within mini-batch to stabilize and accelerate training
<b>Batch Size</b>	Number of training samples processed before updating model parameters
<b>Bayesian Neural Network</b>	Neural network with probability distributions over weights, enabling uncertainty estimation
<b>Beam Search</b>	Search algorithm maintaining top-k candidates at each step for sequence generation
<b>BERT (Bidirectional Encoder Representations)</b>	Bidirectional transformer pre-trained on masked language modeling. Foundation for NLU tasks
<b>Bias-Variance Tradeoff</b>	Balance between model simplicity (high bias) and complexity (high variance)
<b>Bidirectional LSTM</b>	LSTM processing sequences in both forward and backward directions
<b>Binary Cross-Entropy</b>	Loss function for binary classification measuring divergence from true labels
<b>BLEU Score</b>	Metric for evaluating machine translation quality based on n-gram precision
<b>Catastrophic Forgetting</b>	Phenomenon where learning new tasks causes model to forget previously learned tasks
<b>Causal Attention</b>	Attention variant where each position only attends to previous positions. Used in decoder models
<b>Chain-of-Thought (CoT)</b>	Prompting technique encouraging LLMs to show step-by-step reasoning
<b>Checkpoint</b>	Saved model state during training, enabling resumption or evaluation
<b>Class Imbalance</b>	Dataset condition where some classes have significantly more samples than others
<b>Classification</b>	Supervised learning task of assigning inputs to predefined categories
<b>CLIP (Contrastive Language-Image Pre-training)</b>	Model learning joint image-text representations through contrastive learning
<b>Clustering</b>	Unsupervised learning task of grouping similar data points together
<b>Codeformer</b>	Face restoration model using transformer architecture for blind face restoration
<b>Compute Optimal Scaling</b>	Research on optimal allocation of compute between model size and training data (Chinchilla)
<b>Concept Drift</b>	Change in statistical properties of target variable over time, degrading model performance
<b>Conditional Generation</b>	Generating output conditioned on specific input or constraint
<b>Confusion Matrix</b>	Table showing true/false positives/negatives for classification evaluation
<b>Constitutional AI (CAI)</b>	Training approach using principles to guide AI behavior, developed by Anthropic
<b>Contrastive Learning</b>	Self-supervised method learning representations by contrasting positive and negative pairs
<b>Convolutional Neural</b>	Neural network using convolution operations, excelling at spatial pattern

<b>Network (CNN)</b>	recognition
<b>Cross-Attention</b>	Attention between two different sequences. Used in encoder-decoder models
<b>Cross-Entropy Loss</b>	Loss function measuring divergence between predicted and true probability distributions
<b>CUDA</b>	NVIDIA's parallel computing platform enabling GPU acceleration for deep learning
<b>Curriculum Learning</b>	Training strategy presenting samples in order of increasing difficulty
<b>Data Augmentation</b>	Artificially expanding training data through transformations (rotation, cropping, noise)
<b>Data Leakage</b>	Information from outside training set inappropriately influencing model, causing overfitting
<b>Data Parallelism</b>	Distributed training strategy splitting data across multiple devices
<b>Decision Boundary</b>	Hypersurface separating different classes in feature space
<b>Decoder</b>	Network component generating output from encoded representation. Used in seq2seq models
<b>Deep Q-Network (DQN)</b>	Reinforcement learning algorithm combining Q-learning with deep neural networks
<b>DeepSpeed</b>	Microsoft's distributed training library with ZeRO optimization for large models
<b>Denoising Diffusion (DDPM)</b>	Generative model learning to reverse gradual noise addition process
<b>Dense Retrieval</b>	Information retrieval using learned dense embeddings rather than sparse keyword matching
<b>Depth-wise Convolution</b>	Efficient convolution applying separate filter per input channel. Used in MobileNet
<b>Differentiable Programming</b>	Programming paradigm where all operations are differentiable for gradient-based optimization
<b>Diffusion Model</b>	Generative model learning to denoise data through iterative refinement. Stable Diffusion, DALL-E
<b>Direct Preference Optimization (DPO)</b>	Simpler alternative to RLHF, directly optimizing policy from preference data
<b>Discriminator</b>	GAN component that distinguishes real from generated samples
<b>Distillation</b>	Transferring knowledge from large teacher model to smaller student model
<b>Distributed Training</b>	Training models across multiple GPUs or nodes for scalability
<b>Domain Adaptation</b>	Adapting model trained on source domain to perform well on different target domain
<b>Dropout</b>	Regularization randomly deactivating neurons during training to prevent overfitting
<b>Early Stopping</b>	Halting training when validation performance stops improving to prevent

	overfitting
<b>Edge AI</b>	Running AI models on edge devices (phones, IoT) rather than cloud servers
<b>Embedding</b>	Dense vector representation of discrete data (words, entities) in continuous space
<b>Emergent Abilities</b>	Capabilities that appear in large models but are absent in smaller ones
<b>Encoder</b>	Network component that compresses input into latent representation
<b>Encoder-Decoder Architecture</b>	Model with encoder processing input and decoder generating output. T5, BART
<b>Ensemble Learning</b>	Combining multiple models to improve prediction accuracy and robustness
<b>Entity Recognition (NER)</b>	NLP task identifying and classifying named entities in text
<b>Epoch</b>	One complete pass through the entire training dataset
<b>Explainability (XAI)</b>	Methods for understanding and interpreting model decisions. SHAP, LIME, attention visualization
<b>Exploration vs Exploitation</b>	Reinforcement learning dilemma between trying new actions and using known good ones
<b>F1 Score</b>	Harmonic mean of precision and recall, balancing both metrics
<b>Feature Extraction</b>	Deriving informative attributes from raw data for model input
<b>Feature Map</b>	Output of convolutional layer representing detected features in input
<b>Feature Store</b>	Centralized repository for storing and serving ML features in production
<b>Federated Learning</b>	Training on decentralized data without transferring it to central server
<b>Few-Shot Learning</b>	Learning to perform tasks from very few examples
<b>Fine-Tuning</b>	Adapting pre-trained model to specific task through additional training
<b>Flash Attention</b>	Memory-efficient attention algorithm reducing GPU memory usage significantly
<b>FLOPS (Floating Point Operations)</b>	Measure of computational cost, used to quantify model training requirements
<b>Foundation Model</b>	Large pre-trained model serving as base for many downstream applications
<b>Fully Connected Layer</b>	Neural network layer where each neuron connects to all neurons in previous layer
<b>Function Calling</b>	LLM capability to generate structured output for invoking external tools/APIs
<b>GAN (Generative Adversarial Network)</b>	Generative model with competing generator and discriminator networks
<b>Gated Recurrent Unit (GRU)</b>	RNN variant with reset and update gates, simpler than LSTM
<b>Gaussian Process</b>	Non-parametric model defining probability distribution over functions
<b>GELU (Gaussian Error Linear Unit)</b>	Activation function combining properties of dropout and ReLU. Used in transformers

<b>Generalization</b>	Model's ability to perform well on unseen data, not just training data
<b>Generative AI</b>	AI systems that create new content (text, images, audio, video, code)
<b>Generative Pre-trained Transformer (GPT)</b>	Autoregressive transformer architecture for text generation. OpenAI's GPT series
<b>Gibbs Sampling</b>	MCMC algorithm for sampling from joint probability distribution
<b>Gradient Accumulation</b>	Accumulating gradients over multiple batches before updating, simulating larger batch size
<b>Gradient Checkpoint</b>	Trading compute for memory by recomputing activations during backward pass
<b>Gradient Descent</b>	Optimization algorithm minimizing loss by iteratively adjusting parameters
<b>Graph Neural Network (GNN)</b>	Neural network operating on graph-structured data. Node classification, link prediction
<b>Greedy Decoding</b>	Text generation selecting highest probability token at each step
<b>Ground Truth</b>	Correct/actual labels in supervised learning dataset
<b>Group Normalization</b>	Normalization dividing channels into groups, alternative to batch norm for small batches
<b>Grounding</b>	Connecting model outputs to real-world knowledge or external data sources
<b>Guardrails</b>	Safety mechanisms limiting AI model behavior to acceptable boundaries
<b>Hallucination</b>	LLM generating false or fabricated information presented as fact
<b>Hidden State</b>	Internal representation passed between time steps in RNNs
<b>Hierarchical Clustering</b>	Clustering method building tree of nested clusters (dendrogram)
<b>High-Rank Adaptation (HiRA)</b>	PEFT method allowing higher-rank updates than standard LoRA
<b>Hugging Face</b>	Platform and library ecosystem for sharing and using transformer models
<b>Hyperparameter</b>	Model configuration set before training (learning rate, layers, batch size)
<b>Image Classification</b>	Computer vision task assigning images to predefined categories
<b>Image Segmentation</b>	Classifying each pixel in image to specific class. Semantic, instance, panoptic
<b>Imitation Learning</b>	Learning by observing and mimicking expert demonstrations
<b>In-Context Learning (ICL)</b>	LLM ability to learn tasks from examples provided in prompt without fine-tuning
<b>Inference</b>	Using trained model to make predictions on new data
<b>Information Retrieval (IR)</b>	Finding relevant documents from large collections based on queries
<b>Input Embedding</b>	Converting discrete inputs (tokens) to continuous vectors
<b>Instance Normalization</b>	Normalization within single sample, used in style transfer
<b>Instruction Tuning</b>	Fine-tuning models on instruction-following datasets for better alignment

<b>Interpretability</b>	Degree to which humans can understand model's decision process
<b>IoU (Intersection over Union)</b>	Metric measuring overlap between predicted and ground truth regions
<b>Jailbreaking</b>	Techniques to bypass AI model safety restrictions through adversarial prompts
<b>Knowledge Distillation</b>	Training smaller model to mimic larger model's outputs
<b>Knowledge Graph (KG)</b>	Graph structure representing entities and their relationships
<b>KV Cache</b>	Caching key-value pairs in transformer attention for efficient autoregressive generation
<b>Label Smoothing</b>	Regularization technique softening hard labels to prevent overconfidence
<b>LangChain</b>	Framework for building applications with LLMs, including chains and agents
<b>Language Model (LM)</b>	Model trained to predict probability of word sequences
<b>Large Language Model (LLM)</b>	Language model with billions of parameters trained on massive text corpora
<b>Latent Space</b>	Lower-dimensional space where data representations exist
<b>Layer Normalization</b>	Normalizing across features within single sample, standard in transformers
<b>Learning Rate Schedule</b>	Strategy for adjusting learning rate during training. Warmup, cosine decay
<b>Leaky ReLU</b>	ReLU variant allowing small gradient for negative inputs to prevent dead neurons
<b>Linear Probe</b>	Evaluation method training only linear classifier on frozen representations
<b>LLaMA</b>	Meta's open-weight large language model family. LLaMA 2, LLaMA 3
<b>LLMops</b>	Practices for deploying, monitoring, and maintaining LLM applications in production
<b>Logistic Regression</b>	Linear model for binary classification using sigmoid function
<b>Long Short-Term Memory (LSTM)</b>	RNN architecture with gates controlling information flow, handling long sequences
<b>LoRA (Low-Rank Adaptation)</b>	Efficient fine-tuning method adding low-rank matrices to frozen weights
<b>Loss Function</b>	Function measuring model's prediction error, optimized during training
<b>Loss Landscape</b>	Visualization of loss function across parameter space
<b>MAE (Mean Absolute Error)</b>	Average absolute difference between predictions and actual values
<b>Masked Language Modeling (MLM)</b>	Pre-training task predicting masked tokens in input sequence. BERT's approach
<b>Max Pooling</b>	Downsampling by selecting maximum value in each region
<b>Mean Pooling</b>	Aggregating embeddings by taking element-wise mean
<b>Meta-Learning</b>	Learning to learn - acquiring knowledge that transfers across tasks
<b>Mini-Batch</b>	Subset of training data used for single parameter update

<b>Mistral</b>	Open-weight LLM known for efficiency. Mistral 7B, Mixtral MoE
<b>Mixed Precision Training</b>	Using FP16/BF16 with FP32 selectively for faster training
<b>Mixture of Experts (MoE)</b>	Architecture with multiple expert networks, routing inputs to relevant experts
<b>MLflow</b>	Platform for ML lifecycle management including experiment tracking and deployment
<b>MLOps</b>	Practices combining ML and DevOps for reliable model deployment and monitoring
<b>Model Card</b>	Documentation describing model's capabilities, limitations, and intended use
<b>Model Parallelism</b>	Distributing model layers across multiple devices for large models
<b>Model Pruning</b>	Removing unnecessary weights or neurons to reduce model size
<b>Model Quantization</b>	Reducing numerical precision (FP32 → INT8) for efficient inference
<b>Momentum</b>	Optimization technique using moving average of past gradients for smoother updates
<b>Monte Carlo Dropout</b>	Using dropout at inference for uncertainty estimation
<b>Multi-Agent System</b>	System with multiple AI agents collaborating or competing on tasks
<b>Multi-Head Attention</b>	Parallel attention mechanisms learning different aspects of relationships
<b>Multimodal AI</b>	AI processing and integrating multiple data types (text, images, audio)
<b>Multi-Task Learning</b>	Training single model on multiple related tasks simultaneously
<b>N-gram</b>	Contiguous sequence of n items (words, characters) from text
<b>Named Entity Recognition (NER)</b>	NLP task identifying and classifying named entities (people, places, organizations)
<b>Natural Language Generation (NLG)</b>	Generating human-readable text from structured data or meaning representations
<b>Natural Language Processing (NLP)</b>	AI field enabling computers to understand and generate human language
<b>Natural Language Understanding (NLU)</b>	Understanding meaning and intent in human language
<b>Negative Sampling</b>	Training technique using sampled negative examples rather than all possible negatives
<b>Neural Architecture Search (NAS)</b>	Automated design of neural network architectures
<b>Neural Machine Translation (NMT)</b>	Translating text between languages using neural networks
<b>Neural Network</b>	Computational model inspired by biological neurons, organized in layers
<b>Next Token Prediction</b>	Core LLM training objective predicting next token given context
<b>Noise Contrastive Estimation</b>	Training by distinguishing real data from noise samples

<b>Normalization</b>	Scaling data to standard range or distribution for stable training
<b>Nucleus Sampling (Top-p)</b>	Text generation sampling from smallest set of tokens exceeding probability $p$
<b>Object Detection</b>	Computer vision task locating and classifying objects in images
<b>Objective Function</b>	Function model optimizes during training (minimize loss, maximize reward)
<b>Offline Reinforcement Learning</b>	RL learning from fixed dataset without environment interaction
<b>ONNX</b>	Open format for representing deep learning models, enabling framework interoperability
<b>Open-Source Model</b>	Model with publicly available weights and often code. LLaMA, Mistral
<b>Optimizer</b>	Algorithm updating model parameters to minimize loss. Adam, AdamW, SGD
<b>Out-of-Distribution (OOD)</b>	Data samples significantly different from training distribution
<b>Overfitting</b>	Model memorizing training data, performing poorly on new data
<b>Padding</b>	Adding values to sequences to achieve uniform length for batch processing
<b>Parameter-Efficient Fine-Tuning (PEFT)</b>	Fine-tuning updating only small subset of parameters. LoRA, QLoRA, Adapters
<b>Perplexity</b>	Measure of language model uncertainty; lower is better
<b>Pipeline Parallelism</b>	Distributing model layers across devices with micro-batch pipelining
<b>Policy Gradient</b>	RL method directly optimizing policy through gradient ascent on expected reward
<b>Pooling</b>	Downsampling operation reducing spatial dimensions while preserving important features
<b>Position Embedding</b>	Encoding token position in sequence for transformers lacking inherent order
<b>Post-Training Quantization (PTQ)</b>	Quantizing model after training without retraining
<b>Pre-Training</b>	Initial training phase on large general dataset before task-specific fine-tuning
<b>Precision</b>	Fraction of positive predictions that are actually positive ( $TP/(TP+FP)$ )
<b>Prediction Head</b>	Task-specific layer(s) on top of model backbone for specific outputs
<b>Principal Component Analysis (PCA)</b>	Dimensionality reduction finding orthogonal axes of maximum variance
<b>Prompt</b>	Input text provided to LLM to elicit desired response
<b>Prompt Engineering</b>	Designing effective prompts to achieve desired LLM outputs
<b>Prompt Injection</b>	Attack attempting to override model instructions through malicious input
<b>Proximal Policy Optimization (PPO)</b>	RL algorithm with clipped objective for stable policy updates
<b>QLoRA</b>	Quantized LoRA combining 4-bit quantization with low-rank adaptation
<b>Quantization</b>	Reducing numerical precision of model weights for efficiency

<b>Quantization-Aware Training (QAT)</b>	Training with simulated quantization for better quantized performance
<b>Query-Key-Value (QKV)</b>	Three projections in attention mechanism computing weighted representations
<b>RAG (Retrieval-Augmented Generation)</b>	Combining retrieval of relevant documents with LLM generation for grounded responses
<b>Random Forest</b>	Ensemble of decision trees using bagging for robust predictions
<b>Recall</b>	Fraction of actual positives correctly identified (TP/(TP+FN))
<b>Receptive Field</b>	Region of input influencing particular CNN neuron's output
<b>Recommendation System</b>	AI system suggesting items based on user preferences and behavior
<b>Recurrent Neural Network (RNN)</b>	Neural network with loops for processing sequential data
<b>Red Teaming</b>	Adversarial testing to identify vulnerabilities in AI systems
<b>Regularization</b>	Techniques preventing overfitting by constraining model complexity
<b>Reinforcement Learning (RL)</b>	Learning through environment interaction, maximizing cumulative reward
<b>Reinforcement Learning from Human Feedback (RLHF)</b>	Aligning models using human preference feedback. Key for ChatGPT, Claude
<b>ReLU (Rectified Linear Unit)</b>	Activation function outputting $\max(0, x)$ , standard in deep learning
<b>Representation Learning</b>	Learning useful features/representations from raw data automatically
<b>Residual Connection</b>	Skip connection adding input to output, enabling deep network training
<b>ResNet</b>	CNN architecture with residual connections enabling very deep networks
<b>Reward Hacking</b>	RL agent exploiting reward function in unintended ways
<b>Reward Model</b>	Model predicting human preferences for RLHF training
<b>RMSprop</b>	Optimizer adapting learning rates using moving average of squared gradients
<b>RNN (Recurrent Neural Network)</b>	Network with recurrent connections for sequential data processing
<b>ROC-AUC</b>	Area under ROC curve, measuring classifier performance across thresholds
<b>RoPE (Rotary Position Embedding)</b>	Efficient position encoding using rotation matrices. Used in LLaMA, Mistral
<b>Safety Filter</b>	Content moderation system filtering harmful AI outputs
<b>Sampling Temperature</b>	Parameter controlling randomness in text generation. Higher = more diverse
<b>Scaling Laws</b>	Predictable relationships between model size, data, compute, and performance
<b>Self-Attention</b>	Attention where query, key, value all come from same sequence
<b>Self-Supervised Learning</b>	Learning representations from unlabeled data using pretext tasks

<b>Semantic Search</b>	Search based on meaning rather than keyword matching, using embeddings
<b>Semantic Similarity</b>	Measuring meaning similarity between texts using vector representations
<b>Semi-Supervised Learning</b>	Learning from both labeled and unlabeled data
<b>Sentence Embedding</b>	Fixed-length vector representation of entire sentence
<b>Sentiment Analysis</b>	NLP task determining emotional tone or opinion in text
<b>Sequence-to-Sequence (Seq2Seq)</b>	Model architecture mapping input sequence to output sequence
<b>SGD (Stochastic Gradient Descent)</b>	Gradient descent using random sample subsets for updates
<b>SHAP (SHapley Additive exPlanations)</b>	Explainability method using game theory to attribute feature importance
<b>Sigmoid Function</b>	Activation squashing output to (0,1) range, used for probabilities
<b>Similarity Search</b>	Finding nearest neighbors in embedding space. Used in RAG, recommendations
<b>Skip Connection</b>	Direct path bypassing layers, enabling gradient flow in deep networks
<b>Softmax</b>	Function converting logits to probability distribution summing to 1
<b>Sparse Attention</b>	Attention mechanism attending to subset of positions for efficiency
<b>Speculative Decoding</b>	Acceleration technique using small model to draft tokens verified by large model
<b>Speech Recognition (ASR)</b>	Converting spoken audio to text. Whisper, DeepSpeech
<b>Speech Synthesis (TTS)</b>	Generating spoken audio from text. ElevenLabs, Bark
<b>Stable Diffusion</b>	Open-source latent diffusion model for image generation
<b>State-Space Model (SSM)</b>	Sequence model using state space formulation. Mamba architecture
<b>Stochastic Gradient Descent</b>	Gradient descent with random sampling for computational efficiency
<b>Stop Token</b>	Special token signaling end of generated sequence
<b>Stride</b>	Step size for convolution or pooling operations
<b>Structured Prediction</b>	Predicting complex outputs with interdependent components
<b>Style Transfer</b>	Applying artistic style of one image to content of another
<b>Subword Tokenization</b>	Tokenizing into subword units. BPE, WordPiece, SentencePiece
<b>Supervised Fine-Tuning (SFT)</b>	Fine-tuning on labeled instruction-response pairs before RLHF
<b>Supervised Learning</b>	Learning from labeled data with known input-output pairs
<b>Support Vector Machine (SVM)</b>	Classifier finding maximum margin hyperplane between classes
<b>Synthetic Data</b>	Artificially generated data for training when real data is scarce

<b>System Prompt</b>	Hidden instructions defining AI assistant's behavior and persona
<b>t-SNE</b>	Dimensionality reduction for visualizing high-dimensional data in 2D/3D
<b>Tanh</b>	Activation function outputting values in (-1, 1) range
<b>Temperature (Softmax)</b>	Parameter controlling output distribution sharpness in generation
<b>Tensor</b>	Multi-dimensional array, fundamental data structure in deep learning
<b>TensorFlow</b>	Google's open-source deep learning framework
<b>Test Set</b>	Held-out data for final model evaluation after training complete
<b>Text Classification</b>	Categorizing text into predefined classes. Sentiment, spam, topic
<b>Text Generation</b>	Producing coherent text given prompt or context
<b>Token</b>	Basic unit of text processing. Word, subword, or character
<b>Tokenizer</b>	Tool converting text to tokens and vice versa. BPE, SentencePiece
<b>Top-k Sampling</b>	Text generation sampling from k highest probability tokens
<b>Training Loop</b>	Iterative process of forward pass, loss computation, backpropagation, update
<b>Training Set</b>	Data used to train model parameters
<b>Transfer Learning</b>	Applying knowledge from one task/domain to improve performance on another
<b>Transformer</b>	Architecture using self-attention for parallel sequence processing. Foundation of modern LLMs
<b>Transformer Block</b>	Basic unit with multi-head attention and feed-forward network
<b>Tree of Thoughts (ToT)</b>	Prompting technique enabling exploration of multiple reasoning paths
<b>Triton</b>	Language and compiler for writing efficient GPU kernels. NVIDIA's inference server
<b>Truncation</b>	Cutting sequences exceeding maximum length for model processing
<b>UMAP</b>	Dimensionality reduction preserving global and local structure better than t-SNE
<b>Underfitting</b>	Model too simple to capture data patterns, performing poorly on training and test
<b>Unsupervised Learning</b>	Learning patterns from unlabeled data without ground truth
<b>Upsampling</b>	Increasing spatial resolution. Transposed convolution, interpolation
<b>Validation Set</b>	Data for hyperparameter tuning and model selection during training
<b>Vanishing Gradient</b>	Gradient becoming near-zero in deep networks, preventing learning
<b>Variational Autoencoder (VAE)</b>	Autoencoder learning probabilistic latent space for generation
<b>Vector Database</b>	Database optimized for storing and searching embedding vectors. Pinecone, Weaviate

<b>Vector Embedding</b>	Dense numerical representation of data in continuous space
<b>Vision Language Model (VLM)</b>	Model processing both images and text. GPT-4V, LLaVA, Claude
<b>Vision Transformer (ViT)</b>	Transformer architecture applied to image patches for computer vision
<b>Vocabulary</b>	Set of all tokens known to tokenizer/model
<b>Warmup</b>	Gradually increasing learning rate at training start for stability
<b>Weight Decay</b>	Regularization adding penalty proportional to weight magnitudes
<b>Weight Initialization</b>	Strategy for setting initial parameter values. Xavier, He initialization
<b>Weight Sharing</b>	Using same parameters across different parts of model
<b>Whisper</b>	OpenAI's speech recognition model supporting multiple languages
<b>Word2Vec</b>	Algorithm learning word embeddings from co-occurrence patterns. Skip-gram, CBOW
<b>XGBoost</b>	Gradient boosting library known for performance on tabular data
<b>Zero-Shot Learning</b>	Performing tasks without any task-specific training examples
<b>Zero-Shot Prompting</b>	Using LLM without examples, relying only on instruction

AI 기술은 빠르게 진화하고 있습니다.  
이 용어집이 AI 시대를 이해하는 데 도움이 되길 바랍니다.

---