# How LoRA works

LoRA (Low-Rank Adaptation) works by **efficiently fine-tuning large pretrained models** without updating all their parameters.

**Core idea:**

Instead of modifying the full weight matrix $W \in \mathbb{R}^{d \times k}$, LoRA adds two small trainable matrices $A$ and $B$ of much lower rank $r \ll \min(d, k)$:

$$W' = W + \Delta W \quad \text{where} \quad \Delta W = BA$$

Here:

- $A \in \mathbb{R}^{r \times k}$ and $B \in \mathbb{R}^{d \times r}$
- $W$ (original weights) stays **frozen**
- Only $A$ and $B$ are **trained**, drastically reducing trainable parameters

**During inference:**

You merge $\Delta W$ with $W$ or keep them separate—both yield the same output.

**Intuition:**

- LoRA assumes the needed weight updates lie in a **low-dimensional subspace**.
- It's like approximating a large matrix update with a small-rank correction.

**Benefits:**

- Uses <1% of parameters of full fine-tuning
- Memory- and compute-efficient
- Easily swappable—different tasks can load different LoRA adapters on the same base model

**Example (for an attention layer):**

In a Transformer, LoRA is often applied to the query/key projection matrices:

$$h = (W_q + B_q A_q)x$$

This lets you adapt model behavior (e.g., to a new domain or language) with minimal storage.

Let's unpack LoRA's mathematics from a linear algebra perspective.

---

## 1. Base setup

In a standard neural layer (say a linear/fully connected one):

$$h = Wx$$

where

- $W \in \mathbb{R}^{d \times k}$ are the pretrained weights
- $x \in \mathbb{R}^k$ is the input
- $h \in \mathbb{R}^d$ is the output

In full fine-tuning, we update $W$ entirely.

LoRA instead modifies $W$ by adding a *low-rank* matrix.

---

## 2. Low-rank reparameterization

LoRA defines:

$$W' = W + \Delta W, \quad \Delta W = BA$$

where

- $A \in \mathbb{R}^{r \times k}$
- $B \in \mathbb{R}^{d \times r}$
- $r \ll \min(d, k)$ (typically 4–16)

This means the update to $W$ is **rank ≤ r** (since rank of product ≤ min(ranks)).

So instead of learning $d \times k$ parameters, we only learn $r(d + k)$.

---

## 3. Mathematical interpretation

In matrix algebra, any rank-r matrix $M$ can be factored as:

$$M = UV$$

where $U \in \mathbb{R}^{d \times r}, V \in \mathbb{R}^{r \times k}$.

This is exactly what LoRA does — it assumes the *effective parameter update* lies in a **low-dimensional subspace** of the full parameter space.

So LoRA approximates:

$$\text{argmin}_{\Delta W} \ L(f_{W + \Delta W}(x), y) \quad \text{subject to} \quad \text{rank}(\Delta W) \leqslant r$$

This is a **low-rank constrained optimization problem**, solved indirectly by optimizing $A$ and $B$.

---

## 4. Scaling

To stabilize training, LoRA often includes a scaling factor $\alpha / r$:

$$\Delta W = \frac{\alpha}{r} BA$$

This keeps the effective update magnitude roughly invariant with respect to $r$.

---

## 5. Forward pass

During training:

$$h = (W + \frac{\alpha}{r} BA)x = Wx + \frac{\alpha}{r} B(Ax)$$

- $Wx$: base model output
- $Ax$: projection to low-dim subspace
- $B(Ax)$: projection back to output space

Thus, LoRA adds a *low-rank perturbation* to the layer's linear transformation.

---

## 6. Geometric intuition

- $A$: defines a low-dimensional *direction* in input space.
- $B$: lifts that direction into the output space.
- Combined, $BA$ spans a **rank-r subspace** of all possible parameter updates.

So LoRA fine-tunes within a **restricted linear subspace** rather than the full high-dimensional weight space.

## 7. Relation to SVD

If you performed an SVD on the full fine-tuned $W_{FT} - W$:

$$W_{FT} - W = U\Sigma V^{T}$$

then LoRA is essentially learning an approximation:

$$\Delta W \approx U_r\Sigma_r V_r^{T}$$

where $r$ is small — but instead of decomposing it after the fact, LoRA learns it *directly* in low-rank form.

---

Would you like me to visualize this geometrically (how LoRA projects and reconstructs updates in low-dimensional space)?

> So basically, the update matrix we are simply refactoring as product of 2 matrix, which are lower in rank right?

Exactly — that's the essence of it. ✅

LoRA assumes that the **weight update** $\Delta W$ doesn't need to be full-rank. So instead of directly learning a large $\Delta W \in R^{d \times k}$, it **factorizes** it as:

$$\Delta W = BA$$

with

- $A \in R^{r \times k}$,
- $B \in R^{d \times r}$,
  and $r \ll \min(d, k)$.

This makes $\Delta W$ *at most rank-r*.

The two smaller matrices capture the **direction** and **magnitude** of changes in a compressed form — giving nearly the same expressive power but with far fewer trainable parameters.

> Ok lets go to fundamental level. Rank of a matrix is basically the number of basis vectors which are defining the linear space right?

Exactly right — that's the fundamental geometric meaning. ✅

The **rank** of a matrix $A \in \mathrm{R}^{m \times n}$ is the **dimension of the vector space spanned by its columns (or rows)**.
Equivalently:

- It's the **number of linearly independent columns** (column space basis), or
- The **number of linearly independent rows** (row space basis).

So:

- If $\mathrm{rank}(A) = r$, then all columns of $A$ live in an $r$-dimensional subspace of $\mathrm{R}^m$.
- The matrix can be thought of as **mapping inputs through an $r$-dimensional bottleneck** — only $r$ independent directions of transformation actually matter.

That's exactly why LoRA constrains updates to be low-rank: it says, "the new directions we need to adapt the model lie in a much smaller subspace than the full parameter space."

> Got it, now the rank of A and B matrix in LoRA, do we choose it or is it learnable?

Good question — the **rank $r$** in LoRA is a **hyperparameter**, not a learnable quantity.

You **choose $r$** before training (e.g., 4, 8, 16).
Then during training, the **entries** of $A$ and $B$ are learned, but their **shapes** (and hence rank) stay fixed.

Formally:

- $A \in \mathrm{R}^{r \times k}$
- $B \in \mathrm{R}^{d \times r}$
  so the effective update $\Delta W = BA$ can have rank at most $r$.

Choosing $r$ controls a trade-off:

- Smaller $r$: fewer trainable parameters, cheaper, less expressive.
- Larger $r$: more expressive but closer to full fine-tuning.

In practice, $r$ is tuned empirically — many tasks work surprisingly well even with $r = 4$ or $8$.

> But why should we choose the 'r'? Shouldn't it be fixed/constant for a fixed dimensional matrix?

Excellent question — that's the kind of subtle but deep thinking that really gets at what LoRA *is doing under the hood*.

Let's unpack this carefully.

---

### 1. Rank is not an intrinsic property we're matching

The **rank of the LoRA update** $\Delta W = BA$ is not determined by the size of $W$; it's a *modeling choice* about **how expressive we want the update to be**.

- $W \in \mathbb{R}^{d \times k}$ has a *maximum possible* rank of $\min(d, k)$.
- But the "true" update we need (to adapt the model to a new task) likely lies in a *smaller subspace*.

So LoRA assumes:

> The needed adaptation lives in a low-dimensional subspace of the full parameter space.

That's why we **pick $r$** — it defines the *capacity* of that subspace.

---

### 2. Analogy: choosing model capacity

Think of $r$ like the **hidden dimension** in an autoencoder:

- Larger hidden dimension → more capacity to represent complex patterns.
- Smaller → forces the model to find more compact, generalizable features.

Same logic here — we choose $r$ to control how flexible the adaptation is.

---

### 3. Why it can't be "fixed"

Even for the same base matrix dimension ($d \times k$):

- Different tasks need different levels of adaptation.
  - Small domain shift → low $r$ may suffice.
  - Big domain shift → higher $r$ needed.
- Memory and speed constraints also affect what's practical.

So $r$ is like a **regularization knob**, not a mathematically dictated constant.

---

## 4. In principle

You *could* imagine learning $r$ automatically (e.g., via pruning or low-rank approximation during training), but LoRA itself doesn't — it fixes $r$ to keep training simple and predictable.

---

**Summary:**

$r$ isn't fixed by the base weight dimensions — it's a *design choice* that trades off parameter efficiency and representational power.

> rank and dimension of a matrix are the same thing in general?

Not quite — they're related but **not the same thing**.

Here's the distinction:

---

## 1. Matrix dimension

- Refers to the **shape** of the matrix.
- For $A \in \mathbb{R}^{m \times n}$:
  - $m$ = number of rows
  - $n$ = number of columns
  - So the *dimension* of the matrix is "$m \times n$" — it tells you how big the matrix is, not how much information it actually contains.

---

## 2. Matrix rank

- Refers to the **intrinsic dimensionality of the information** inside the matrix.
- It's the number of **linearly independent columns (or rows)** — i.e., the dimension of the *column space* or *row space*.
- Denoted $\text{rank}(A) = r \leqslant \min(m, n)$.

---

## 3. Intuition

Think of the matrix as a linear transformation:

$$A : \mathbb{R}^n \rightarrow \mathbb{R}^m$$

- The *dimension* says: it maps an $n$-dimensional input to an $m$-dimensional output.
- The *rank* says: how many independent directions in the input actually produce distinct outputs.

If the rank is smaller than $n$, some input directions are **collapsed** (mapped to the same output).

---

## 4. Example

$$A = \begin{bmatrix} 1 & 2 & 3 \\ 2 & 4 & 6 \end{bmatrix}$$

Here:

- Dimension = $2 \times 3$
- Rank = 1 (since the second row is 2× the first → only one independent direction)

So:

**Dimension** = physical size of matrix

**Rank** = number of independent information-carrying directions

> How to decide the 'r' in LoRA? I mean is it a trial n error way or there are systematic way of doing it?

Great question — and a subtle one, because **deciding the LoRA rank $r$** connects both theory and practice.

Let's break it down cleanly:

## 1. In practice — mostly empirical (trial & tuning)

Today, most people still **choose** *r* through experimentation:

- Try a few small values like $r = 4, 8, 16, 32$
- Compare validation performance and resource usage
- Pick the smallest *r* that gives near–full fine-tuning accuracy

Why this works: many NLP and vision tasks show diminishing returns beyond $r \approx 8$ or $16$.

---

## 2. Theoretical reasoning — low-rank structure of weight updates

Research shows that **fine-tuning gradients and updates themselves are low-rank** in many large models.
That means the true $\Delta W$ often lives in a low-dimensional subspace.
So if you could compute the full fine-tuned update and take its **SVD decomposition**:

$$\Delta W = U \Sigma V^T$$

then the top *r* singular values (and vectors) often capture most of the useful variation.
LoRA effectively learns that low-rank approximation directly.

In theory, you could:

1. Fine-tune a small portion of the model fully,
2. Measure singular value decay of the weight deltas,
3. Choose *r* so that cumulative explained variance $\geqslant$ some threshold (e.g., 95%).

That's a **systematic**, data-driven way — but rarely done because it's expensive.

---

## 3. Practical heuristics

Some simple rules of thumb:

| Layer type | Typical *r* | Notes |
| --- | --- | --- |
| Attention (Q,K,V,O) | 4–16 | Smaller usually fine |

| Layer type | Typical $r$ | Notes |
| --- | --- | --- |
| MLP (dense) | 8–32 | Higher if domain shift is large |
| Small models | lower $r$ | They're already overparameterized |
| Large domain adaptation | higher $r$ | More flexibility helps |

---

## 4. Emerging systematic methods

Recent research explores **adaptive or learned-rank LoRA**:

- *AdaLoRA* (Hu et al., 2023): dynamically adjusts per-layer rank during training by tracking importance via singular values.
- *Dynamic Rank Adaptation*: starts with high $r$ and prunes unimportant directions over time.

So, systematic approaches exist, but most practitioners still start with empirical tuning and use dynamic-rank variants only if needed.

---

**In short:**

- $r$ controls expressive power vs. efficiency.
- Start with small $r$, scale up if performance drops.
- For rigorous setups, estimate $r$ from singular value spectra or use adaptive methods like AdaLoRA.

> What is the benefit of matrix decomposition?

Excellent question — matrix decomposition (or **factorization**) is one of the most fundamental tools in linear algebra and machine learning. Its power lies in revealing the **structure**, **relationships**, and **efficient representations** hidden inside a matrix.

Here's a breakdown of the **key benefits**:

---

## 1. Dimensionality Reduction

Decomposition expresses a large matrix in terms of a few key basis directions.
Example: **SVD**

$$A = U\Sigma V^T$$

If singular values in $\Sigma$ decay fast, you can approximate:

$$A \approx U_r \Sigma_r V_r^T$$

This keeps only the top $r$ dimensions — a low-rank approximation that captures most of the important information while removing noise.
Used in **PCA, LoRA, and compression**.

---

## 2. Reveals Linear Structure

Decompositions show:

- How vectors transform under the matrix (stretch, rotate, reflect, scale).
- Which directions (eigenvectors/singular vectors) are *invariant* or *dominant*.
  This helps understand stability, correlations, and intrinsic geometry.

---

## 3. Computational Efficiency

Decomposition converts hard problems into simpler subproblems:

- Linear systems $Ax = b$ → solve easily via **LU** or **QR**.
- Eigenvalue problems → diagonalize for **fast exponentiation or powers of** $A$.
- LoRA exploits this idea by storing only low-rank factors $A$ and $B$ instead of full $W$.

---

## 4. Noise Filtering and Regularization

In data with redundancy or noise, low-rank decomposition acts like a **filter**:

- Removes small singular values → eliminates directions dominated by noise.
- Keeps dominant subspace → captures essential signal.

That's why decompositions underpin PCA, denoising autoencoders, and collaborative filtering.

## 5. Interpretability and Insight

Each decomposition type reveals a different perspective:

| Type | Insight |
| --- | --- |
| **SVD** | Energy (variance) along orthogonal directions |
| **Eigendecomposition** | Principal modes of transformation |
| **LU / QR** | Numerical stability & efficient solving |
| **NMF** | Additive, interpretable parts (e.g., topics, components) |

## 6. Storage & Transmission Efficiency

Low-rank factorization drastically reduces storage:

$$A \in \mathbb{R}^{m \times n} \quad \text{vs.} \quad (U_r, \Sigma_r, V_r^T)$$

For $r \ll \min(m, n)$, this is a massive compression — exactly the same idea LoRA uses to store weight updates compactly.

**In short:**

Matrix decomposition transforms a complex, high-dimensional operator into simpler, interpretable, and often smaller components — enabling **understanding, compression, and computation** all at once.

> The decomposed matrices ideally have lower dimensions, right?

Exactly — that's the key idea. ✅

When you decompose a matrix, the resulting **factor matrices usually have lower dimensions** that capture the **essential structure** of the original high-dimensional matrix.

Let's make this concrete:

---

## 1. Example: Singular Value Decomposition (SVD)

For $A \in \mathbb{R}^{m \times n}$,

$$A = U\Sigma V^T$$

where

- $U \in \mathbb{R}^{m \times m}$,
- $\Sigma \in \mathbb{R}^{m \times n}$,
- $V \in \mathbb{R}^{n \times n}$.

Now, if $A$ has **rank** $r$, you can approximate it by truncating:

$$A \approx U_r \Sigma_r V_r^T$$

where

- $U_r \in \mathbb{R}^{m \times r}$
- $\Sigma_r \in \mathbb{R}^{r \times r}$
- $V_r \in \mathbb{R}^{n \times r}$

These are **much smaller matrices** that reconstruct the main structure of $A$.

So yes — decomposition exposes a **low-dimensional representation** of the same linear mapping.

---

## 2. Intuitive view

You can think of decomposition as finding a **bottleneck representation**:

$$\text{input (n-dim)} \xrightarrow{V_r} \text{r-dim latent} \xrightarrow{U_r \Sigma_r} \text{output (m-dim)}$$

It compresses data into $r$ dimensions — exactly what LoRA does by learning two smaller matrices $A$ and $B$ of rank $r$.

---

### 3. Why this matters

- Lower-dimensional factors reduce **storage** and **computation**.
- They often remove **redundant** or **noisy** information.
- They reveal **core patterns or directions** that dominate the data.

---

So yes — decomposed matrices (or truncated ones) are smaller and represent the **intrinsic dimensionality** of the data or transformation, not its full ambient size.

> Then Autoencoder is also a kind of matrix decomposition technique right

Beautifully reasoned — yes, you're absolutely right. ✅

At its **core**, an **autoencoder** *is* a **nonlinear generalization of matrix decomposition** (specifically, low-rank factorization). Let's make that connection precise.

---

### 1. Linear autoencoder = matrix factorization

If you build an autoencoder with:

- a single linear encoder layer $h = W_e x$
- a single linear decoder layer $\hat{x} = W_d h$
- no activation function
- and train it with MSE loss

then the optimal solution learns:

$$W_e, W_d \text{ such that } X \approx W_d W_e X$$

That's **exactly** a rank-$r$ approximation of $X$.
In fact, it can be shown (Baldi & Hornik, 1989) that:

> a linear autoencoder with a bottleneck of size $r$ learns the same subspace as PCA.

So mathematically:

$$X \approx U_r \Sigma_r V_r^T$$

and the encoder–decoder pair corresponds to $V_r$ and $U_r \Sigma_r$.

Hence, **linear autoencoder = learned SVD/PCA decomposition**.

---

## 2. Nonlinear autoencoder = nonlinear factorization

When you add nonlinear activations (ReLU, tanh, etc.):

$$h = f(W_e x + b_e), \quad \hat{x} = g(W_d h + b_d)$$

the model can now represent **nonlinear manifolds** instead of linear subspaces.

That's why:

- PCA captures linear correlations
- Autoencoders can capture **nonlinear latent structures** in the data

In this sense, nonlinear autoencoders are like **nonlinear matrix decomposers** — they still compress data through a bottleneck but in a *curved manifold* rather than a flat subspace.

---

## 3. Summary: conceptual equivalence

| Concept | Linear Case | Nonlinear Case |
|---|---|---|
| **Goal** | Approximate $X$ by low-rank reconstruction | Approximate $X$ via low-dim nonlinear manifold |
| **Math form** | $X \approx W_d W_e X$ | $X \approx g(W_d f(W_e X))$ |
| **Equivalent to** | PCA / SVD | Nonlinear generalization |
| **Latent space** | Linear subspace | Nonlinear manifold |

So yes — the **autoencoder is a matrix decomposition framework**,
→ **linear** ones recover SVD/PCA,
→ **nonlinear** ones generalize it to more complex data manifolds.