

AI Applications Lecture 18

Parameter Efficient Fine Tuning and LoRA

SUZUKI, Atsushi

Jing WANG

2025-11-18

Contents

1	Introduction	3
1.1	Learning Outcomes	4
2	Preparation: Mathematical Notations Revisited	4
3	The Idea of PEFT: Representing Model Differences with Low-Dimensional Parameters	6
3.1	Review of PEFT Motivation	6
3.2	General Definition of PEFT Adapters	6
3.3	A Simple Example of PEFT	7
4	The Idea of Low Rank Adaptation (LoRA)	9
4.1	Low-Rank Constraint on the Difference Representation	9
4.2	Parameterization of the Difference Matrix in LoRA	12
4.3	General Definition of LoRA	14
4.4	Training LoRA Parameters	15
4.5	Controlling LoRA Strength After Training	15
5	Specific Application of LoRA: Example with Stable Diffusion 1.5	16
5.1	Overview of Stable Diffusion 1.5 and LoRA Application Points	16
5.2	LoRA in Self-Attention and Cross-Attention within U-Net	17
5.2.1	Mathematical Representation of the CrossAttention Layer	17
5.2.2	Applying LoRA to the CrossAttention Layer	18
5.3	LoRA in the Text Encoder	19
5.3.1	Mathematical Representation of Self-Attention in CLIPTextModel	19
5.3.2	LoRA for CLIP Text Self-Attention	20
5.4	Implementation Handling of LoRA at Inference Time	20

6 Summary and Next Lecture Preview	21
6.1 Answers to Learning Outcomes	21
6.2 Next Lecture Preview	21

1 Introduction

The machine learning framework used in much of generative AI is divided into two steps: the "learning" step, which appropriately determines the parameters of a parametric function, and the "inference" step, which actually executes the desired task using the function determined by the parameters set in the learning step. This lecture, like the previous one, focuses on "learning".

The neural networks currently used in practice in the field of generative AI have an enormous number of parameters. Since machine learning is a framework for determining the parameters of a parametric function from training data, an amount of training data corresponding to the number of parameters to be determined is required. For example, if one attempts to minimize an objective function using a gradient method, the objective function itself must contain enough information to sufficiently judge the quality of the parameters for the task one wants to solve, and it is easy to understand that this requires an amount of training data corresponding to the number of parameters.

Collecting a large amount of training data is accompanied by difficulties, and even if it is collected, learning on it requires a large amount of computation. Therefore, if possible (i.e., if there are no major performance problems), one would like to reduce the amount of training data used. For learning to succeed even with less training data, the number of parameters changed during learning needs to be small (of course, this is a necessary condition for successful learning, not a sufficient condition).

For the reasons mentioned above, **one would like to reduce the number of parameters changed during learning, if possible**. Of course, in a situation where one is building a model by training from scratch, if the task itself is complex and the desired input-output relationship is complicated, a high-dimensional checkpoint (a large number of parameters) and a large-scale architecture are necessary to represent that complex input-output relationship, and data for training it is also necessary. Thus, the policy of reducing the number of parameters to be learned is inherently unreasonable.

On the other hand, if one is in a situation where one **already has** a model f_θ that gives a good input-output relationship, and it is expected that a better input-output relationship can be obtained by modifying it slightly (in some sense), then a scheme of **learning the minimum number of parameters just enough to represent the part corresponding to that slight modification** might be practically feasible. The idea of starting from an existing model and modifying it slightly to obtain a better model is called **fine-tuning**, and the idea of realizing this by learning a small number of parameters is called **PEFT (parameter-efficient fine-tuning)**. In this lecture, we will learn about **LoRA (low-rank adaptation)** as a typical example of PEFT [1, 2].

1.1 Learning Outcomes

After completing this lecture, students should be able to do the following:

- Control the behavior of generative AI using LoRA.
- Explain the advantages of LoRA.
- Distinguish and explain what LoRA does mathematically and what it does in implementation.

2 Preparation: Mathematical Notations Revisited

- **Definition:**

- (LHS) := (RHS): Indicates that the left-hand side is defined by the right-hand side.
For example, $a := b$ indicates that a is defined by b .

- **Set (Set):**

- Sets are often denoted by uppercase calligraphic letters. Example: \mathcal{A} .
 - $x \in \mathcal{A}$: Indicates that the element x belongs to the set \mathcal{A} .
 - $\{\}$: Empty set.
 - $\{a, b, c\}$: The set consisting of elements a, b, c (extensional notation).
 - $\{x \in \mathcal{A} \mid P(x)\}$: The set of elements in set \mathcal{A} for which the proposition $P(x)$ is true (intensional notation).
 - \mathbb{R} : The set of all real numbers.
 - $\mathbb{R}_{>0}$: The set of all positive real numbers.
 - $\mathbb{R}_{\geq 0}$: The set of all non-negative real numbers.
 - \mathbb{Z} : The set of all integers.
 - $\mathbb{Z}_{>0}$: The set of all positive integers.
 - $\mathbb{Z}_{\geq 0}$: The set of all non-negative integers.
 - $[1, k]_{\mathbb{Z}} := \{1, 2, \dots, k\}$: For a positive integer k , the set of integers from 1 to k .

- **Function (Function):**

- $f : \mathcal{X} \rightarrow \mathcal{Y}$: Indicates that the function f is a map that takes elements from set \mathcal{X} as input and outputs elements from set \mathcal{Y} .
 - $y = f(x)$: Indicates that the output is $y \in \mathcal{Y}$ when $x \in \mathcal{X}$ is input to the function f .

- **Vector (Vector):**

- In this course, a vector refers to a column of numbers.
- Vectors are denoted by bold italic lowercase letters. Example: \mathbf{v} .
- $\mathbf{v} \in \mathbb{R}^n$: Indicates that the vector \mathbf{v} is an n -dimensional real vector.
- The i -th element of vector \mathbf{v} is denoted as v_i .

$$\mathbf{v} = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix}. \quad (1)$$

• **Matrix (Matrix):**

- Matrices are denoted by bold italic uppercase letters. Example: \mathbf{A} .
- $\mathbf{A} \in \mathbb{R}^{m,n}$: Indicates that the matrix \mathbf{A} is an $m \times n$ real matrix.
- The element in the i -th row and j -th column of matrix \mathbf{A} is denoted as $a_{i,j}$.

$$\mathbf{A} = \begin{bmatrix} a_{1,1} & a_{1,2} & \cdots & a_{1,n} \\ a_{2,1} & a_{2,2} & \cdots & a_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m,1} & a_{m,2} & \cdots & a_{m,n} \end{bmatrix}. \quad (2)$$

- The transpose of matrix \mathbf{A} is denoted as \mathbf{A}^\top . If $\mathbf{A} \in \mathbb{R}^{m,n}$, then $\mathbf{A}^\top \in \mathbb{R}^{n,m}$, and

$$\mathbf{A}^\top = \begin{bmatrix} a_{1,1} & a_{2,1} & \cdots & a_{m,1} \\ a_{1,2} & a_{2,2} & \cdots & a_{m,2} \\ \vdots & \vdots & \ddots & \vdots \\ a_{1,n} & a_{2,n} & \cdots & a_{m,n} \end{bmatrix} \quad (3)$$

- A vector is also a matrix with 1 column, and its transpose can also be defined.

$$\mathbf{v}^\top = \begin{bmatrix} v_1 & v_2 & \cdots & v_n \end{bmatrix} \in \mathbb{R}^{1,n} \quad (4)$$

• **Tensor (Tensor):**

- In this lecture, the term tensor simply refers to a multi-dimensional array. A vector can be regarded as a 1st-order tensor, and a matrix as a 2nd-order tensor. Tensors of 3rd-order or higher are denoted by underlined bold italic uppercase letters, like $\underline{\mathbf{A}}$.

- Students who have already learned about abstract tensors in mathematics or physics may feel resistant to calling a mere multi-dimensional array a tensor. If one considers that the basis is always fixed to the standard basis and the mathematical tensor is identified with its component representation (which becomes a multi-dimensional array), then there is (arguably) consistency in the terminology.

3 The Idea of PEFT: Representing Model Differences with Low-Dimensional Parameters

3.1 Review of PEFT Motivation

Let's review the motivation for PEFT. We assume a situation where we already have a model f_θ that provides a good input-output relationship, and we expect to obtain a better input-output relationship by modifying it slightly. In this situation, the idea of PEFT is to try to represent the part corresponding to that slight modification with a small number of parameters (or a low-dimensional parameter vector).

Remark 3.1. Here, "slight modification" is not something to be defined strictly mathematically. Although there are many mathematical criteria for measuring the closeness of functions, it does not necessarily mean "slight" in that sense here. "Slight" here, in the end, can only be defined as a modification that can be represented by a small number of parameters, which becomes a tautology. To be more precise, it is more accurate to say that methods have been empirically found that can achieve good performance improvements with specific modifications represented by a small number of parameters, and these are called PEFT.

3.2 General Definition of PEFT Adapters

We will strictly define PEFT as a function. Here, we let the parameter space of the original model f_θ be

$$\boldsymbol{\theta} \in \mathbb{R}^{d_{\text{Original}}} \quad (5)$$

and the input and output spaces be \mathcal{X} and \mathcal{Y} , respectively.

Definition 3.1 (PEFT Adapter). Let \mathcal{X} be the input space and \mathcal{Y} be the output space. Fix $d_{\text{Original}} \in \mathbb{Z}_{>0}$, and suppose we are given a family of functions with parameters on the parameter vector space $\mathbb{R}^{d_{\text{Original}}}$

$$f_{(\cdot)} : \mathbb{R}^{d_{\text{Original}}} \rightarrow \mathcal{Y}^{\mathcal{X}} \quad (6)$$

. Here, $f_{(\cdot)}$ is a "map from parameters to models", and given $\boldsymbol{\theta} \in \mathbb{R}^{d_{\text{Original}}}$, a specific function $f_\theta : \mathcal{X} \rightarrow \mathcal{Y}$ is determined.

Fix $d_{\text{PEFT}} \in \mathbb{Z}_{>0}$, and let the **PEFT parameter** be $\boldsymbol{\psi} \in \mathbb{R}^{d_{\text{PEFT}}}$. At this time, a **PEFT adapter**

is a map

$$\text{PEFTAdapter}_{f(\cdot), \psi} : \mathbb{R}^{d_{\text{Original}}} \rightarrow \mathcal{Y}^X \quad (7)$$

. The original model $\theta \in \mathbb{R}^{d_{\text{Original}}}$ is transformed into

$$\text{PEFTAdapter}_{f(\cdot), \psi}(\theta) \quad (8)$$

, depending on the small number of parameters ψ . The specific way this transformation is determined depends on the PEFT method.

Remark 3.2. Definition 3.1 defines PEFT very abstractly, but the essence is as follows. We fix a base model f_θ and construct a new function $\text{PEFTAdapter}_{f(\cdot), \psi}(\theta)$ through the PEFTAdapter. At this time, $\text{PEFTAdapter}_{f(\cdot), \psi}(\theta)$ is controlled by the low-dimensional parameter ψ . That is, it is a framework that allows access to various models in the vicinity of f_θ by changing ψ .

Under Definition 3.1, **PEFT is a framework for determining the PEFT parameter ψ in the PEFT adapter by learning, so that it becomes a "good" function for the desired task**. More specifically, the learning in PEFT involves defining some objective function (e.g., a loss function) as a function of ψ and minimizing it.

3.3 A Simple Example of PEFT

Since Definition 3.1 is abstract, let's consider the simplest specific example. Consider the case where the original parameter vector θ is explicitly divided into a part that is not updated during learning and a part that is updated.

Example 3.1 (Simple PEFT by Modifying Only Some Parameters). For $d_{\text{fixed}}, d_{\text{tuned}} \in \mathbb{Z}_{>0}$, let

$$d_{\text{Original}} := d_{\text{fixed}} + d_{\text{tuned}} \quad (9)$$

and represent the parameter vector by block partitioning as

$$\theta = \begin{bmatrix} \theta_{\text{fixed}} \\ \theta_{\text{tuned}} \end{bmatrix} \in \mathbb{R}^{d_{\text{Original}}} \quad (10)$$

. Here, suppose $\theta_{\text{fixed}} \in \mathbb{R}^{d_{\text{fixed}}}$ are the parameters not changed during learning (fixed parameters), and $\theta_{\text{tuned}} \in \mathbb{R}^{d_{\text{tuned}}}$ are the parameters changed during learning.

At this time, we define the **simple PEFT adapter** $\text{SimplePEFTAdapter}_{f(\cdot), \Delta\theta_{\text{tuned}}}$ as follows:

$$\text{SimplePEFTAdapter}_{f(\cdot), \Delta\theta_{\text{tuned}}} \left(\begin{bmatrix} \theta_{\text{fixed}} \\ \theta_{\text{tuned}} \end{bmatrix} \right) := f \begin{bmatrix} \theta_{\text{fixed}} \\ \theta_{\text{tuned}} + \Delta\theta_{\text{tuned}} \end{bmatrix}. \quad (11)$$

Here, $\Delta\theta_{\text{tuned}} \in \mathbb{R}^{d_{\text{tuned}}}$ is the difference parameter determined by learning, and it corre-

sponds to the PEFT parameter in this example. That is, matching the notation of Definition 3.1, we have

$$\psi = \Delta\theta_{\text{tuned}}, \quad d_{\text{PEFT}} = d_{\text{tuned}} \quad (12)$$

Remark 3.3. Let's make some observations from Example 3.1.

- First, the number of parameters to be determined is $d_{\text{PEFT}} = d_{\text{tuned}}$, which can be set freely by the designer. In particular, it can be designed such that $d_{\text{tuned}} \ll d_{\text{Original}}$. This is a great advantage in terms of learning stability when the available training data is scarce.
- Next, consider the perspective of memory for storage. If one only wants to obtain the result of PEFT, one may discard the original θ and save only

$$\begin{bmatrix} \theta_{\text{fixed}} \\ \theta_{\text{tuned}} + \Delta\theta_{\text{tuned}} \end{bmatrix} \quad (13)$$

- . The number of parameters to save in this case is d_{Original} , which is the same as saving the original model.
- On the other hand, suppose that both the original f_θ and the PEFT-modified $f_{\begin{bmatrix} \theta_{\text{fixed}} \\ \theta_{\text{tuned}} + \Delta\theta_{\text{tuned}} \end{bmatrix}}$ are valuable for applications. In this case, naively thinking, one would save two sets of parameters, θ and (13), resulting in storing a total of $2d_{\text{Original}}$ floating-point values.
- However, in reality, the PEFT-modified (13) can be recovered from the original θ and $\Delta\theta_{\text{tuned}}$ as

$$\begin{bmatrix} \theta_{\text{fixed}} \\ \theta_{\text{tuned}} + \Delta\theta_{\text{tuned}} \end{bmatrix} = \begin{bmatrix} \theta_{\text{fixed}} \\ \theta_{\text{tuned}} \end{bmatrix} + \begin{bmatrix} \mathbf{0} \\ \Delta\theta_{\text{tuned}} \end{bmatrix} \quad (14)$$

, so the number of floating-point values that need to be saved is

$$d_{\text{Original}} + d_{\text{PEFT}} = d_{\text{Original}} + d_{\text{tuned}} \quad (15)$$

- . Usually, $d_{\text{tuned}} \ll d_{\text{Original}}$, so PEFT can significantly reduce the number of floating-point values to be saved.
- In particular, if PEFT is performed separately for multiple purposes, for example, if K PEFTs are performed and each PEFT result is valuable for applications. Naive storage would require saving

$$(K + 1)d_{\text{Original}} \quad (16)$$

floating-point values. On the other hand, when using PEFT, one only needs to save the original θ and each task-specific $\Delta\theta_{\text{tuned}}^{(k)}$ ($k = 1, \dots, K$), so the number of floating-point values to be saved is

$$d_{\text{Original}} + Kd_{\text{PEFT}} = d_{\text{Original}} + Kd_{\text{tuned}} \quad (17)$$

. When $d_{\text{tuned}} \ll d_{\text{Original}}$ and K is large, this difference is practically very significant.

Thus, PEFT is advantageous from both the perspective of "number of parameters to update during learning" and "number of parameters to save". In the next chapter, we will introduce LoRA, which achieves even higher parameter efficiency by imposing further structural constraints on parameters represented as matrices.

4 The Idea of Low Rank Adaptation (LoRA)

4.1 Low-Rank Constraint on the Difference Representation

In the method described above, if $d_{\text{PEFT}} \ll d_{\text{Original}}$, the parameter efficiency of PEFT is good, but a concern is whether it can provide practically good differences when d_{PEFT} is small. Of course, mathematically considering the worst case, it's impossible for a small number of parameters to substitute for many parameters. The issue here is not the mathematical worst case, but whether, when considering practical utility in natural or social sciences, a better difference can be represented with a small number of parameters.

In fact, for the neural networks of generative AI in practical use, it is empirically known that when a part of the parameters can be regarded as a matrix representing a linear transformation, it is better to limit the difference to be represented by a **low-rank matrix** using fewer parameters, rather than representing the difference directly. This can represent practically useful differences despite the small number of parameters [1, 2]. **LoRA (low-rank adaptation)** is based on this idea, representing the difference for matrix-form parameters as a product of a matrix with few rows (a matrix with a small number of rows) and a matrix with few columns (a matrix with a small number of columns).

As a basic matter of linear algebra, we will strictly define the rank of a matrix based on matrix factorization.

Definition 4.1 (Definition of Matrix Rank by Factorization). Let $d_{\text{out}}, d_{\text{in}} \in \mathbb{Z}_{>0}$ and $\mathbf{W} \in \mathbb{R}^{d_{\text{out}}, d_{\text{in}}}$. Consider the set

$$S(\mathbf{W}) := \{r \in \mathbb{Z}_{\geq 0} \mid \exists \mathbf{B} \in \mathbb{R}^{d_{\text{out}}, r}, \mathbf{A} \in \mathbb{R}^{r, d_{\text{in}}} \text{ exist such that } \mathbf{W} = \mathbf{BA}\} \quad (18)$$

. At this time, the **rank** of \mathbf{W} is defined as

$$\text{rank}(\mathbf{W}) := \min S(\mathbf{W}) \quad (19)$$

. That is, $\text{rank}(\mathbf{W})$ is the minimum value of the number of column vectors r required when representing \mathbf{W} as

$$\mathbf{W} = \mathbf{B}\mathbf{A}, \quad \mathbf{B} \in \mathbb{R}^{d_{\text{out}},r}, \quad \mathbf{A} \in \mathbb{R}^{r,d_{\text{in}}} \quad (20)$$

, and it represents the "minimum number of vectors" necessary to represent \mathbf{W} . In particular, if $\mathbf{W} = \mathbf{0}$, then $\text{rank}(\mathbf{W}) = 0$.

Proposition 4.1 (Equivalent Characterizations of Matrix Rank). For $\text{rank}(\mathbf{W})$ defined in Definition 4.1, the following hold.

- (1) $\text{rank}(\mathbf{W}) = \dim C(\mathbf{W})$. Here, $C(\mathbf{W})$ is the subspace spanned by the column vectors of \mathbf{W} , called the **column space**.
- (2) $\text{rank}(\mathbf{W}) = \dim \mathcal{R}(\mathbf{W})$. Here, $\mathcal{R}(\mathbf{W})$ is the subspace spanned by the row vectors of \mathbf{W} , called the **row space**.

In particular, the dimensions of the column space and row space of \mathbf{W} are always equal.

Proof. First, we show (1). Let $\text{rank}(\mathbf{W}) = r$. By Definition 4.1, there exist $\mathbf{B} \in \mathbb{R}^{d_{\text{out}},r}$ and $\mathbf{A} \in \mathbb{R}^{r,d_{\text{in}}}$ such that

$$\mathbf{W} = \mathbf{B}\mathbf{A} \quad (21)$$

. Let the column vectors of \mathbf{B} be $\mathbf{b}_1, \dots, \mathbf{b}_r$ and the column vectors of \mathbf{W} be $\mathbf{w}_1, \dots, \mathbf{w}_{d_{\text{in}}}$. If we write the j -th column of \mathbf{A} as $\mathbf{a}^{(j)} \in \mathbb{R}^r$, then from equation (21), for each $j \in [1, d_{\text{in}}]_{\mathbb{Z}}$,

$$\mathbf{w}_j = \mathbf{B}\mathbf{a}^{(j)} = \sum_{i=1}^r a_i^{(j)} \mathbf{b}_i \quad (22)$$

. Therefore, any column vector \mathbf{w}_j can be expressed as a linear combination of $\{\mathbf{b}_1, \dots, \mathbf{b}_r\}$, so

$$C(\mathbf{W}) \subset \text{span}\{\mathbf{b}_1, \dots, \mathbf{b}_r\} \quad (23)$$

holds. Thus,

$$\dim C(\mathbf{W}) \leq r = \text{rank}(\mathbf{W}) \quad (24)$$

To show the reverse inequality, let $s := \dim C(\mathbf{W})$ and let $\mathbf{u}_1, \dots, \mathbf{u}_s$ be a basis for $C(\mathbf{W})$. Define the matrix $\mathbf{B}' \in \mathbb{R}^{d_{\text{out}},s}$ having these as columns:

$$\mathbf{B}' := [\mathbf{u}_1 \quad \mathbf{u}_2 \quad \cdots \quad \mathbf{u}_s] \quad (25)$$

. Since each column vector \mathbf{w}_j is an element of the column space, there exists a coefficient vector $\mathbf{c}^{(j)} \in \mathbb{R}^s$ such that

$$\mathbf{w}_j = \sum_{i=1}^s c_i^{(j)} \mathbf{u}_i \quad (26)$$

. By arranging these coefficients into columns,

$$\mathbf{A}' := \begin{bmatrix} \mathbf{c}^{(1)} & \mathbf{c}^{(2)} & \cdots & \mathbf{c}^{(d_{\text{in}})} \end{bmatrix} \in \mathbb{R}^{s, d_{\text{in}}} \quad (27)$$

, it follows from equations (25) and (27) that

$$\mathbf{W} = \mathbf{B}' \mathbf{A}' \quad (28)$$

. That is, $s \in S(\mathbf{W})$, and by Definition 4.1,

$$\text{rank}(\mathbf{W}) \leq s = \dim C(\mathbf{W}) \quad (29)$$

holds. Combining equations (24) and (29),

$$\text{rank}(\mathbf{W}) = \dim C(\mathbf{W}) \quad (30)$$

follows. This proves (1).

Next, we show (2). First, from Definition 4.1, we show that for any \mathbf{W} ,

$$\text{rank}(\mathbf{W}) = \text{rank}(\mathbf{W}^T) \quad (31)$$

holds. In fact, if $r \in S(\mathbf{W})$, then there exist \mathbf{B}, \mathbf{A} such that equation (21) holds. Taking the transpose,

$$\mathbf{W}^T = \mathbf{A}^T \mathbf{B}^T \quad (32)$$

. Here, $\mathbf{A}^T \in \mathbb{R}^{d_{\text{in}}, r}$ and $\mathbf{B}^T \in \mathbb{R}^{r, d_{\text{out}}}$, so $r \in S(\mathbf{W}^T)$. The reverse can be shown similarly, so $S(\mathbf{W}) = S(\mathbf{W}^T)$, and thus equation (31) holds.

On the other hand, the row space $\mathcal{R}(\mathbf{W})$ is the subspace spanned by the row vectors of \mathbf{W} , which is identical to the column space $C(\mathbf{W}^T)$ of \mathbf{W}^T . Therefore,

$$\mathcal{R}(\mathbf{W}) = C(\mathbf{W}^T) \quad (33)$$

. Applying (1) to \mathbf{W}^T ,

$$\text{rank}(\mathbf{W}^T) = \dim C(\mathbf{W}^T) \quad (34)$$

, and combining equations (31), (33), and (34), we get

$$\text{rank}(\mathbf{W}) = \text{rank}(\mathbf{W}^T) = \dim C(\mathbf{W}^T) = \dim \mathcal{R}(\mathbf{W}) \quad (35)$$

. This shows (2) and completes the proof of the proposition. \square

Remark 4.1. From Proposition 4.1, $\text{rank}(\mathbf{W}) = r$ is equivalent to the dimension of the column space of \mathbf{W} being r , i.e., "the number of linearly independent column vectors needed to generate all column vectors of \mathbf{W} is r ". On the other hand, according to Definition 4.1, when

$\text{rank}(\mathbf{W}) = r$, \mathbf{W} can be expressed using $\mathbf{B} \in \mathbb{R}^{d_{\text{out}},r}$ and $\mathbf{A} \in \mathbb{R}^{r,d_{\text{in}}}$ as in equation (21). In this case, the number of elements in \mathbf{B} is $d_{\text{out}}r$, and the number of elements in \mathbf{A} is rd_{in} , so \mathbf{W} can be described by specifying a total of

$$d_{\text{rank}} := d_{\text{out}}r + rd_{\text{in}} \quad (36)$$

scalars. In particular, when $r \ll \min\{d_{\text{out}}, d_{\text{in}}\}$,

$$d_{\text{rank}} \ll d_{\text{out}}d_{\text{in}} \quad (37)$$

, so a low-rank matrix is a class of matrices that can be described with far fewer parameters compared to handling the original full-rank matrix directly. This is perfectly consistent with the motivation of PEFT, which is to represent models and differences as much as possible with a small number of parameters.

LoRA improves parameter efficiency by artificially constraining the difference matrix $\Delta\mathbf{W}$ to be represented in this form.

4.2 Parameterization of the Difference Matrix in LoRA

In LoRA, the difference matrix $\Delta\mathbf{W} \in \mathbb{R}^{d_{\text{out}},d_{\text{in}}}$ for a matrix parameter $\mathbf{W} \in \mathbb{R}^{d_{\text{out}},d_{\text{in}}}$ is parameterized as follows.

Definition 4.2 (Difference Matrix in LoRA). Fix $d_{\text{out}}, d_{\text{in}} \in \mathbb{Z}_{>0}$, and let the original matrix parameter be $\mathbf{W} \in \mathbb{R}^{d_{\text{out}},d_{\text{in}}}$. Let $r \in \mathbb{Z}_{>0}$ be a parameter called **rank**, and $\alpha \in \mathbb{R}$ be a parameter called **scaling factor**.

At this time, the difference matrix $\Delta\mathbf{W} \in \mathbb{R}^{d_{\text{out}},d_{\text{in}}}$ by LoRA is defined as

$$\Delta\mathbf{W} := \frac{\alpha}{r} \mathbf{BA} \quad (38)$$

. Here,

$$\mathbf{A} \in \mathbb{R}^{r,d_{\text{in}}}, \quad \mathbf{B} \in \mathbb{R}^{d_{\text{out}},r} \quad (39)$$

are matrices determined by learning, and are the main learning targets of LoRA.

The matrix parameter after applying LoRA is defined as

$$\tilde{\mathbf{W}} := \mathbf{W} + \Delta\mathbf{W} = \mathbf{W} + \frac{\alpha}{r} \mathbf{BA} \quad (40)$$

Remark 4.2. As an implementation note regarding LoRA, the α/r appearing in equation (38) is redundant in a strict sense. That is, the same expressive power can be maintained by setting α to 1 and incorporating the scalar into \mathbf{A} or \mathbf{B} .

In reality, the LoRA paper [1], after including α and r in the formulation as above, explicitly

states that α is set equal to the rank r , i.e., $\alpha = r$, as an experimental setup. In this case, the coefficient α/r in equation (38) is always 1, which is mathematically equivalent to adding $\Delta W = BA$. The inclusion of α and r in the formulation is for the convenience of theoretical explanation.

Also, in the LoRA training script `train_text_to_image_lora.py` for Stable Diffusion in Hugging Face diffusers, the setting for LoRA added to the UNet is `LoraConfig(r=args.rank, lora_alpha=args.rank, ...)`, which is designed such that $\alpha = r$ in implementation as well^a. Therefore, although α and r are mathematically redundant, the form $\frac{\alpha}{r}$ is preserved in both the original paper and representative implementations, and on top of that, it is typically used with $\alpha = r$.

^aSee the code example in the Hugging Face diffusers LoRA training guide. At <https://huggingface.co/docs/diffusers/training/lora>, the setting `unet_lora_config = LoraConfig(r=args.rank, lora_alpha=args.rank, init_lora_weights="gaussian", target_modules=["to_k", "to_q", "to_v", "to_out.0"])` is shown, from which it is clear that $\alpha = r$.

Let's count the number of parameters in the difference matrix. Whereas the original full degrees of freedom are $d_{\text{out}}d_{\text{in}}$, in LoRA, the degrees of freedom for A, B are

$$d_{\text{LoRA}} := rd_{\text{in}} + d_{\text{out}}r = r(d_{\text{in}} + d_{\text{out}}) \quad (41)$$

and one for the scalar α , totaling $r(d_{\text{in}} + d_{\text{out}}) + 1$. Usually, we choose $r \ll \min\{d_{\text{in}}, d_{\text{out}}\}$, so $d_{\text{LoRA}} \ll d_{\text{out}}d_{\text{in}}$.

Example 4.1 (Comparison of Specific Parameter Counts). As specific numbers, consider the case $d_{\text{in}} = 768$, $d_{\text{out}} = 3072$, $r = 32$. This is a size typically appearing in the linear transformations of intermediate layers in Transformer-based models [3, 4].

First, the number of elements in the original matrix parameter $W \in \mathbb{R}^{3072, 768}$ is

$$d_{\text{full}} = d_{\text{out}}d_{\text{in}} = 3072 \times 768 \quad (42)$$

. Calculating this gives

$$d_{\text{full}} = 3072 \times 768 = 2,359,296 \quad (43)$$

scalars.

On the other hand, the number of elements in A, B representing the difference by LoRA is

$$d_{\text{LoRA}} = r(d_{\text{in}} + d_{\text{out}}) = 32(768 + 3072) = 32 \times 3840 = 122,880 \quad (44)$$

, and even adding the scalar α , it is 122,881.

Therefore, compared to representing the difference matrix at full size, the number of parameters required is

$$\frac{d_{\text{LoRA}}}{d_{\text{full}}} \approx \frac{1.23 \times 10^5}{2.36 \times 10^6} \approx 0.052 \quad (45)$$

, which is reduced to about 1/19. This is for representing the difference for a single matrix,

and in practice, LoRA is applied to many matrices in the network, but a similar reduction effect is obtained for each matrix.

4.3 General Definition of LoRA

Next, we define the framework for applying LoRA to a general neural network containing multiple matrix parameters.

Definition 4.3 (LoRA for General Neural Networks). Let $f_{W_1, W_2, \dots, W_K, \theta_{\text{others}}}$ be a neural network consisting of K matrix parameters and other parameters. That is, $W_k \in \mathbb{R}^{d_{\text{out},k}, d_{\text{in},k}}$ ($k = 1, \dots, K$) and

$$\theta_{\text{others}} \in \mathbb{R}^{d_{\text{others}}} \quad (46)$$

. Let the entire set of parameters be

$$\Theta = (W_1, \dots, W_K, \theta_{\text{others}}) \quad (47)$$

, and this is the parameter of the model f .

For each $k \in [1, K]_{\mathbb{Z}}$, fix a rank $r_k \in \mathbb{Z}_{>0}$, and let

$$A_k \in \mathbb{R}^{r_k, d_{\text{in},k}}, \quad B_k \in \mathbb{R}^{d_{\text{out},k}, r_k} \quad (48)$$

be the LoRA matrix parameters, and $\alpha_k \in \mathbb{R}$ be the LoRA scaling factor. Also, let $\eta \in \mathbb{R}$ be a parameter that controls the overall scaling of LoRA. η is usually a hyperparameter set at inference time and is often not a learning target.

At this time, the **LoRA adapter** is defined as

$$\begin{aligned} & \text{LoRA}_{f(\cdot), \eta, \alpha_1, A_1, B_1, \dots, \alpha_K, A_K, B_K}(W_1, \dots, W_K, \theta_{\text{others}}) \\ &:= f_{W_1 + \eta \frac{\alpha_1}{r_1} B_1 A_1, W_2 + \eta \frac{\alpha_2}{r_2} B_2 A_2, \dots, W_K + \eta \frac{\alpha_K}{r_K} B_K A_K, \theta_{\text{others}}}. \end{aligned} \quad (49)$$

Here, for $k = 1, 2, \dots, K$, r_k is a predetermined positive integer and is not changed during learning. Also, α_k is usually a fixed positive number.

Remark 4.3. In Definition 4.3, η is a parameter that adjusts the "strength" of LoRA. In the learning phase, learning is often performed with $\eta = 1$, and by changing η in the inference phase,

$$\tilde{W}_k(\eta) = W_k + \eta \frac{\alpha_k}{r_k} B_k A_k \quad (50)$$

, the contribution of the modification by LoRA can be strengthened or weakened. Setting $\eta = 0$ returns to the original model $f_{W_1, \dots, W_K, \theta_{\text{others}}}$.

4.4 Training LoRA Parameters

Under Definition 4.3, training LoRA means determining $\mathbf{A}_k, \mathbf{B}_k$ ($k = 1, \dots, K$) based on the training data. Here, we assume that the original matrix parameters \mathbf{W}_k and other parameters θ_{others} are fixed and not updated.

Given m input-output pairs $(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})$ for $i = 1, \dots, m$, define the objective function (loss function) as

$$\mathcal{L}(\alpha_1, \mathbf{A}_1, \mathbf{B}_1, \dots, \alpha_K, \mathbf{A}_K, \mathbf{B}_K) := \frac{1}{m} \sum_{i=1}^m \ell(\mathbf{y}^{(i)}, f^{\text{LoRA}}(\mathbf{x}^{(i)})) \quad (51)$$

. Here,

$$f^{\text{LoRA}}(\mathbf{x}) := \text{LoRA}_{f(\cdot), \eta=1, \alpha_1, \mathbf{A}_1, \mathbf{B}_1, \dots, \alpha_K, \mathbf{A}_K, \mathbf{B}_K}(\mathbf{W}_1, \dots, \mathbf{W}_K, \theta_{\text{others}})(\mathbf{x}) \quad (52)$$

$$= f_{\mathbf{W}_1 + \frac{\alpha_1}{r_1} \mathbf{B}_1 \mathbf{A}_1, \dots, \mathbf{W}_K + \frac{\alpha_K}{r_K} \mathbf{B}_K \mathbf{A}_K, \theta_{\text{others}}}(\mathbf{x}) \quad (53)$$

, and ℓ is a task-dependent loss function, such as mean squared error or cross-entropy.

Definition 4.4 (LoRA Training Problem). The LoRA training problem is, given the fixed original parameters $\mathbf{W}_1, \dots, \mathbf{W}_K, \theta_{\text{others}}$, to solve

$$\underset{\mathbf{A}_1, \mathbf{B}_1, \dots, \mathbf{A}_K, \mathbf{B}_K}{\text{Minimize}} \quad \mathcal{L}(\mathbf{A}_1, \mathbf{B}_1, \dots, \mathbf{A}_K, \mathbf{B}_K) \quad (54)$$

Remark 4.4. In actual training, a gradient method (stochastic gradient descent, AdamW, etc.) is used to numerically solve the minimization problem in Definition 4.4. That is, at each step,

$$\frac{\partial \mathcal{L}}{\partial \mathbf{A}_k}, \quad \frac{\partial \mathcal{L}}{\partial \mathbf{B}_k} \quad (55)$$

are obtained by backpropagation, and $\mathbf{A}_k, \mathbf{B}_k$ are updated based on them. At this time, gradients with respect to \mathbf{W}_k and θ_{others} are calculated but not used for updates (they are not updated). This allows training only the LoRA parameters without changing the original model.

4.5 Controlling LoRA Strength After Training

After the LoRA parameters $\alpha_k, \mathbf{A}_k, \mathbf{B}_k$ are determined by Definition 4.4, the strength of the LoRA contribution can be controlled using η from Definition 4.3 in the inference phase.

That is, at inference time, by specifying an arbitrary $\eta \in \mathbb{R}$ and defining

$$f^{\text{LoRA}, \eta}(\mathbf{x}) = f_{\mathbf{W}_1 + \eta \frac{\alpha_1}{r_1} \mathbf{B}_1 \mathbf{A}_1, \dots, \mathbf{W}_K + \eta \frac{\alpha_K}{r_K} \mathbf{B}_K \mathbf{A}_K, \theta_{\text{others}}}(\mathbf{x}) \quad (56)$$

, one can return to the original model with $\eta = 0$, use the same strength as during training

with $\eta = 1$, or strengthen or weaken the LoRA contribution more than usual with $|\eta| > 1$.

In implementation, by allowing the user to explicitly specify η , the "strength of the LoRA flavor" can be easily adjusted. This is frequently used for purposes such as adjusting style transfer or the tendency of prompt responses.

5 Specific Application of LoRA: Example with Stable Diffusion 1.5

In this section, we will outline the application of LoRA in Stable Diffusion 1.5 [5, 6] as a specific example. Here, we will mathematically describe which matrix parameters LoRA is applied to, referring to the layer names in the Python libraries `diffusers`¹ and `transformers`².

5.1 Overview of Stable Diffusion 1.5 and LoRA Application Points

Stable Diffusion 1.5 is broadly composed of the following components:

- **Text encoder**: Usually a `CLIPTextModel` from `transformers` or a similar model, which generates embeddings (condition vectors) from the input text token sequence [7].
- **U-Net denoiser** (U-Net style noise estimator): Implemented as `UNet2DConditionModel` in `diffusers`, it estimates noise from a noisy latent variable. It includes self-attention and cross-attention blocks.
- **VAE (variational auto-encoder)**: Implemented as `AutoencoderKL` or similar, it handles the mutual conversion between images and latent representations.

In practice, the main places where LoRA is applied are the following two:

- Targeting only the U-Net: Apply LoRA to the linear layers related to self-attention and cross-attention in the noise estimator (especially `to_q`, `to_k`, `to_v`, `to_out`).
- Also applying to the text encoder: In addition to the above, apply LoRA to the self-attention layers (`self.attn.q_proj`, `k_proj`, `v_proj`, `out_proj`) or MLP parts of the `CLIPTextModel`.

In the following, we will define the functions represented by each layer, separating input variables and parameters, and strictly describe in mathematical formulas which matrices LoRA is applied to.

¹Official repository for the `diffusers` library: <https://github.com/huggingface/diffusers>

²Official repository for the `transformers` library: <https://github.com/huggingface/transformers>

5.2 LoRA in Self-Attention and Cross-Attention within U-Net

The attention blocks in the U-Net of Stable Diffusion 1.5 are often implemented as the `diffusers.models.attention.CrossAttention` class. Here, we will abstractly define the function determined by this attention block.

5.2.1 Mathematical Representation of the CrossAttention Layer

Let d_{model} be the dimension of the feature vectors, d_{head} be the dimension of each head, and n_{head} be the number of heads, such that

$$d_{\text{model}} = n_{\text{head}} d_{\text{head}} \quad (57)$$

Definition 5.1 (CrossAttention Layer in U-Net). Let the query input sequence and the key/-value input sequences be

$$\mathbf{X}_q \in \mathbb{R}^{L_q, d_{\text{model}}}, \quad \mathbf{X}_k \in \mathbb{R}^{L_k, d_{\text{model}}}, \quad \mathbf{X}_v \in \mathbb{R}^{L_k, d_{\text{model}}} \quad (58)$$

. Here, L_q, L_k are the query sequence length and the key/value sequence length, respectively.

The CrossAttention layer of the U-Net has the following matrix parameters:

$$\mathbf{W}_q \in \mathbb{R}^{d_{\text{model}}, d_{\text{model}}}, \quad (59)$$

$$\mathbf{W}_k \in \mathbb{R}^{d_{\text{model}}, d_{\text{model}}}, \quad (60)$$

$$\mathbf{W}_v \in \mathbb{R}^{d_{\text{model}}, d_{\text{model}}}, \quad (61)$$

$$\mathbf{W}_o \in \mathbb{R}^{d_{\text{model}}, d_{\text{model}}}. \quad (62)$$

These correspond to `to_q`, `to_k`, `to_v`, `to_out` in the `diffusers` implementation.

First, by linear transformation,

$$\mathbf{Q} := \mathbf{X}_q \mathbf{W}_q^\top \in \mathbb{R}^{L_q, d_{\text{model}}}, \quad (63)$$

$$\mathbf{K} := \mathbf{X}_k \mathbf{W}_k^\top \in \mathbb{R}^{L_k, d_{\text{model}}}, \quad (64)$$

$$\mathbf{V} := \mathbf{X}_v \mathbf{W}_v^\top \in \mathbb{R}^{L_k, d_{\text{model}}} \quad (65)$$

are determined. $\mathbf{Q}, \mathbf{K}, \mathbf{V}$ are split for each head, and written as

$$\mathbf{Q} \leftrightarrow (\mathbf{Q}^{(h)})_{h=1}^{n_{\text{head}}}, \quad \mathbf{Q}^{(h)} \in \mathbb{R}^{L_q, d_{\text{head}}}, \quad (66)$$

$$\mathbf{K} \leftrightarrow (\mathbf{K}^{(h)})_{h=1}^{n_{\text{head}}}, \quad \mathbf{K}^{(h)} \in \mathbb{R}^{L_k, d_{\text{head}}}, \quad (67)$$

$$\mathbf{V} \leftrightarrow (\mathbf{V}^{(h)})_{h=1}^{n_{\text{head}}}, \quad \mathbf{V}^{(h)} \in \mathbb{R}^{L_k, d_{\text{head}}} \quad (68)$$

. Scaled dot-product attention is calculated for each head. That is,

$$\mathbf{S}^{(h)} := \frac{1}{\sqrt{d_{\text{head}}}} \mathbf{Q}^{(h)} \mathbf{K}^{(h)\top} \in \mathbb{R}^{L_q, L_k} \quad (69)$$

, and applying softmax row-wise,

$$\mathbf{A}^{(h)} := \text{softmax}_{\text{row}}(\mathbf{S}^{(h)}) \in \mathbb{R}^{L_q, L_k} \quad (70)$$

is determined. Then

$$\mathbf{O}^{(h)} := \mathbf{A}^{(h)} \mathbf{V}^{(h)} \in \mathbb{R}^{L_q, d_{\text{head}}} \quad (71)$$

is obtained. The outputs for $h = 1, \dots, n_{\text{head}}$ are concatenated in the feature dimension direction

$$\mathbf{O} := \text{Concat}_{h=1}^{n_{\text{head}}} \mathbf{O}^{(h)} \in \mathbb{R}^{L_q, d_{\text{model}}} \quad (72)$$

, and finally

$$\mathbf{Y} := \mathbf{O} \mathbf{W}_o^{\top} \in \mathbb{R}^{L_q, d_{\text{model}}} \quad (73)$$

is the output of the CrossAttention layer.

Remark 5.1. Definition 5.1 writes down the standard multi-head attention of the Transformer [3], adapted for the attention blocks of the U-Net. Although batch dimensions, positional encodings, etc., are included in the implementation, the important point is that the linear layers (63)–(65), (73) targeted by LoRA have been explicitly shown in mathematical formulas.

5.2.2 Applying LoRA to the CrossAttention Layer

A typical method for applying LoRA is to add separate LoRAs to each of $\mathbf{W}_q, \mathbf{W}_k, \mathbf{W}_v, \mathbf{W}_o$. For example, consider LoRA for \mathbf{W}_q .

Prepare rank r_q , scaling factor α_q , matrices $\mathbf{A}_q \in \mathbb{R}^{r_q, d_{\text{model}}}$, $\mathbf{B}_q \in \mathbb{R}^{d_{\text{model}}, r_q}$, and

$$\tilde{\mathbf{W}}_q := \mathbf{W}_q + \eta \frac{\alpha_q}{r_q} \mathbf{B}_q \mathbf{A}_q \quad (74)$$

is determined. Similarly,

$$\tilde{\mathbf{W}}_k := \mathbf{W}_k + \eta \frac{\alpha_k}{r_k} \mathbf{B}_k \mathbf{A}_k, \quad (75)$$

$$\tilde{\mathbf{W}}_v := \mathbf{W}_v + \eta \frac{\alpha_v}{r_v} \mathbf{B}_v \mathbf{A}_v, \quad (76)$$

$$\tilde{\mathbf{W}}_o := \mathbf{W}_o + \eta \frac{\alpha_o}{r_o} \mathbf{B}_o \mathbf{A}_o \quad (77)$$

are defined. Here, r_q, r_k, r_v, r_o are the ranks for each respective matrix, and in implementation, a common value is often used.

The CrossAttention layer after applying LoRA is defined as in Definition 5.1, but replacing

$\mathbf{W}_q, \mathbf{W}_k, \mathbf{W}_v, \mathbf{W}_o$ with $\tilde{\mathbf{W}}_q, \tilde{\mathbf{W}}_k, \tilde{\mathbf{W}}_v, \tilde{\mathbf{W}}_o$ respectively in (63)–(73).

5.3 LoRA in the Text Encoder

The text encoder of Stable Diffusion 1.5 is the CLIP text model [7], implemented as `CLIPTextModel` in `transformers`. Here too, LoRA is often applied to the linear transformations of the self-attention layers.

5.3.1 Mathematical Representation of Self-Attention in `CLIPTextModel`

Let $d_{\text{model}}^{\text{text}}$ be the dimension of the text embedding, $n_{\text{head}}^{\text{text}}$ be the number of heads, and $d_{\text{head}}^{\text{text}}$ be the dimension of each head.

$$d_{\text{model}}^{\text{text}} = n_{\text{head}}^{\text{text}} d_{\text{head}}^{\text{text}} \quad (78)$$

. Let the input sequence be

$$\mathbf{X}_{\text{text}} \in \mathbb{R}^{L_{\text{text}}, d_{\text{model}}^{\text{text}}} \quad (79)$$

Definition 5.2 (CLIP Text Self-Attention Layer). The self-attention layer of the CLIP text model has the following matrix parameters:

$$\mathbf{W}_q^{\text{text}} \in \mathbb{R}^{d_{\text{model}}^{\text{text}}, d_{\text{model}}^{\text{text}}}, \quad (80)$$

$$\mathbf{W}_k^{\text{text}} \in \mathbb{R}^{d_{\text{model}}^{\text{text}}, d_{\text{model}}^{\text{text}}}, \quad (81)$$

$$\mathbf{W}_v^{\text{text}} \in \mathbb{R}^{d_{\text{model}}^{\text{text}}, d_{\text{model}}^{\text{text}}}, \quad (82)$$

$$\mathbf{W}_o^{\text{text}} \in \mathbb{R}^{d_{\text{model}}^{\text{text}}, d_{\text{model}}^{\text{text}}}. \quad (83)$$

Query, key, and value are defined as

$$\mathbf{Q}^{\text{text}} := \mathbf{X}_{\text{text}} \mathbf{W}_q^{\text{text}\top}, \quad (84)$$

$$\mathbf{K}^{\text{text}} := \mathbf{X}_{\text{text}} \mathbf{W}_k^{\text{text}\top}, \quad (85)$$

$$\mathbf{V}^{\text{text}} := \mathbf{X}_{\text{text}} \mathbf{W}_v^{\text{text}\top}, \quad (86)$$

, and then multi-head self-attention is calculated similarly to Definition 5.1, and the output is

$$\mathbf{Y}^{\text{text}} := \mathbf{O}^{\text{text}} \mathbf{W}_o^{\text{text}\top} \quad (87)$$

5.3.2 LoRA for CLIP Text Self-Attention

The application of LoRA is similar to the CrossAttention in the U-Net, and for each matrix parameter,

$$\tilde{\mathbf{W}}_q^{\text{text}} = \mathbf{W}_q^{\text{text}} + \eta \frac{\mathbf{a}_q^{\text{text}}}{r_q^{\text{text}}} \mathbf{B}_q^{\text{text}} \mathbf{A}_q^{\text{text}}, \quad (88)$$

$$\tilde{\mathbf{W}}_k^{\text{text}} = \mathbf{W}_k^{\text{text}} + \eta \frac{\mathbf{a}_k^{\text{text}}}{r_k^{\text{text}}} \mathbf{B}_k^{\text{text}} \mathbf{A}_k^{\text{text}}, \quad (89)$$

$$\tilde{\mathbf{W}}_v^{\text{text}} = \mathbf{W}_v^{\text{text}} + \eta \frac{\mathbf{a}_v^{\text{text}}}{r_v^{\text{text}}} \mathbf{B}_v^{\text{text}} \mathbf{A}_v^{\text{text}}, \quad (90)$$

$$\tilde{\mathbf{W}}_o^{\text{text}} = \mathbf{W}_o^{\text{text}} + \eta \frac{\mathbf{a}_o^{\text{text}}}{r_o^{\text{text}}} \mathbf{B}_o^{\text{text}} \mathbf{A}_o^{\text{text}} \quad (91)$$

are determined. This modifies the self-attention of the text encoder by LoRA. In practice, by combining the LoRA of the text encoder with the LoRA of the U-Net, both the tendency of prompt interpretation and the image generation process are controlled.

5.4 Implementation Handling of LoRA at Inference Time

Finally, we describe two typical methods for handling LoRA in implementation at inference time.

Remark 5.2 (Fusing LoRA vs. On-the-fly Application). There are at least the following two options for how to specifically use the LoRA-applied matrix $\tilde{\mathbf{W}} = \mathbf{W} + \frac{\alpha}{r} \mathbf{BA}$ in calculations.

- **Pre-fusing method:** Before performing inference, $\mathbf{W}_{\text{LoRA}} := \mathbf{W} + \frac{\alpha}{r} \mathbf{BA}$ is calculated once and saved, and thereafter \mathbf{W}_{LoRA} is always used to calculate $\mathbf{W}_{\text{LoRA}} \mathbf{x}$. In the diffusers library, this is implemented as `fuse_lora()`. This method is advantageous in terms of computation, because once \mathbf{W}_{LoRA} is calculated, there is no need to calculate \mathbf{BAx} every time during numerous inferences (image generation steps).
- **Applying LoRA during each forward pass method:** Even at inference time, the original matrix \mathbf{W} is kept as is, and for the input \mathbf{x} ,

$$\mathbf{W}\mathbf{x} + \frac{\alpha}{r} \mathbf{BAx} \quad (92)$$

is calculated every time. This method is highly flexible when one wants to switch multiple LoRAs at once or turn LoRA on/off during inference. It is also useful in situations where one absolutely does not want to change the original \mathbf{W} .

In general, the on-the-fly method is convenient in situations where LoRA application and removal are frequently repeated, while the fuse method is advantageous in terms of computational efficiency when a certain LoRA is fixed for a long time to perform a large amount of inference.

6 Summary and Next Lecture Preview

6.1 Answers to Learning Outcomes

We briefly review the Learning Outcomes set forth in this lecture.

- LoRA is a type of parameter-efficient fine-tuning (PEFT) that can efficiently realize a better model for solving a task by determining a small number of parameters that represent the difference from the original model.
- LoRA requires changing only a small number of parameters, yet possesses sufficient expressive power for applications. This brings practical advantages such as improved learning stability with small data, reduction in storage usage, and ease of managing numerous task-specific adapters.
- Mathematically, LoRA restricts the difference matrix to one represented by the product of a matrix with few rows and a matrix with few columns, and determines those matrices with few rows and few columns.

6.2 Next Lecture Preview

Next time, we will explain the **relationship between the number of parameters or the size of the function space, and the amount of data required for learning**. That is, we will outline the relationship between model capacity, generalization performance, and the required sample size from the perspective of statistics and learning theory, and consider how methods like PEFT and LoRA can contribute to data efficiency.

References

- [1] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, L. Wang, and W. Chen, “Lora: Low-rank adaptation of large language models,” *arXiv preprint arXiv:2106.09685*, 2021.
- [2] T. Dettmers, A. Pagnoni, A. Holtzman, and L. Zettlemoyer, “Qlora: Efficient finetuning of quantized llms,” *Advances in Neural Information Processing Systems*, 2023. NeurIPS 2023.
- [3] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, “Attention is all you need,” *Advances in Neural Information Processing Systems*, vol. 30, 2017.
- [4] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever, “Language models are unsupervised multitask learners,” *OpenAI Technical Report*, 2019. Technical report.
- [5] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, “High-resolution image synthesis with latent diffusion models,” *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10684–10695, 2022.
- [6] S. AI, “Stable diffusion v1-5 model card.” Technical report, 2022. Model card and weights.
- [7] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, G. Krueger, and I. Sutskever, “Learning transferable visual models from natural language supervision,” *Proceedings of the 38th International Conference on Machine Learning*, pp. 8748–8763, 2021.