

Deep Learning System Design



Engineering and Service Architectures

Han Cheol Moon
School of Computer Science and Engineering
Nanyang Technological University

Singapore
hancheol001@e.ntu.edu.sg
February 4, 2025

Contents

I	Introduction	1
1	Introduction	2
1.1	Complexity of Matrix Multiplication	2
II	Transformers	4
1.2	Flash Attention	5
2	Positional Embeddings	6
2.0.1	Rotary Positional Embeddings	6
3	Tokenization	11
4	Model Compression	12

Part I

Introduction

Chapter 1

Introduction

1.1 Complexity of Matrix Multiplication

Matrix multiplication is a fundamental operation in many computational tasks, including neural networks. The complexity of multiplying two matrices depends on their dimensions. Let's dive into the specifics.

- Let A be a matrix of size $m \times k$.
- Let B be a matrix of size $k \times n$.
- The result C will be a matrix of size $m \times n$.

Standard Matrix Multiplication: For each element c_{ij} in the resulting matrix C :

$$c_{ij} = \sum_{l=1}^k a_{il} \cdot b_{lj}$$

This involves:

- Multiplications: k multiplications for each element c_{ij} .
- Additions: $k - 1$ additions for each element c_{ij} .

Complexity

- The total number of elements in C is $m \times n$.
- Therefore, the total number of multiplications is $m \times n \times k$.
- The total number of additions is $m \times n \times (k - 1)$.

Thus, the total complexity is $O(m \times n \times k)$.

Even though there are several advanced methods, the standard $O(m \times n \times k)$ complexity is often used in practice, due to the simplicity and efficiency of implementation on modern hardware. Optimized libraries (like BLAS, cuBLAS for GPUs) leverage hardware-specific optimizations to improve practical performance.

Part II

Transformers

1.2 Flash Attention

Chapter 2

Positional Embeddings

Rather than focusing on a token's absolute position in a sentence, relative positional embeddings concentrate on the distances between pairs of tokens. This method doesn't add a position vector to the word vector directly. Instead, it alters the attention mechanism to incorporate relative positional information.

2.0.1 Rotary Positional Embeddings

Chapter 3

Tokenization

Chapter 4

Model Compression