# Deep Learning System Design

Engineering and Service Architectures

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

Singapore
hancheol001@edu.ntu.edu.sg
June 23, 2024

# Contents

# 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.

### 1.1.1   Complexity in Neural Networks

In the context of neural networks:

- Input Matrices: Weight matrices and input feature vectors.

- Typical Sizes:

  - Weight matrix: $d \times d_{in}$ for RNNs, $d \times d$ for Transformers.
  - Input/Output vectors: Usually batch-processed, leading to sizes like $batch\_size \times sequence\_length \times feature\_size$.

# Bibliography

[1] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015.

[2] Rachit Singh. . `https://rachitsingh.com/elbo_surgery/`, 2017. Online; accessed 29 January 2014.

[3] Hany Hassan, Anthony Aue, Chang Chen, Vishal Chowdhary, Jonathan Clark, Christian Federmann, Xuedong Huang, Marcin Junczys-Dowmunt, William Lewis, Mu Li, Shujie Liu, Tie-Yan Liu, Renqian Luo, Arul Menezes, Tao Qin, Frank Seide, Xu Tan, Fei Tian, Lijun Wu, Shuangzhi Wu, Yingce Xia, Dongdong Zhang, Zhirui Zhang, and Ming Zhou. Achieving human parity on automatic chinese to english news translation. *CoRR*, abs/1803.05567, 2018.