

Lu Sun, and many more.

A Notebook on Machine Learning Updates



*To my family, friends and communities members who
have been dedicating to the presentation of this
notebook, and to all students, researchers and faculty
members who might find this notebook helpful.*



Contents

Foreword	vii
Preface	ix
List of Figures	xi
List of Tables	xiii
I Techniques	1
1 Transformer	3
1.1 It Seems That Attention is All You Need	3
1.2 CNN, RNN, and Their Limitations	4
1.2.1 Brief Review of CNN	5
1.2.2 Brief Review of RNN	6
1.3 Attention Mechanism	8
1.4 Transformer	8
1.5 Research Trends	9
II Products	11
2 OpenAI	15
2.1 General Introduction	15
2.1.1 Background	15
2.1.2 Business Model	15
2.2 ChatGPT	15
2.3 OpenAI API	16
2.4 DELL-E-2	16
3 Meta AI	17
3.1 General Introduction	17
3.1.1 Background	17
3.1.2 Business Model	17
3.2 PyTorch	17
3.3 LLaMA	18
3.3.1 LLaMA: Open and Efficient Foundation Language Mod- els	18

Foreword

If software or e-books can be made completely open-source, why not a notebook?

This brings me back to the summer of 2009 when I started my third year as a high school student in Harbin No. 3 High School. In around the end of August when the results of Gaokao (National College Entrance Examination of China, annually held in July) are released, people from photocopy shops would start selling notebooks photocopies that they claim to be from the top scorers of the exam. Much curious as I was about what these notebooks look like, never have I expected myself to actually learn anything from them, mainly for the following three reasons.

First of all, some (in fact many) of these notebooks were more difficult to understand than the textbooks. I guess we cannot blame the top scorers for being so smart that they sometimes make things extremely brief or overwhelmingly complicated.

Secondly, why would I want to adapt to notebooks of others when I had my own notebooks which in my opinion should be just as good as theirs.

And lastly, as a student in the top-tier high school myself, I knew that the top scorers of the coming year would probably be a schoolmate or a classmate. Why would I want to pay that much money to a complete stranger in a photocopy shop for my friend's notebook, rather than requesting a copy from him or her directly?

However, things had changed after my becoming an undergraduate student in 2010. There were so many modules and materials to learn in a university, and as an unfortunate result, students were often distracted from digging deeply into a module (For those who were still able to do so, you have my highest respect). The situation became even worse as I started pursuing my Ph.D. in 2014. As I had to focus on specific research areas entirely, I could hardly split much time on other irrelevant but still important and interesting contents.

This motivated me to start reading and taking notebooks for selected books and articles, just to force myself to spent time learning new subjects out of my comfort zone. I used to take hand-written notebooks. My very first notebook was on *Numerical Analysis*, an entrance level module for engineering background graduate students. Till today I still have on my hand dozens of these notebooks. Eventually, one day it suddenly came to me: why not digitalize them, and make them accessible online and open-source, and let everyone read and edit it?

As most of the open-source software, this notebook (and it applies to the other notebooks in this series as well) does not come with any “warranty” of any kind, meaning that there is no guarantee for the statement and knowledge in this notebook to be absolutely correct as it is not peer reviewed. **Do NOT cite this notebook in your academic research paper or book!** Of course, if you find anything helpful with your research, please trace back to the origin of the citation and double confirm it yourself, then on top of that determine whether or not to use it in your research.

This notebook is suitable as:

- a quick reference guide;
- a brief introduction for beginners of the module;
- a “cheat sheet” for students to prepare for the exam (Don’t bring it to the exam unless it is allowed by your lecturer!) or for lecturers to prepare the teaching materials.

This notebook is NOT suitable as:

- a direct research reference;
- a replacement to the textbook;

because as explained the notebook is NOT peer reviewed and it is meant to be simple and easy to read. It is not necessary brief, but all the tedious explanation and derivation, if any, shall be “fold into appendix” and a reader can easily skip those things without any interruption to the reading experience.

Although this notebook is open-source, the reference materials of this notebook, including textbooks, journal papers, conference proceedings, etc., may not be open-source. Very likely many of these reference materials are licensed or copyrighted. Please legitimately access these materials and properly use them.

Some of the figures in this notebook is drawn using Excalidraw, a very interesting tool for machine to emulate hand-writing. The Excalidraw project can be found in GitHub, [excalidraw/excalidraw](https://github.com/excalidraw/excalidraw).

Preface

Artificial intelligence (AI) and specifically machine learning were originally a research topic under the scope of control systems, and it is primarily used for system identification (“machine learning was used to be called system identification”, as some professors say). Its artificial neuron network (ANN) structure is a promising approach for building highly nonlinear self-tuning functions.

Due to the limitations of computational power and data storage in the past years, training large-size ANN with massive data points was hardly possible. With the advancement of computer and material science in the beginning of the 21st century, in particular the development of GPU and TPU, nowadays we can manage networks with dozens of layers, each containing hundreds of neurons. This is referred as the deep learning (DL) network structure. With DL as the building block, effective convolutional neural networks (CNN) and recurrent neural networks (RNN) have been proposed which are effective at finding trends in spacial and sequential data. The proposition of transformers, yet another attention-based deep learning structure, has taken natural language processing to the next level.

AI has grown so significantly in the past decades that it is now considered a separate area apart from classic control systems, and it is drawing more and more attentions. To put it in perspective, while a well-known textbook on control systems may have 10,000+ citations records after years since its publication, a famous conference paper on modern AI can easily have 50,000+ citations after a few months.

This notebook does not focus on the introduction of AI mechanisms, as most of which are already covered in control systems related notebooks. The purpose of this notebook is mainly to keep up-to-date with the latest technology in machine learning, starting from transformer and its applications in chatbot, the most famous of which being ChatGPT developed by OpenAI. The notebook cares not only the techniques, but also the business models used in these technologies and companies.

Some part of this notebook is written by ChatGPT-4.



List of Figures

1.1	A demonstration of CNN kernel. The input is given by the white box (3D), and the kernel by the red box.	6
1.2	A example of RNN.	7



List of Tables





Part I

Techniques



1

Transformer

CONTENTS

1.1	It Seems That Attention is All You Need	3
1.2	CNN, RNN, and Their Limitations	4
1.2.1	Brief Review of CNN	5
1.2.2	Brief Review of RNN	6
1.3	Attention Mechanism	8
1.4	Transformer	8
1.5	Research Trends	8

This chapter introduces the development of “transformer” in its main branch. Transformer is an AI framework initially proposed in 2017 in a paper “Attention is All You Need”, and it has magically revolutionized how we train AI for natural language processing and interpreting.

1.1 It Seems That Attention is All You Need

It started in 2017 when “transformer” is proposed in a paper “Attention is all you need”, by Vaswani, et al., a bunch of guys mainly from Google.

“Attention” was not a new concept previous to this paper. It was proposed a few years earlier than the transformer, and had been integrated into several CNN and RNN models to boost their performances to capture long-range dependencies. In these researches, attention mechanism is more often a supporting role than the main character. It is quite surprising indeed, that when getting rid of CNN and RNN and purely relying on attention mechanism, everything actually gets better: it is faster to train the network, and its performance gets improved.

Does this mean the RNN, which we had been relying on so much, is in fact the drawback to natural language processing? The following is the answer of ChatGPT-4 of this question.

Q: By getting rid of RNN and using attention mechanism alone, natural language processing performance gets better. Is this correct? Does this mean that RNN is a drawback to natural language processing?

A: Yes, by getting rid of RNNs and using the attention mechanism alone, the Transformer model demonstrated improved performance in natural language processing tasks. However, it's important to note that RNNs are not inherently a drawback to NLP; they were successfully used for many tasks before the advent of Transformers.

RNNs, particularly LSTM (Long Short-Term Memory) and GRU (Gated Recurrent Unit) variants, were designed to handle sequences and capture dependencies in the data. They have been effective in many NLP tasks, such as language modeling, machine translation, and sentiment analysis. However, RNNs do have certain limitations, such as difficulty in capturing long-range dependencies and sequential processing, which can lead to slower training and inference times.

The attention mechanism, as employed in the Transformer model, addresses some of these limitations. It enables parallel processing of input sequences, leading to faster training and inference, and effectively captures long-range dependencies without the vanishing gradient problem that can affect RNNs.

While Transformers have shown superior performance in many NLP tasks, RNNs still have their place in certain scenarios, especially when working with resource-constrained environments or smaller datasets. They can be more lightweight and computationally efficient compared to large-scale Transformer models.

In summary, RNNs are not a drawback to NLP per se, but their limitations have been mitigated by the introduction of the attention mechanism in Transformer models, which has led to improved performance in a wide range of NLP tasks.

(April 14, 2023, ChatGPT-4)

In the remaining part of the chapter, a quick review of CNN and RNN is given and their limitations introduced. Attention mechanism is illustrated. The original transformer model is explained. Finally, research trends of transformer as of today are briefly summarized.

1.2 CNN, RNN, and Their Limitations

This section does not focus on the detailed introduction of CNN and RNN mechanisms. Instead, only their main features are reviewed, and their limita-

tions discussed. This section does not reflect on the state-of-art of researches in CNN and RNN.

The conventional feedforward network (network without cycles) with dense ANN is widely used. It can be trained systematically using back-propagation methods such as stochastic gradient descent. However, dense ANN is not very good at interpreting correlations among the inputs. At least, it is not good at doing it efficiently.

For example, consider the input to be pixels of an image of size $3 \times 128 \times 128$, where 3 corresponds to the RGB of the image, and 128×128 its pixel size. It is obvious to a human that the pixel at (1, 1) should definitely has a closer relationship with the one at (1, 2) (they are probably forming an object in the image together) than (128, 128). The same applies to sequential inputs. Consider a signal sampled continuously. The first and second samples are very likely to be more strongly correlated than the first and last samples.

In a conventional dense ANN, the correlation cannot be captured efficiently. In contrast, a dense ANN would treat all the inputs “equally” in a symmetric manner. To enforce the ANN to “memorize” the correlation, it would require a lot more layers and nodes (which is often considered low-efficient), and requires more data points during the training. Other problems of conventional dense ANN include, for example, the lack of ability in handling data with arbitrary length.

CNN and RNN try to tackle the above problems by implementing a “pre-processing” stage, where the correlation of the spacial and sequential data is first abstracted using some mechanism, and the correlation information is sent as (additional) inputs to the followed dense ANN.

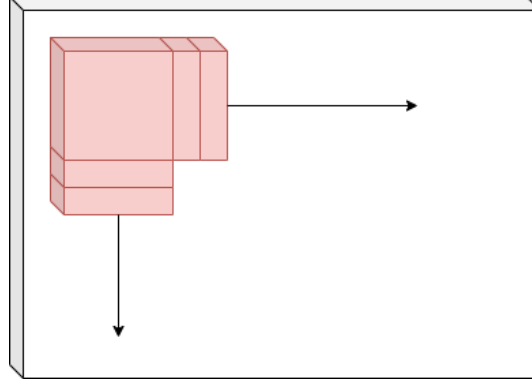
The details of CNN and RNN mechanisms are not discussed in this notebook. Brief reviews are given as follows.

1.2.1 Brief Review of CNN

CNN is a type of ANN structure designed to handle grid-like data. It is effective when dealing when data spatially correlated, thus becomes very popular in computer vision. It defines “kernel” that aggregates nearby pixels information before sending it to a dense network.

A CNN kernel, also known as a filter, is a small matrix of weights that slides across the input image or feature map to perform a mathematical operation called convolution. An demonstration of CNN kernel is given in Fig. 1.1. Multiple kernels can be defined on the same layer to handle the same feature map, each kernel associated with an output channel. In practice, each kernel or channel is designed to detect a specific features in the feature map. For example, there might be a kernel detecting edges, while a second kernel detect color codes.

Notice that CNN differs quite largely from transformer in the problems they are expected to address. CNN is more for spatial data processing such

**FIGURE 1.1**

A demonstration of CNN kernel. The input is given by the white box (3D), and the kernel by the red box.

as image processing, while transformer targets more on sequential data processing such as natural language processing and machine translation.

1.2.2 Brief Review of RNN

RNN is a connectionist model with the ability to selectively pass information across sequence steps [1]. It is good at handling sequence of data such as voice message, text contents, or a flow of images (videos). It is worth mentioning that the “sequence” does not necessarily mean a time sequence. Nevertheless, without losing generality, we will consider time sequence in the review for simplicity and convenience.

Denote inputs $x(1), x(2), \dots, x(k), \dots$ where $x(k)$ is a vector sampled at time instant k . The length of the sequence may be finite or infinite. In the case of finite sequence, its maximum sample index is denoted by T . For example, in the context of natural language processing, each input might be a word in a dictionary. For example, $x(1) = \text{“Pandas”}$, $x(2) = \text{“are”}$, $x(3) = \text{“so”}$, $x(4) = \text{“cute”}$, $x(5) = \text{“!”}$. The corresponding target output sequence is given by $y(1), y(2), \dots, y(k), \dots$, respectively.

RNN differs from the conventional dense ANN by introducing “recurrent edges”, which allows the output of hidden layers at $k - 1$ be used as additional inputs to the system at k . This means, at any time k , the input of the system includes both $x(k)$ and also selected $h(k - 1)$, where $h(\cdot)$ is the outputs of hidden layers. We can think of the “weights” of a trained RNN the “long-term memory” that does not change with specific sequence of inputs, while the information passing through recurrent edges the “short-term memory” that links previous inputs with future inputs.

A demonstration is given in Fig. 1.2. Notice that each hidden layer is a

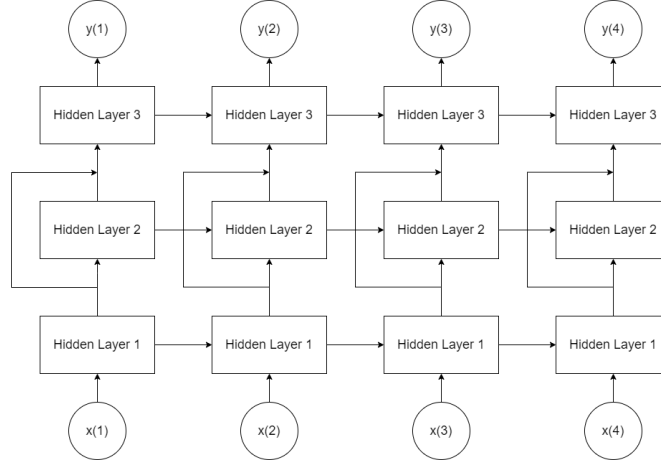


FIGURE 1.2
A example of RNN.

multi-input-multi-output subsystem containing multiple nodes. Different from a conventional dense ANN, each hidden layer takes additional inputs from its corresponding hidden layer in the previous instant. It is also common to see “bypass” (this is widely used in different ANN structures, not unique to RNN) for better performance of the system.

RNN has some limitations. One of the major problem is that it is difficult to train an RNN even for the basic standard feedforward networks. The optimization of RNN is NP-complete. It is especially difficult for RNN to learn long-range dependencies due to the vanishing and exploding gradients problem that could occur when backpropagating errors across many time-steps (long sequence) [1]. This is one of the main challenges why RNN has difficulties building on long-range dependencies. The vanishing and exploding gradient problems are caused by the structure of the system as well as the backpropagation-based training methods.

Different approaches have been proposed to prevent vanishing and exploding gradient problems. Famous ones among these approaches include strategic weight initialization, long short term memory (LSTM), gated recurrent units (GRUs), skip connections, and more. Many of these approaches try to reduce the effect of vanishing and exploding gradient problems by carefully design the ANN structures. For example, both LSTM and GRUs introduce “memory cells” with built-in “gates” that balance and control the flow of information from previous cells versus current inputs. The cells are used to replace the traditional perceptron nodes. These approaches have made the training of RNN a feasible problem. LSTM and GRUs are almost certainly used in modern RNNs.

Bidirectional RNN (BRNN) is proposed at about the same time with

LSTM. It allows information to travel not only from previous hidden layers to future layers, but also from future hidden layers to previous layers. LSTM and BRNN can be used together to boost the RNN performance. Notice that BRNN cannot run continuously as it requires fixed endpoints in both the future and the past. It is useful for prediction over a sequence of fixed length, such as part-of-speech tagging in natural language processing.

Another problem that people have found during the training of RNN is local optima. However, recent studies have shown that local optima is not as serious issue as we might thought when the network is large, since many critical points are actually saddle points rather than local minima.

Successful implementations of the above RNN structures include natural language translation such as [2] where an encoder-decoder structure is used, each is a LSTM. Another example is image captioning, where the AI tries to explain what is in an image using texts. A solution to this is to use CNN to encode the image, and use LSTM to decode to generate texts. Following similar ideas is hand-writing reorganization.

1.3 Attention Mechanism

“nobreak

1.4 Transformer

While CNN is mainly used for spatial data processing such as computer vision, both RNN and transformer are mainly used for sequential data processing. It is worth mention in the very beginning that transformer does not guarantee superior performance in all scenarios comparing with traditional RNN-based natural language processing models such as LSTM.

For one thing, transformer consumes larger computational capabilities. Transformer processes data in batch, while RNN does them in sequence, meaning that RNN can response faster in some real-time applications. It is difficult for the transformer to process super-long text due to the computational burden (quadratic computational complexity with respect to sequence length), while RNN can process arbitrarily long text, although the later suffers from capturing long-range dependencies.

With the above been said, the transformer does demonstrated superior performance than RNN in many tasks. The mechanism and the reasons why it performs better in these occasions are introduced as follows.

1.5 Research Trends





Part II

Products



Many data companies, including Microsoft, Google, Meta (formerly known as Facebook) and Amazon, provide their own data models and AI engines of different kinds. Famous and AI engines and models include TensorFlow by Google, PyTorch by Meta, and ChatGPT by OpenAI, which is a company supported by Microsoft.

These data companies usually open (part of) the source code of these models and engines in the community and encourage individual developers to learn and use them. At the same time, they provide enterprise level services such as model fine-tuning to make profit from the tools.

Some state-of-art models, engines, as well as the business models of these companies are reviewed in this part of the notebook.



2

OpenAI

CONTENTS

2.1	General Introduction	15
2.1.1	Background	15
2.1.2	Business Model	15
2.2	ChatGPT	15
2.3	OpenAI API	15
2.4	DELL-E-2	16

“nobreak

2.1 General Introduction

“nobreak

2.1.1 Background

“nobreak

2.1.2 Business Model

“nobreak

2.2 ChatGPT

“nobreak

2.3 OpenAI API

“nobreak

2.4 DELL-E-2

3

Meta AI

CONTENTS

3.1	General Introduction	17
3.1.1	Background	17
3.1.2	Business Model	17
3.2	PyTorch	17
3.3	LLaMA	17
3.3.1	LLaMA: Open and Efficient Foundation Language Models	18

Meta, formerly known as Facebook, provides variety of solutions for AI, including computer vision, stream data processing, natural language processing, and many more.

3.1 General Introduction

“nobreak

3.1.1 Background

“nobreak

3.1.2 Business Model

“nobreak

3.2 PyTorch

“nobreak

3.3 LLaMA

LLaMA is the large language model published by Meta AI. The associated paper that introduced LLaMA for the first time was published in February 2023 on ArXiv as [3]. The initial LLaMA is trained on publicly available datasets exclusively.

Different levels of complexity of the models are provided. Like many other LLMs, Meta AI provides charged service for fine-tuning the model. There is also a re-write of LLaMA, namely Lit-LLaMA, that is completely open source under Apache 2.0 license. Both LLaMA and Lit-LLaMA can be fine-tuned via tools such as LoRA and Stanford Alpaca.

Comparing with other models such as ChatGPT-3, LLaMA is claimed to use smaller models (with less number of parameters) to achieve the same performance.

3.3.1 LLaMA: Open and Efficient Foundation Language Models

Bibliography

- [1] Zachary C Lipton, John Berkowitz, and Charles Elkan. A critical review of recurrent neural networks for sequence learning. *arXiv preprint arXiv:1506.00019*, 2015.
- [2] Ilya Sutskever, Oriol Vinyals, and Quoc V Le. Sequence to sequence learning with neural networks. *Advances in neural information processing systems*, 27, 2014.
- [3] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models, 2023.