

**Imperial College
London**

IMPERIAL COLLEGE LONDON

DEPARTMENT OF MATHEMATICS

**Optimizing Transformer Neural Network
for Real-Time Outlier Detection**

Author: Ilia Sobakinskikh (CID: 02273801)

A thesis submitted for the degree of

MSc in Mathematics and Finance, 2022-2023

Declaration

The work contained in this thesis is my own work unless otherwise stated.

Acknowledgements

I would like to thank my supervisor, Dr Paul Bilokon, for his guidance and support throughout the project. He kept me motivated and helped me to stay on track.

I would like to thank the Department of Mathematics at Imperial College London for providing me with the opportunity to study at the institution.

Abstract

In this work, we explore how the inference time of a Transformer Neural Network can be efficiently optimized with applications to real-time anomaly detection in financial time series. The financial time series are price series such as asset prices. Unfortunately, the data is often with errors or outliers that make the downstream data processing tasks useless, unstable or even harmful [1] [2]. Moreover, the amount of financial time-series data has been significantly increasing [3]. Hence, there is a need for better data-cleaning methods in terms of accuracy and in terms of processing speed.

Transformers as a neural network architecture have achieved superior performances in many tasks such as Natural Language Processing and Computer Vision [4]. Time series modelling and especially anomaly detection tasks can benefit from the features of transformers architecture in multiple ways, including the capacity to capture long-range dependencies and interactions [5].

Increasingly powerful hardware, such as field-programmable gate arrays (FPGAs), have seen increasing usage in recent years due to their reconfigurability and high performance [6]. They can be efficiently utilized to speed up the computations of the Transformer architecture.

We explore different Transformer architectures for time series modelling and how they can be efficiently implemented on an FPGA board (PYNQ-Z2). In particular, we examine the application of Transformers to detect anomalies in time series and we show how they can be efficiently implemented on an FPGA board to minimize latency.

The code is available at https://github.com/thxi/icl_thesis

Contents

1	Problem statement	6
1.1	Definitions and examples	6
1.2	Different Anomalies Types	7
2	Transformers	9
2.1	General architecture	9
2.1.1	Positional encoding	9
2.1.2	Encoder Block	10
2.2	Attention mechanism	11
2.2.1	Dot-Product Attention and Multi-Head Attention	11
2.2.2	Linear attention	12
2.3	Transformers for time series modelling	13
3	FPGA design	14
3.1	Introduction to FPGA programming and development	14
3.1.1	Common Terms	14
3.1.2	Common optimizations	16
3.1.3	Example of optimizing matrix multiplication	18
3.2	Transformer architecture optimizations	21
3.2.1	Vanilla Transformer	21
3.2.2	Linear transformer	22
4	Experiments	23
4.1	Architecture and hyperparameters	23
4.1.1	Model Fitting	23
4.2	Datasets	27
4.2.1	Numenta Anomaly Benchmark (NAB)	27
4.2.2	KPI Anomaly Detection Dataset	27
4.2.3	FI2010	28
4.3	Accuracy	28
4.3.1	Metrics	29
4.3.2	Comparison	31
4.4	Performance/Speed	31
5	Conclusion	33
A	Training plots	34
B	Empirical latency	40
	Bibliography	45

List of Figures

1.1	Sample mutli-variate time-series of Apple stock price. The blue line are the close prices and the red line is the differenced close prices.	6
1.2	Apple stock modified time series with injected Point anomalies . Top panel: price time series and its distribution. Bottom panel: differenced price time series and its distribution.	8
2.1	Model architecture of the Transformer [7]	10
2.2	Positional encoding example embeddings for a sequence of max length of 100 and embedding dimension of 48. The embeddings are generated using the sinusoidal positional encoding technique. The columns represent the embeddings that we would add to the input vector x_i	11
2.3	Scaled Dot-Product Attention and Multi-Head Attention [7]	12
3.1	Latency vs Initiation Interval illustration. Source: [8]	15
3.2	Loop pipelining illustration. Source: [9]	16
3.3	Array partitioning illustration.	18
3.4	Naive matrix multiplication synthesis report	19
3.5	Matrix multiplication with pipelined loop <code>loop_k</code>	19
3.6	Matrix multiplication with pipelined loop <code>loop_j</code>	20
3.7	Matrix multiplication with pipelined loop <code>loop_i</code>	20
3.8	Matrix multiplication with loop reordering	21
4.1	Example of gradient explosion at around epoch 40 which leads to the loss function diverging for a few following epochs.	26
4.2	NYC Taxi demand - anomalies highlighted in red. Left: Time series of observations, Right: Distribution of observations.	27
4.3	Sensor data from a machine in a data center. The red dots indicate the anomalies. Left: Time series of observations, Right: Distribution of observations.	28
4.4	Example of the injected outliers in the FI2010 dataset.	29
4.5	Confusion matrix description	30
A.1	Training metrics for different epochs for Transformer model on KPI dataset	34
A.2	Training metrics for different epochs for Transformer model on NAB dataset . . .	35
A.3	Training metrics for different epochs for Linear Transformer model on KPI dataset	36
A.4	Training metrics for different epochs for Linear Transformer model on NAB dataset	37
A.5	Training metrics for different epochs for Linear Regression model on KPI dataset .	38
A.6	Training metrics for different epochs for Linear Regression model on NAB dataset	39
B.1	Empirical latency for Transformer model on CPU.	40
B.2	Empirical latency for Transformer model on GPU.	41
B.3	Empirical latency for Linear Transformer model on CPU.	42
B.4	Empirical latency for Linear Transformer model on GPU.	43

List of Tables

4.1	Hyperparameters	23
4.2	Final metrics. The reported value are for the train and validation set.	31
4.3	Inference time for the models on different devices. For CPU and GPU we report the mean and standard deviation. For FPGA, the deterministic latency is reported (see Section 3.2).	31
4.4	Relative resource utilization for the PYNQ-Z2 board.	32

Introduction

The rapid growth of financial time series data has given rise to significant challenges in maintaining data quality and processing speed. In this study, we address these challenges by proposing an implementation of Transformer-based Neural Networks which achieve good performance in anomaly detection tasks while having a low inference time.

Financial time series data often suffer from data inaccuracies, errors, and outliers, rendering downstream data processing tasks ineffective or even detrimental to decision-making processes. The urgency of this issue is underscored by the ever-increasing volume of financial time-series data available. Consequently, there is a pressing demand for more accurate and faster data-cleaning methods.

Transformers, as a powerful neural network architecture, have consistently demonstrated exceptional performance across various domains, including Natural Language Processing (NLP) and Computer Vision. Leveraging the capabilities of Transformer architectures for time series modeling, especially in the realm of anomaly detection, offers numerous advantages, including the ability to capture long-range dependencies and intricate interactions within the data.

The evolution of hardware, for example, the growing utilization of field-programmable gate arrays (FPGAs), has brought about increased computational capabilities due to their reconfigurability, high performance and deterministic latency. These FPGAs can be effectively used to accelerate computations associated with Transformer architectures.

Outline of the thesis.

In Chapter 1, we will describe the problem statement of anomaly detection task. In Chapter 2, we will describe the main concepts and ideas used in the Transformer architecture. In Chapter 3, we will describe the main concepts of programming an FPGA and the specific optimizations that can be applied to speed up the computations. In Chapter 4, we will describe the experiments that we conducted to evaluate the performance of the proposed architectures. Chapter 5 will conclude the thesis and will provide some ideas for future work.

Contributions

The main contributions of this paper are as follows:

1. Implementation of vanilla Transformer architecture on an FPGA board.
2. Implementation of a Linear Attention Transformer on an FPGA board.
3. Analysis of the latency/resource utilization tradeoff for the proposed architectures.
4. Comparative analysis of the performance, related to accuracy, of the proposed architectures on the anomaly detection task.

Chapter 1

Problem statement

In this chapter, we will describe the problem statement of anomaly detection task.

1.1 Definitions and examples

Definition 1.1.1. We consider a **time-series** \mathcal{T} which is simply a timestamped sequence of observations $x_i \in R^n$, potentially with some metadata y_i .

Most of the times we will consider univariate case, i.e. $n = 1$. An example of this is a price time-series of a single stock. However, the multivariate case is also important and we will consider it in the experiments. For example, one can consider a time-series of prices of multiple stocks to get a multivariate time series or one can also create features (for example, by calculating the **differenced time series**) on the original time-series to get a multivariate time-series.

Definition 1.1.2. The **differenced time series** is a time-series \mathcal{T} where each observation \hat{x}_i is simply the difference between the current and the previous observation. That is, $\hat{x}_i = x_{i-1} - x_i$.

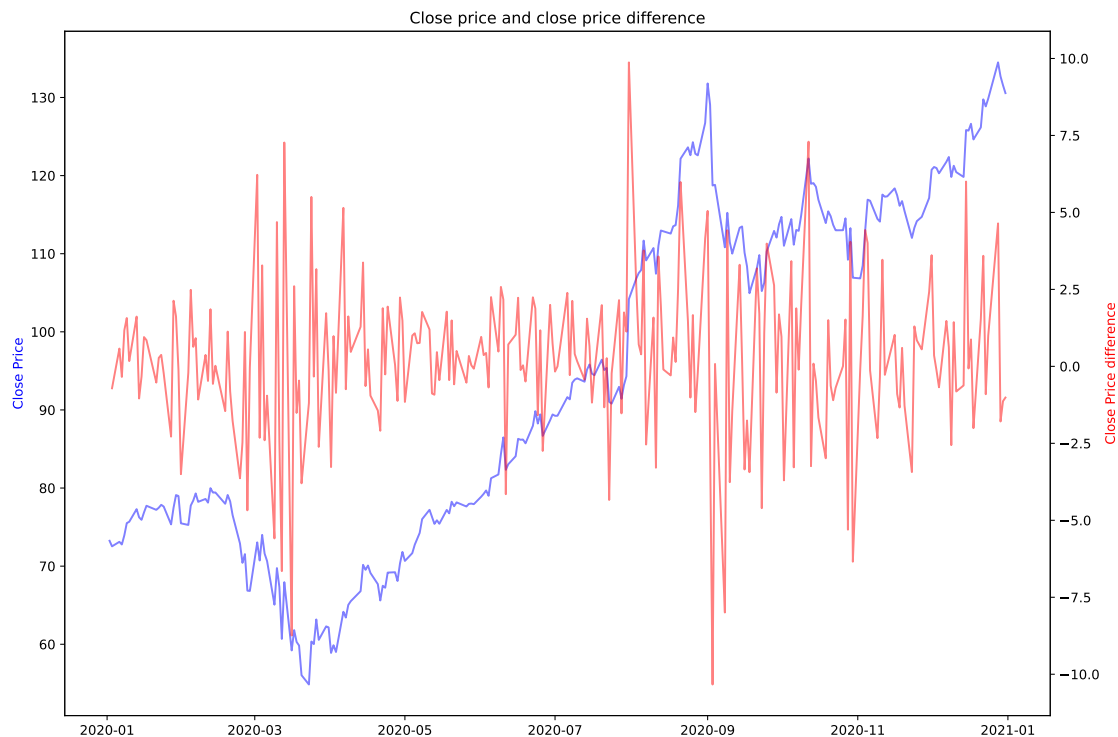


Figure 1.1: Sample mutli-variate time-series of Apple stock price. The blue line are the close prices and the red line is the differenced close prices.

Definition 1.1.3. The **Anomaly Detection** task: for any time-series $\hat{\mathcal{T}}$ of length n , we need to predict $\mathcal{Y} = \{y_1, \dots, y_n\}, y_i \in \{0, 1\}$, whether the datapoint at the i -th timestamp anomalous (where by convention we will use 1 as anomaly and 0 as not an anomaly).

In this work, we will restrict ourselves to the **supervised case** where the labels y_i are known for the *seen* (or training) part of the dataset.

Remark 1.1.4. One can also consider an unsupervised task. However, one issue with the unsupervised task is that it is hard to evaluate the performance (i.e., accuracy) of the model’s predictions [10].

1.2 Different Anomalies Types

Although there is no universally accepted definition of an anomaly, there exists a common classification of *synthetic* anomalies into different types. For a broad comprehensive review, reader is referred to [11], [12].

In this work, we will only consider **Point anomalies** types since the main focus of this work is not to show the superior performance of the Transformer architecture on different types of anomalies but to show that it can be efficiently implemented on an FPGA board.

Definition 1.2.1. **Point anomalies** are individual data points that are considered anomalous with respect to the rest of the data. For example, a point anomaly can be a sudden spike in the price of a stock. See Figure 1.1.

Mathematically, we could define a point anomaly as follows:

$$x_i = \begin{cases} x_i^{\text{original}} + S & y_i = 1 \\ x^{\text{original}} & \text{otherwise} \end{cases} \quad (1.2.1)$$

where x^{original} is the original time-series (as if there was no outliers) and S is the spike.

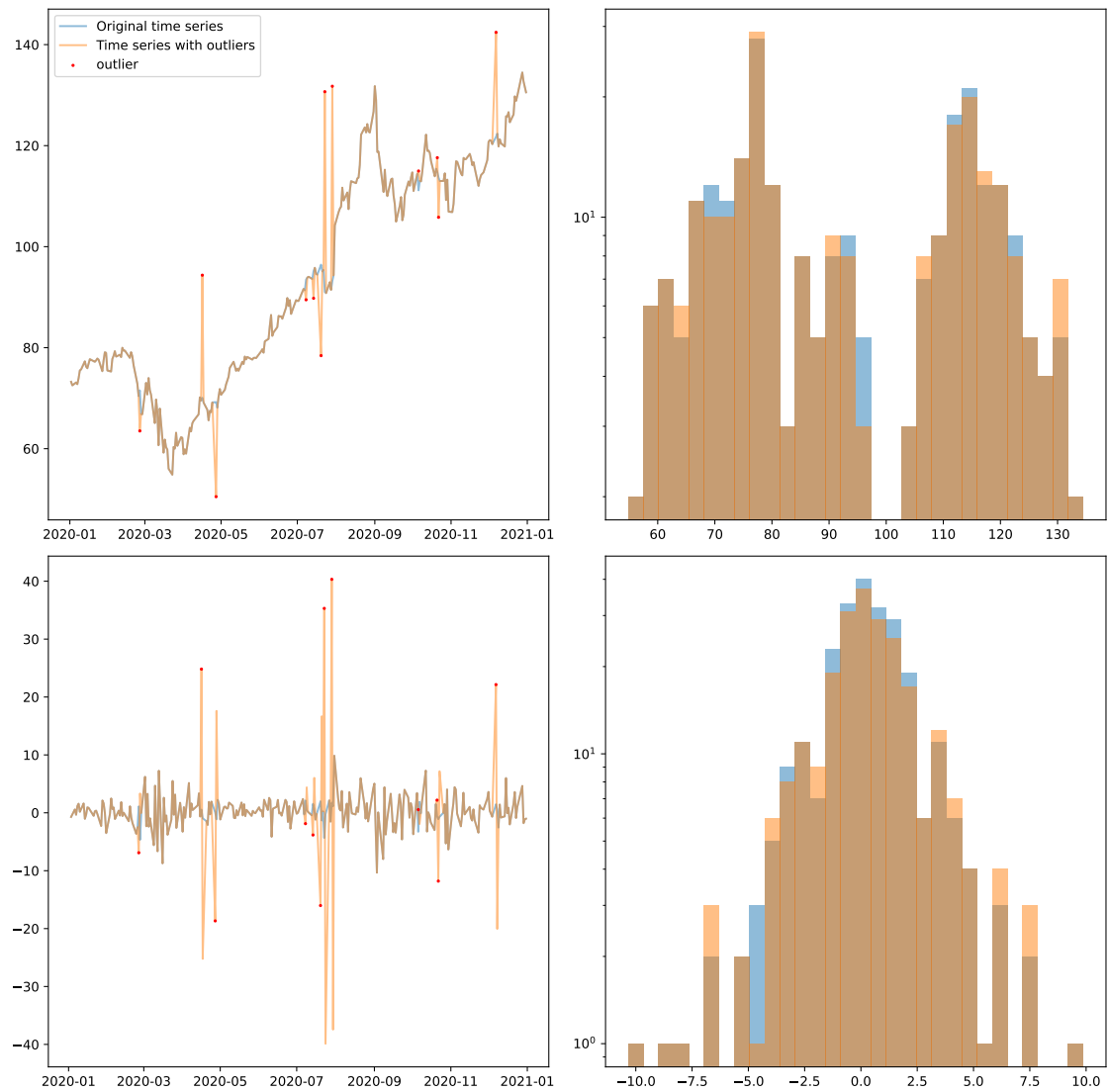


Figure 1.2: Apple stock modified time series with injected **Point anomalies**. Top panel: price time series and its distribution. Bottom panel: differenced price time series and its distribution.

Chapter 2

Transformers

In this chapter, we will describe the main concepts and ideas used in the Transformer architecture. and how to use them for anomaly detection time series task. We will describe the main building blocks of the Transformer architecture and will give a special treatment to the attention mechanism firstly introduced in [13].

2.1 General architecture

In [7], authors introduced the *vanilla* Transformer architecture which a neural network architecture which is the architecture that is dominantly used in Natural Language Processing tasks. The architecture's main feature was reliance on the attention mechanism and the complete elimination of recurrent and convolutional layers.

Figure 2.1 presents the main architecture of the Transformer. The architecture consists of an **encoder** and a **decoder** blocks. For the purpose of the thesis we will only consider the **encoder** part of the architecture (and in that regard the idea is similar to BERT model for language processing [14] which only uses a transformer encoder and this work which uses BERT for anomaly detection [15]). The **encoder** is preceded by a **positional encoding** layer which is used to *inject* the positional information to the input vectors x_i because the attention mechanism is permutation invariant, this will be explained in Section 2.2.

2.1.1 Positional encoding

Since the attention mechanism is permutation invariant (i.e., it does not take into account the order of the input, refer to Section 2.2), we need to inject the positional information into the input vectors x_i . In short, a positional encoding is usually a vector of the same dimension as the input vectors x_i which can either be fixed (for example, we have the same vector for different positions) or learned.

For the comprehensive overview of different positional encoding techniques, reader is referred to [16] and [17].

The vanilla Transformer architecture uses a sinusoidal positional encoding technique which is a fixed positional encoding scheme. Suppose we have a time-series of length T . Then the positional encoding for the i -th input vector x_i of dimension d is defined as follows:

$$\begin{aligned} PE_{i,2j} &= \sin\left(\frac{i}{10000^{2j/d}}\right) \\ PE_{i,2j+1} &= \cos\left(\frac{i}{10000^{2j/d}}\right) \end{aligned} \tag{2.1.1}$$

Refer to Figure 2.2 for an example of the embedding vectors that are obtained using this method.

This combination of sine and cosine functions allows the positional encoding to generate the embedding of a position in a unique way. The embedding is then added to the input vector x_i , however, it can also be simply stacked with the input vector x_i to get a vector of dimension $2d$. In this work, we will use the former approach.

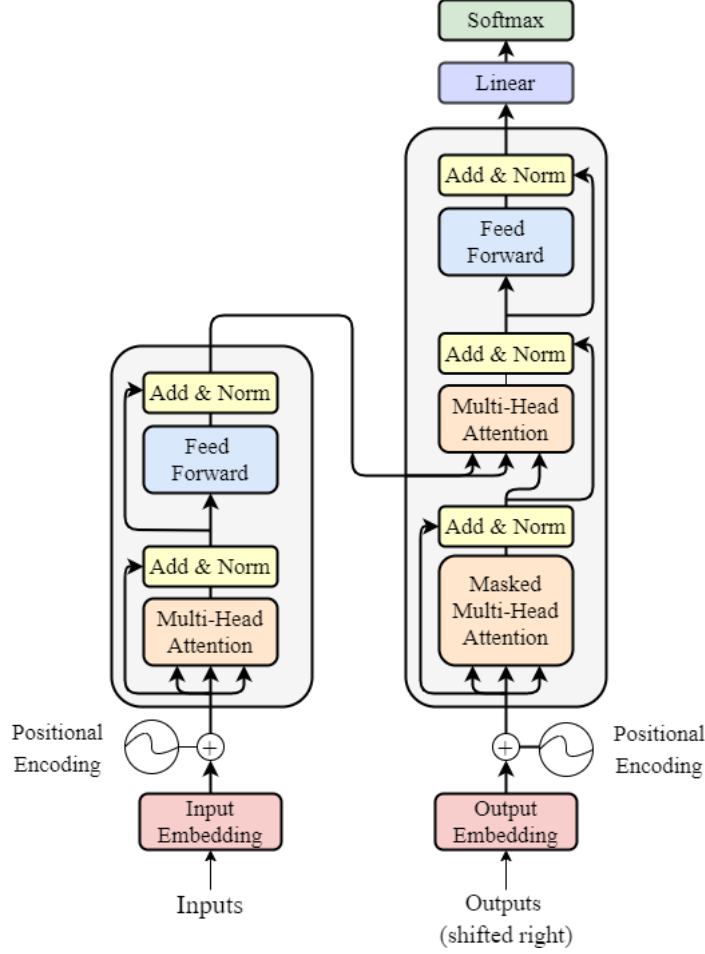


Figure 2.1: Model architecture of the Transformer [7]

2.1.2 Encoder Block

The encoder consists of N identical layers. Each layer has two sub-layers which are a **multi-head self-attention** layer (described in detail in Section 2.2.1) and a **feed-forward (FF)** layer. The **feed-forward (FF)** layer $\text{FFN}(\cdot)$ is a simple neural network which comprises of 2 fully-connected (FC) linear layers and an activation function in between them. Specifically, authors of [7] used an FC layer with the ReLU activation function

$$\text{ReLU}(x) = \max(0, x)$$

followed by another FC layer without activation function, i.e.

$$\text{FFN}(x) = W_2 \cdot \text{ReLU}(W_1 \cdot x + b_1) + b_2$$

where W_1, W_2, b_1, b_2 are the weight matrices and bias vectors respectively for the FC layers.

The **Add & Norm** layer is a residual connection [18] followed by a layer normalization layers [19]. That is,

$$\text{Add \& Norm}(x) = \text{LayerNorm}(x + \text{Sublayer}(x))$$

where $\text{Sublayer}(x)$ is either the multi-head self-attention layer or the feed-forward layer (refer to Figure 2.1).

Residual connections (i.e., adding the input to the output of the layer), also known as skip connections, are commonly used to avoid the vanishing gradient problem [18]. In this case, the residual connection is used to avoid the degradation of the model performance when the model is deep (i.e., when the number of sequential layers N is large). While we do not use very deep models in this work, we still use residual connections as in the original architecture for more stable training (refer to Section 4.1.1). Moreover, the residual connection is used to propagate

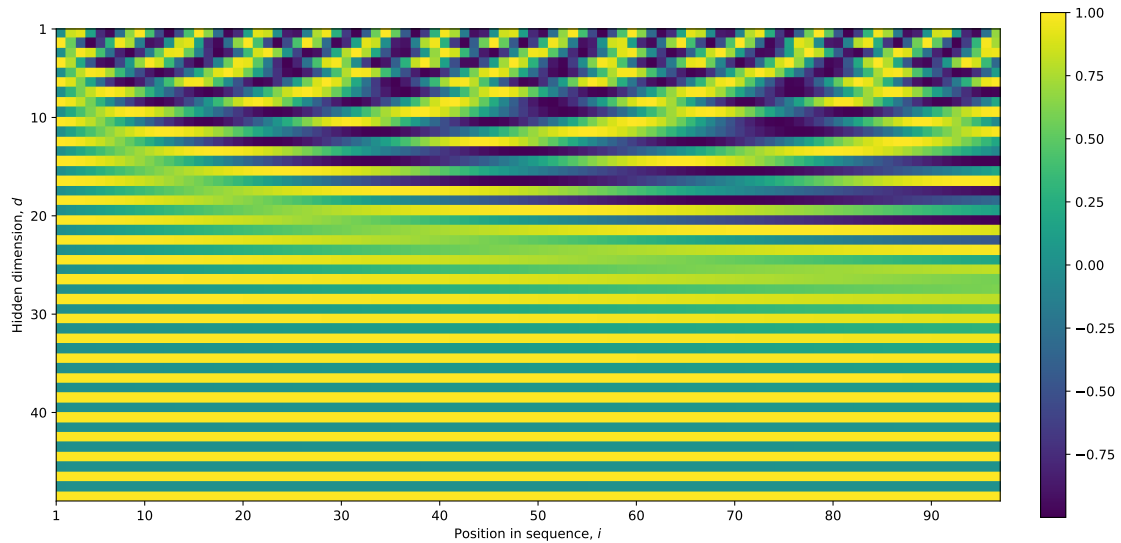


Figure 2.2: Positional encoding example embeddings for a sequence of max length of 100 and embedding dimension of 48. The embeddings are generated using the sinusoidal positional encoding technique. The columns represent the embeddings that we would add to the input vector x_i .

the positional information to the next layers as the attention mechanism is permutation invariant (refer to Section 2.2).

The layer normalization layer [19], LayerNorm is used to normalize the output of the residual connection. The normalization is used to stabilize the training process (however, it is not strictly necessary, refer to Section 4.1.1 for our empirical findings).

This constitutes the main building block of the Transformer encoder architecture. The next section will describe the attention mechanism in detail.

2.2 Attention mechanism

This section will describe the attention mechanism, its variations and the intuition behind it. Moreover, we will compare different attention mechanism implementations in terms of their computational complexity which is important for fast calculations. Most importantly, we will describe the **multi-head attention** mechanism which is used in the vanilla Transformer architecture as a continuation of Section 2.1.

2.2.1 Dot-Product Attention and Multi-Head Attention

The attention mechanism introduced in the Transformer architecture [7] used a **scaled dot-product attention**.

The main idea of the **dot-product attention** mechanism is to compute the mapping of a query q_i for each input vector x_i to a set of key-value pairs (k_j, v_j) . The query q_i , key k_i and value v_i vectors are simply linear transformations of the input vectors x_i , i.e., $q_i = W_Q \cdot x_i$, $k_i = W_K \cdot x_i$, $v_i = W_V \cdot x_i$ where W_Q , W_K and W_V are the weight matrices. The attention mechanism is a weighted sum of the values v_j where the weights are computed as a function of the query q_i and the key k_j . That is, $Attention(x_i) = \sum_j \alpha_{ij} v_j$ where $\alpha_{ij} = \text{softmax}(q_i \cdot v_i)$ is the weight of the j -th value v_j . In practice, the attention mechanism is computed for all the queries q_i at the same

time by utilizing the following expression in matrix form:

$$\begin{aligned} Q &= W_Q \cdot X, \\ K &= W_K \cdot X, \\ V &= W_V \cdot X \end{aligned} \tag{2.2.1}$$

$$\text{Attention}(Q, K, V) = \text{softmax}(QK^T \cdot \frac{1}{\sqrt{d_k}})V \tag{2.2.2}$$

Figure 2.1 visualizes the attention mechanism introduced in [7].

Remark 2.2.1. In [7], authors additionally scaled the weights α_i by the square root of the dimension of the key vectors d_k . This is, however, not strictly necessary and is done for numerical stability reasons.

The **multi-head attention** mechanism is simply a concatenation of multiple attention mechanisms. That is, we can compute h different attention mechanisms in parallel and then concatenate the results. The main idea behind this is that different attention mechanisms can learn different features of the input vectors. After the concatenation, we apply a linear transformation (i.e., a fully-connected layer) to the concatenated vectors to get the final output with a desired dimensionality which is usually the same as the input vectors x_i .

However, in the main paper [7], authors implemented the multi-head attention mechanism in a more efficient way compared to the naive concatenation of multiple attention heads of size d_k . They, instead, linearly transformed the input vectors to the dimension of $d_k h$ where h is the number of attention heads and applied the attention mechanism to the transformed vectors. This is more efficient because we require less parameters and can benefit from the increased accuracy of having multiple attention heads.

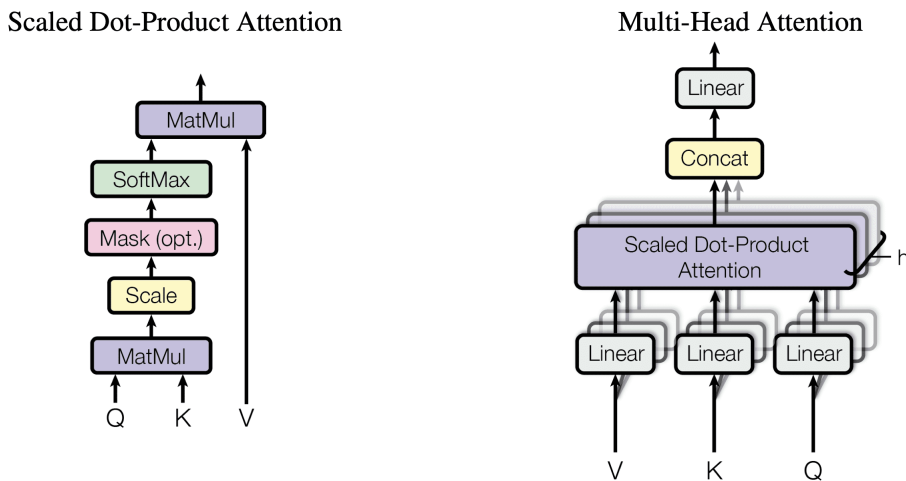


Figure 2.3: Scaled Dot-Product Attention and Multi-Head Attention [7]

2.2.2 Linear attention

In [20], authors propose an extension to the dot-product attention mechanism called **linear attention**. This significantly reduces the computational complexity of the attention mechanism by eliminating the need to compute the softmax function.

Notice that in 2.2.2, the softmax function is applied rowwise to the matrix QK^T . The softmax function can be substituted with a general similarity function $\text{sim}(\cdot, \cdot)$ between a query q_i and a key k_j . The equation 2.2.2 for output value v'_i can then be rewritten as follows:

$$v'_i = \text{Attention}(Q, K, V)_i = \frac{\sum_j \text{sim}(q_i, k_j) v_j}{\sum_j \text{sim}(q_i, k_j)} \tag{2.2.3}$$

The only constraint imposed on the similarity function $\text{sim}(\cdot, \cdot)$ is that it should be non-negative for it to define an attention function. This conveniently includes all kernels. That is $\text{sim}(q_i, k_j) = \phi(q_i)^T \phi(k_j)$ where $\phi(\cdot)$ is a feature map.

So that given a kernel with a feature map $\phi(\cdot)$, the attention mechanism can be computed as follows:

$$v'_i = \text{Attention}(Q, K, V)_i = \frac{\sum_j \phi(q_i)^T \phi(k_j) v_j}{\sum_j \phi(q_i)^T \phi(k_j)} \quad (2.2.4)$$

And we can rewrite the attention mechanism in matrix form as follows:

$$\text{Attention}(Q, K, V) = \frac{\phi(Q)^T \phi(K) V}{\phi(Q)^T \phi(K)} \quad (2.2.5)$$

Regrouping the terms, we get the following expression for the attention mechanism:

$$\text{Attention}(Q, K, V) = \phi(Q)^T \frac{\phi(K) V}{\phi(Q)^T \phi(K)} \quad (2.2.6)$$

which makes it evident that we can compute $\sum_j \phi(k_j) v_j$ once and reuse them for all the queries q_i which reduces the computational complexity from $O(N^2)$ to $O(N)$ where N is the number of input vectors in the attention layer.

Remark 2.2.2. In [20], authors used the $\phi(x) = \text{elu}(x) + 1$ feature map where $\text{elu}(x) = \max(0, x) + \min(0, \alpha(\exp(x) - 1))$ is the exponential linear unit activation function. This feature map is used to ensure that the attention mechanism is non-negative and hence defines a valid attention function. Moreover, $\text{elu}(\cdot)$ is used instead of $\text{ReLU}(\cdot)$ to ensure the differentiability when x is negative.

2.3 Transformers for time series modelling

Originally, the Transformer architecture was introduced for Natural Language Processing tasks [7] but has since found applications in other domains and time-series modelling specifically [5]. The motivation behind using the Transformer architecture for time series modelling is that they show good performance in many sequence modelling tasks. While many specialized Transformer architectures were proposed for time series modelling (e.g., [21], [22]) and specifically deployed in asset management setting [23], in the current work, we will use the least tuned architecture, i.e., the vanilla Transformer architecture and the architecture with linear attention [20]. The reason for using the not so specialized architectures is that we want to show that the Transformer can achieve good out of the box performance on the anomaly detection task with minimal tuning and that it can be efficiently implemented on an FPGA board.

Chapter 3

FPGA design

In this section, the main concepts of programming an FPGA will be introduced and the specific optimizations that can be applied to speed up the computations.

Readers will be introduced to the common optimization techniques and how they are achieved. A matrix multiplication example will be provided to illustrate the concepts and how the optimizations affect the latency and resource utilization.

Lastly, the analysis of the specific optimizations that can be applied to the Transformer architecture will be provided.

3.1 Introduction to FPGA programming and development

In this section, we will introduce the main concepts of FPGA programming and development tools.

3.1.1 Common Terms

In this section, common terms will be introduced. The terms will be used throughout Chapter 3. It is not required to read all of them at once, but it is recommended to refer to this section when a term is not clear (i.e., only when necessary).

Definition 3.1.1. LUT (Look-Up Table) is a small, fast memory that stores the output of a Boolean function of its inputs. The LUT is the basic building block of an FPGA and is capable of implementing any logic function of N Boolean variables.

Definition 3.1.2. BRAM (Block RAM) is a dedicated two-port memory that can be used to store data.

Definition 3.1.3. DSP (Digital Signal Processing) is a specialized hardware unit that is optimized for performing mathematical operations.

Definition 3.1.4. Clock cycle is the time between two consecutive rising edges of the clock signal. It is the amount of time between two pulses of an oscillator, a single increment of the central processing unit (CPU) clock during which the smallest unit of processor activity is carried out.

Definition 3.1.5. Latency is the time between the start of an operation and the moment its results become available or the number of clock cycles required to complete an operation. **Latency of a loop** is the number of clock cycles required to complete one iteration of the loop.

Definition 3.1.6. Throughput is the number of operations that can be completed in a given amount of time.

Definition 3.1.7. Initiation Interval (II) is the number of clock cycles between the start of two consecutive iterations of a loop. That is, it is the maximum rate (in clock cycles) at which loop iterations can be initiated. In the ideal case, the II is equal to 1 so that we can start a new iteration of the loop every clock cycle. Initiation interval is different from latency. The reason for this is pipelining which will be described in Section 3.1.2. For a visual explanation, see Figure 3.1.

Definition 3.1.8. Trip count is simply the number of iterations of a loop.

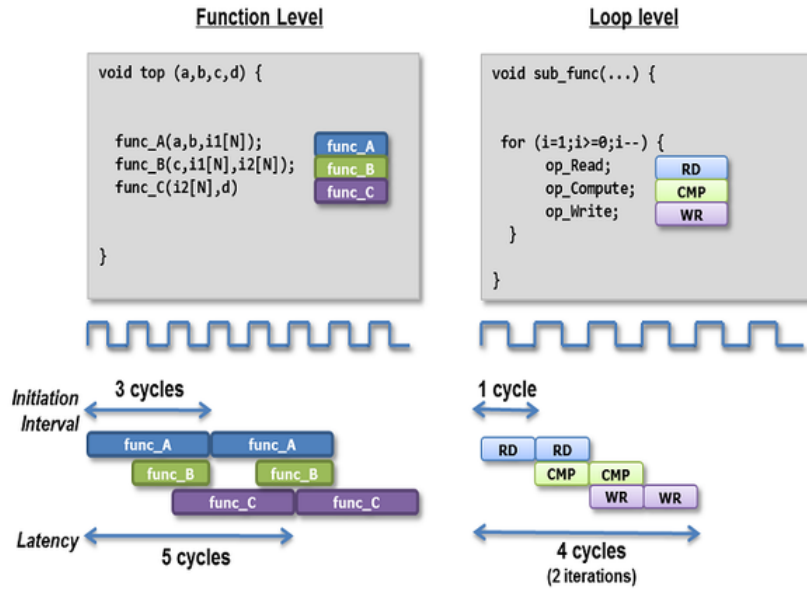


Figure 3.1: Latency vs Initiation Interval illustration. Source: [8]

HLS synthesis

In this section, HLS synthesis will be described [24]. It is now the common workflow in the FPGA development because it significantly improves the productivity when working with design.

HLS synthesis is a methodology that bridges the gap between high-level programming languages, such as C, C++, or OpenCL, and the low-level hardware description languages typically used in FPGA design, like Verilog or VHDL. This approach enables developers to describe their algorithms and functionality at a higher level of abstraction, allowing them to focus on the problem-solving aspect rather than the intricacies of hardware implementation.

The HLS process starts with a high-level description of the desired functionality. This description can be written in familiar programming languages (in this case, C++ HLS), taking advantage of their abstractions and concise syntax. The HLS tool then performs a series of transformations on this high-level description to generate an optimized hardware implementation that can be deployed on an FPGA. The high-level description is transformed by the HLS tool into a RTL (Register Transfer Level) representation, which is the low-level hardware description that defines the behavior of the FPGA.

Simulation, Cosimulation

In this section, the processes of **simulation** and **cosimulation** will be described. In the context of developing for Field-Programmable Gate Arrays (FPGAs), simulation and cosimulation are two crucial techniques for verifying and testing the functionality of your design before actually programming it onto the FPGA hardware.

Definition 3.1.9. Simulation is the process of running a software-based model of your FPGA design on a computer to simulate its behavior. The process is very similar to running a software program on a computer for the purpose of unit testing certain parts of functionality of your code [25]. That is, the simulation does not involve any RTL code and is simply a software simulation of the high-level description.

Remark 3.1.10. Simulation might not always capture all aspects of hardware behavior, such as timing delays, which can be critical on FPGAs.

Definition 3.1.11. Cosimulation is a technique that combines simulation of the high-level description with simulation of the generated RTL description. This means that the simulations of both the original high-level code and the RTL representation in parallel, comparing their behavior. The purpose of cosimulation is to ensure that the high-level synthesis tool accurately transformed the high-level description into the desired RTL behavior. [25].

3.1.2 Common optimizations

In this section, common optimization techniques and how they are achieved will be introduced.

It is quite common to process data blocks (for example, a sequence of samples in anomaly detection) using for loops. For loops are usually the main bottleneck in the performance of the design and it is the area where most of the optimizations are applied first [24].

Loop Pipelining

Loop pipelining is a technique used in FPGAs programming to optimize the performance of sequential operations within a loop. It improves the throughput of loops by breaking them down into multiple stages that can execute concurrently. That is, it allows to start the next iteration of a loop before the current iteration has finished [9]. Refer to Figure 3.2 and Figure 3.1 for a visual illustration.

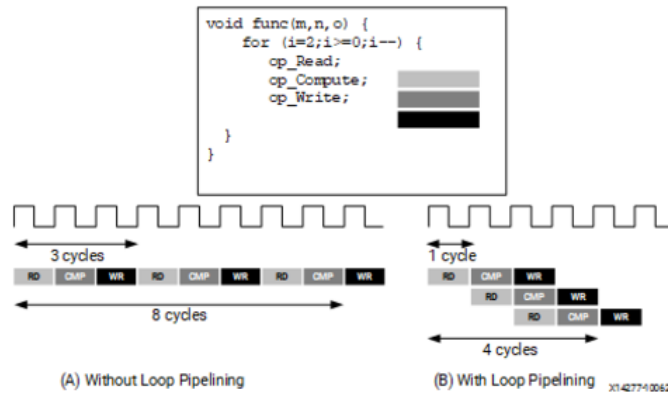


Figure 3.2: Loop pipelining illustration. Source: [9]

Consider example of pipelining a simple for loop which adds 2 vectors ([9]):

```
void toplevel(int* a, int* b, int* c, int len) {
    vadd: for(int i = 0; i < len; i++) {
#pragma HLS PIPELINE
        c[i] = a[i] + b[i];
    }
}
```

Taking the reference numbers from [9], assume that len is 20 and that one loop iteration takes 3 clock cycles, then the total latency of the loop is 60 clock cycles without pipelining. The pipelining pragma `#pragma HLS PIPELINE` allows to start the next iteration of the loop before the current iteration has finished. By default, the pipelining pragma will try to achieve an II of 1 but this can be specified manually by using the II parameter. So, the pipelining pragma reduces the latency of the loop to 22 clock cycles.

In general, the latency of a loop with pipelining is given by the following formula [26]:

$$\text{Latency} = \text{II} \cdot (\text{Trip Count} - 1) + \text{Loop Body Latency} \quad (3.1.1)$$

Remark 3.1.12. It is not always possible to achieve an II of 1. This is because of the dependencies between the iterations of the loop. For example, if the loop body depends on the result of the previous iteration, then the II cannot be 1 and we have to wait until the previous iteration has finished before starting the next one.

Loop Unrolling

Loop unrolling is an optimization technique that involves expanding or unwinding loops in code to potentially enable better utilization of hardware resources and/or to minimize control flow (branching) in loop iterations. This technique is not limited to FPGA development but can be particularly useful in optimizing code for FPGA implementations.

Loop unrolling works by duplicating loop iterations (read as `copy-pasting`). This allows utilizing more hardware as the loop body is duplicated multiple times and loop iterations will

utilize different hardware resources. This increase in performance (i.e., throughput) comes at the cost of increased resource utilization [9].

Consider a simple example of a function which multiplies the input vector of length 4 by a constant, 2:

```
void toplevel(int* a, int* b) {
    smult: for(int i = 0; i < 4; i++) {
#pragma UNROLL
        b[i] = 2 * a[i];
    }
}
```

The unroll pragma `#pragma UNROLL` allows to unroll the loop and execute the loop iterations in parallel.

The unrolled version of the loop will be equivalent to the following code:

```
void toplevel(int* a, int* b) {
    b[0] = 2 * a[0];
    b[1] = 2 * a[1];
    b[2] = 2 * a[2];
    b[3] = 2 * a[3];
}
```

Remark 3.1.13. It might not be possible to unroll a loop. For example, if the trip count is not known at compile time then the loop cannot be unrolled. Sometimes it is possible to unroll a loop partially. The unroll pragma allows to specify the factor of unrolling, i.e., how many iterations to unroll. For example, `#pragma UNROLL factor=2` will only duplicate the loop body so that there are 2 of them.

Loop Reordering

Loop reordering optimization is an optimization used to improve the performance by changing the order in which loops are executed. This optimization is not strictly related to FPGA development and it involves altering the nesting order of loops in a way that improves data locality, cache utilization (for example, on modern CPU) and enables usage of SIMD resources.

In the context of FPGA, loop reordering can be used achieve better II when we are dealing with pipelining and nested loops (see Section 3.1.2 and Section 3.1.3 for an example).

Function Inlining

Function inlining optimization technique is not specific to FPGA development only. The technique is used to improve the performance of a program by reducing the overhead associated with function calls. Calling a function incurs some overhead in terms of memory and execution time due to the need to set up the function call stack, pass arguments, and jump to the function's code. Function inlining aims to eliminate this overhead by replacing a function call with the actual body of the function at the call site. In other words, the compiler takes the contents of the called function and inserts it directly into the location where the function is called. This, however, increases the size of the code and, specifically in the case of FPGAs, the resource utilization (LUTs and FF).

Array Partitioning and Reshaping

Partitioning arrays in an FPGA (Field-Programmable Gate Array) refers to the process of dividing a large memory block (e.g., one array) into smaller sections, often referred to as memory banks or partitions.

Partitioning allows accessing different parts of the array in parallel so that bottlenecks caused by a single memory interface being overwhelmed with requests can be avoided.

There are 3 different types of partitioning that can be applied to an array (Refer to Figure 3.3):

1. **Cyclic.** In a cyclic partition, the array is divided blocks of interleaved elements of the original array.
2. **Block.** In a block partition, the array is divided into non-overlapping blocks of sequential elements in the original array. Each block is assigned to a separate memory bank.

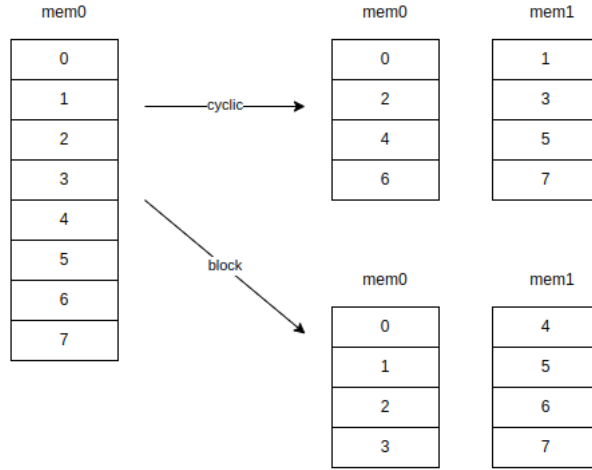


Figure 3.3: Array partitioning illustration.

3. Complete. The array is split into individual elements

The disadvantage of partitioning is that it increases the number of memory interfaces which leads to the increased resource utilization (i.e., more FFs, LUTs are used because each memory block requires separate control logic). This can be partly mitigated by using the **array reshaping**. The difference between partitioning and reshaping is that partitioning creates multiple memory interfaces while reshaping still uses a single memory interface (i.e., all partitions are merged to a single physical memory).

3.1.3 Example of optimizing matrix multiplication

In this section, we will describe the process of optimizing a matrix multiplication using the techniques described in Section 3.1.2. This section can be treated as a tutorial on how to optimize a simple matrix multiplication.

Full source code for with all the files can be located in `vitis_hls/matmul_naive`.

Naive implementation

The naive implementation of matrix multiplication is simply a triple for loop which directly implements the definition of $C = A \cdot B$ where A, B and C are the matrices and $C_{i,j}$ are defined as in Equation 3.1.2.

$$C_{i,j} = \sum_{k=0}^N A_{i,k} \cdot B_{k,j} \quad (3.1.2)$$

```
#include "matrixmul.h"

void matmul(mat_a_t a[MAT_A_ROWS][MAT_A_COLS],
            mat_b_t b[MAT_B_ROWS][MAT_B_COLS],
            result_t res[MAT_A_ROWS][MAT_B_COLS]) {
loop_i:
  for (int i = 0; i < MAT_A_ROWS; i++) {
loop_j:
    for (int j = 0; j < MAT_B_COLS; j++) {
      res[i][j] = 0;
loop_k:
      for (int k = 0; k < MAT_B_ROWS; k++) {
#pragma HLS PIPELINE off
        res[i][j] += a[i][k] * b[k][j];
      }
    }
  }
}
```

The Vitis HLS Synthesis Report for the naive implementation is presented in Figure 3.4. The data types for matrices used are int32_t. Matrix A is of size 3x4 and matrix B has size 4x3. For this baseline implementation, we are reaching a latency of 205 clock cycles, using 327 FF and 282 LUTs.

Modules & Loops	Issue Type	Violation Type	Distance	Slack	Latency(cycles)	Latency(ns)	Iteration Latency	Interval	Trip Count	Pipelined	BRAM	DSP	FF	LUT	URAM
matmul					205	2.050E3	-	206	-	no	0	3	327	282	0
loop_j					204	2.040E3	68	-	3	no	-	-	-	-	-
loop_j					66	660.000	22	-	3	no	-	-	-	-	-
loop_k					20	200.000	5	-	4	no	-	-	-	-	-

Figure 3.4: Naive matrix multiplication synthesis report

Loop pipelining and unrolling

The baseline code can be optimized by pipelining. There are three loops that can be pipelined (loop_i, loop_j and loop_k).

Case 1: Pipelining loop_k. A simple pipelining of the innermost loop leads to the pipelining of Multiply and Accumulate operation (MAC) which is the main operation in the loop. There is no need to partition arrays a and b as memory in arrays a and b only need to supply 1 element per cycle.

The pipelining leads to decrease in latency to 42 clock cycles at the cost of using 526 FF and 501 LUT which is almost twice as many resources as in the naive implementation.

Modules & Loops	Issue Type	Violation Type	Distance	Slack	Latency(cycles)	Latency(ns)	Iteration Latency	Interval	Trip Count	Pipelined	BRAM	DSP	FF	LUT	URAM
matmul					42	420.000	-	43	-	no	0	3	526	501	0
loop_i_loop_j_loop_k					40	400.000	6	1	36	yes	-	-	-	-	-

Figure 3.5: Matrix multiplication with pipelined loop loop_k

Case 2: Pipelining loop_j. The pipelining of the middle loop requires partitioning of arrays a by MAT_A_COLS and b by MAT_B_ROWS. There are MAT_A_COLS (or MAT_B_ROWS) MAC operations per cycle so that memory of a and b needs to be sufficiently divided to supply MAT_B_ROWS elements per cycle. Array a is partitioned along the second dimension and array b is partitioned along the first dimension because of the access patterns. In this example, we use the complete partitioning which divides the memory into individual registers.

```
#include "matrixmul.h"

void matmul(mat_a_t a[MAT_A_ROWS][MAT_A_COLS],
            mat_b_t b[MAT_B_ROWS][MAT_B_COLS],
            result_t res[MAT_A_ROWS][MAT_B_COLS]) {
#pragma HLS ARRAY_PARTITION variable = a complete dim = 2
#pragma HLS ARRAY_PARTITION variable = b complete dim = 1
loop_i:
  for (int i = 0; i < MAT_A_ROWS; i++) {
#pragma HLS PIPELINE off
    loop_j:
      for (int j = 0; j < MAT_B_COLS; j++) {
        int tmp = 0;
#pragma HLS PIPELINE
        loop_k:
          for (int k = 0; k < MAT_B_ROWS; k++) {
#pragma HLS UNROLL
            tmp += a[i][k] * b[k][j];
          }
        res[i][j] = tmp;
      }
    }
  }
}
```

With this addition, the latency is down to 15 clock cycles and the resource usage has increased to 1221 FF and 511 LUTs (Figure 3.6).

Case 3: Pipelining loop_i. Pipelining the outermost loop requires additional partitioning of array res by MAT_B_COLS and full partitioning of array b as loop_k and loop_j are getting unrolled.

Modules & Loops	Issue Type	Violation Type	Distance	Slack	Latency(cycles)	Latency(ns)	Iteration Latency	Interval	Trip Count	Pipelined	BRAM	DSP	FF	LUT	URAM
matmul					15	150.000	-	16	-	no	0	12	1231	511	0
loop_i_loop_j					13	130.000	6	1	9	yes	-	-	-	-	-

Figure 3.6: Matrix multiplication with pipelined loop loop_j

```
#include "matrixmul.h"

void matmul(mat_a_t a[MAT_A_ROWS][MAT_A_COLS],
            mat_b_t b[MAT_B_ROWS][MAT_B_COLS],
            result_t res[MAT_A_ROWS][MAT_B_COLS]) {
#pragma HLS ARRAY_PARTITION variable = a complete dim = 2
#pragma HLS ARRAY_PARTITION variable = b complete dim = 0
#pragma HLS ARRAY_PARTITION variable = res complete dim = 2
loop_i:
  for (int i = 0; i < MAT_A_ROWS; i++) {
#pragma HLS PIPELINE
  loop_j:
    for (int j = 0; j < MAT_B_COLS; j++) {
#pragma HLS UNROLL
      int tmp = 0;
      loop_k:
        for (int k = 0; k < MAT_B_ROWS; k++) {
#pragma HLS UNROLL
          tmp += a[i][k] * b[k][j];
        }
      res[i][j] = tmp;
    }
  }
}
```

The latency is down to 9 clock cycles and the resource usage is up to 3819 FF and 1383 LUTs (Figure 3.7).

Modules & Loops	Issue Type	Violation Type	Distance	Slack	Latency(cycles)	Latency(ns)	Iteration Latency	Interval	Trip Count	Pipelined	BRAM	DSP	FF	LUT	URAM
matmul					9	90.000	-	10	-	no	0	36	3819	1383	0
loop_i					7	70.000	6	1	3	yes	-	-	-	-	-

Figure 3.7: Matrix multiplication with pipelined loop loop_i

Loop reordering

We can also compare how the loop reordering affects the performance of the matrix multiplication and specifically how it compares to the pipelining of the loop_j in terms of latency and resource usage.

The loop reordering technique allows us to avoid partitioning matrix a compared to the loop_j pipelining solution. However, we now have to partition the output matrix by MAT_B_COLS.

```
#include "matrixmul.h"

void matmul(mat_a_t a[MAT_A_ROWS][MAT_A_COLS],
            mat_b_t b[MAT_B_ROWS][MAT_B_COLS],
            result_t res[MAT_A_ROWS][MAT_B_COLS]) {
  int temp_sum[MAT_B_COLS];
#pragma HLS ARRAY_PARTITION variable = b dim = 2 complete
#pragma HLS ARRAY_PARTITION variable = res dim = 2 complete
#pragma HLS ARRAY_PARTITION variable = temp_sum dim = 1 complete
loop_i:
  for (int i = 0; i < MAT_A_ROWS; i++) {
  loop_k:
    for (int k = 0; k < MAT_B_ROWS; k++) {
#pragma HLS PIPELINE
    loop_j:
      for (int j = 0; j < MAT_B_COLS; j++) {
#pragma HLS UNROLL
        int result = (k == 0) ? 0 : temp_sum[j];
        result += a[i][k] * b[k][j];
        temp_sum[j] = result;
        if (k == MAT_B_ROWS - 1) {
```

```

        res[i][j] = result;
    }
}
}
}
}

```

From Figure 3.8, the latency is 17 clock cycles and the resource usage is 1031 FF and 611 LUTs. The performance is worse than the pipelined loop `loop_j` but the resource usage is lower for FF and higher for LUTs.

Modules & Loops	Issue Type	Violation Type	Distance	Slack	Latency(cycles)	Latency(ns)	Iteration Latency	Interval	Trip Count	Pipelined	BRAM	DSP	FF	LUT	URAM
matmul					17	170.000	-	18	-	no	0	9	1031	621	0
loop_j_loop_k					15	150.000	5	1	12	yes	-	-	-	-	-

Figure 3.8: Matrix multiplication with loop reordering

Summary The pipelining of the outermost loop leads to the best performance with the lowest latency. However, this comes at the big cost of using the most resources as array `b` needs to be fully partitioned. In this work, we will use the pipelined loop `loop_j` as it provides a good trade-off between latency and resource usage.

3.2 Transformer architecture optimizations

In this section, we present the specific optimizations that were applied to the Transformer architecture. For HLS code (i.e., implementation), see `vitis_hls` directory. The comparison of the performance and resource utilization of the FPGA optimized architecture and the CPU implementation is presented in Section 4.4.

3.2.1 Vanilla Transformer

In this section, we will describe the specific optimizations that can be applied to the Transformer architectures and present the HLS synthesis reports for the optimized architecture.

For the purpose of the analysis, we will use the architectures with the hyperparameters as described in Section 4.1. The models achieve reasonable performance while still being small enough to fit the board.

For all **matrix multiplications** performed in the architecture, we will use the pipelined loop `loop_j` as described in Section 3.1.3 because it provides a reasonable performance-resource utilization tradeoff.

We utilize `float data type` for all the matrices in the architecture. We found that using `double` data type does not affect the accuracy performance of the model, while it significantly increases the resource utilization (DSP) and makes the model not fit into the board.

We also use the following numerical optimization technique applied to **softmax** to make it numerically stable. Notice that

$$\begin{aligned}
 \text{softmax}(x_i) &= \frac{e^{x_i}}{\sum_{j=1}^N e^{x_j}} \\
 &= \frac{e^{x_i - \max(x)}}{\sum_{j=1}^N e^{x_j - \max(x)}} \\
 &= \text{softmax}(x_i - \max(x))
 \end{aligned}$$

Since e_i^x term can be very large, subtracting the maximum value from all the elements bounds the values of e_i^x (the largest element will be 0) so that the sum of exponents will never become large.

For the **ReLU** activation function, we which needs to be applied to a 2 dimensional matrix (see Feed Forward layer of the encoder block in Section 2.1.2) fully unroll the loops and use the **complete** partitioning of the matrix.

All the loops in the architecture are automatically pipelined unless it is impossible to do so (e.g., because of the dependencies between the iterations).

The optimized architecture reaches latency of 3714 cycles (37.14 microseconds) utilizes 5% BRAM, 30% FFs and 90% LUTs.

See full report in `vitis_hls/transformer/csynth.rpt`.

For comparison, the architecture with all optimization turned off reaches latency of 34745 cycles (347.45 microseconds, 10 times as many) and utilizes only 2% BRAM, 2% FFs and 10% LUTs.

See full report in `vitis_hls/transformer_slow/csynth.rpt`.

3.2.2 Linear transformer

We also implement the linear transformer architecture (i.e., linear attention) as described in Section 2.2.2 on the board.

Similarly to the vanilla transformer implementation (described in Section 3.2.1), we also utilize `float` data type for all the matrices in the architecture, use `loop_j` pipelining for matrices, fully unroll the loops for ReLU activation function.

Remark 3.2.1. As an implementation detail, the matrices had to be manually partitioned with the usage of `pragma HLS ARRAY_PARTITION` directive because Vitis HLS could not infer the partitioning automatically.

The optimized version achieves 2986 cycles (29.86 microseconds) latency and utilizes 18% BRAM, 43% FFs and 99% LUTs

The full report is in `vitis_hls/linear_transformer/csynth.rpt`.

For comparison, the architecture with all optimization turned off reaches latency of 68894 cycles (688.94 microseconds) and utilizes 17% BRAM, 6% FFs and 24% LUTs. The full report is in `vitis_hls/linear_transformer_slow/csynth.rpt`.

Chapter 4

Experiments

In this section, we will describe the experiments that were performed to evaluate the performance of the proposed architecture. While the main focus of this work is the inference performance of the models (refer to Section 4.4), we will also describe the training procedure and the accuracy performance of the models.

4.1 Architecture and hyperparameters

Here we will describe the model architecture and the hyperparameters used for the experiments.

In the experiments 3 models were used for comparison:

- **Linear Regression** - a simple linear regression model on handcrafted features
- **Transformer Encoder** - a transformer encoder model on raw time series data
- **Linear Transformer Encoder** - a linear transformer model on raw time series data

The reason for using the linear regression model is to provide a baseline for the performance of the transformer models to show that even simple untuned models can outperform the linear regression model.

For both transformer models, we used the encoder architectures as described in Section 2.2 with a linear layer on top of the output of the transformer encoder to get the final prediction (i.e., if a sample is an anomaly).

Remark 4.1.1. The encoder part has positional encoding and layer normalization disabled. We found that the positional encoding does not improve the performance of the model and leads to more unstable training (see Section 4.1.1).

The transformer hyperparameters used for the experiments are presented in Table 4.1.

Parameter, Model	Transformer Encoder	Linear Transformer Encoder
Window Size	8	8
Number of heads	2	-
Dim. of FeedForward network	16	16
Number of blocks	2	1

Table 4.1: Hyperparameters

The learning rate was chosen using the learning rate finder [27] (see Section 4.1.1) and the batch size was chosen to be the maximum value that fitted the dataset in memory (2^{13} samples), we also used **Adam** optimizer [28] (see Section 4.1.1).

4.1.1 Model Fitting

In this section, we will describe the training procedure, the main tools used and issues that we encountered and how they were addressed.

General procedure

The General process of training a neural network involves iteratively adjusting its parameters to minimize a specified loss function. This is typically achieved through an optimization algorithm, such as gradient descent.

That is, at each iteration, the parameters of the model θ are updated as follows:

$$\theta_{t+1} = \theta_t - \alpha \nabla_{\theta} \mathcal{L}(\theta_t) \quad (4.1.1)$$

where α is the learning rate and $\mathcal{L}(\theta_t)$ is the loss function at the t -th iteration.

While the general procedure is simple there are multiple methods to improve the convergence of the optimization algorithm, which will be described in the following sections.

Optimizer

The optimizer is the algorithm that is used to update the parameters of the model. The general procedure of updating the parameters as described in Equation 4.1.1, can be improved by using more sophisticated optimization algorithms.

The **Adam** optimizer is a popular optimization algorithm used in training artificial neural networks and is a reasonable baseline choice for many problems. It stands for "Adaptive Moment Estimation" and combines ideas from two other optimization techniques: RMSprop (Root Mean Square Propagation) and Momentum. The Adam optimizer is known for its efficiency and fast convergence, and robustness to various types of neural network architectures and problem domains.

The optimization procedure performed by Adam is as follows:

1. Initialize 2 moving average accumulators, m and v with 0. They will store the exponentially decaying average of past gradients and squared gradients respectively for each parameter θ
2. At each iteration, compute the gradient of the loss function with respect to the parameters, $\nabla_{\theta} \mathcal{L}(\theta_t)$
3. Update the moving averages m and v as follows:

$$\begin{aligned} m &= \beta_1 m + (1 - \beta_1) \nabla_{\theta} \mathcal{L}(\theta_t) \\ v &= \beta_2 v + (1 - \beta_2) \nabla_{\theta} \mathcal{L}(\theta_t)^2 \end{aligned}$$

where β_1 and β_2 are the hyperparameters that control the decay rate of the moving averages. The usual values for β_1 and β_2 are 0.9 and 0.999 respectively.

4. Compute the bias-corrected moving averages \hat{m} and \hat{v} as follows:

$$\begin{aligned} \hat{m} &= \frac{m}{1 - \beta_1^t} \\ \hat{v} &= \frac{v}{1 - \beta_2^t} \end{aligned}$$

The bias correction is necessary because the moving averages are initialized with 0 (see original paper for derivations [28]).

5. Update the network parameters as follows:

$$\theta_{t+1} = \theta_t - \alpha \frac{\hat{m}}{\sqrt{\hat{v}} + \epsilon}$$

where ϵ is a small constant used for numerical stability (i.e., to avoid division by zero).

The update is similar to the gradient descent update rule (Equation 4.1.1) except that the gradient is divided by the square root of the moving average of squared gradients.

We found that the optimizer with default parameters (β_1 and β_2) works well for our task, and we did not experiment with other optimizers.

Learning rate finder

The learning rate α is one of the most important hyperparameters that determines the step size in parameter space during gradient descent optimization. An appropriate learning rate is essential for model convergence. The problem of choosing the learning rate is a well-known problem in machine learning and badly chosen learning rate can lead to either underfitting where the model learns too slowly or it can lead to divergence where the parameters are updated too abruptly.

In [27], authors proposed a simple method to find an appropriate learning rate automatically by plotting the loss function against the learning rate.

The procedure is performed as follows:

1. Start with a very small learning rate α and increase it at each iteration
2. At each iteration, train the model for a few epochs and compute the loss function
3. Plot the loss function against the learning rate. This plot is crucial in identifying the "sweet spot" in the learning rate range where the loss is decreasing effectively (i.e., in the most steeper way).
4. Choose the point on the learning rate vs. loss curve where the loss starts to decrease most steeply. This point indicates that the model is making the most significant progress towards convergence.

Remark 4.1.2. While there are no guarantees that the learning rate finder will find the optimal learning rate, it provides a good empirical estimate of the optimal learning rate.

Instead of manually plotting the loss function against the learning rate, we used the implementation provided in Pytorch Lightning [29].

Gradient explosion and Gradient clipping

A different challenge of training neural networks is the phenomenon known as the "gradient explosion problem." We have found that the gradient explosion problem is especially prevalent in the proposed architectures.

The gradient explosion problem occurs when the gradient of the loss function with respect to the parameters becomes too large and the parameters are updated too abruptly (e.g., as in Equation 4.1.1). This can lead to the model diverging and the loss function increasing instead of decreasing. An example of the loss function diverging is presented in Figure 4.1.

To solve the problem of gradient explosion, we used the gradient clipping technique introduced in [30]. The idea behind gradient clipping is to clip the gradient to a maximum value g_{max} . This remediates the problem of gradient explosion because the gradient is bounded and the parameters are updated more smoothly.

Class imbalance and loss functions

In anomaly detection task the dataset is often imbalanced, i.e., the number of normal samples is much larger than the number of anomalous samples (see Section 4.2).

This poses a problem for the training of the model because the model can simply learn to predict the majority class (i.e., normal samples) and achieve a high accuracy without having good recall (refer to Section 4.3.1 for the description to the metrics).

The loss functions that we use for training is the binary cross-entropy loss function [31] which is computed as follows:

$$\mathcal{L}(\theta) = -\frac{1}{N} \sum_{i=1}^N y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i) \quad (4.1.2)$$

where y_i is the true label and \hat{y}_i is the predicted label.

While the binary cross-entropy loss function is a good choice for the anomaly detection task, it does not take into account the class imbalance problem.

A way to solve the class imbalance problem is to use a weigh positive samples (i.e., anomalous samples) more than negative samples in the loss function.

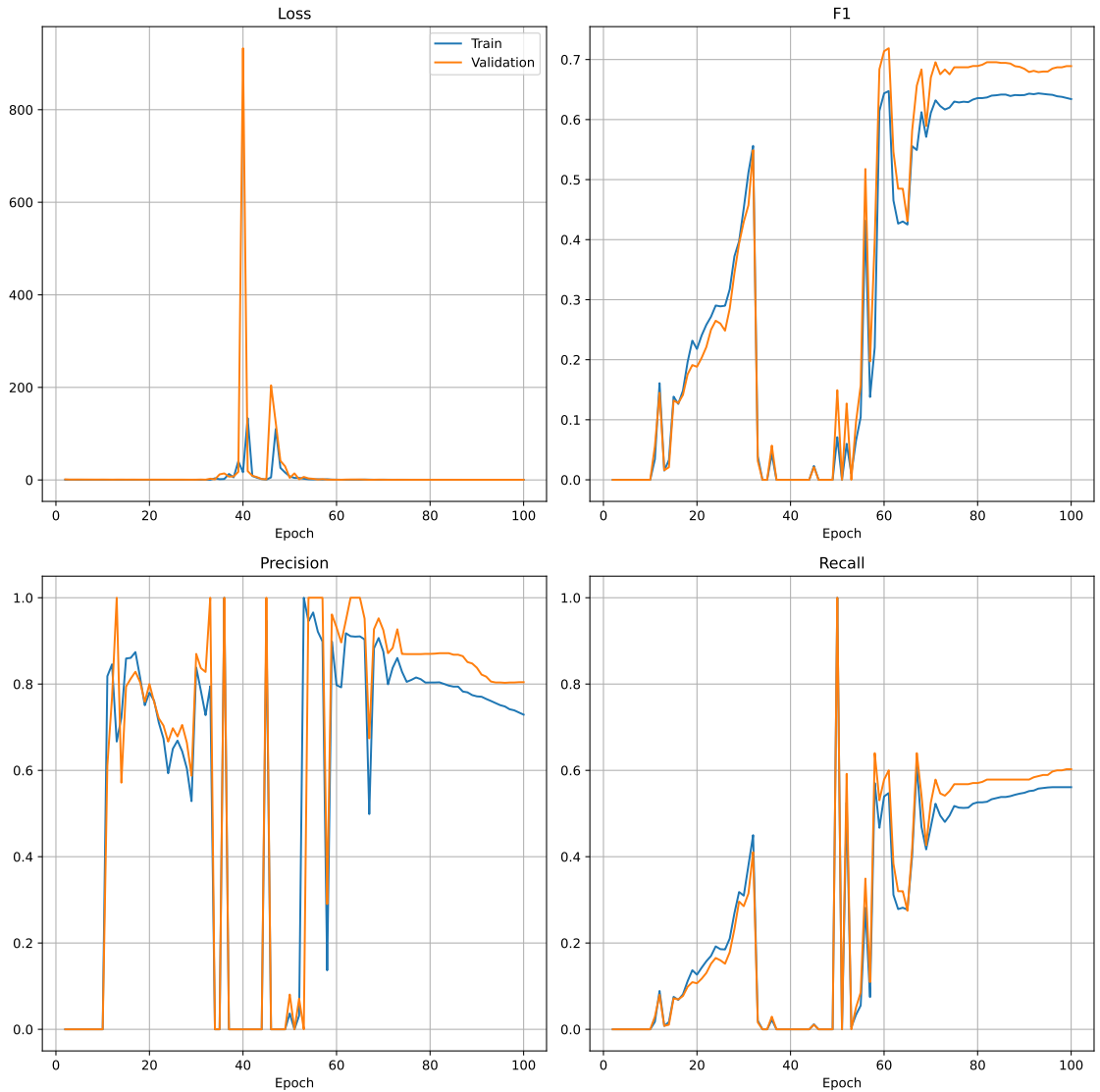


Figure 4.1: Example of gradient explosion at around epoch 40 which leads to the loss function diverging for a few following epochs.

So we can modify the loss function as follows:

$$\mathcal{L}(\theta) = -\frac{1}{N} \sum_{i=1}^N w_i (y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)) \quad (4.1.3)$$

where w_i is the weight of the i -th sample. In the experiments, we found that weighting the positive samples 5 times more than the negative samples yielded good results for most of the datasets which we used throughout the experiments.

Training stability: Layer normalization and positional encoding

While the base transformer encoder architecture uses layer normalization and positional encoding layers, we found that they lead to unstable training and worse performance of the model on most of the datasets. Hence, we disabled them for the experiments and replaced them with identity layers.

Data Preprocessing

Although, this is not a part of the training procedure, before fitting a neural network, it is usually a good idea to preprocess the data by scaling it.

For the experiments, we used the **MinMaxScaler** from the `scikit-learn` library [32], to scale the input features to the range of $[0, 1]$ for a more stable training.

4.2 Datasets

In this section, the datasets used for model training and performance evaluation will be described.

4.2.1 Numenta Anomaly Benchmark (NAB)

To assess the accuracy of predictions, we use the Numenta Anomaly Benchmark [33] dataset, which contains various real-world labeled time series of temperature sensor readings, CPU utilization of cloud machines, service request latencies, and taxi demands in New York City. It is commonly used to assess the performance of anomaly detection models on time-series data.

The reason why we use this dataset is that it is a standard benchmark dataset for anomaly detection in time series and because it has a large number of labeled time series.

A sample time series of NYC taxi demand is presented in Figure 4.2.

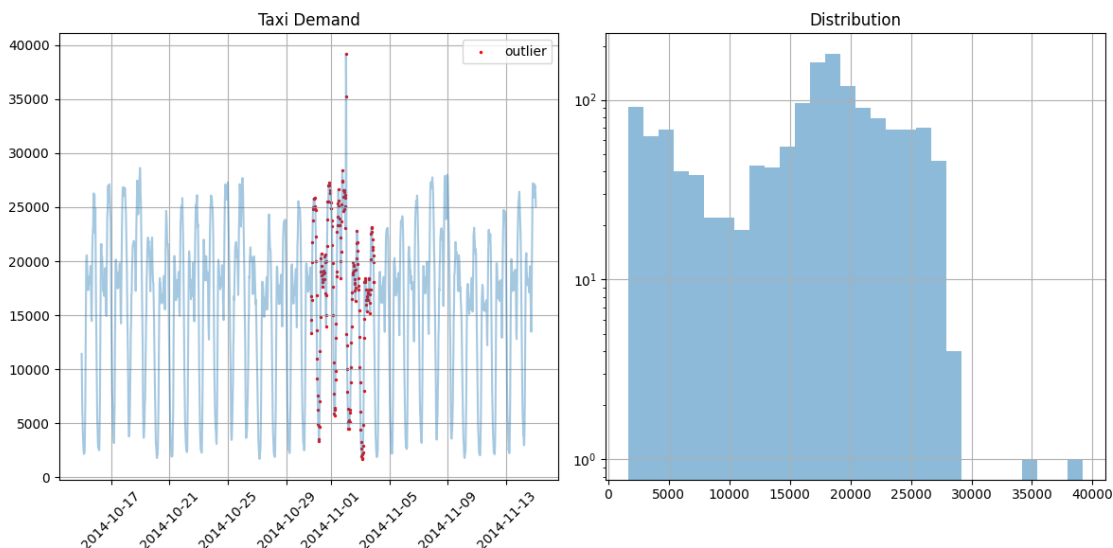


Figure 4.2: NYC Taxi demand - anomalies highlighted in red. Left: Time series of observations, Right: Distribution of observations.

In this work, we only used the NYC taxi demand dataset as in the [21]. The dataset consists of 10320 samples with 10% of anomalous samples.

4.2.2 KPI Anomaly Detection Dataset

The other labeled dataset that we use is the KPI Anomaly Detection Dataset (KPI AIOps) [34]. This dataset alongside the NAB dataset will be used to evaluate the predictive performance of the anomaly detection models.

The dataset consists of KPI (key performance index) time series data from many real scenarios of Internet companies with ground truth label. KPIs fall into two broad categories: service KPIs and machine KPIs. Service KPIs are performance metrics that reflect the size and quality of a Web service, such as page response time, page views, and number of connection errors. Machine KPIs are performance indicators that reflect the health of the machine (server, router, switch), such as CPU utilization, memory utilization, disk IO and network card throughput.

A sample time series of a sensor readings is presented in Figure 4.3. We can clearly see the outliers for some of the observations (colored in red).

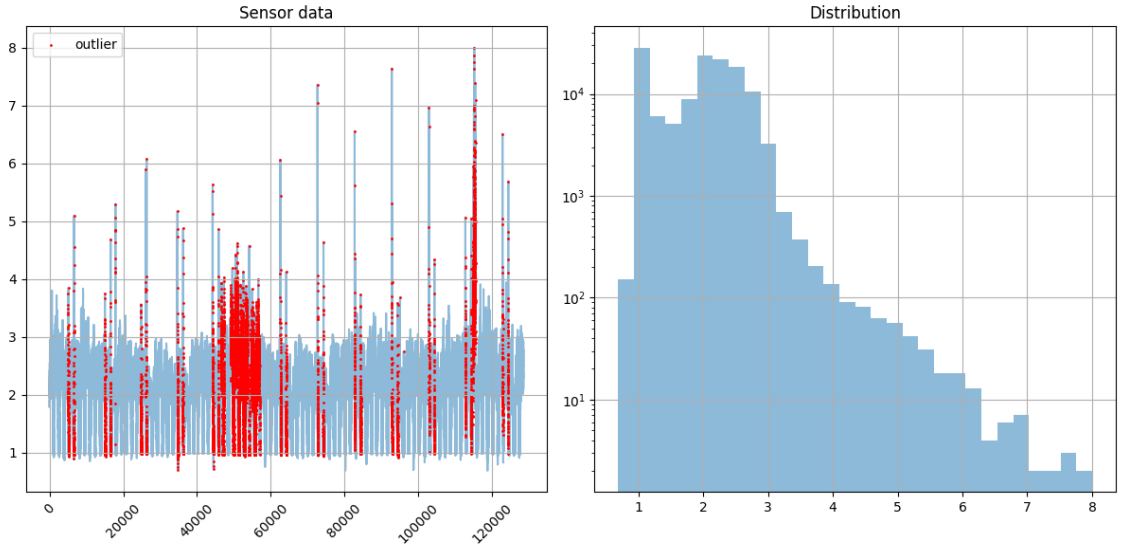


Figure 4.3: Sensor data from a machine in a data center. The red dots indicate the anomalies. Left: Time series of observations, Right: Distribution of observations.

4.2.3 FI2010

In [35], authors described the first publicly available benchmark dataset of high-frequency limit order markets for mid-price prediction. The dataset contains 10-day limit order book data from June 2010 for five stocks that are listed on the Helsinki exchange. Each entry in the time series provides price details and aggregate order sizes for the top ten levels on both the bid and offer sides of the market, totaling forty data points. The time series consists of approximately four million messages, representing incoming buy/sell orders or cancellations. The dataset features order book snapshots taken after every 10 messages, resulting in approximately 400,000 records for the five stocks.

A number of versions of the dataset are available using different normalization schemes. We used the not normalized version of the dataset.

For the purpose of this work, we only extract only the mid price from the dataset which will be used for anomaly detection task.

Synthetic outliers

Since the dataset is not labeled, we have to inject synthetic anomalies into the dataset. We employ the approach similar to [36] with a slight modification. The algorithm can be summarized as follows:

1. Select n samples from the time series which will be contaminated (i.e., anomalous)
2. Replace the sample S_i with $\hat{S}_i = S_i(1 + \delta)$ where δ is the injected outlier in the return space.

Authors model δ as a uniformly distributed random variable $\mathcal{U}[0, \rho]$. We instead use the uniform distribution $\mathcal{U}[-\rho, \rho]$ (i.e., just increase the range) to allow for positive and negative outliers.

An example of the injected outliers is presented in Figure 4.4

4.3 Accuracy

In this section, we will describe the main metrics used to evaluate the performance of the anomaly detection models. The Transformer encoder model will be compared with the simple linear regression model on handcrafted features and with the Linear Transformer model [20]. The inference procedure will be described and the results will be presented.

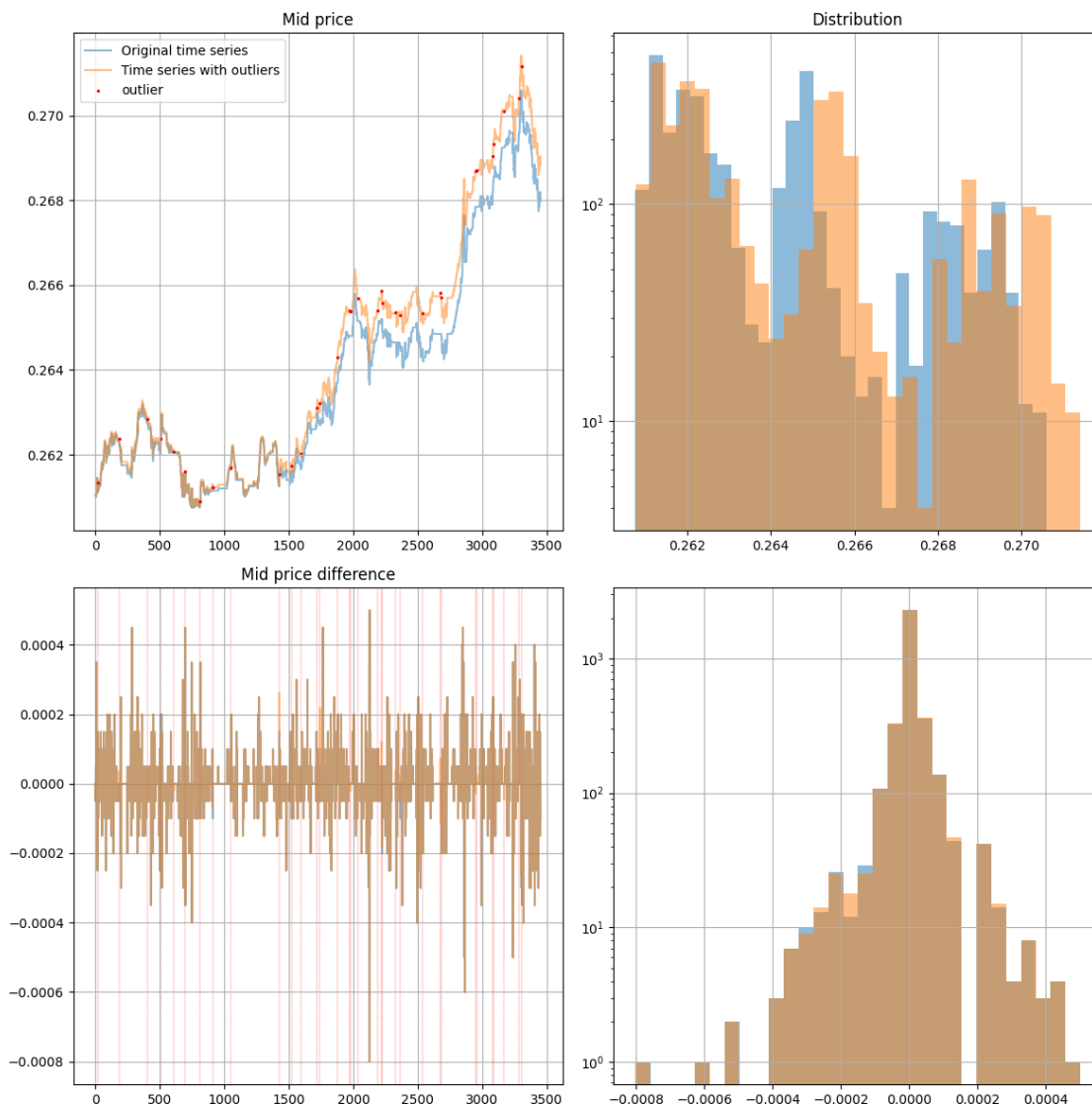


Figure 4.4: Example of the injected outliers in the FI2010 dataset.

4.3.1 Metrics

In this section the metrics used to evaluate the performance of the anomaly detection models will be described.

Before we describe the metrics, we need to introduce the confusion matrix and the following notation:

- **TP** - True Positive calculated as the number of correctly predicted anomalies
- **TN** - True Negative calculated as the number of correctly predicted non-anomalies
- **FP** - False Positive which is the number of incorrectly predicted anomalies
- **FN** - False Negative which is the number of incorrectly predicted non-anomalies
- **P** - Number of positive samples, i.e., $\mathbf{P} = \mathbf{TP} + \mathbf{FN}$
- **N** - Number of negative samples, i.e., $\mathbf{N} = \mathbf{TN} + \mathbf{FP}$

The confusion matrix is a table with two rows and two columns that reports the number of false positives, false negatives, true positives, and true negatives.

The matrix summarizes the predictions from a classification model, i.e., how well the model performed when predicting the class labels for positive and negative samples. While the matrix

		prediction outcome		
		P	n	
actual value	P'	True positive TP	False negative FN	P'
	N'	False positive FP	True negative TN	N'
total		P	N	

Figure 4.5: Confusion matrix description

presents the most informative view of the performance of the model, we still need to summarize the information in the matrix into a single number(s) that can be used to compare different models.

In this paper, we will use the following metrics to compare the performance of the anomaly detection models:

- **Accuracy** is the fraction of predictions that the model got right. It is defined as follows:

$$\text{Accuracy} = \frac{\mathbf{TP} + \mathbf{TN}}{\mathbf{TP} + \mathbf{TN} + \mathbf{FP} + \mathbf{FN}}$$

While this metric is easy to understand, it is not very informative when the dataset is imbalanced which is the case for the anomaly detection task where the proportion of positive samples is low.

- **Precision** is the fraction of positive predictions that were correct.

$$\text{Precision} = \frac{\mathbf{TP}}{\mathbf{TP} + \mathbf{FP}}$$

This metric is useful when the cost of false positives is high. For example, in the case of anomaly detection, we want to have a high precision so that we do not have to manually check many false positives or trigger any downstream filtering task too often.

- **Recall** is the fraction of positive samples that were correctly predicted.

$$\text{Recall} = \frac{\mathbf{TP}}{\mathbf{TP} + \mathbf{FN}}$$

This metric is useful when the cost of false negatives is high. For example, in the case of anomaly detection, we want to minimize the number of false negatives because we do not want to miss any anomalies.

- **F1** is the harmonic mean of precision and recall which ranges from 0 to 1.

$$\text{F1} = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

This metric is useful when we want to balance the precision and recall. For example, in the case of anomaly detection, we want to have a high precision so that we do not have to manually check many false positives or trigger any downstream filtering task too often. At the same time, we want to minimize the number of false negatives because we do not want to miss any anomalies.

The advantage of using this metric instead of the accuracy is that it can be used even when the dataset is highly imbalanced and it would detect if the model performs poorly in terms of precision and/or recall.

4.3.2 Comparison

Dataset	Model	Accuracy	Precision	Recall	F1
KPI	Linear Regression	0.61/0.66	0.07/0.15	0.73/0.69	0.12/0.25
	Transformer	0.98/0.97	0.89/0.95	0.59/0.63	0.71/0.76
	Linear Transformer	0.97/0.95	0.58/0.73	0.62/0.69	0.60/0.71
NAB	Linear Regression	0.09/0.20	0.09/0.20	1.00/1.00	0.16/0.33
	Transformer	0.76/0.75	0.16/0.41	0.39/0.55	0.23/0.47
	Linear Transformer	0.53/0.56	0.13/0.27	0.76/0.71	0.22/0.40
FI2010	Linear Regression	0.14/0.16	0.14/0.16	1.00/1.00	0.24/0.27
	Transformer	0.84/0.80	0.17/0.21	0.03/0.11	0.06/0.14
	Linear Transformer	0.82/0.77	0.15/0.14	0.06/0.09	0.09/0.11

Table 4.2: Final metrics. The reported value are for the train and validation set.

The final metrics for the models are presented in Table 4.2. Refer to Appendix A for the training plots for the models.

The Transformer model outperforms the other models on all datasets. And that is expected because the Transformer model is more complex (i.e., has more parameters) than the other models. However, the Linear Transformer model performs almost as well and it provides better performance-resource utilization tradeoff (see Section 4.4).

In terms of training robustness (i.e., how stable the training is, refer to Section 4.1.1), a simple Linear Regression model is the most stable model, as expected. Both Transformer models suffer from the gradient explosion problem which leads to unstable training which can be partially remediated (see Section 4.1.1).

However, the performance of a simple Linear regression depends on the dataset. For example, the Linear regression is not able to learn the patterns in the NAB dataset and is stuck at predicting only the positive samples (i.e., anomalies). As for the KPI dataset, the performance on the KPI can be tweaked by training the model for a different number of epochs to achieve the required balance between precision and recall (refer to Figure A.5). Both Transformer models, on the other hand, perform equally well on both datasets as long as there is no gradient explosion problem.

In short, Transformer models outperform the simple linear regression but can have unstable training. Linear Transformer model is slightly behind the vanilla model. This which shows that untuned Transformer models provide a better baseline for anomaly detection than the simple linear regression model.

4.4 Performance/Speed

The performance of FPGA implementation of the models will be compared with the pytorch CPU and GPU implementation of the models. The CPU used for the experiments is Intel(R) Core(TM) i7-8750H CPU @ 2.20GHz clocked to 4.20GHz and the GPU is NVIDIA GeForce GTX 1050 Ti. The latency on CPU and GPU was measured empirically by running the inference multiple times (10,000 runs). Refer to the Appendix B for the histograms of inference times.

Model	Device	Inference time (microseconds)
Transformer	CPU	353.35(113.50)
	GPU	626.34(323.18)
	FPGA	37.14
LinearTransformer	CPU	365(114.21)
	GPU	472.44(110.83)
	FPGA	29.86

Table 4.3: Inference time for the models on different devices. For CPU and GPU we report the mean and standard deviation. For FPGA, the deterministic latency is reported (see Section 3.2).

According to the results (Table 4.3), the FPGA implementation of the models outperforms the CPU implementation by a large margin (achieving a speedup of 10x). The reason for this is that the FPGA implementation is highly parallelized compared to the CPU implementation.

Moreover, the benefit of using the FPGA implementation is that it has deterministic latency (i.e., the inference time is constant) which is not the case for the CPU and GPU implementation where the inference time varies depending on the load of the CPU/GPU (see the standard deviation for the inference times of CPU and GPU).

Remark 4.4.1. The GPU implementation of the models is *empirically* slower than the CPU implementation. This is because the models are small and the overhead of transferring the data to the GPU is larger than the time it takes to perform the inference on the CPU for a single sample.

Regarding the resource utilization, the Linear Transformer model uses more resources than the Transformer model (see Table 4.4) and achieves a slightly better runtime performance (see Table 4.3).

Model	BRAM	FF	LUT
Transformer	5%	30%	90%
LinearTransformer	18%	43%	99%

Table 4.4: Relative resource utilization for the PYNQ-Z2 board.

Chapter 5

Conclusion

Summary

In this work, we have examined the performance of the Transformer model for the anomaly detection task and how it can be effectively implemented on FPGA.

We have shown how different Transformer Architectures (and especially the attention module) affect the performance of the model both in terms of accuracy and inference speed. It was shown that the Transformer model outperforms the simple linear regression model without too much hyperparameter tuning. The only downside of the Transformer model that was observed was the unstable training which can be partially remediated by using gradient clipping.

We have also demonstrated how different optimizations affect the performance and resource utilization of the model on FPGA. We also provided an optimized implementation of the Transformer model on FPGA which can be used as a reference.

This work can be used as a starting point for further research on the topic of anomaly detection using Transformer models and for further research on the topic of efficient implementation of Transformer models on FPGA in different domains (e.g., natural language processing).

Future work

For future work, we propose the following:

1. To try bigger FPGA boards to push the limit of maximum performance. The current implementations hit the limit of the PYNQ-Z2 board.
2. To evaluate the performance on more recent **labeled** financial market data. Currently, there is no publicly available labeled financial market dataset with anomalies. The FI2010 dataset is not labeled and we had to inject synthetic anomalies.
3. To evaluate the performance Transformer on FPGA board for other tasks (e.g., natural language processing). The current implementation is only for the anomaly detection task. The other domains might involve a slightly different architecture (e.g., different number of layers, different number of heads, etc.), it can be valuable to have those models on edge devices because of their recent surge in popularity and increased use. Currently, the language models are too big to fit on the FPGA board (especially, PYNQ-Z2), however, the smaller models can still be implemented and evaluated on bigger boards..

Appendix A

Training plots

Figure A.1: Training metrics for different epochs for Transformer model on KPI dataset

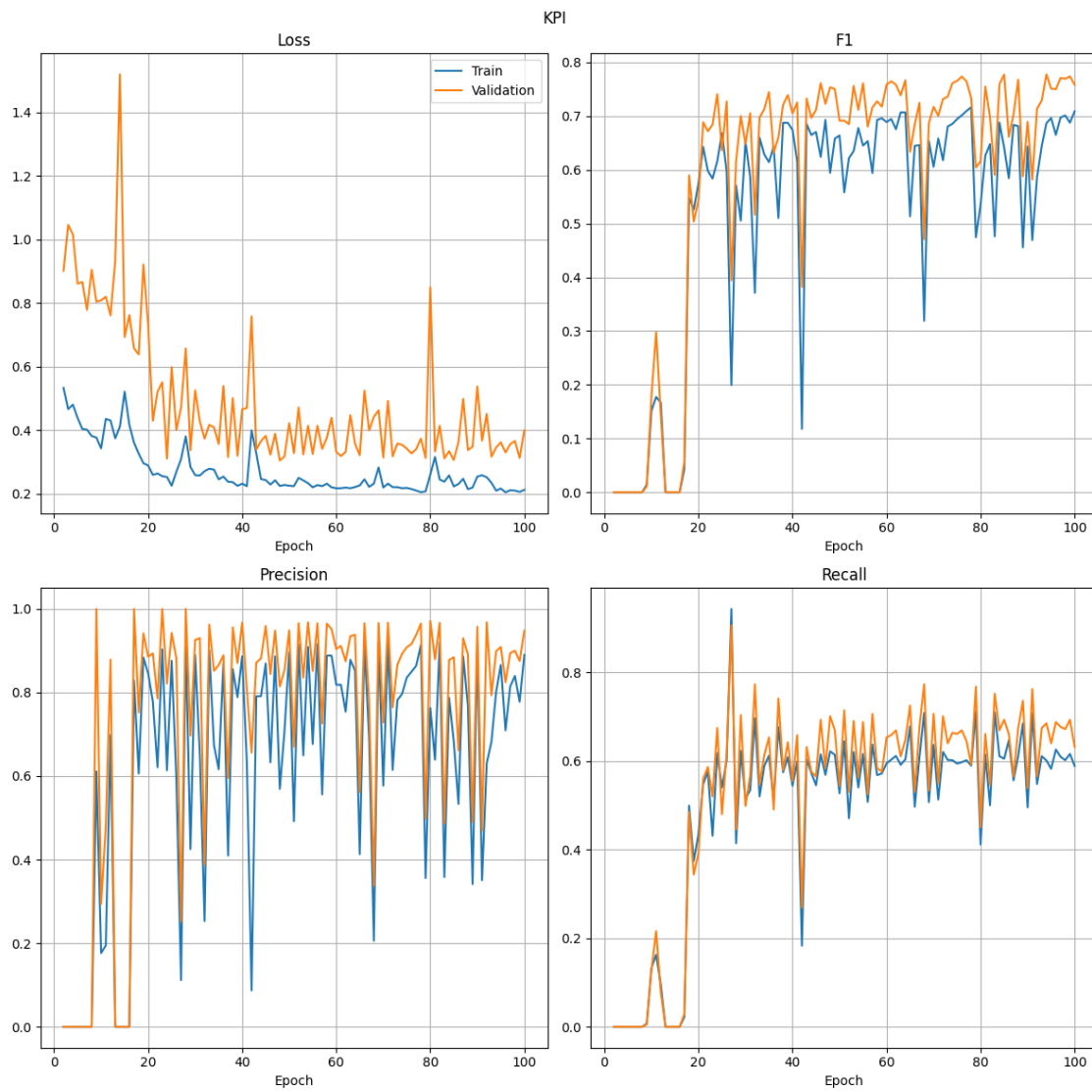


Figure A.2: Training metrics for different epochs for Transformer model on NAB dataset

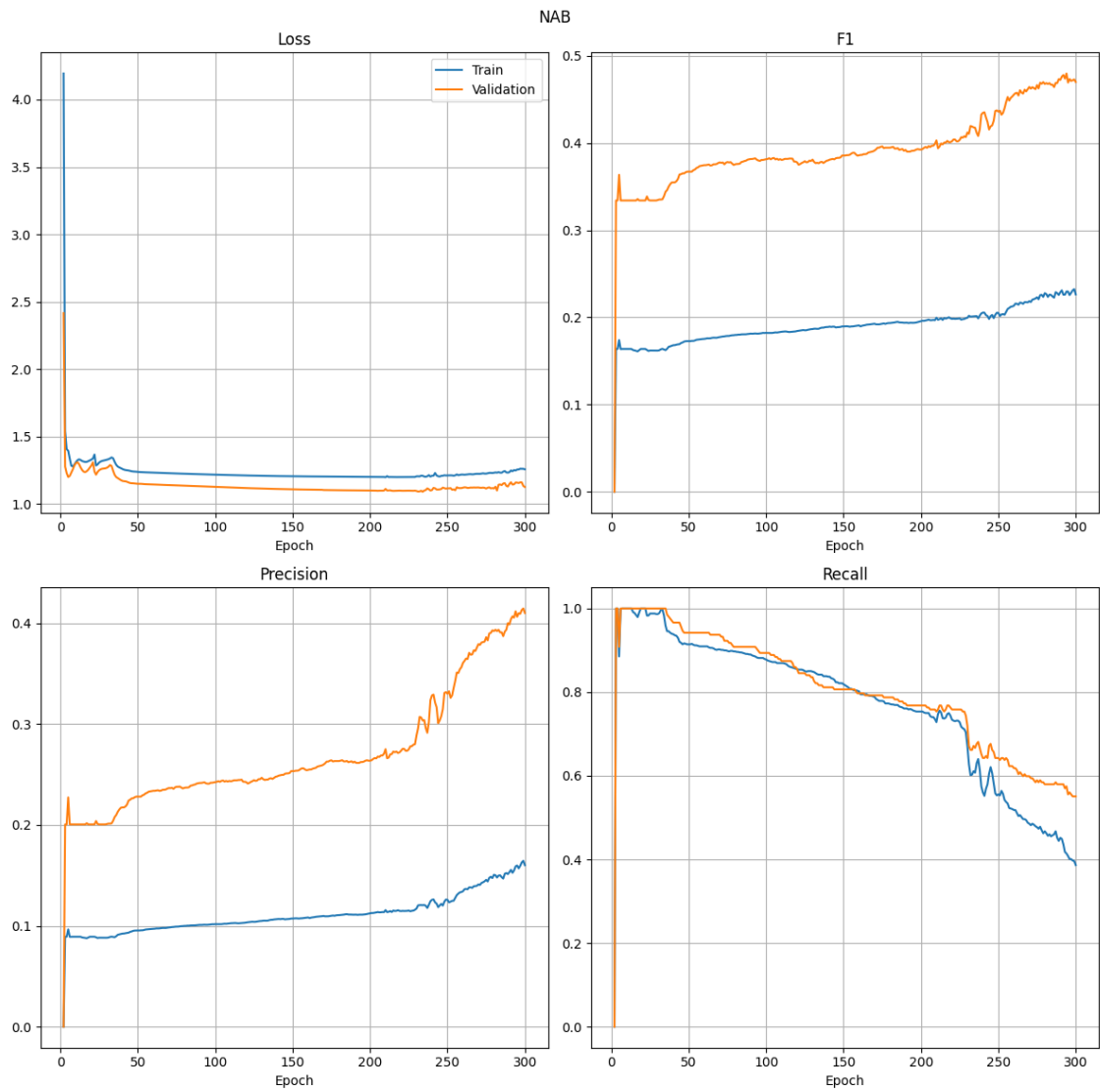


Figure A.3: Training metrics for different epochs for Linear Transformer model on KPI dataset

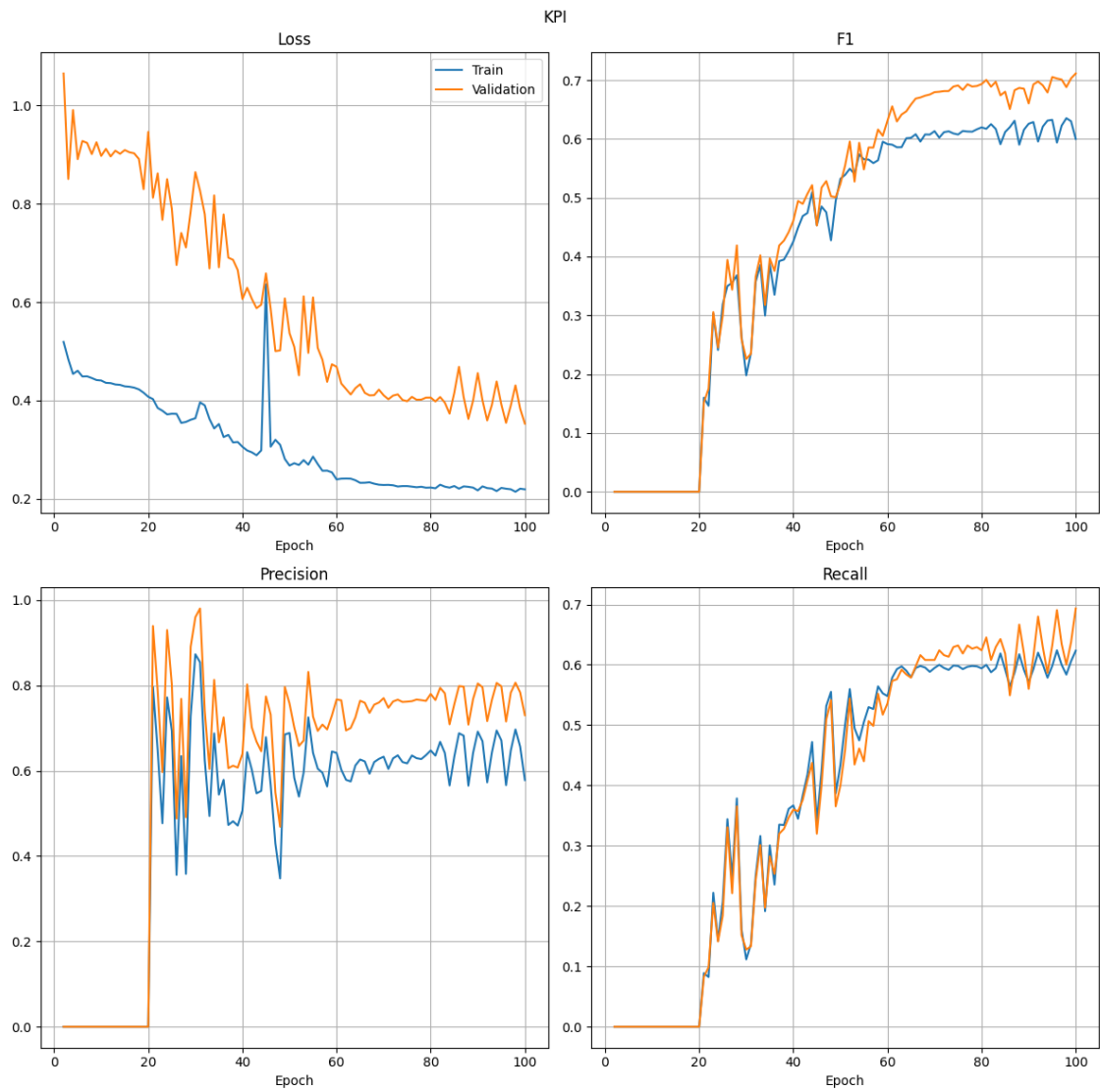


Figure A.4: Training metrics for different epochs for Linear Transformer model on NAB dataset

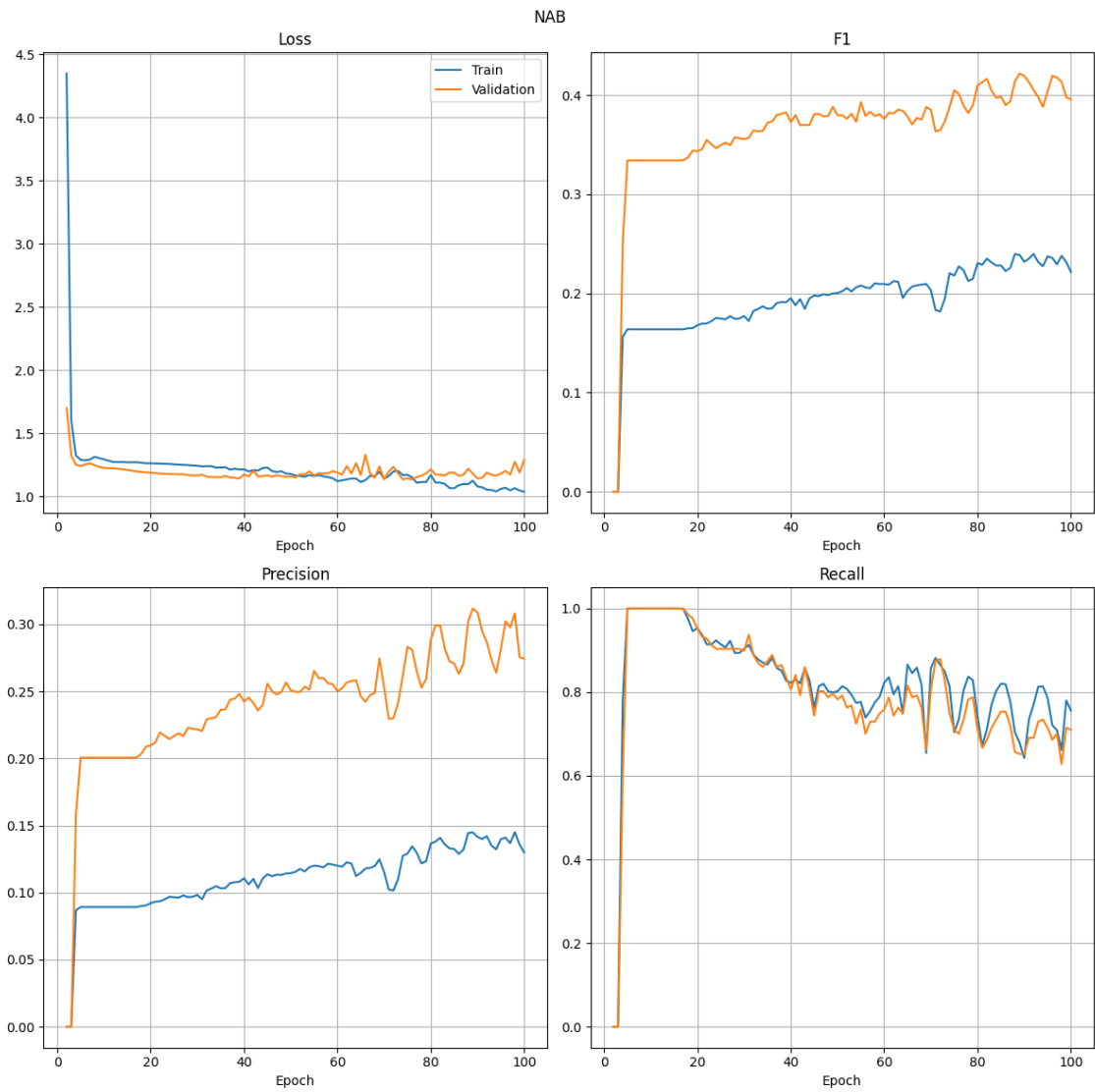


Figure A.5: Training metrics for different epochs for Linear Regression model on KPI dataset

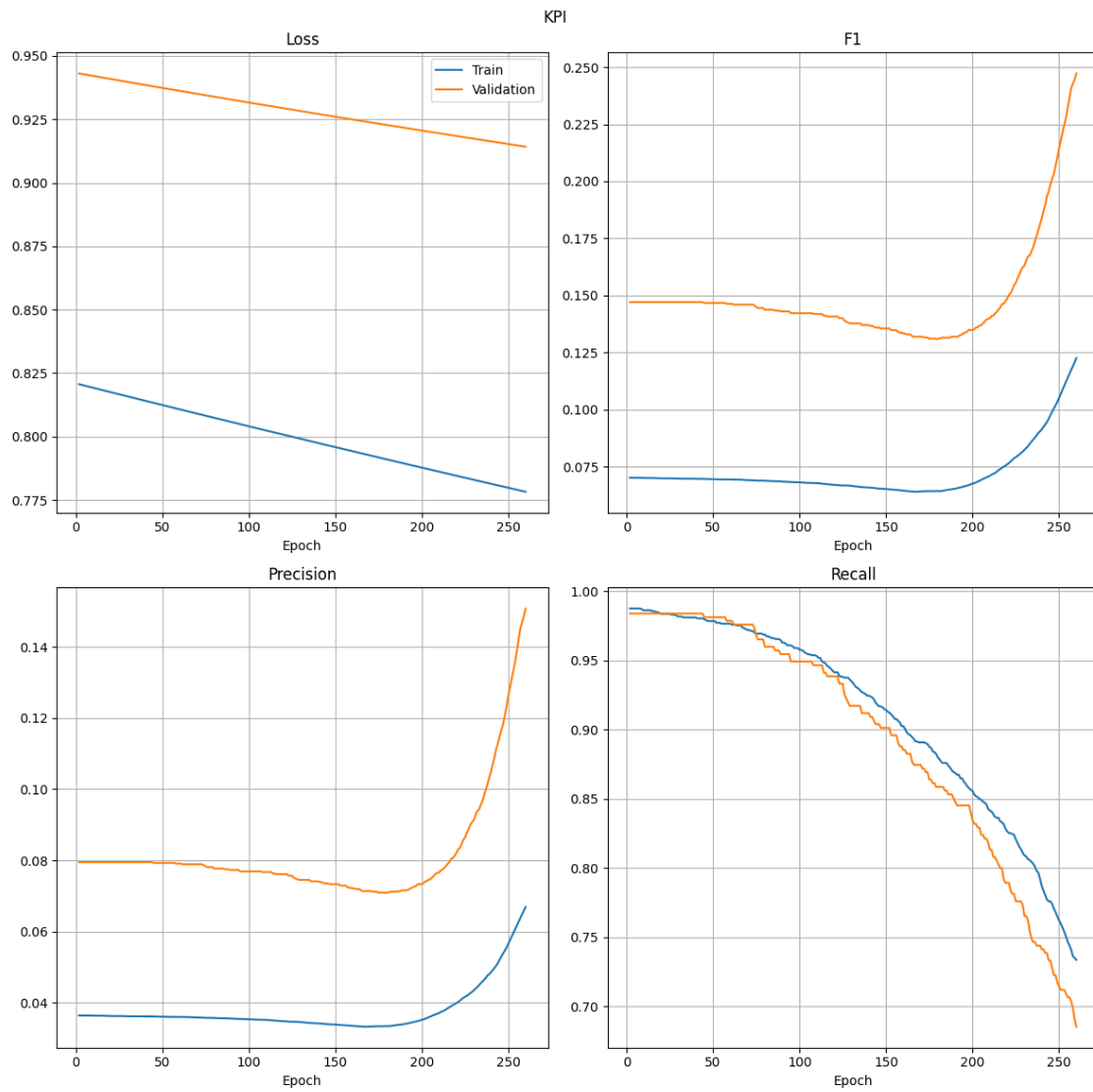
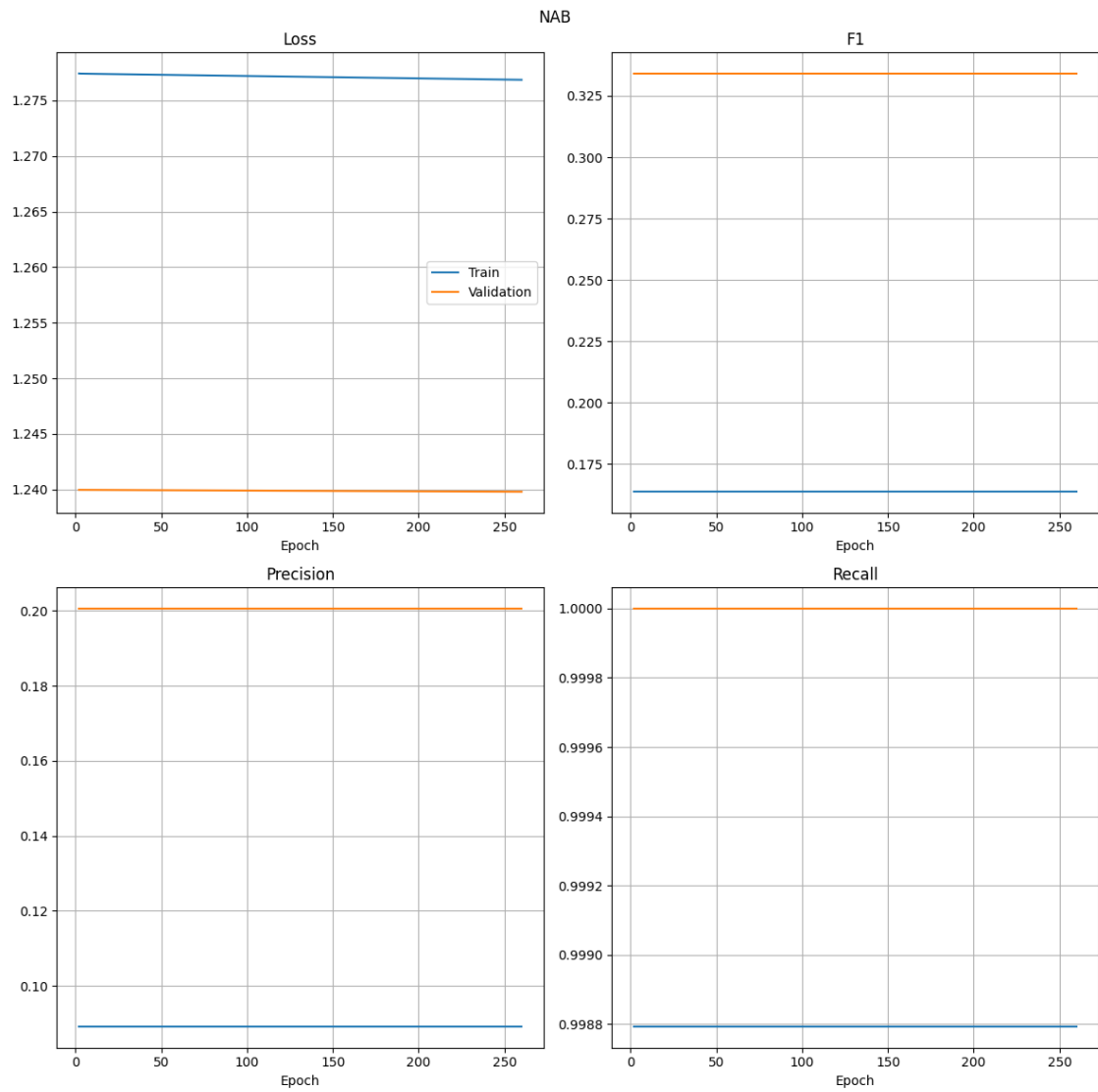


Figure A.6: Training metrics for different epochs for Linear Regression model on NAB dataset



Appendix B

Empirical latency

Figure B.1: Empirical latency for Transformer model on CPU.

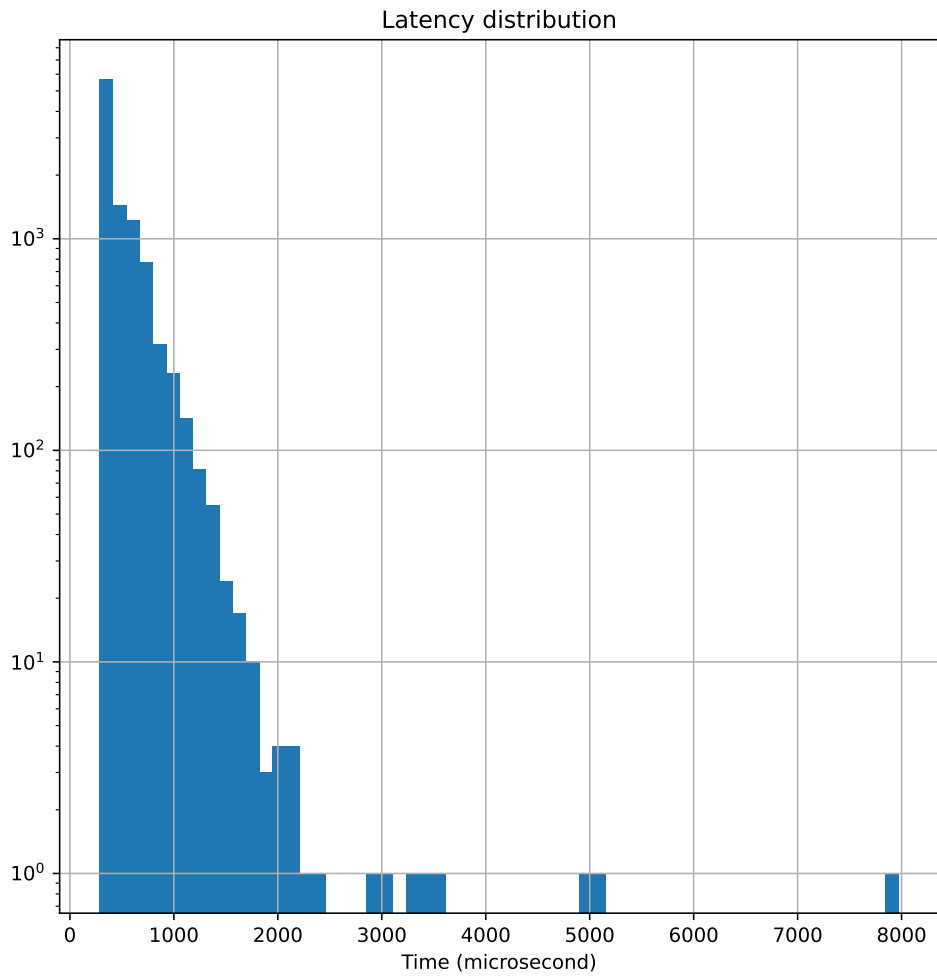


Figure B.2: Empirical latency for Transformer model on GPU.

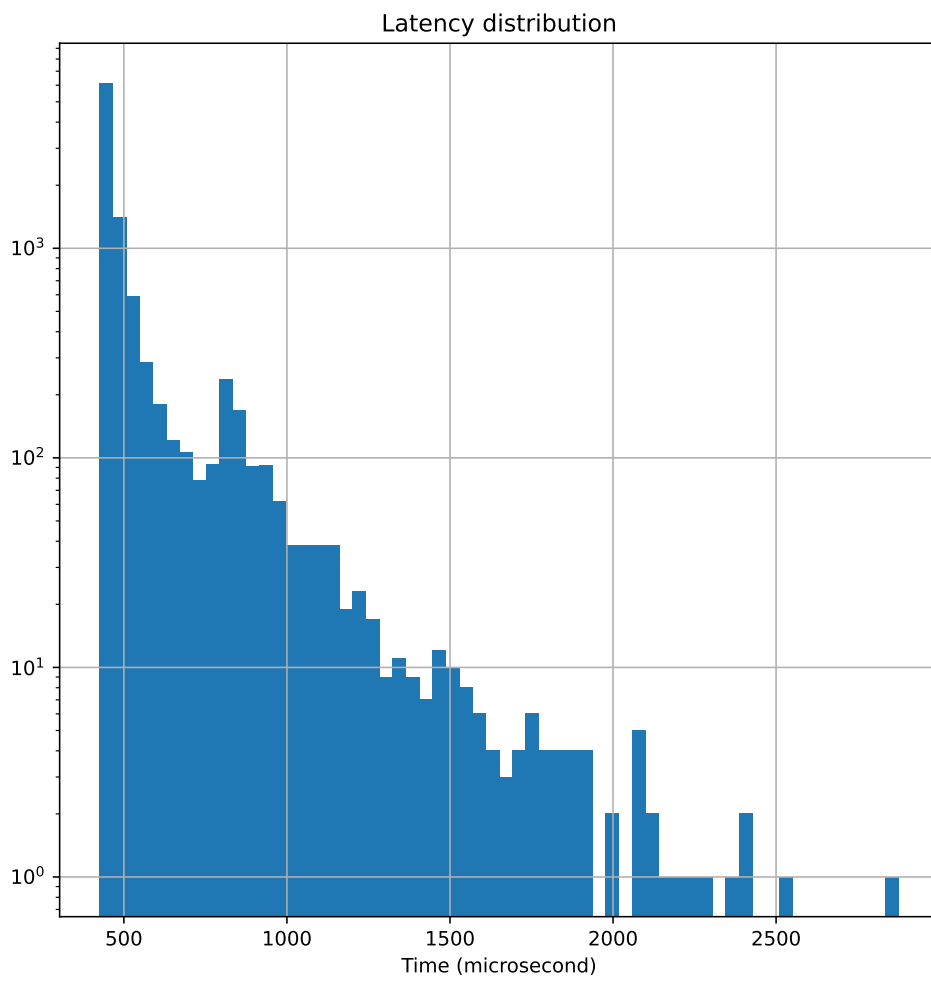


Figure B.3: Empirical latency for Linear Transformer model on CPU.

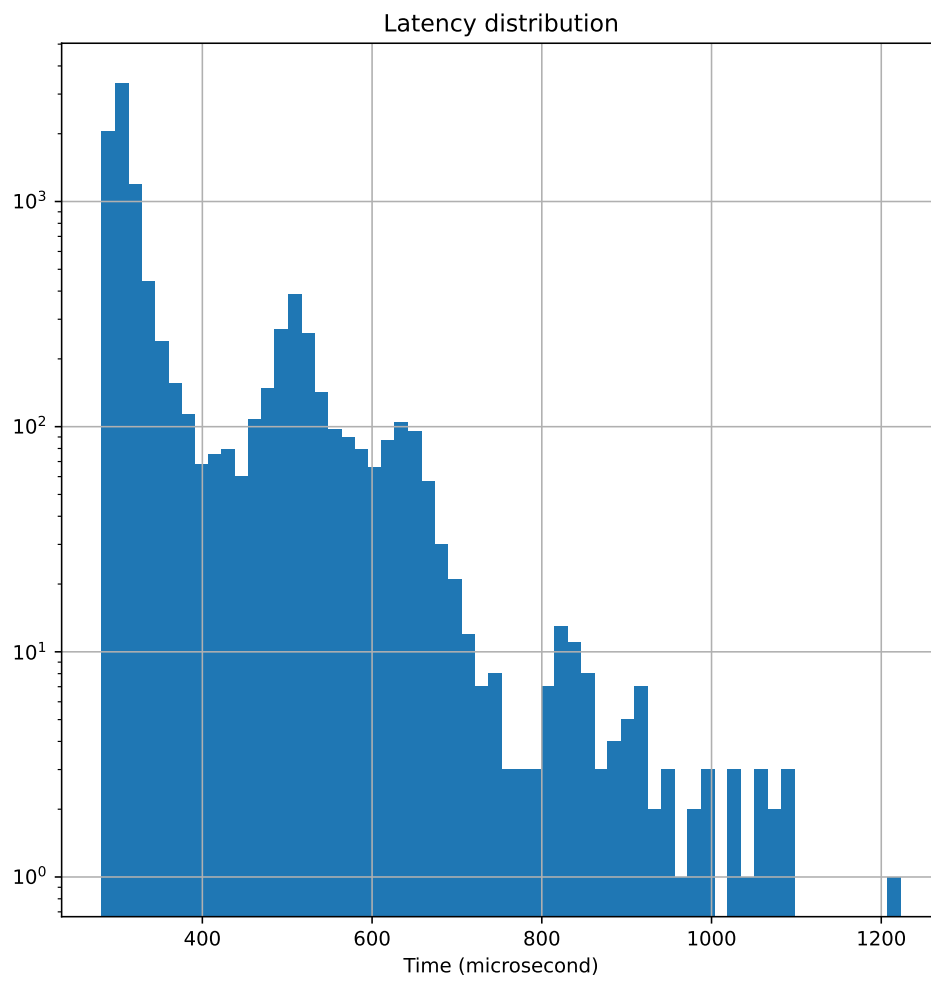
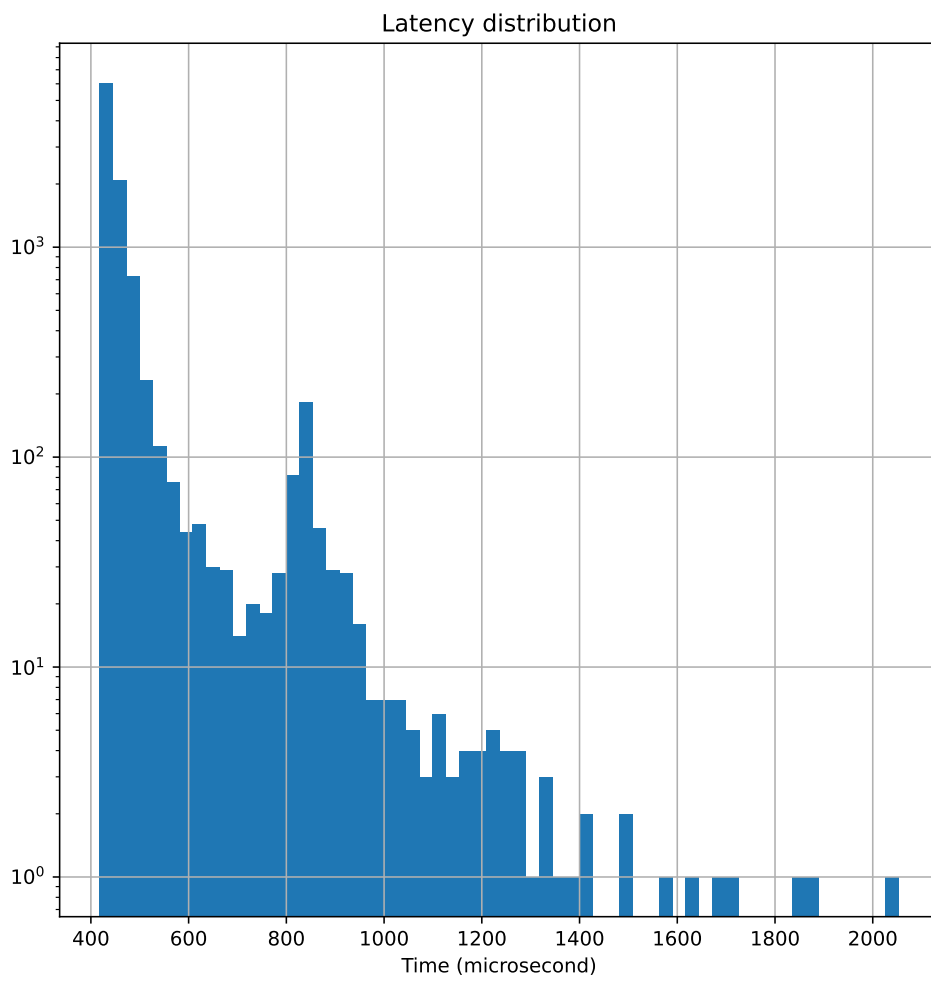


Figure B.4: Empirical latency for Linear Transformer model on GPU.



Bibliography

- [1] Thomas Neil Falkenberry CFA. High frequency data filtering. <https://www.tickdata.com/whitepaper/high-frequency-data-filtering>, Sep 2008.
- [2] Owen Vallis, Jordan Hochenbaum, and Twitter. Introducing practical and robust anomaly detection in a time series. https://blog.twitter.com/engineering/en_us/a/2015/introducing-practical-and-robust-anomaly-detection-in-a-time-series.
- [3] Mohiuddin Ahmed, Nazim Choudhury, and Shahadat Uddin. Anomaly detection on big data in financial markets. In *2017 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM)*, pages 998–1001, 2017.
- [4] Anthony Gillioz, Jacky Casas, Elena Mugellini, and Omar Abou Khaled. Overview of the transformer-based models for nlp tasks. In *2020 15th Conference on Computer Science and Information Systems (FedCSIS)*, pages 179–183, 2020.
- [5] Qingsong Wen, Tian Zhou, Chaoli Zhang, Weiqi Chen, Ziqing Ma, Junchi Yan, and Liang Sun. Transformers in time series: A survey. 2022.
- [6] Andreea-Ingrid Funie, Liucheng Guo, Xinyu Niu, Wayne Luk, and Mark Salmon. Custom framework for run-time trading strategies. In Stephan Wong, Antonio Carlos Beck, Koen Bertels, and Luigi Carro, editors, *Applied Reconfigurable Computing*, pages 154–167, Cham, 2017. Springer International Publishing.
- [7] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need, 2017.
- [8] University of York. Vitis HLS Knowledge Base - Real-Time Systems - York Wiki Service. <https://wiki.york.ac.uk/display/RTS/Vitis+HLS+Knowledge+Base>, July 2020.
- [9] Xilinx Inc. Vitis hls: Pipelining loops. <https://docs.xilinx.com/r/en-US/ug1399-vitis-hls/Design-Principles>, May 2023.
- [10] Julio-Omar Palacio-Niño and Fernando Berzal. Evaluation metrics for unsupervised learning algorithms, 2019.
- [11] Georg Steinbuss and Klemens Böhm. Generating artificial outliers in the absence of genuine ones — a survey. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 15:1 – 37, 2020.
- [12] Sebastian Schmidl, Phillip Wenig, and Thorsten Papenbrock. Anomaly detection in time series: A comprehensive evaluation. *Proc. VLDB Endow.*, 15:1779–1797, 2022.
- [13] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate, 2014.
- [14] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *ArXiv*, abs/1810.04805, 2019.
- [15] Yungi Jeong, Eu-Hui Yang, Jung Hyun Ryu, Imseong Park, and Myung joo Kang. Anomalybert: Self-supervised transformer for time series anomaly detection using data degradation scheme. *ArXiv*, abs/2305.04468, 2023.

- [16] Philipp Dufter, Martin Schmitt, and Hinrich Schütze. Position information in transformers: An overview. *Computational Linguistics*, 48:733–763, 2021.
- [17] Lilian Weng. The transformer family version 2.0. <https://lilianweng.github.io/posts/2023-01-27-the-transformer-family-v2>, Jan 2023.
- [18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition, 2015.
- [19] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer normalization, 2016.
- [20] Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are rnns: Fast autoregressive transformers with linear attention, 2020.
- [21] Shreshth Tuli, Giuliano Casale, and Nicholas R Jennings. TranAD: Deep Transformer Networks for Anomaly Detection in Multivariate Time Series Data. *Proceedings of VLDB*, 15(6):1201–1214, 2022.
- [22] Ailing Zeng, Muxi Chen, Lei Zhang, and Qiang Xu. Are transformers effective for time series forecasting?, 2022.
- [23] Edmond Lezmi and Jiali Xu. Time series forecasting with transformer models and application to asset management. *SSRN Electronic Journal*, 2023.
- [24] Xilinx Inc. Vitis High-Level Synthesis User Guide. <https://docs.xilinx.com/r/en-US/ug1399-vitis-hls/Design-Principles>, May 2023.
- [25] Xilinx Inc. C/RTL Co-Simulation in Vitis HLS. <https://docs.xilinx.com/r/en-US/ug1399-vitis-hls/C/RTL-Co-Simulation-in-Vitis-HLS>, July 2023.
- [26] Intel. Intel high level synthesis: Best practices guide. <https://www.intel.com/content/www/us/en/docs/programmable/683152/21-3/pipeline-loops.html>.
- [27] Leslie N. Smith. Cyclical learning rates for training neural networks, 2015.
- [28] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2014.
- [29] lightning.ai. Learningratefinder — pytorch lightning 2.0.7 documentation. <https://lightning.ai/docs/pytorch/stable/api/lightning.pytorch.callbacks.LearningRateFinder.html>, Aug 2023.
- [30] Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. On the difficulty of training recurrent neural networks, 2012.
- [31] I. J. Good. Rational decisions. *Journal of the Royal Statistical Society: Series B (Methodological)*, 14(1):107–114, Jan 1952.
- [32] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [33] Subutai Ahmad, Alexander Lavin, Scott Purdy, and Zuha Agha. Unsupervised real-time anomaly detection for streaming data. *Neurocomputing*, 262:134–147, 11 2017. [Online; accessed 2023-07-19].
- [34] Zeyan Li, Nengwen Zhao, Shenglin Zhang, Yongqian Sun, Pengfei Chen, Xidao Wen, Minghua Ma, and Dan Pei. Constructing large-scale real-world benchmark datasets for aiops, 2022.
- [35] Adamantios Ntakaris, Martin Magris, Juho Kannianen, Moncef Gabbouj, and Alexandros Iosifidis. Benchmark dataset for mid-price forecasting of limit order book data with machine learning methods. 2017.
- [36] Stéphane Crépey, Noureddine Lehdili, Nisrine Madhar, and Maud Thomas. Anomaly Detection in Financial Time Series by Principal Component Analysis and Neural Networks. *Algorithms*, 15(10):385, oct 19 2022. [Online; accessed 2023-07-19].