

# **Hardware Accelerated High-Frequency Trading System**

## **Design Document**

TA:

Mingjia Huo

Team Members:

Siyi Yu(siyiyu2)

Kevin Lim (wzlim2)

Richard Deng (ruichao4)

September 29, 2022

# Contents

## **1 Introduction**

- 1.1 Problem
- 1.2 Solution
- 1.3 Visualization
- 1.4 High Level Requirements

## **2 Design**

- 2.1 Block Diagram
- 2.2 Subsystem Overview
- 2.3 Subsystem Requirements
- 2.4 Tolerance Analysis

## **3 Cost & Schedule**

- 3.1 Cost Analysis
- 3.2 Schedule

## **4 Ethics & Safety**

- 4.1 Ethics
- 4.2 Safety

## **5 References**

# 1 Introduction

## 1.1 Problem

Modern electronic markets have been characterized by a relentless drive towards faster decision making. Significant technological investments have led to dramatic improvements in latency, the delay between a trading decision and the resulting trade execution. We describe a theoretical model for the quantitative valuation of latency. A low latency for the communication, which refers to sending the order to exchange and receiving the order from exchange, is critical since sometimes a delay of milliseconds could result in a loss of millions dollars.

“Low latency” in a contemporary electronic market would be qualified as under 10 milliseconds, “ultra low latency” as under 1 millisecond. This change represents a dramatic reduction by five orders of magnitude. To put this in perspective, human reaction time is thought to be in the hundreds of milliseconds.

In the financial market today there is a lot of need to optimize the trading/execution latency to support automated quantitative trading systems. While most of the computer software runs on generic operating systems and the CPU executing the logic has many parts of unnecessary instructions/procedures, the automated trading strategy can be highly optimized with hardware to achieve low-latency and high-frequency.

## 1.2 Solution

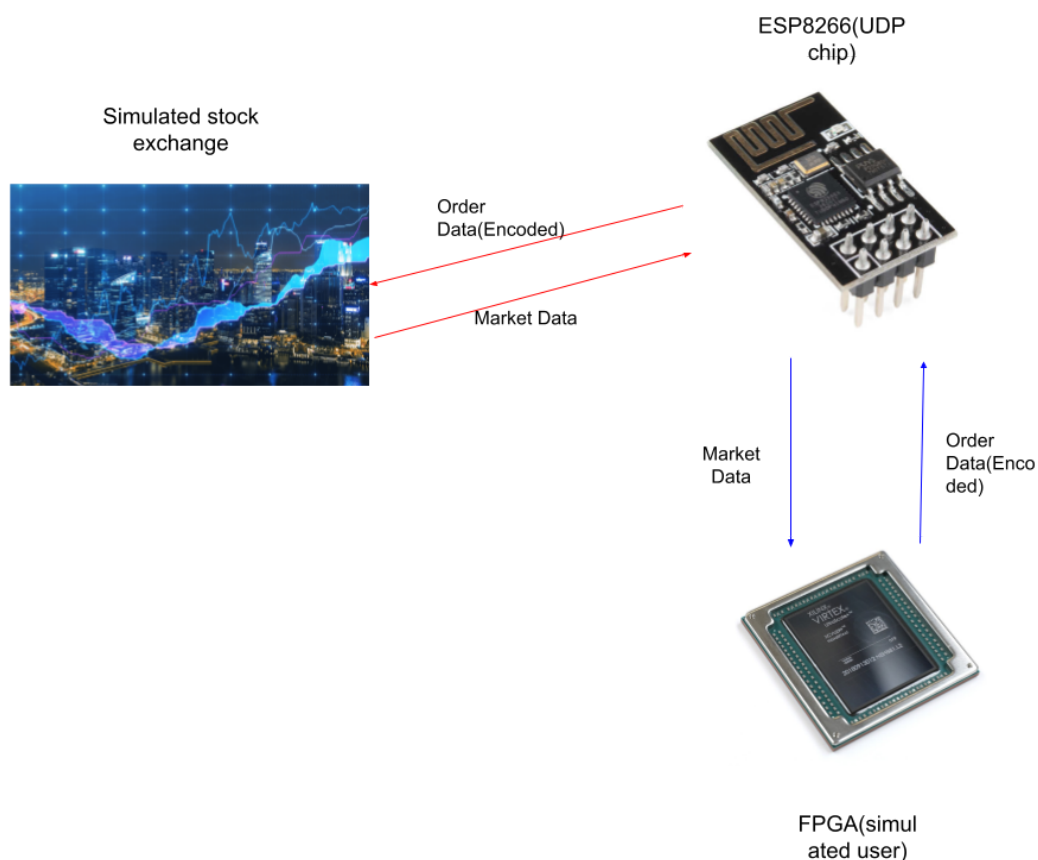
The purpose of this project is to prototype a hardware based trading tool that can perform such a high-frequency and low-latency trading. while being able to communicate with the exchange using basic network protocol.

We plan to build a trading system that uses one PCB to connect to a fake exchange and get/send binary market data and use highly optimized FPGA to consume this market data and make decisions based on the data. Since we will not be able to directly connect to the exchanges without approval, we are also going to simulate an exchange using software. Most of the existing low-latency trading strategy will be implemented using FPGA solely. However, our project will utilize an ESP8266 chip for networking purposes.

The ESP8266 is a low-cost Wi-Fi microchip, with built-in TCP/IP networking software, and microcontroller capability. FPGAs which could accomplish this task are pretty high end and relatively expensive. FPGAs tend to offer vastly greater flexibility which we do not need for this project. It will help us bring down the cost while maintaining a relatively good performance of the system.

### 1.3 Visualization

Figure 1 illustrates the overview of our project design.



**Figure 1: User diagram**

- Simulated stock exchange will be implemented fully using software using C++ code
- We are planning on to use USB for the data transmission between FPGA and ESP8266 chip

- We will use UDP protocol for the communication between ESP8266 and simulated stock exchange

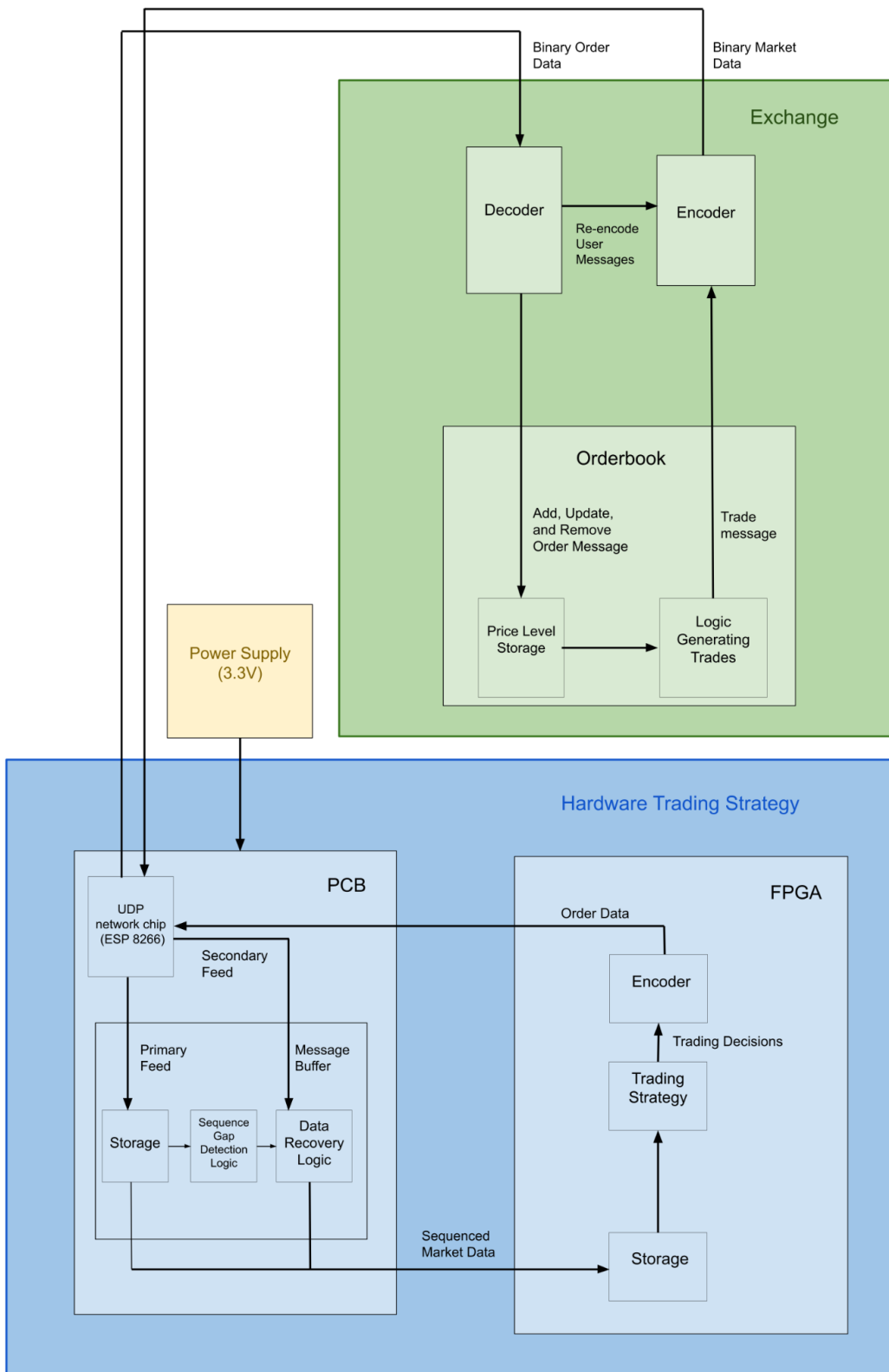
## **1.4 High-Level Requirements**

For this project, we have 3 high-level requirements that we aim to achieve:

1. 80% orders being successfully sent to exchange and received by users. We will use about 100MB of data which we download from the NYSE as a dataset.
2. 30% faster than using purely software implementation of this system. We will be implement a python software for reference
3. 80% of orders are correctly processed by the exchange. We will use about 100MB of data which we download from the NYSE as a dataset.

## 2 Design

### 2.1 Block Diagram



## 2.2 Subsystem Overview

Our project consists of three major systems, simulated exchange, networking hardware, and FPGA trading strategy.

### 2.2.1 Simulated Exchange

This is a simulation of real-world exchange. It can receive orders from market participants and automatically generate trade messages and broadcast it back to market participants. This simulated exchange consists of three parts:

- I. Order Data Decoder. This decoder will decode the binary order data sent by market participants into internal, more organized, human-readable data and feed it to the orderbook.
- II. Orderbook. This is where all the orders from different market participants are stored, and it will generate a trade if the bid side and ask side meets.
- III. Market Data Encoder. This encoder will consume the data from the orderbook and encode it into binary market data and broadcast it to market participants.

The simulated exchange will be a piece of software that simulates an exchange in real-time. This simulated exchange will only have one security to trade and will be receiving three types of binary-encoded messages (add order, cancel order, update order) from the market participants using a certain protocol, building a full-depth limit order book based on the order received, automatically trade two orders when the bid side and ask side meet, and re-broadcasting the trade message with the other three client messages in binary encoded form as market data messages.

To test our strategy we will simulate the market data received from market participants, constantly adding limit ask/bid at the same price level in high frequency. We will expect 100% of the orders to be processed correctly. This means an ask order and a bid order with the same price will not be existing in the order book, since it will generate a trade. This part of the project is designed to be run completely by software so there is no power supply unit involved in the subsystem.

### 2.2.2 Networking Hardware

Networking hardware will be responsible for the communication between users and the stock exchange. It will be made of a PCB with ESP8266 chip on it and other parts/ports to communicate with the simulated exchange through the network. We plan to use UDP since it is what most exchanges will use. This part will be optimized to reduce the connectivity latency between the fake exchange and the trading system and feed the data to the FPGA.

We would utilize ESP8266WiFi.h and WiFiUdp.h these two libraries to enable the UDP communication. The first library ESP8266WiFi.h is required by default if we are using ESP8266's Wi-Fi. The second one WiFiUdp.h is needed specifically for programming of UDP routines.

Without this component, there will be no way for users and stock exchanges to communicate with each other. Since ESP8266 requires a power supply for it to be functional, we would need a power supply for this subsystem. The ESP8266 requires a 3.3V power supply and 3.3V logic levels for communication and a maximum of 170mA current.



### 2.2.3 FPGA Trading Strategy

This part is an FPGA board that consumes market data from the networking hardware and generates decisions based on the market data. This consists of two parts:

- I. Decision maker (Strategy). This part is the core trading strategy that takes in market data in real-time and makes decisions based on some logic that directly manipulates the binary market data to optimize speed.
- II. Order Encoder. This part will take in the decisions that have been made from our strategy, and encode it to binary and send it back to the networking chip, which will eventually be sent to the exchange.

In order to simplify the process(since trading algorithm is not the key part of this project), we plan to implement two simple strategies:

1. High-frequency market-making: This is a liquidity-providing trading strategy that simultaneously generates many bids and asks for a security at ultra-low latency while maintaining a relatively neutral position. When put into implementation, it will make a spread of \$a ~ \$b by sending limit bid at \$a and limit ask at \$b, adjust the spread based on Best Bid Offer (BBO) market data. This part is to test how fast our trading system can be adjusted based on real-time information changes.
2. High-filling rate limit order: A limit order is used to buy or sell a security at a predetermined price and will not execute unless the security's price meets those qualifications. A High-filling rate limit order trading strategy sends a limit bid order whenever there is a limit ask order at \$a. This part is designed to test the latency of our trading system and is extremely useful in options trading where the bid-ask spread is very large and has a low order filling rate.

These two trading algorithms will be sufficient enough for us to mock the real world trading cases and allow us to test the system.

### 2.2.4 Power Supply

We will need a 3.3V power supply for our PCB board.

## 2.3 Subsystem Requirements

### 2.3.1 Simulated Exchange

This simulated exchange will be purely software that is optimized for its speed and correctness. For speed optimization, it will be written in C++ and optimized as much as possible.

- I. Order Data Decoder. The decoder should be able to decode the order data sent by market participants in 100% accuracy, and be able to reject invalid orders. We will design a specific protocol for this decoder, and reject any orders that don't follow this protocol. We will also build in an authentication system by forcing the user to send a key as part of the binary message, and our system should be able to verify whether this user is authorized to trade in our exchange.

Requirements	Verification
80% accuracy for order data decoding	Encode several messages according to the protocol, send it to the exchange and check whether the data we sent is in the system
Reject binary data that doesn't follow the protocol	Send some random binary data to see whether it's rejected by the exchange
Reject unauthorized users	Send order data using random authentication key to see whether it's rejected by the exchange

- II. Orderbook. The orderbook should be able to store valid order data in our system into two sides, the ask(sell) side and bid(buy) side. To make sure to update the same order later, we will generate an Order Reference Number that associates with each individual order, and store this relation into a map. It will also automatically generate trade when the bid side and ask side meets, i.e. whenever there exists some ask order with price  $\$x$  with quantity  $a$  and some bid order with price  $\$y$  with quantity  $b$  that satisfies  $(x \leq y)$ , a trade of quantity  $\min(a, b)$  will be generated.

Requirements	Verification
--------------	--------------

orders are being stored in orderbook and be able to modify	Send orders to the exchange to see whether it's in our system. Modify it later and see whether it's modified in our system.
Trade is generated automatically and properly	Send random orders and make sure best bid price and best ask price never meets

III. Market Data Encoder. The encoder should be able to encode the market data in 100% accuracy and send it back to market participants. To make sure users can fully replicate everything happening in exchange in a time-sequenced manner and also enable users to build sequence gap detection, there will be a sequence number in each market data message. This encoder will also send data in two feeds through different ports/hosts to enable lost packet recovery.

Requirements	Verification
80% accuracy for market data encoding and sequenced	Generate several artificial events in the exchange, and using another test software to receive and decode this data, making sure it's what's happening in the exchange and sequenced correctly
Being able to recover data from backup feed	Receive and check whether both feeds are sending the same data

### 2.3.2 Networking Hardware

This part is generally a PCB that handles networking between the simulated exchange and the FPGA trading strategy. This system consists of two parts:

- I. UDP Networking Chip. This is a ESP8266 chip that processes network packets using UDP protocol.

Requirements	Verification
Keep the average latency time within 5ms	<ul style="list-style-type: none"> <li>● Send about 10,000 packages and develop a simple application to take the average latency.</li> <li>● The code testing code will be like below:</li> </ul> <pre>while (1) { char buf[BUFLLEN]; recvfrom(s, buf, BUFLLEN);</pre>

	<pre>sendto(s, buf, BUFLLEN); }</pre>
Successful receive and process 80% order data sent from FPGA(user)	<ul style="list-style-type: none"> <li>• Feed UDP network chip with 10,000 packages check if the receive and send out process went thoroughly as we planned</li> </ul>

II. Sequence Gap Handler. This part is responsible for detecting the sequence gap when packets are dropped given the unreliable nature of UDP protocol, and it is also responsible for recovering the lost packets from backup data feed.

Requirements	Verification
Successfully detect all lost packages	<ul style="list-style-type: none"> <li>• Feed 100 sequenced packages with 20% of them being intentionally dropped randomly</li> <li>• Verify whether the lost packets detected by the handler are the same with those being dropped</li> </ul>
Successfully recover 100% of the lost packages from backup data feed	<ul style="list-style-type: none"> <li>• Send 100 sequenced packages on primary feed with 20% of them being intentionally dropped on the primary feed, send the dropped packets through secondary feed</li> <li>• Check whether the handler can get all 100 packets in sequence</li> </ul>

### 2.3.3 FPGA Trading Strategy

This FPGA trading strategy simply needs to process the data as fast as possible while maintaining correctness, and make decisions based on the data.

I. Decision maker (Strategy). This strategy needs to process the data fast and correctly. Compared to the same strategy using purely software, the hardware strategy needs to be making the same decision while processing the same data, and process it 30% faster than the software ones.

Requirements	Verification
--------------	--------------

Correctly process the data	<ul style="list-style-type: none"> <li>● Implement same strategy using python, using the same input test data</li> <li>● Verify whether the python strategy outputs the same decision as the FPGA ones</li> </ul>
30% faster than software strategy	<ul style="list-style-type: none"> <li>● Implement same strategy using python, using the same input test data</li> <li>● Verify whether the FPGA strategy is 30% faster than the python one's</li> </ul>

II. Order Data Encoder. The encoder should be able to encode the order decision in 100% accuracy and send it to the networking hardware. This part needs to send authorized orders correctly, following the same protocol as the exchange's order data encoder.

Requirements	Verification
Order is encoded correctly	<ul style="list-style-type: none"> <li>● Send a few test orders from FPGA to the exchange</li> <li>● Verify whether it can be decoded correctly and appear on orderbook</li> </ul>
User is authorized (correct authentication key)	<ul style="list-style-type: none"> <li>● Send a few test orders from FPGA to the exchange</li> <li>● Verify whether it can be received and authorized to trade and appear on exchange's order book</li> </ul>

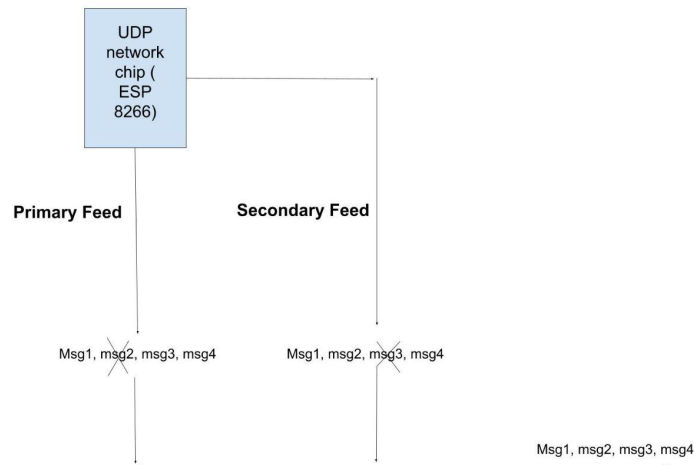
### 2.3.3 3.3V power supply

3.3V power supply will be provided by the ECE department which will be used to supply our PCB board.

Requirements	Verification
Supply +3.3V +-0.3V to the chips and PCB	<ul style="list-style-type: none"> <li>● Use a multimeter to check if voltage output is within the acceptable range</li> </ul>

## 2.4 Tolerance Analysis

Since we are using UDP as our network communication protocol, it might fail to deliver some of the packets. It is very critical for us that every single packet is received, otherwise we cannot process the market data correctly. However, we designed a sequence gap detection and recovery method built in the PCB, only when both primary data feed and secondary data feed fails on the same packet we will experience a failure, thus our system has high tolerance for network failure. A diagram shown below demonstrates when message 2 drops on primary feed and message 3 drops on secondary feed, we can still recover all 4 messages after combining them.



### 3 Cost & Schedule

#### 3.1 Cost Analysis

##### 3.1.1 Labor

Name	Hourly Rate(\$)	Total Hours	Total
Ricahrd Deng	33	144	4752
Kevin Lim	33	144	4752
Siyi Yu	33	144	4752
Total	99	432	14,256

##### 3.1.2 Parts

Parts	Part#	Quantity	Unit Cost(\$)	Cost(\$)
FPGA	DE10-Nano	1	\$135	\$135
PCB	-	1	\$0	\$0
OR gate	SN74HC32N	20	\$0.5	\$10
AND gate	CD4085BE	20	\$0.5	\$10
Register	HCF40100BE	20	\$2.5	\$50.0

#### 3.2 Schedule

Week	Task	Responsibility
9/26	Complete Design Document	Siyi Kevin Richard
	FPGA Trading Algorithm Research	Siyi Kevin

		Richard
10/3	ESP8266 UDP design	Siyi Kevin Richard
	Build Power Supply	Siyi Kevin Richard
10/10	PCB design	Siyi Kevin Richard
	Trading Algorithm Implementation	Siyi Kevin Richard
10/17	Trading Algorithm Implementation	Siyi Kevin Richard
	Testing Trading Algorithm	Siyi Kevin Richard
	PCB design	Siyi Kevin Richard
	Exchange Implementation	Siyi Kevin Richard
10/24	Exchange Implementation and testing	Siyi Kevin Richard
10/31	Integrate PCB to the system	Siyi Kevin Richard
11/7	Testing the system	Siyi Kevin Richard



11/14	Debug and fix any errors which may occur	Siyi Kevin Richard
11/21	Final test for the whole system	Siyi Kevin Richard
11/28	Prepare for final presentation	Siyi Kevin Richard
12/5	Final presentation	Siyi Kevin Richard

## 4 Ethics & Safety

### 4.1 Ethics

We have identified 3 major concerns of ethics from the Association for Computing Machinery (ACM) Code of Ethics and Professional Conduct for our project[3].

Respect privacy (ACM 1.6)

- Our design may enable users to collect and trade personal information, which violates the privacy of the user. We shall protect the privacy of our users and make sure that no one is able to access personal information.

Honor Confidentiality (ACM 1.7)

- Our design may process highly valuable information such as trade secrets, financial information, and client data. We must not share this data with anyone, and we must not access the information ourselves.

Design and Implement Systems that are Robust and Usably Secure (ACM 2.9)

- A fragile or easily breakable design will likely lead to data leakage among other problems. We should make sure that our design is robust and all information is secure.

The Code of Ethics from the Institute of Electrical and Electronic Engineers (IEEE) also provides important guidelines for our project. We shall properly credit the contribution of others to our project (IEEE #5), uphold safety standards and disclose any safety concerns to users (IEEE #1), and to not tolerate any form of discrimination with our project (IEEE #7).[4]

Our design aims to reduce the latency of trading and make financial operations more efficient. As a team, our goal is to make this design safe and accessible while respecting the code of ethics and other moral concerns.

### 4.2 Safety

One concern regarding safety within our project is the power supply. If not handled properly, there is a risk of overcharging which could damage the PCB. We will make sure to test our power supply to not exceed 3.3V.

Our team takes ethical and safety concerns seriously and we will strictly follow all ethics and safety guidelines, including those listed above and other concerns that are not covered by these topics.

## 5 References

[1] Ultra Low Latency Networking with FPGA, 2020.  
[https://www.youtube.com/watch?v=32cK\\_yDcouQ](https://www.youtube.com/watch?v=32cK_yDcouQ)

[2] UDP - ESP8266 Arduino Core, 2017.  
<https://arduino-esp8266.readthedocs.io/en/latest/esp8266wifi/udp-examples.html>

[3] Association for Computing Machinery, “ACM Code of Ethics and Professional Conduct”, 2018. [Online]. Available: <https://www.acm.org/code-of-ethics>.

[4] Institute of Electrical and Electronic Engineers, “IEEE Code of Ethics”, 2017. [Online]. Available: <http://www.ieee.org/about/corporate/governance/p7-8.html>.