

Final Project Report

[EE488] Introduction to Reinforcement Learning

Finance Trading Policy

- **Department:** School of Electrical Engineering
 - **Student ID:** 20200625
 - **Name:** Yejun Joo
-

1. Code Implementation Explanation

- Environment.py
 - self.observation_space
 - Used spaces from gymnasium library to define observation space. Observation state is set to a dictionary type, including information of agent itself and that from market.
 - self.action_dim
 - Also used spaces for action space definition. Action space is 1-dimension array with n-elements. Here n corresponds to action dimension, which is same with the number of tickers. Each element of action means the number of stock to trade, while positive value meaning buying and negative value meaning selling stocks.
 - reset()
 - I designed the code to randomly select the initial balance. However, it is not continuously random. It is discretized in unit of \$2500 to prevent unstable learning with high variance, resulting in long time to converge.
 - It resets the timestep, but for policies that are using state histories, timestep needs to be start after the length of history, meaning that it has to have enough history to recall.
 - step()
 - First it checks the last prices of each stock. Then for every stocks, it computes the maximum number of stock that can buy by comparing the last price and current balance. It also checks the current holding stocks to know the maximum number of stocks to sell.
 - Then it trades the stock based on the action, but here, if the action exceeds the possible trading number, it just clips the action to the maximum value. That's why at the test plot, the agent is keep selling 50 stocks even they don't have any.

- I added trading fee to prevent the agent buying and selling meaninglessly.
- I also check if it outputs impossible number of trading(overflow).
- runner.py
 - make_env()
 - By using gymnasium library, I flattened the observation to easily feed it in to PPO network. Also normalized the observation state to make the network learn well. I later save the observation statistics file and load in tester.py code.
 - Using RescaleAction(), I can make the PPO actor net output values between -1 to 1, but still actually using those value in range of maximum tradable stock numbers.
 - I first store every transition data into PPO using Storage class, then update after every epoch to enable stable learning.
- tester.py
 - Its a code to test the policy. It make the environment similar way with runner does. It loads pre-trained policy and also the pre-saved observation state statistics for normalization.
- ppo.py
 - It has Actor for actor network, Critic for critic network, Storage to save all the transition data in buffer and then rollout later to update ppo.
 - _train_step()
 - It computes deltas(td error) and advantage from transition data in compute_returns() function in Storage class.
 - Then, by rolling out these datas, it computes the ratio of action log probability between old policy and updated policy, which means the importance factor. Then it compute the surrogate loss by multiplying advantage, then clip the loss. Finally adding the entropy loss, we step the optimizer.
 - I added the value loss here, but it is independent with the ppo loss, and it only back propagates to value net.
- dqn.py
 - DQN is only possible for discrete actions, and the number of discrete actions explode if there are many tickers. So first, I discretize the actions into 21 steps, which means with maximum trading limitation of 50 stocks, each step is a 5 unit trade.
 - For multi ticker system, I implemented several action net heads, independent network for each tickers. This makes the DQN algorithm possible for large number of discrete actions.

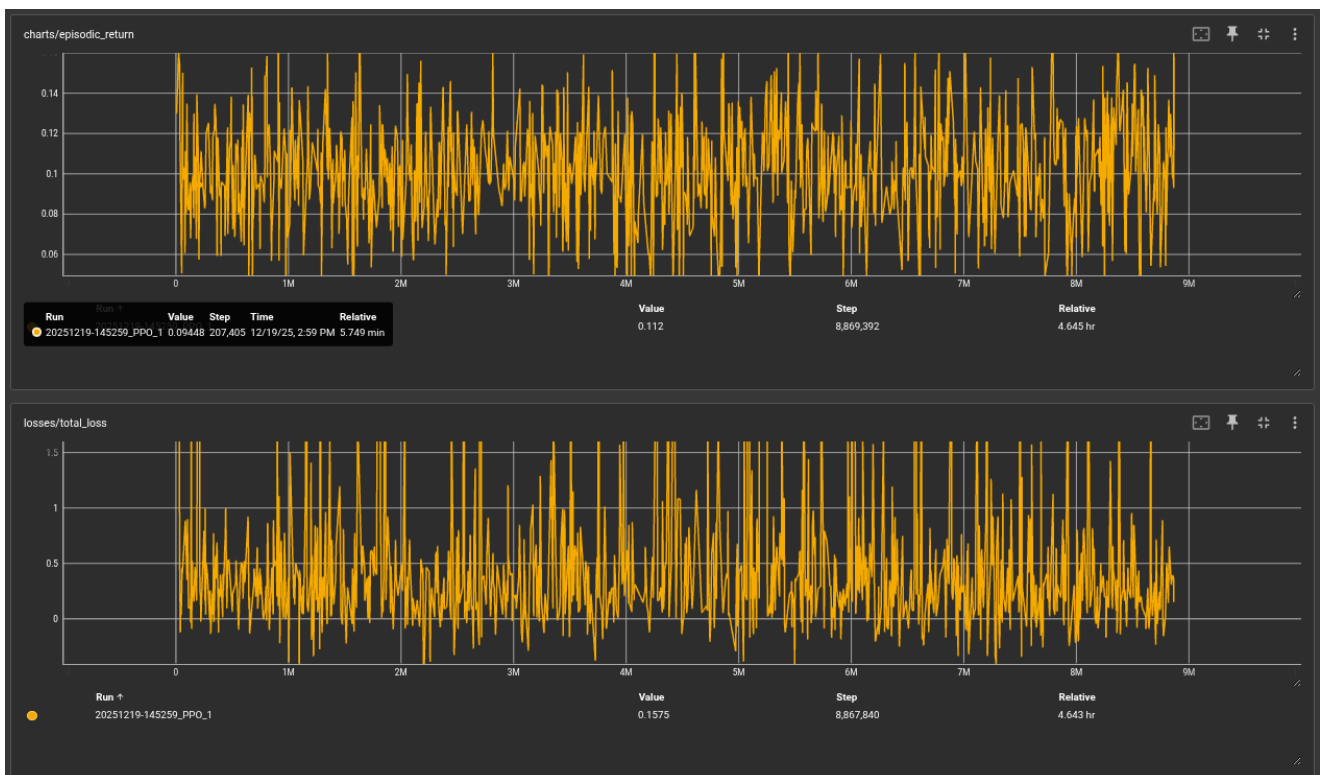
2. Experiments - Performance Upgrade, Single Ticker

Baseline Setup

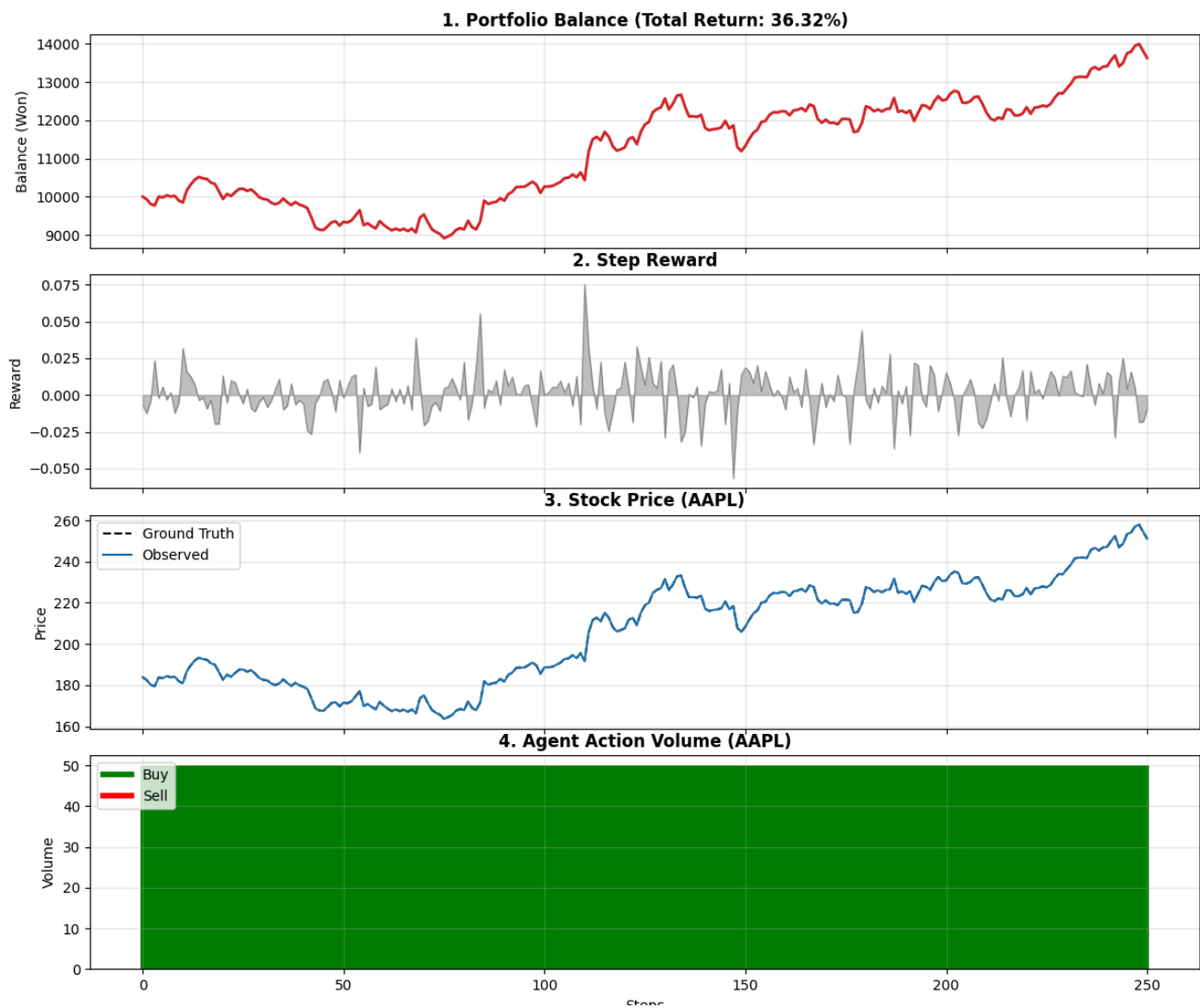
- Observation States
 - agent
 - current balance
 - currently holding stock's number
 - market
 - open price
 - high price
 - low price
 - close price
 - volume
- Reward Function
 - 1 day portfolio difference divided by initial balance
- Dataset
 - 2019-2023 APPL stock
- Algorithm
 - PPO

Results

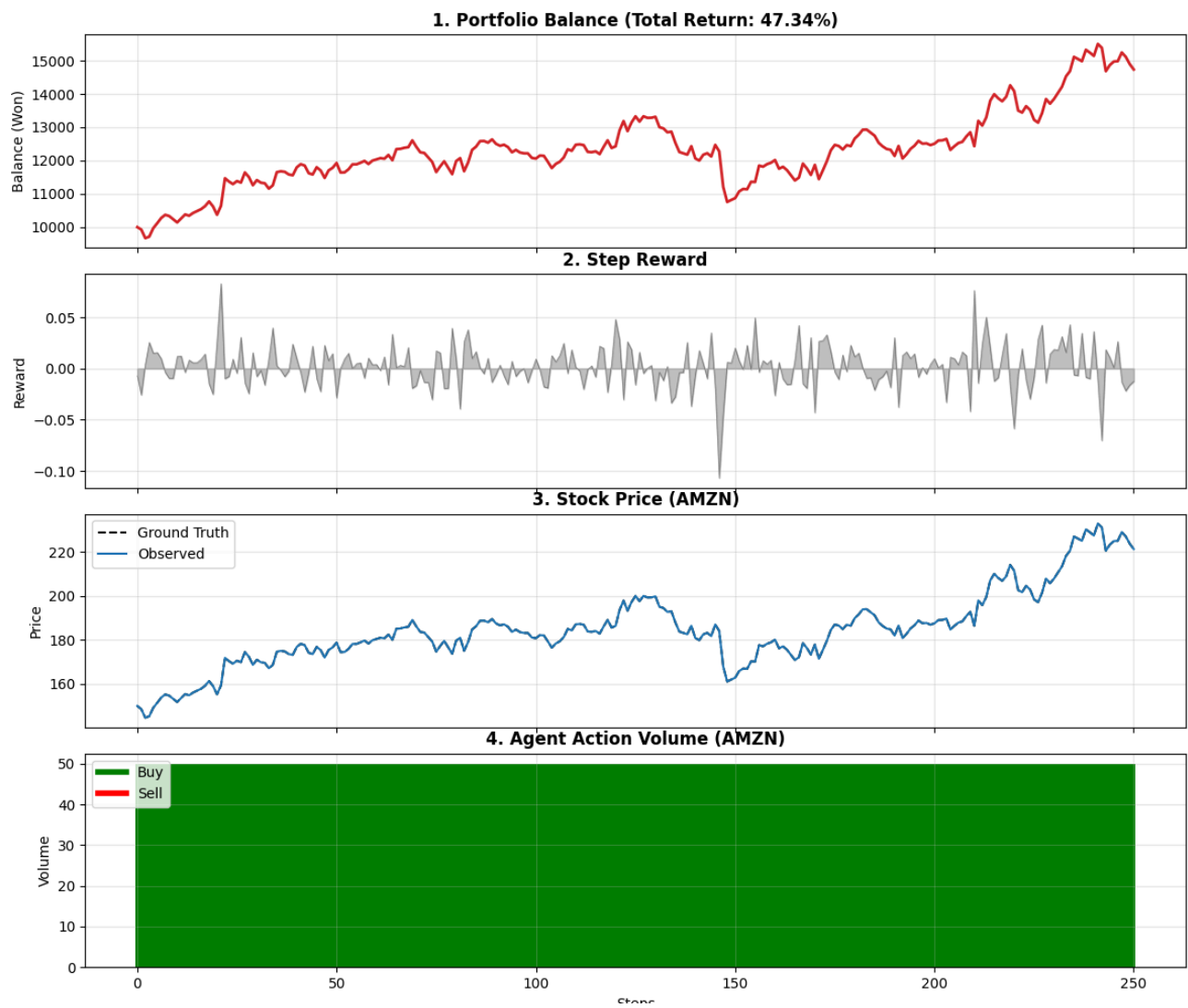
- Learning graph



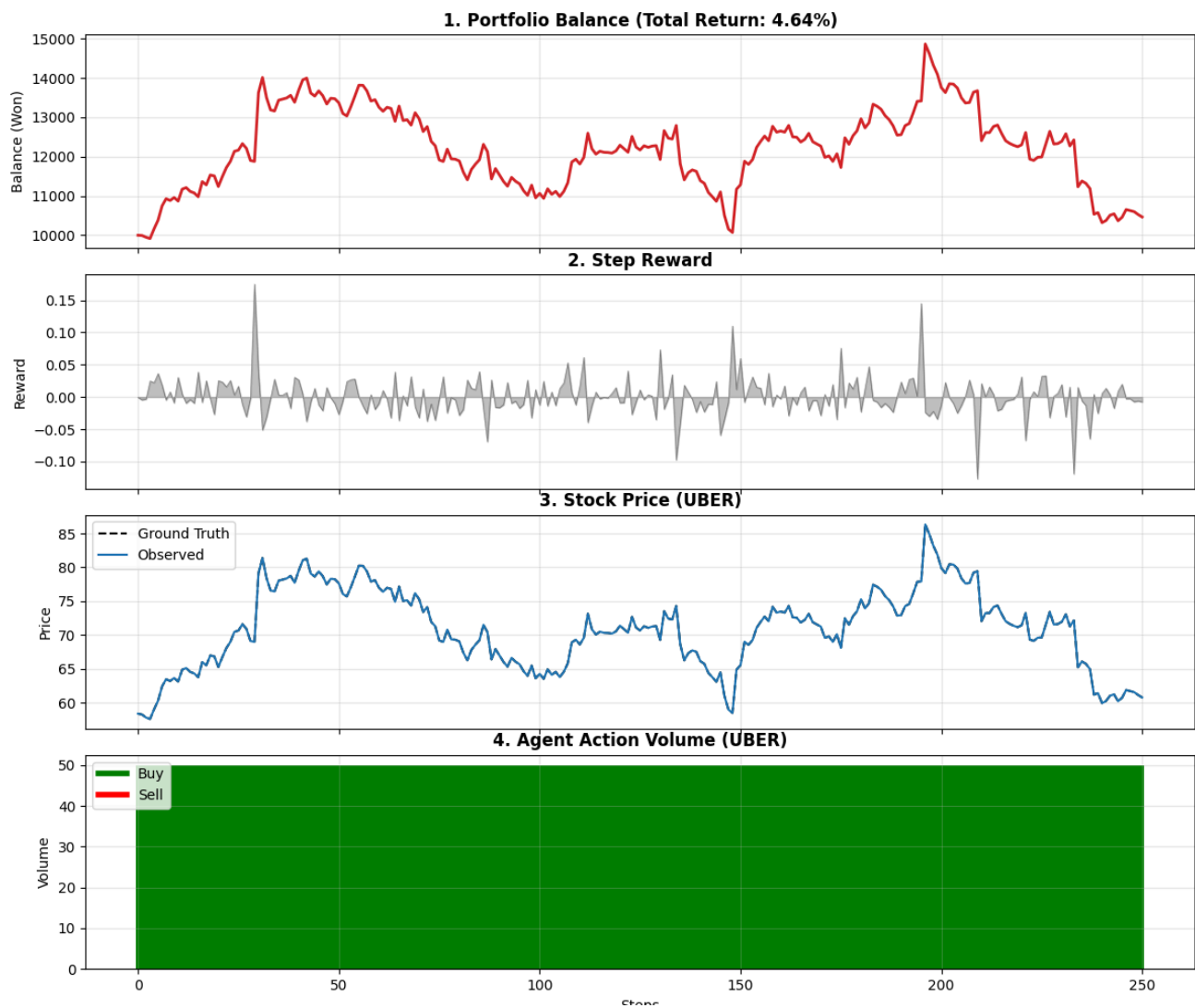
- Same stock; APPL 2024



- **Different stock; AMZN 2024

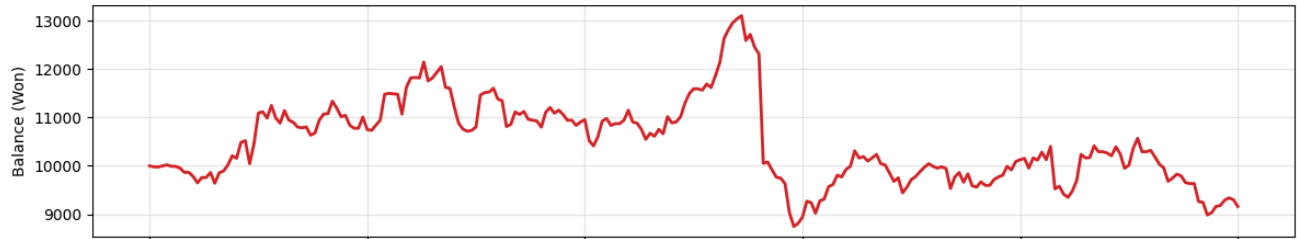


- Different stock; UBER 2024

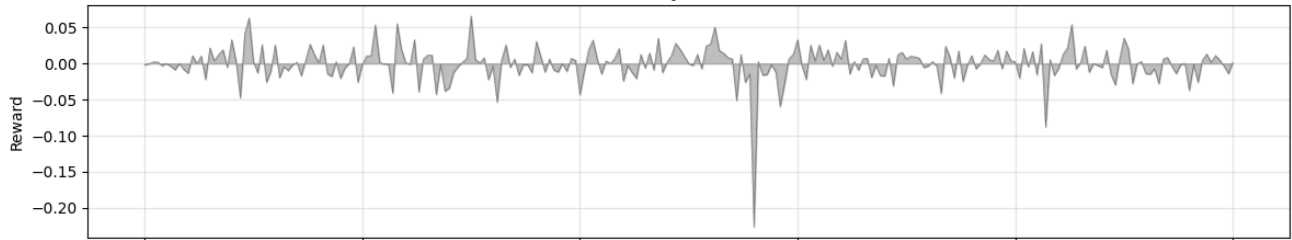


- Different stock; F 2024

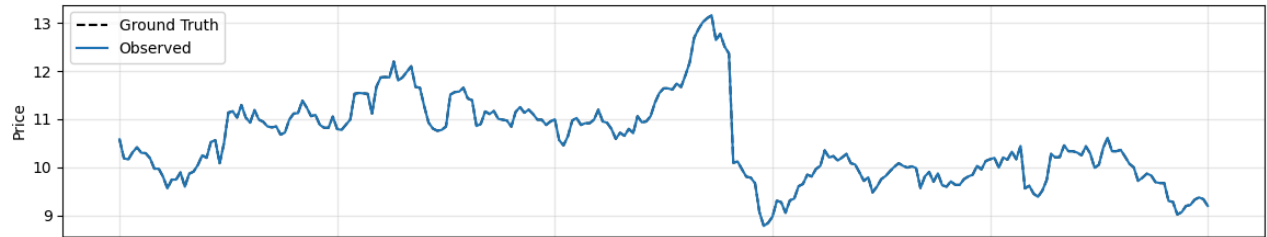
1. Portfolio Balance (Total Return: -8.38%)



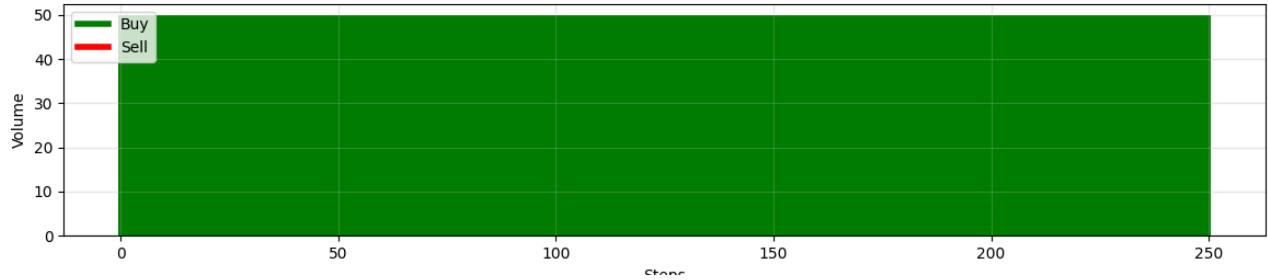
2. Step Reward

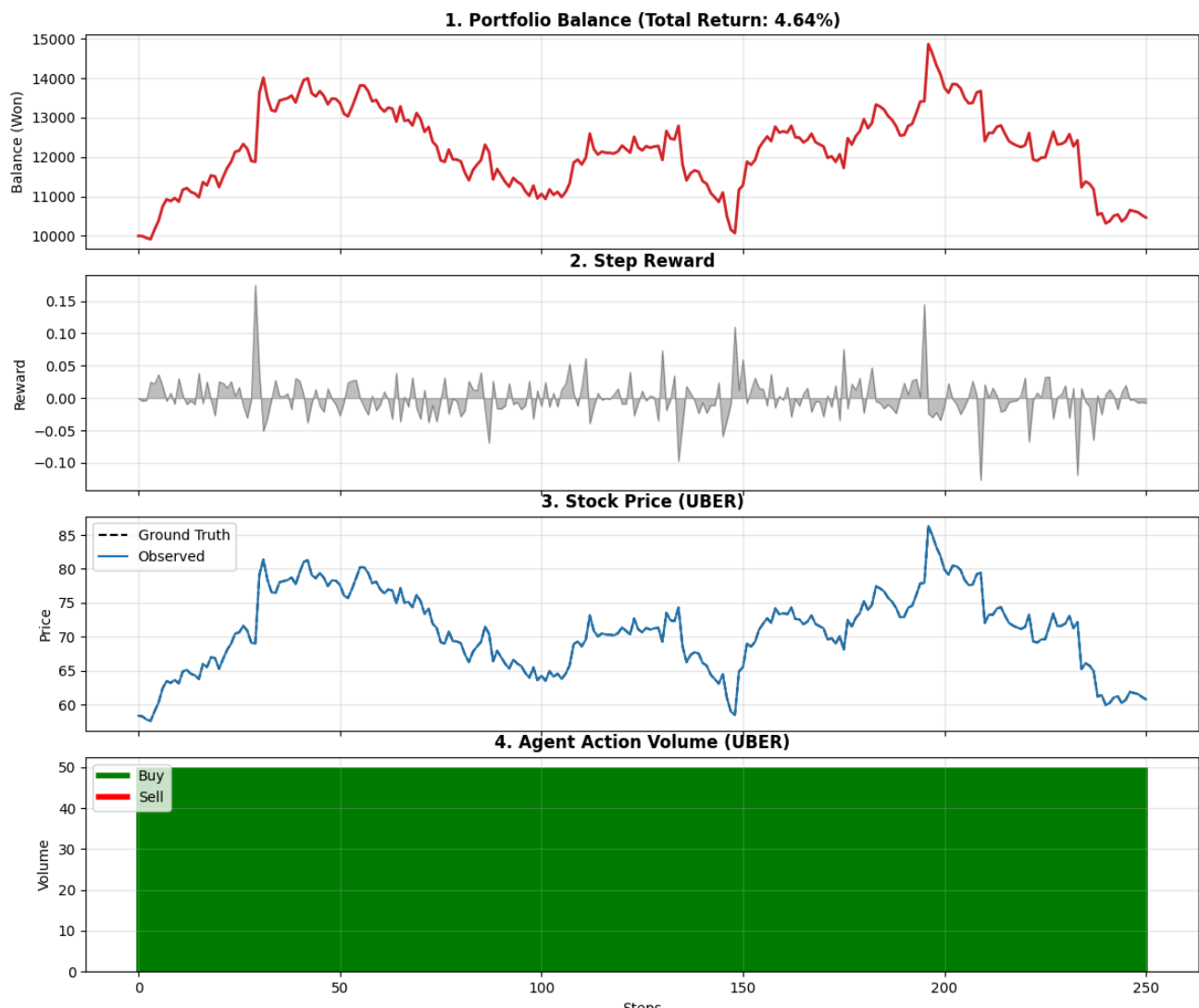


3. Stock Price (F)



4. Agent Action Volume (F)





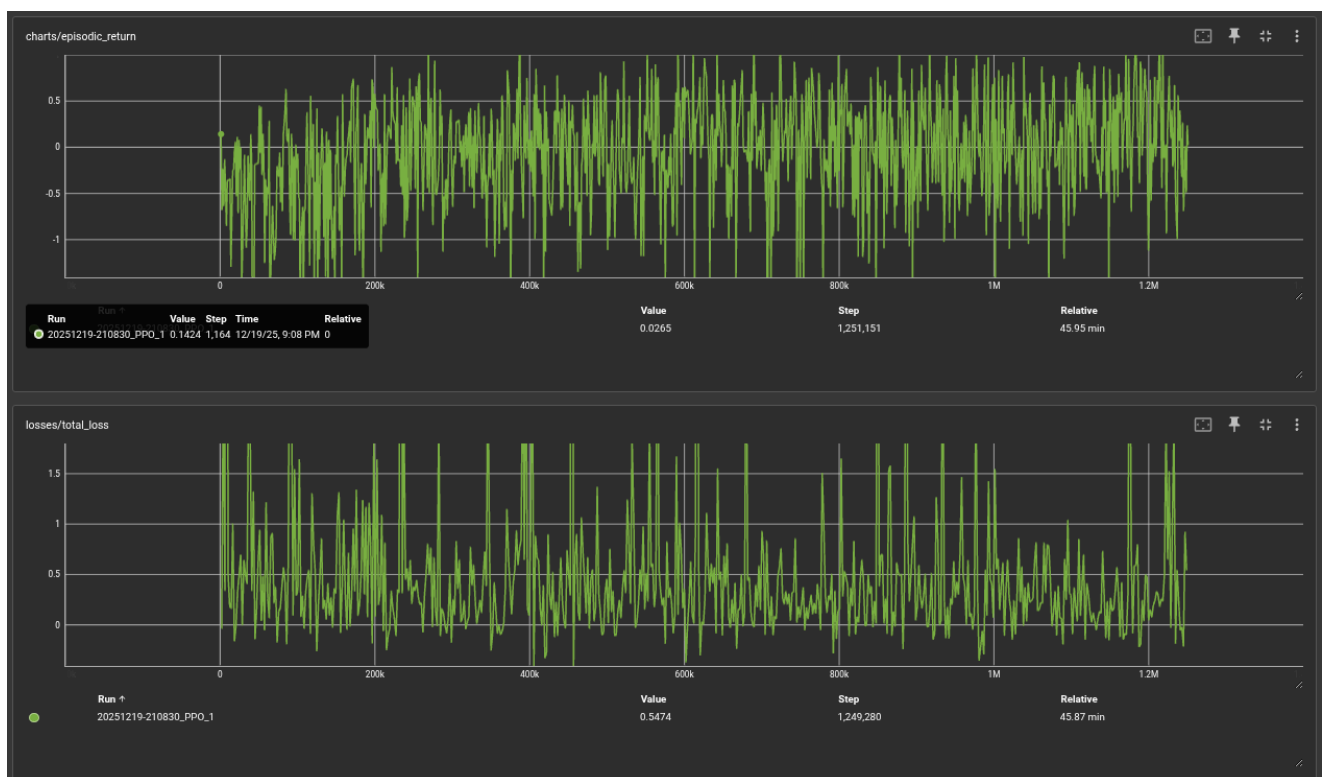
Advanced Setup

- Observation States
 - stack of observation history for 5 timestep
 - agent
 - current balance divided by initial balance
 - number of holding stocks
 - market
 - price over maximum price ever seen
 - target price over actual price
 - open price
 - high price
 - low price
 - close price
 - volume
 - Relative Strength Index(RSI)
 - Trend Placeholder

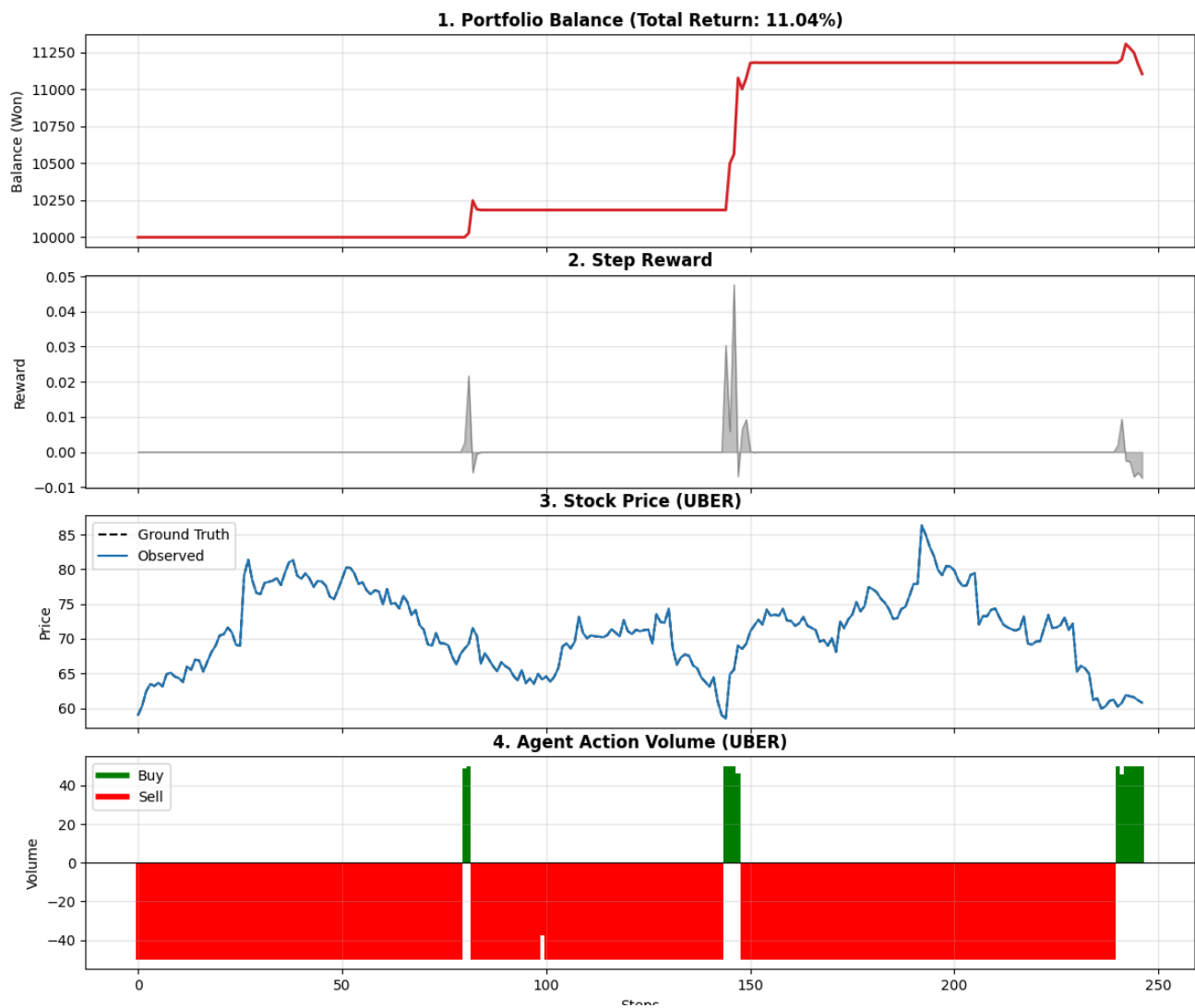
- P/E Ratio
- Dividend Yield
- Reward Function
 - Log return
 - Action penalty
- Dataset
 - 2019-2023 UBER stock

Results

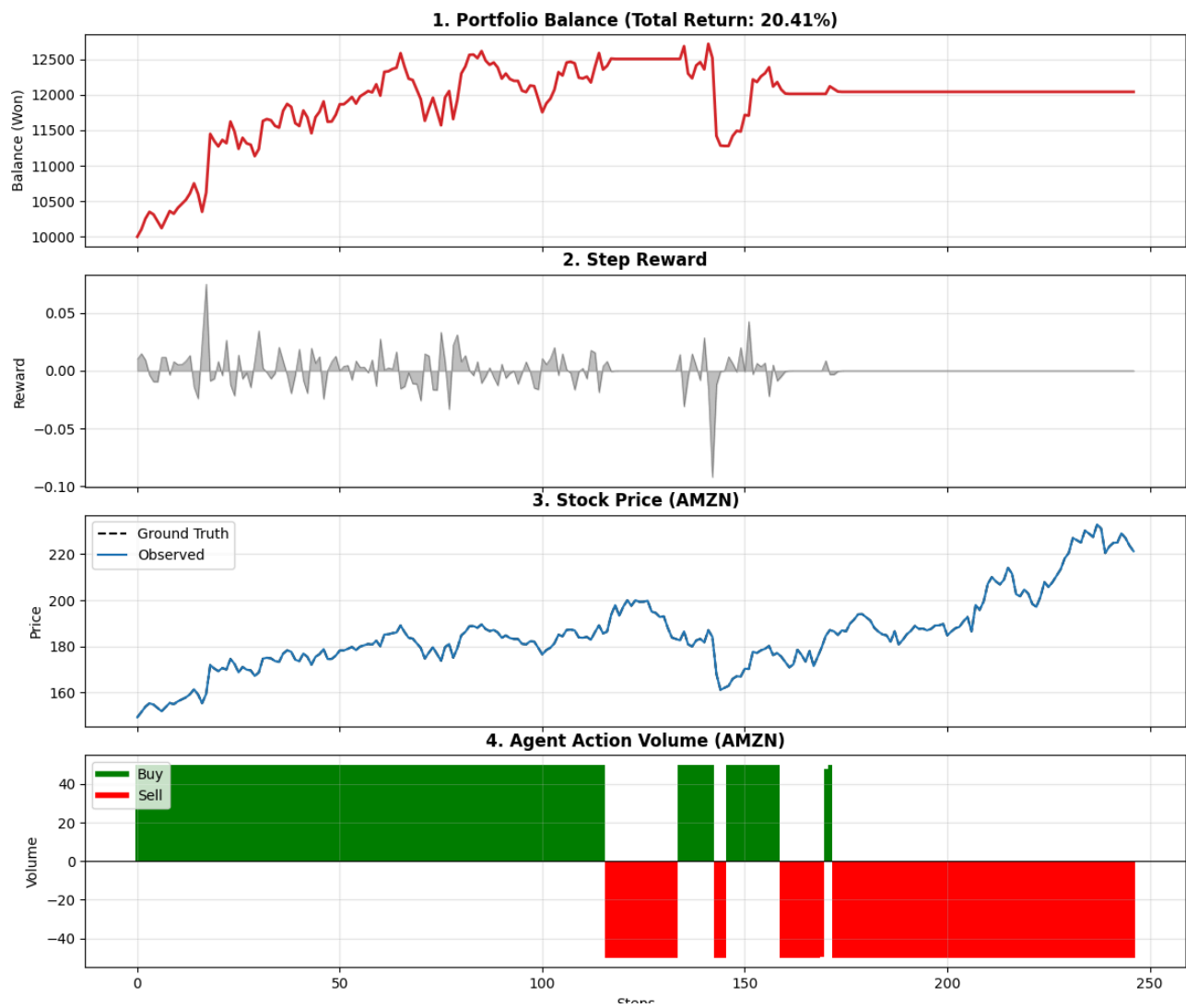
- Learning graph



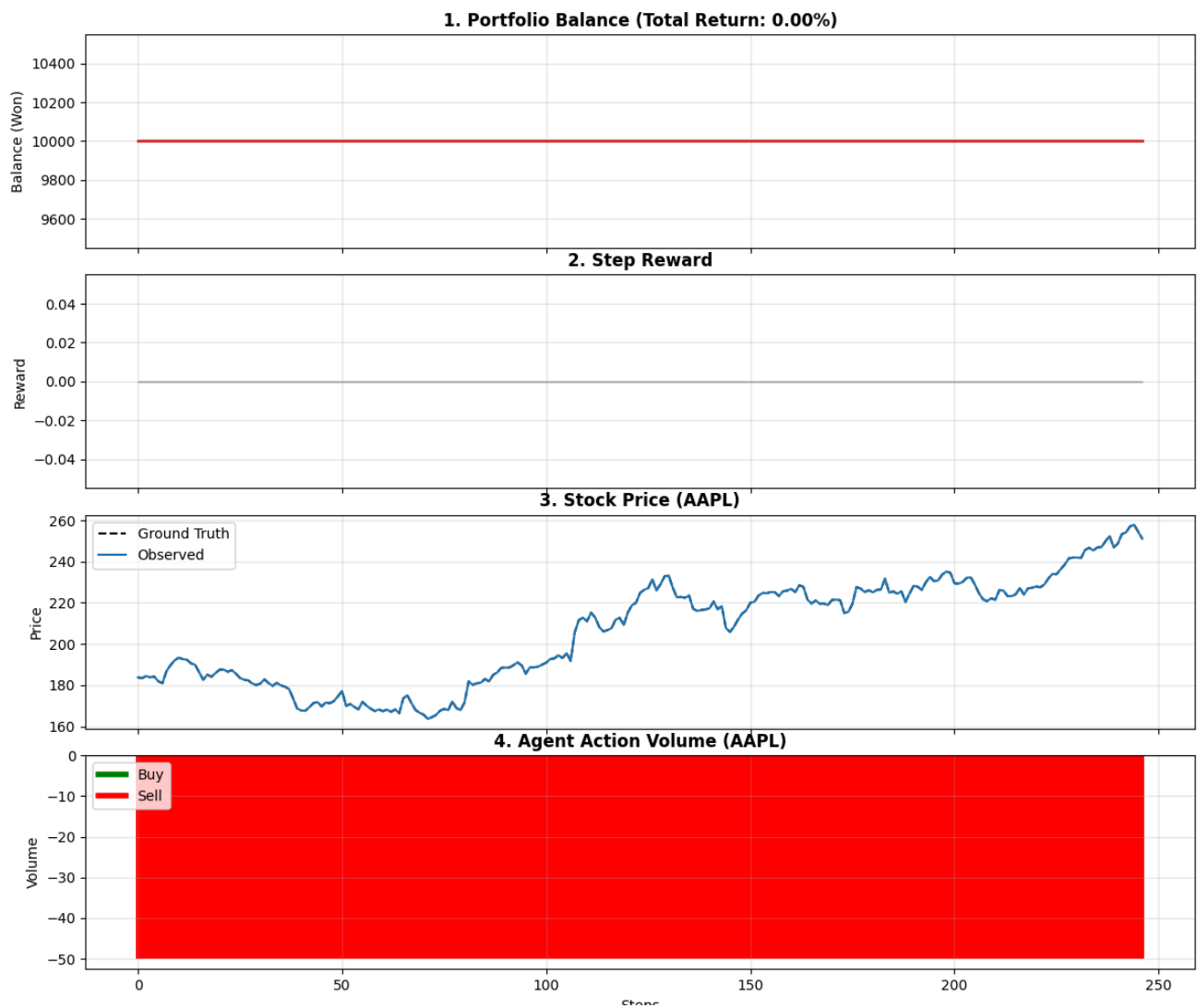
- Same stock; UBER 2024



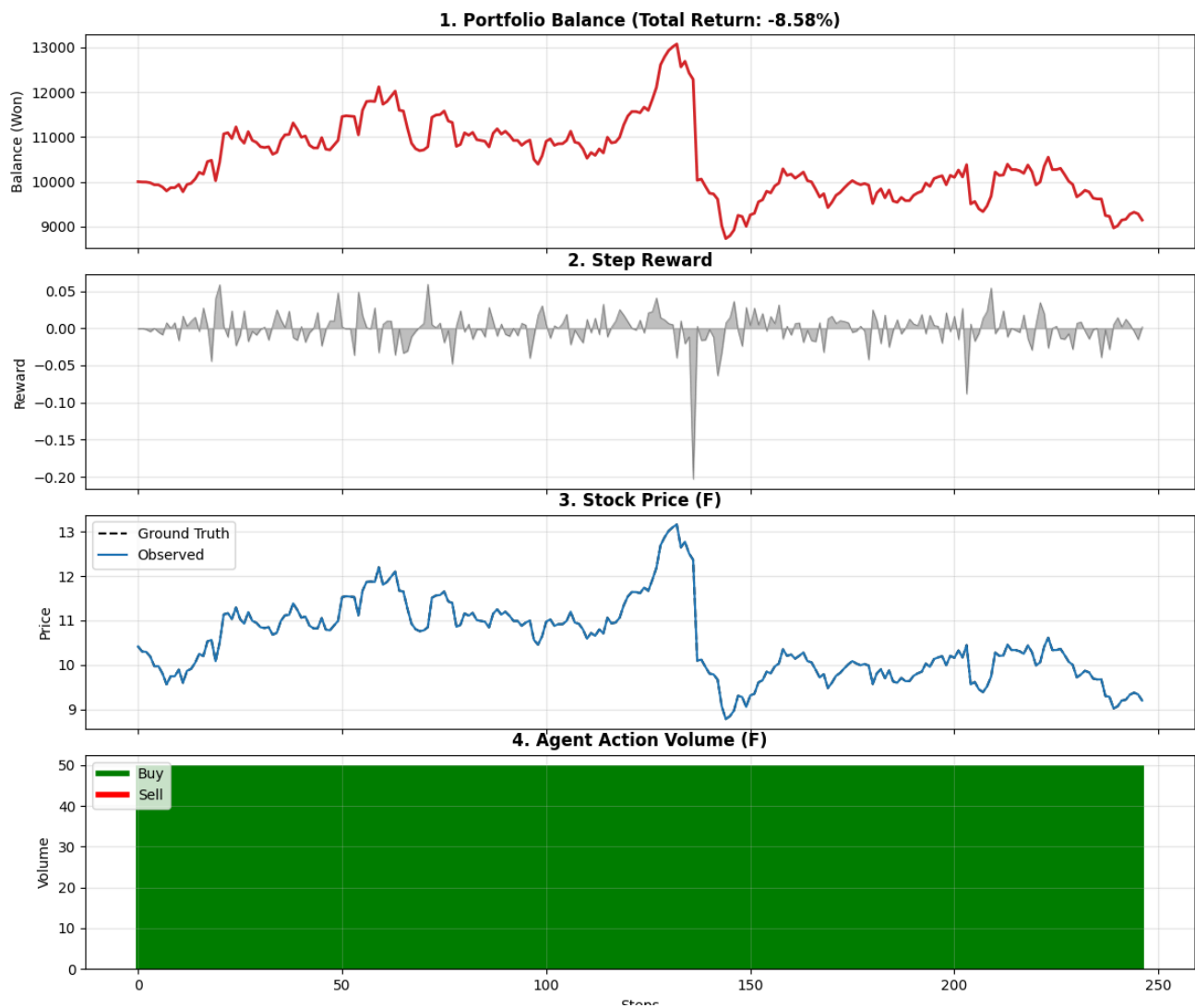
- Different stock; AMZN 2024



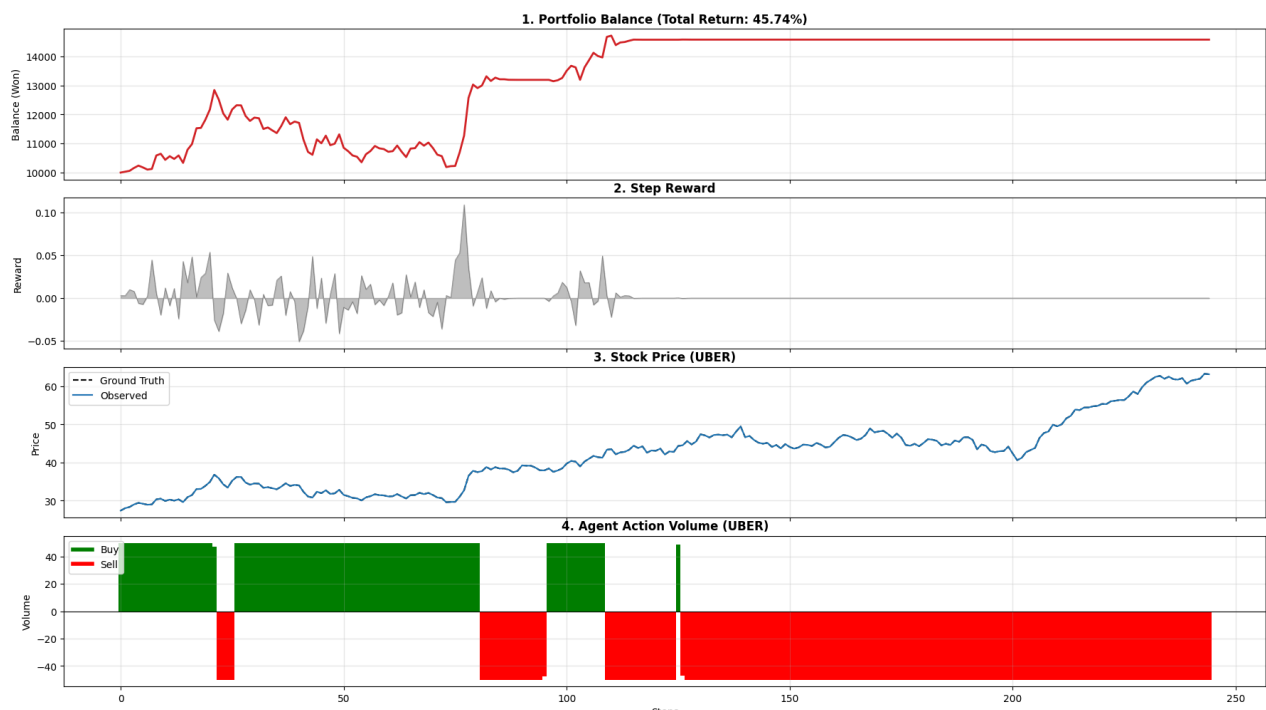
- Different stock; APPL 2024



- Different stock; F 2024



- Same stock, same period; UBER 2023



Discussion

- At first I learned the policy only using tech-related stocks, which are all increasing highly between 2019-2023. The tendency also continued in 2024. So by training the policy using tech-relating monotonously increasing stocks, the policy might not learn about the rule of up and down, but just buy and wait. It can make the policy impossible to handle with non-tech stocks. As you can see in the baseline policy, it works pretty well with tech-related rising stocks, but it fails with Uber and Ford.
 - However, the advanced policy learned more general policy with knowing about the tendency of the stock price, using the history of observation, non-monotonous data, and from various market related values. So it behaves better for Uber, F 2024 datasets.
 - By knowing these values, the agent doesn't just predict the stock by looking at the current price, but actually understanding the value of the company.
 - Analysts' price standard, Relative Strength Index(RSI) and P/E Ratio can also let the policy know about if the stock is over-evaluated or under-evaluated.
 - I add the action penalty term to make the action smoother, in intention of making the policy to do a long-term investment. However, it didn't work well, but it actually somehow increased the tendency to stay at fully buying all stocks, or fully selling all stocks.
 - Here, by the way, even the policy's action is to sell 50 stocks, they can only sell they are holding, so the value of the graph is not directly influencing the environment, but it's being clipped.
-

3. Experiments - Performance Upgrade, Multi Ticker

Baseline Setup

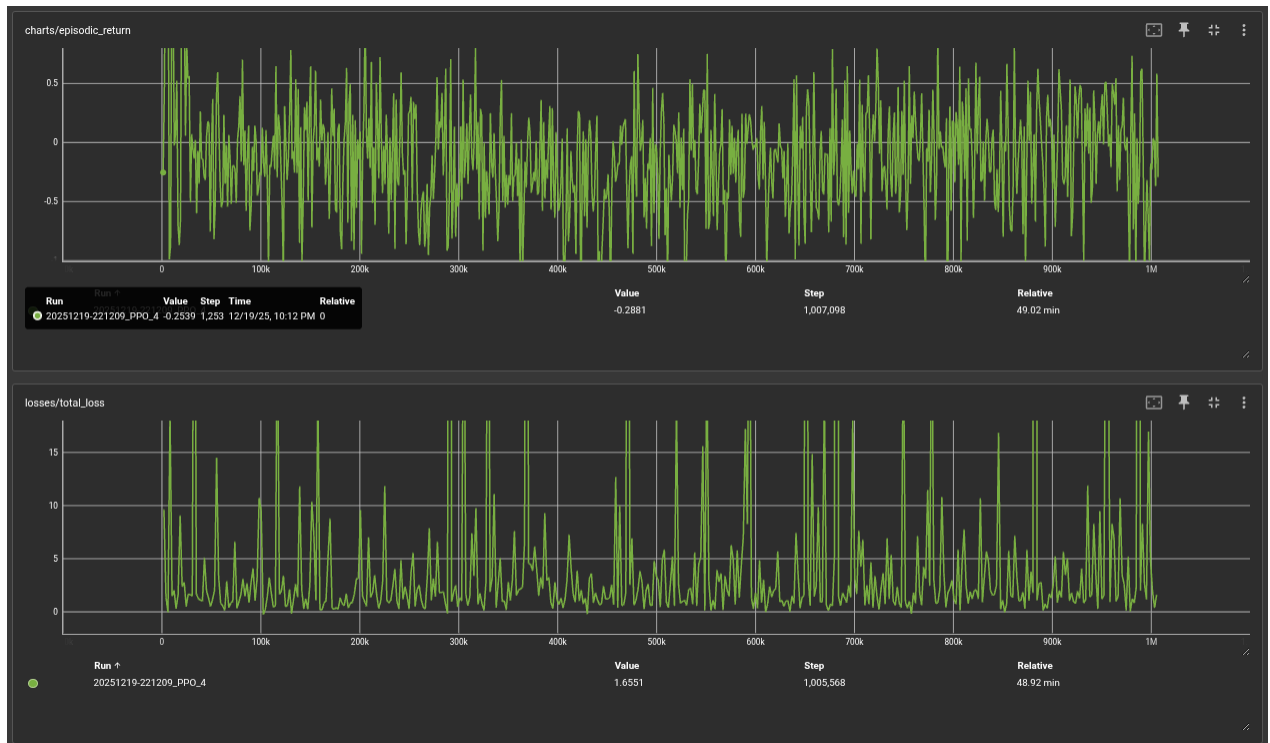
- Same as Single Ticker's

Advanced Setup

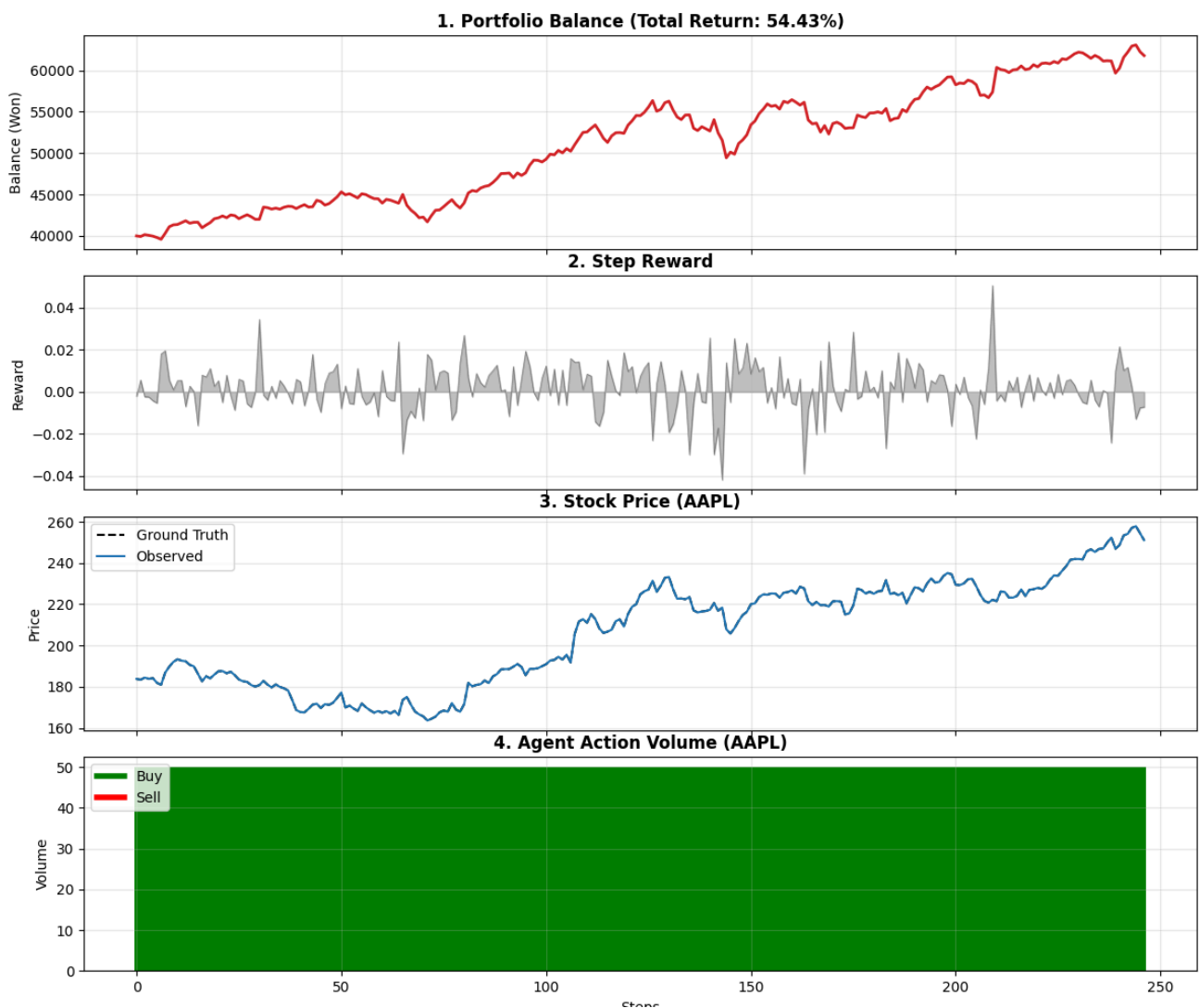
- Same as Single Ticker's

Results

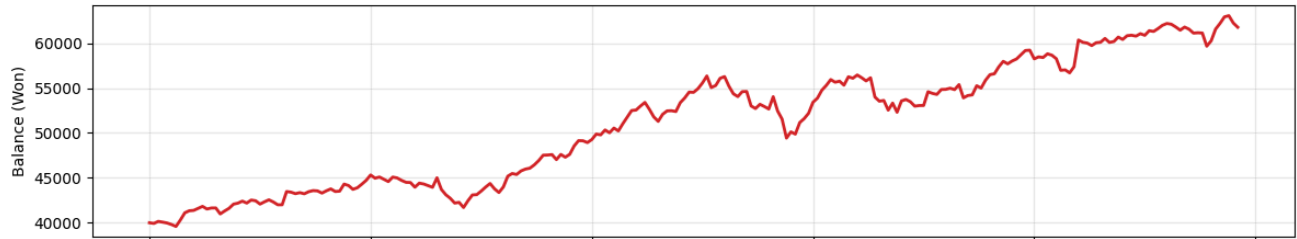
- Learning Graph



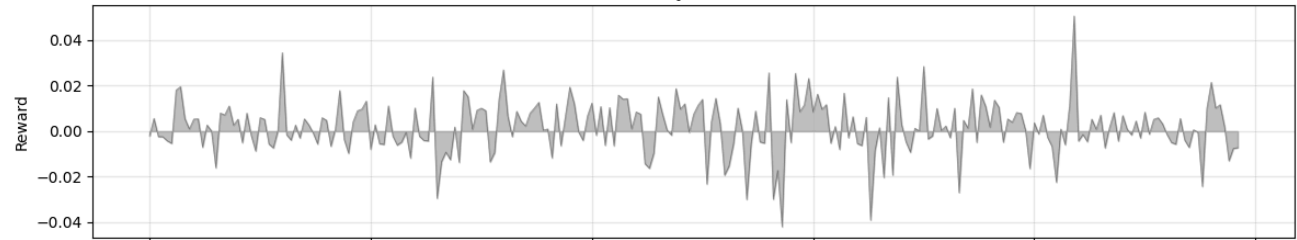
- Exp1. All same stocks



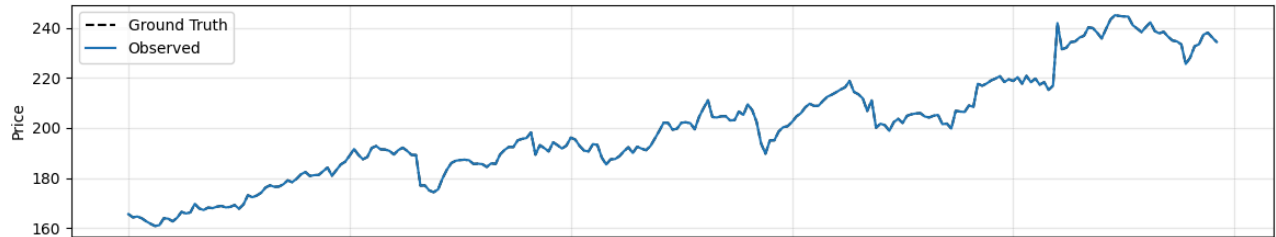
1. Portfolio Balance (Total Return: 54.43%)



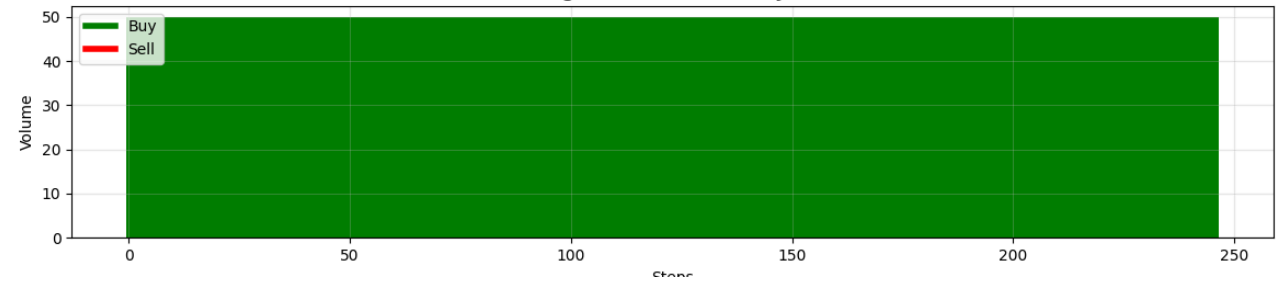
2. Step Reward



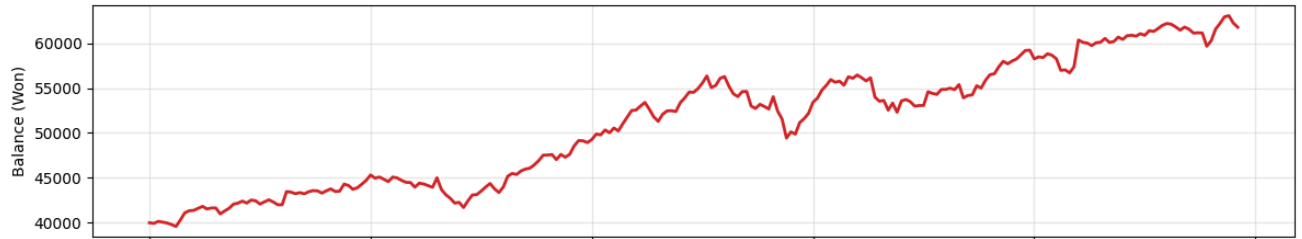
3. Stock Price (JPM)



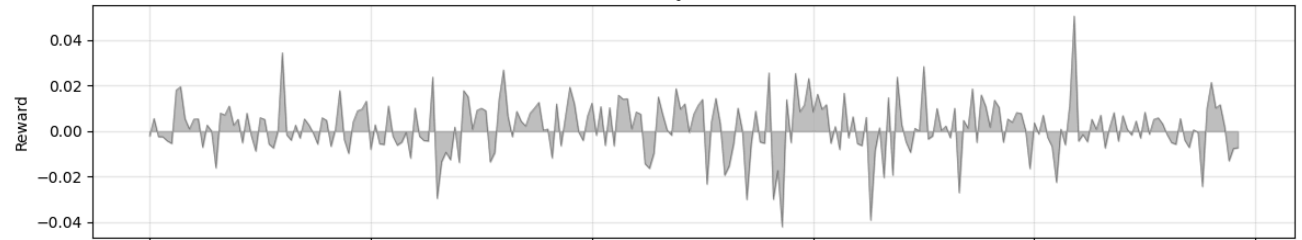
4. Agent Action Volume (JPM)



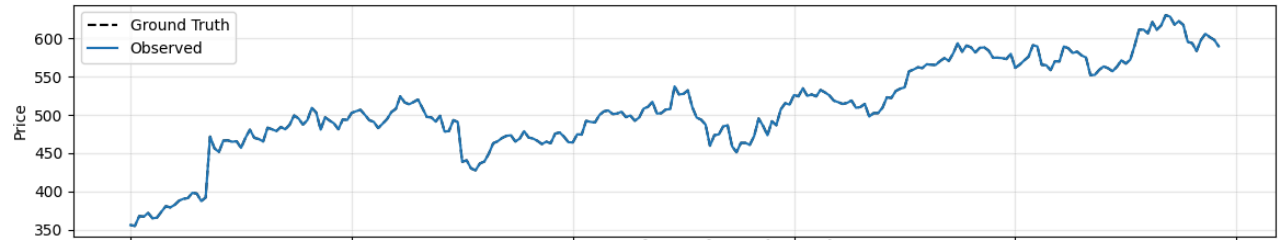
1. Portfolio Balance (Total Return: 54.43%)



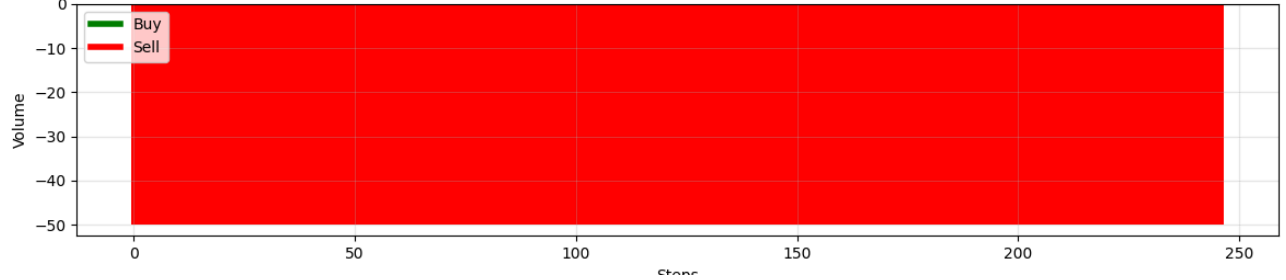
2. Step Reward

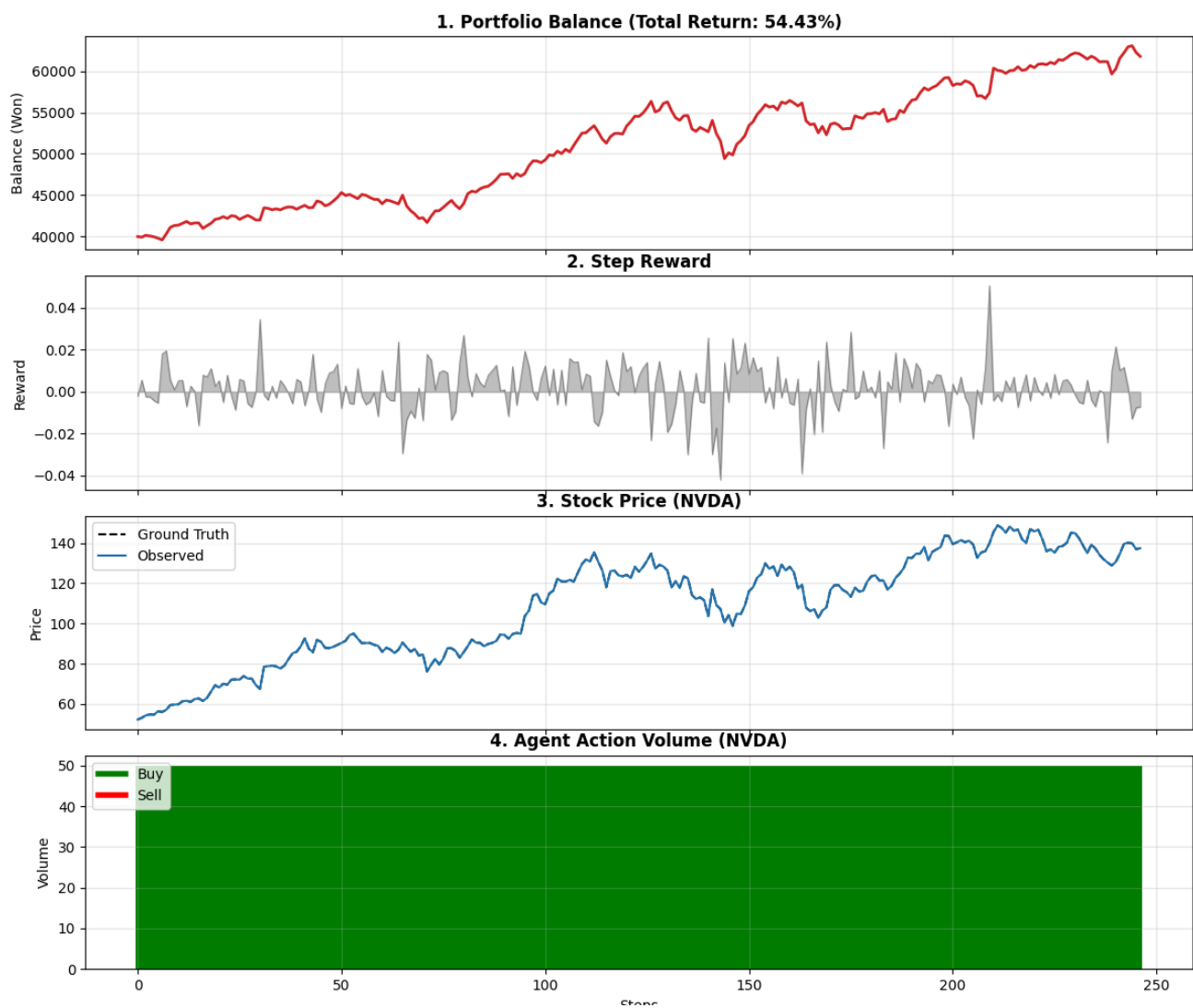


3. Stock Price (META)



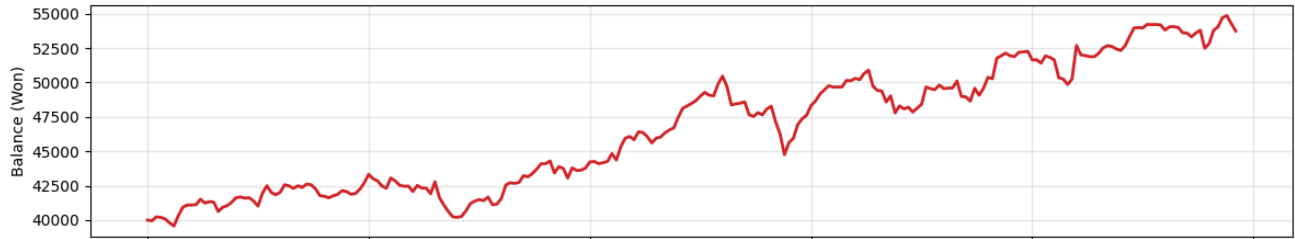
4. Agent Action Volume (META)



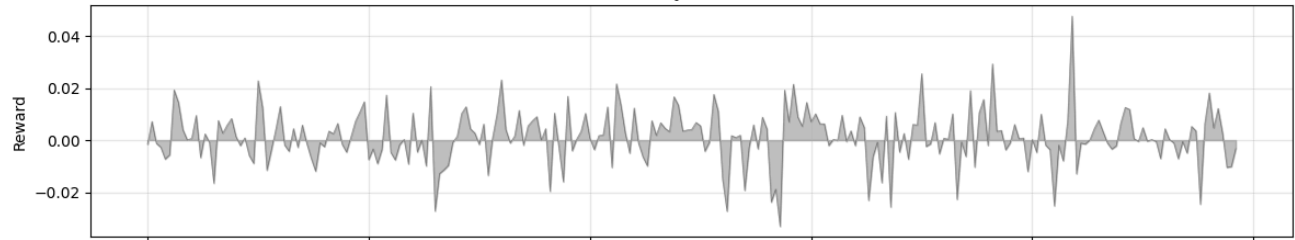


- **Exp2. 2 trained stocks + 2 unseen stocks**

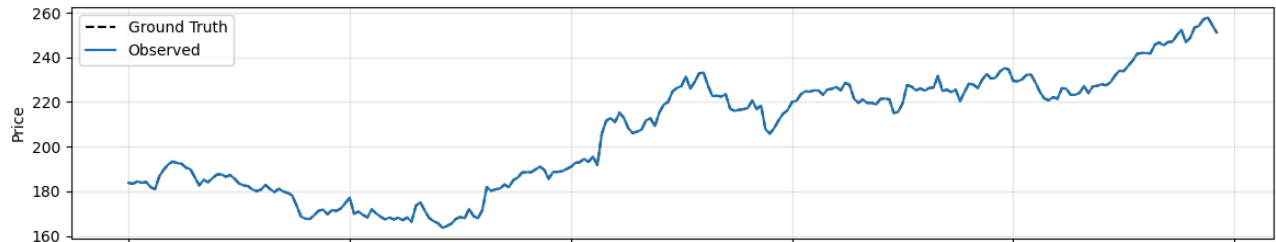
1. Portfolio Balance (Total Return: 34.34%)



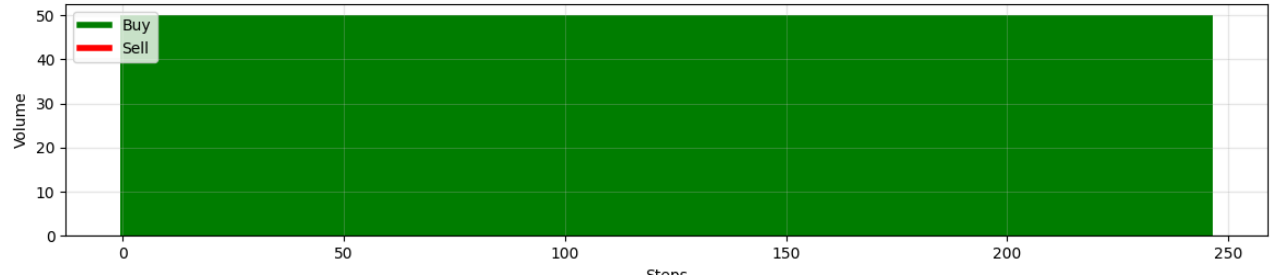
2. Step Reward



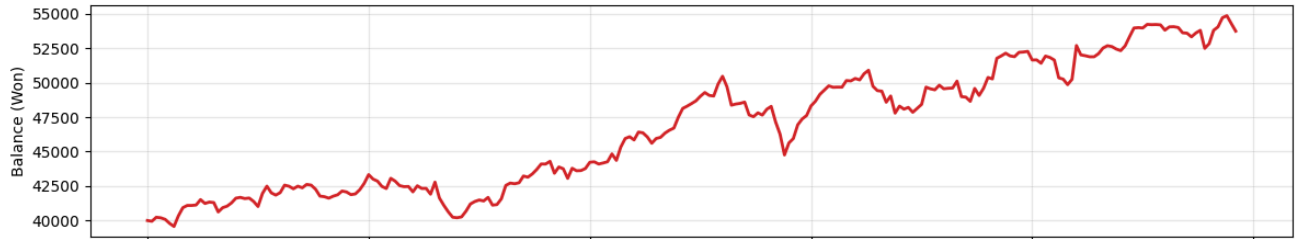
3. Stock Price (AAPL)



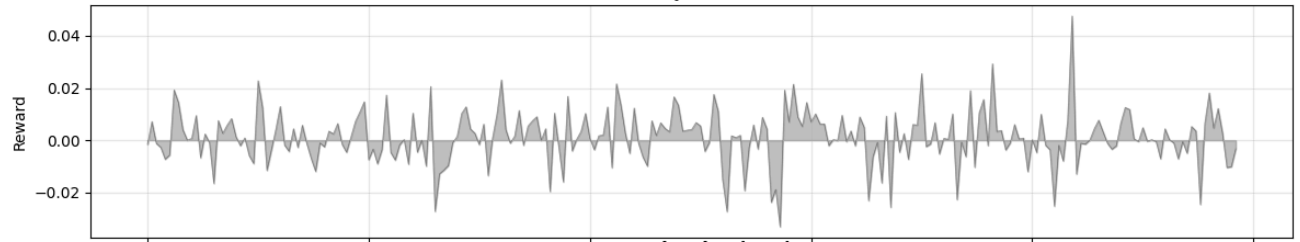
4. Agent Action Volume (AAPL)



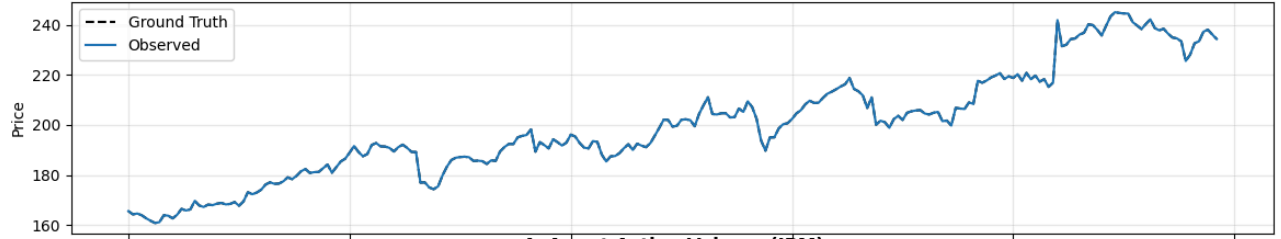
1. Portfolio Balance (Total Return: 34.34%)



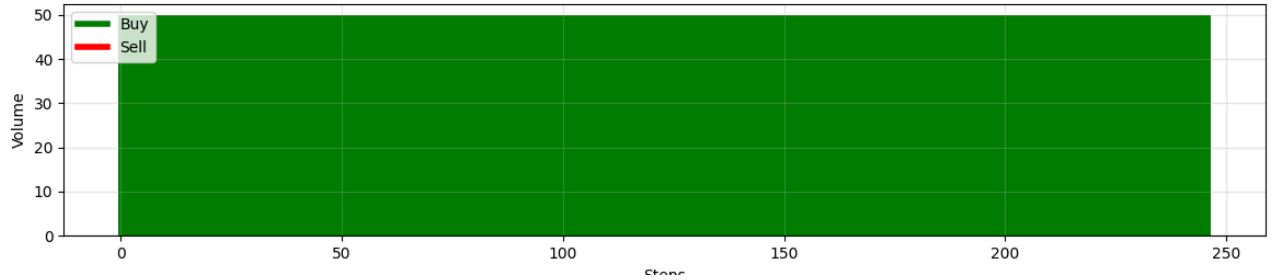
2. Step Reward



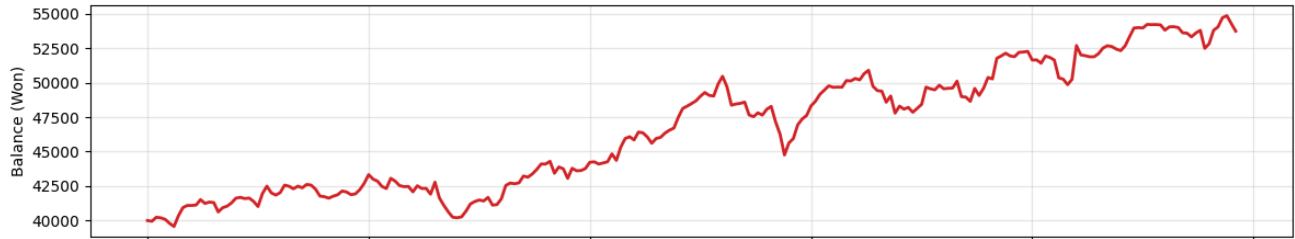
3. Stock Price (JPM)



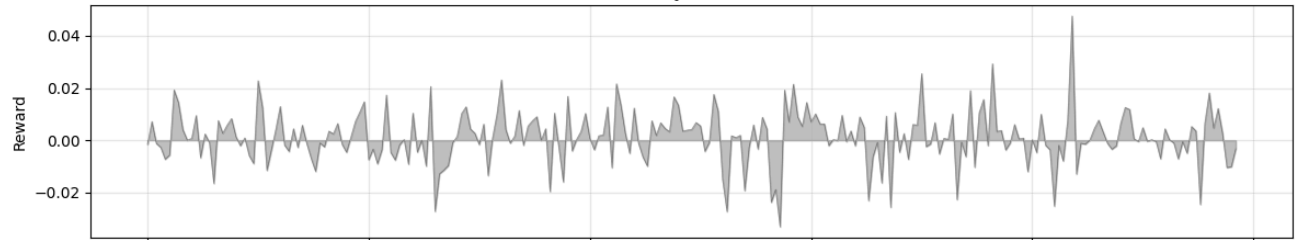
4. Agent Action Volume (JPM)



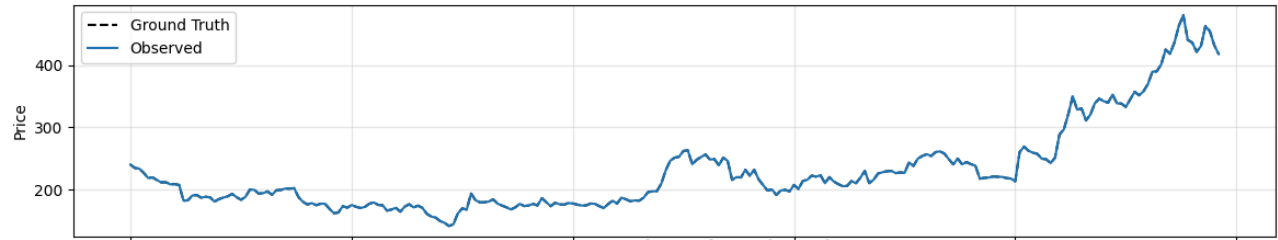
1. Portfolio Balance (Total Return: 34.34%)



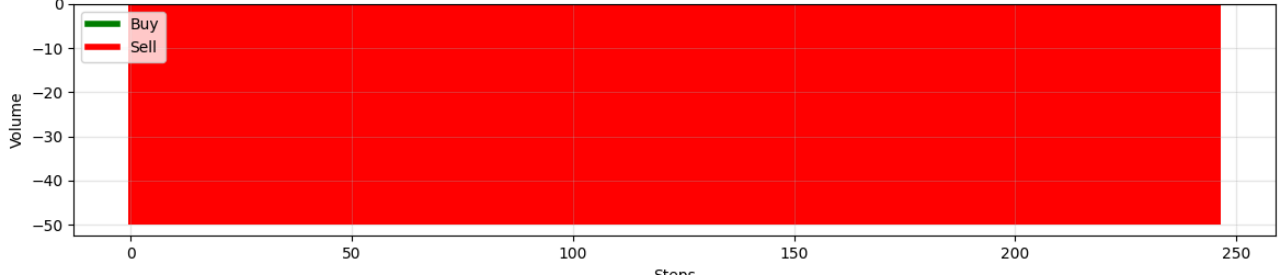
2. Step Reward



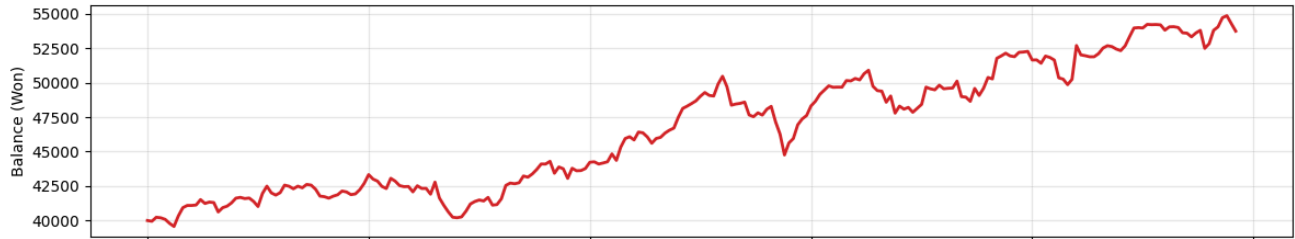
3. Stock Price (TSLA)



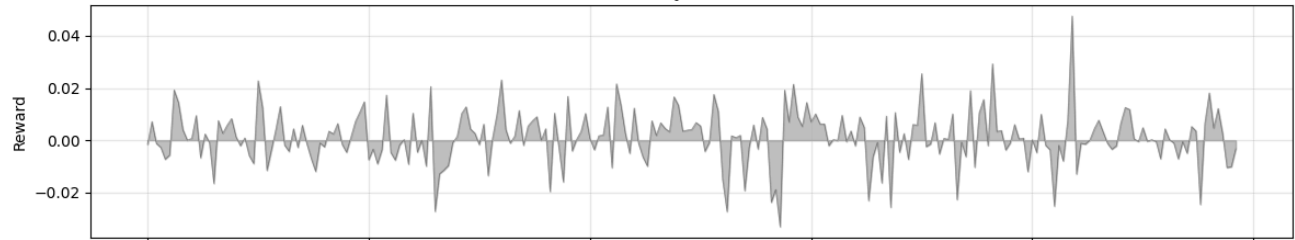
4. Agent Action Volume (TSLA)



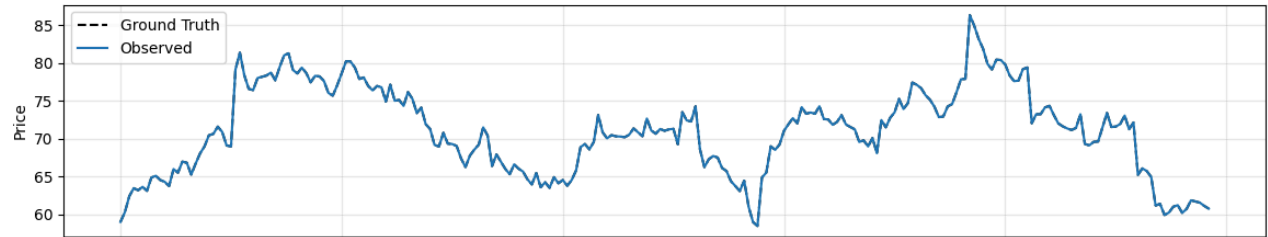
1. Portfolio Balance (Total Return: 34.34%)



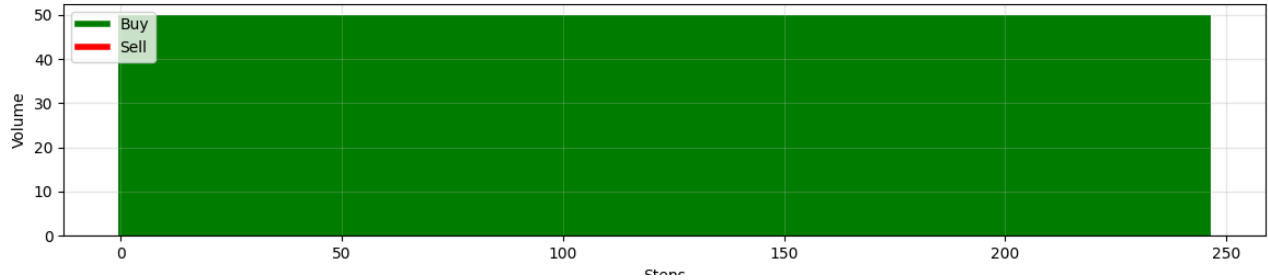
2. Step Reward



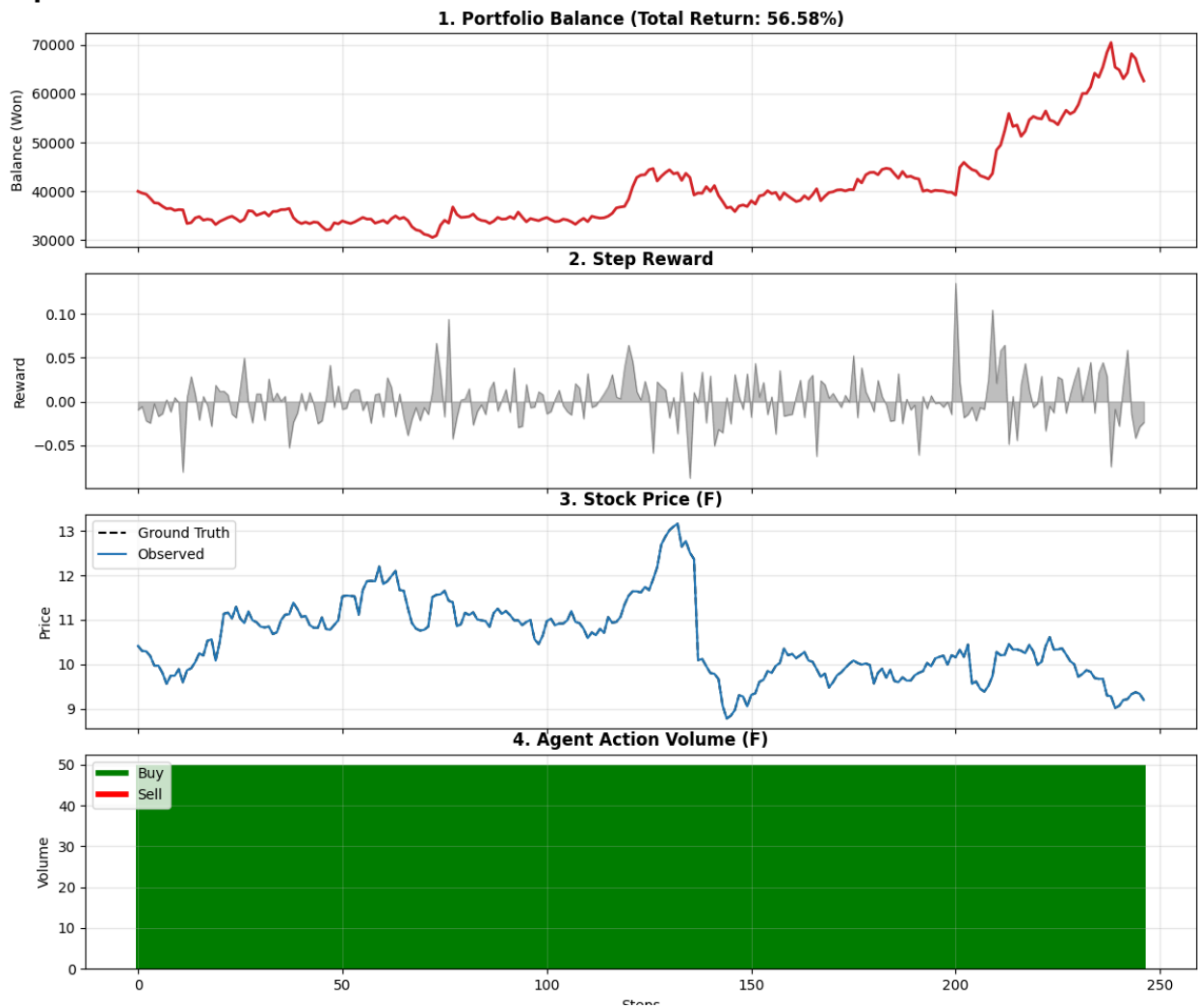
3. Stock Price (UBER)



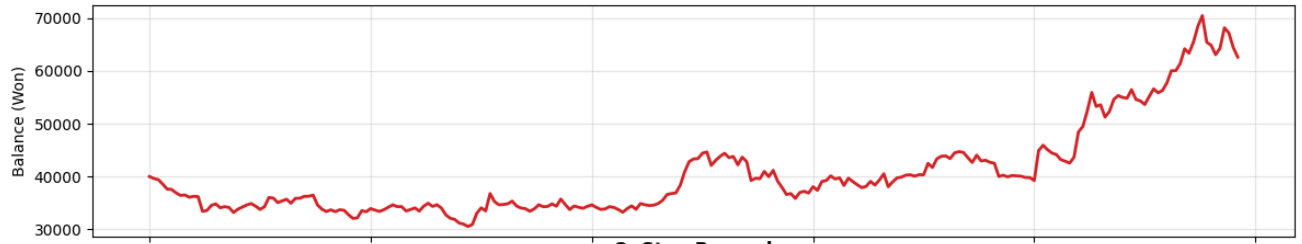
4. Agent Action Volume (UBER)



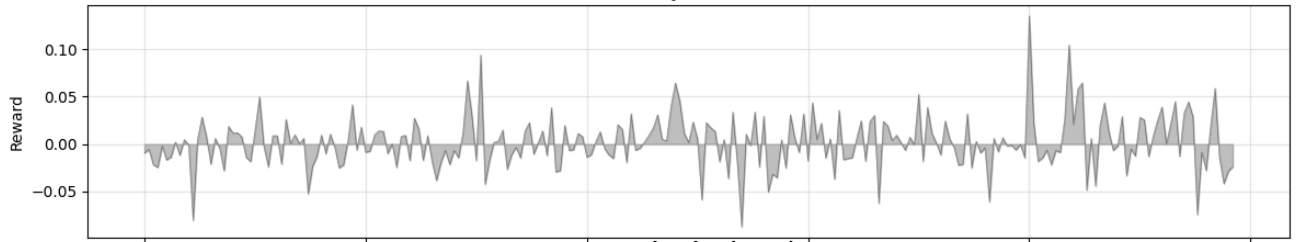
- **Exp3. All unseen stocks**



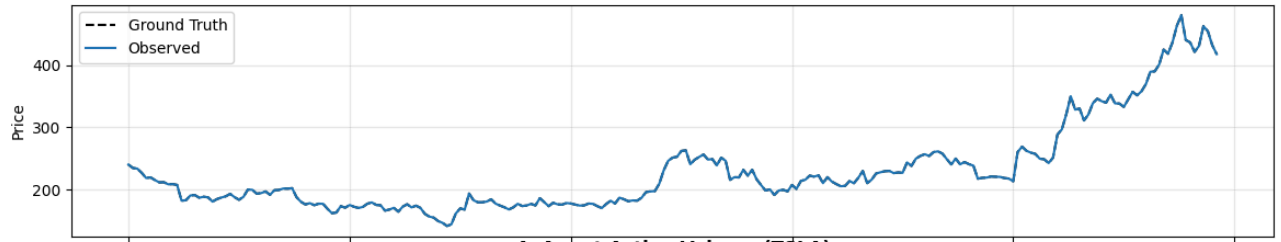
1. Portfolio Balance (Total Return: 56.58%)



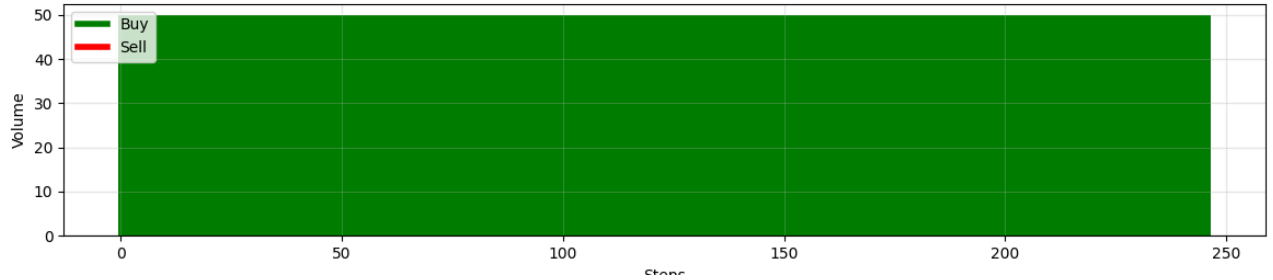
2. Step Reward



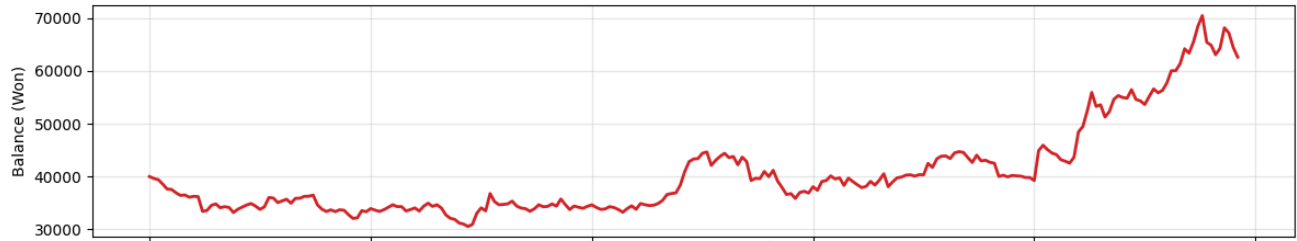
3. Stock Price (TSLA)



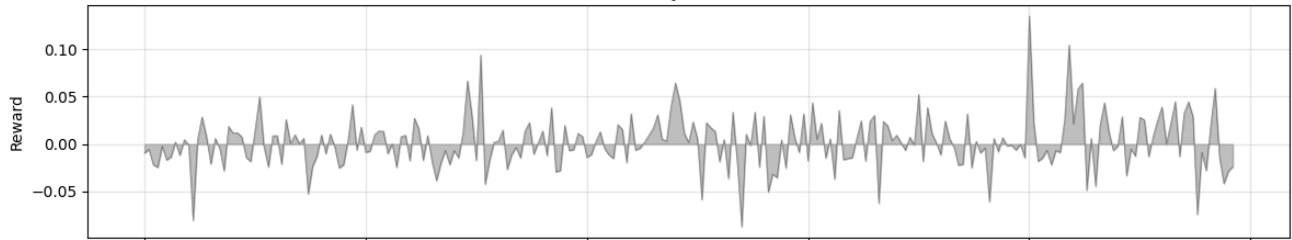
4. Agent Action Volume (TSLA)



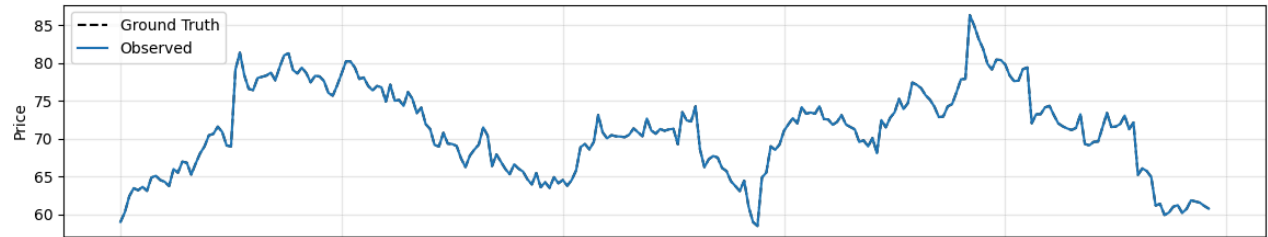
1. Portfolio Balance (Total Return: 56.58%)



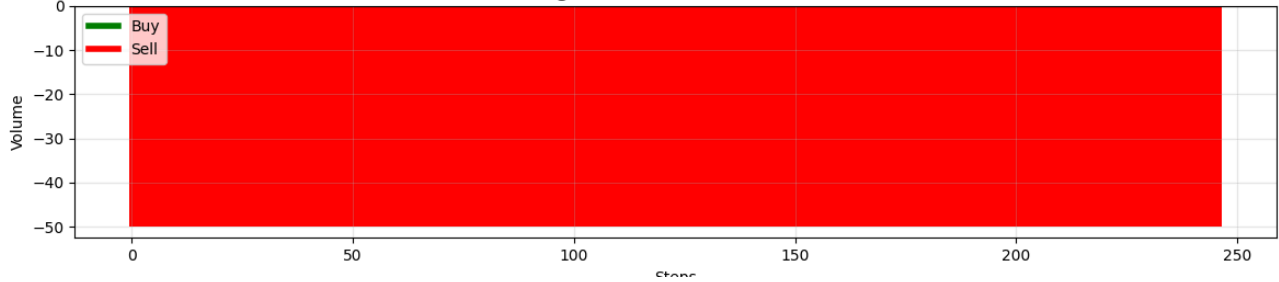
2. Step Reward



3. Stock Price (UBER)



4. Agent Action Volume (UBER)





Discussion

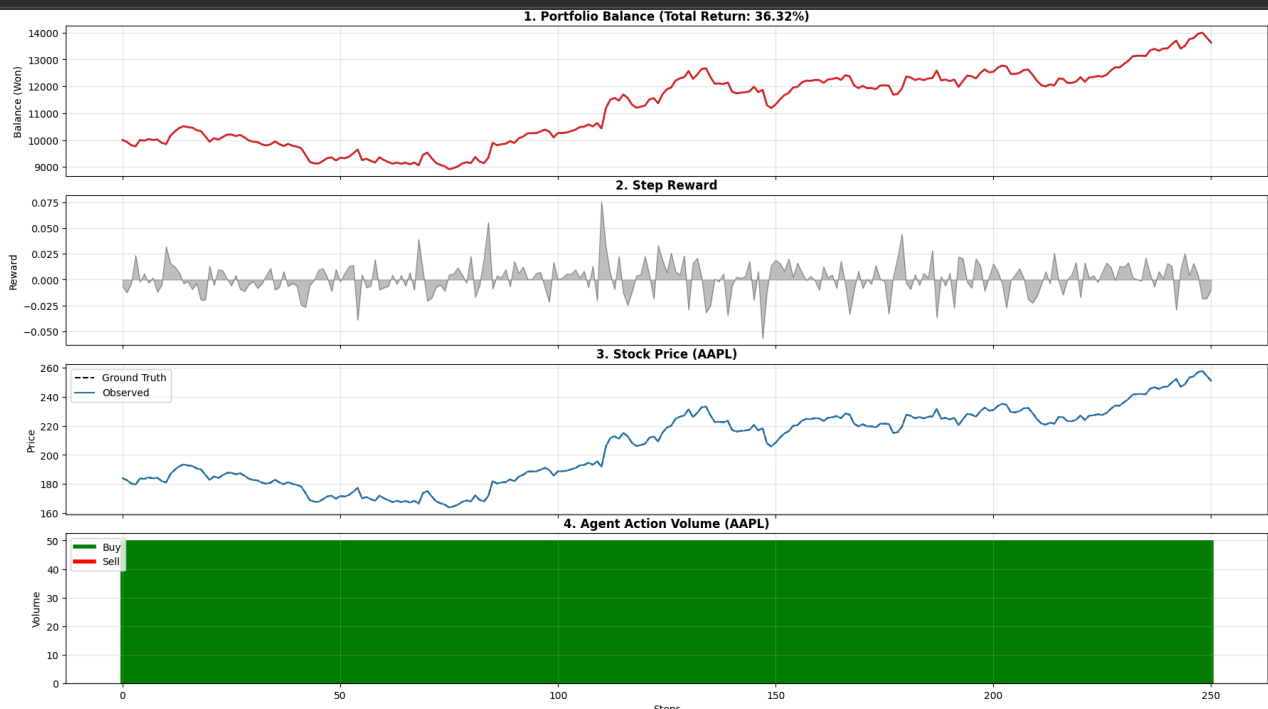
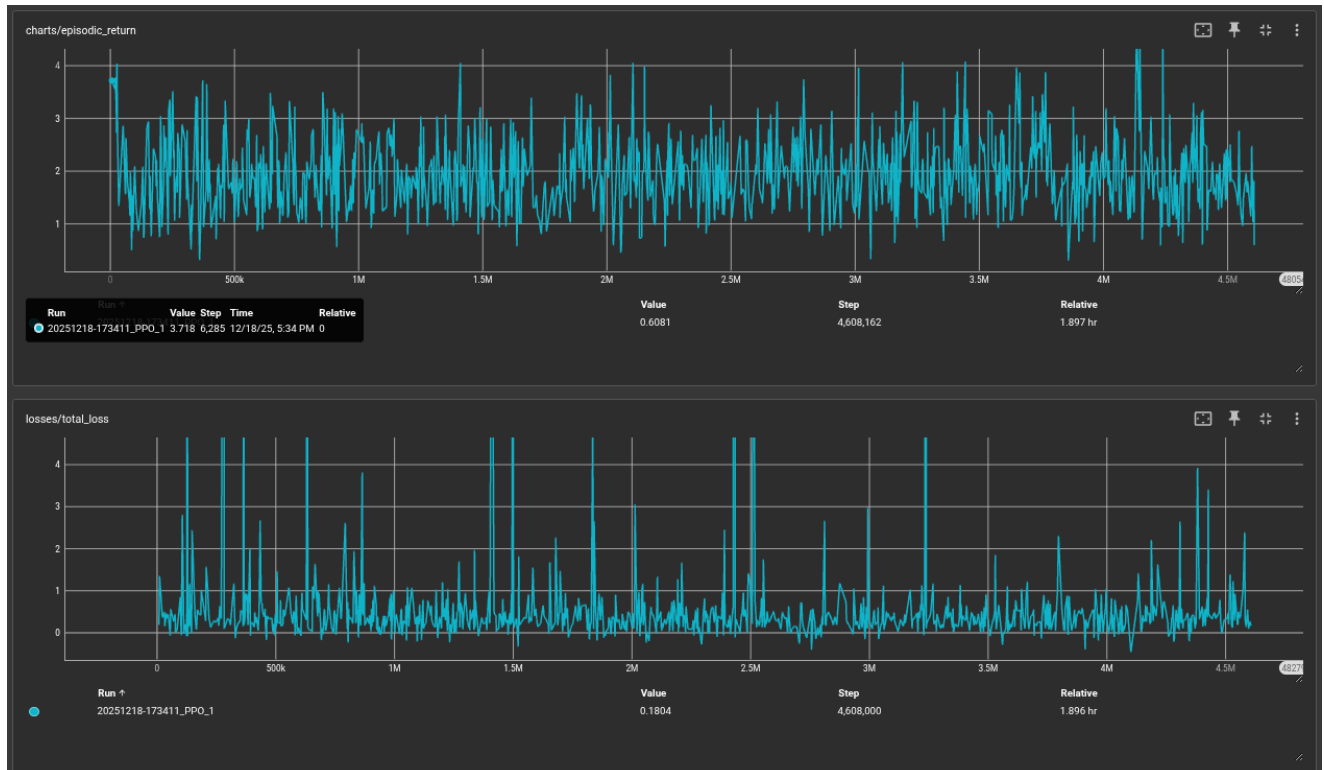
- The improvement of the performance based on the change of states and reward would be similarly discussed as previous discussion. So I want to talk the overall tendency of all policies in this whole project.
- We are giving very small dataset to each policy. Just 1 or up to 4 stocks, and up to 10years (about 2500 trading days). And the datasets that I selected are all tech-related dataset, and recent years, almost all tech-related stocks are growing. Some are even sky-rocketing in several days.
- It means, we are giving too-biased dataset to policies, so the policy learns to just buy all at the first time, then keep it, which might be a fair strategy within recent tech-related stocks.
- However, it fails to learn the actual principal of trading, considering the value of the company, or interpret and predicting the rise or fall of the stock price in long term or even in short term.
- So I would like to leave this as a future work. The main thing to try to improve the performance of the policy, and making the policy to act diversely based on the observation state, the dataset needs to be more larger with various tendencies.

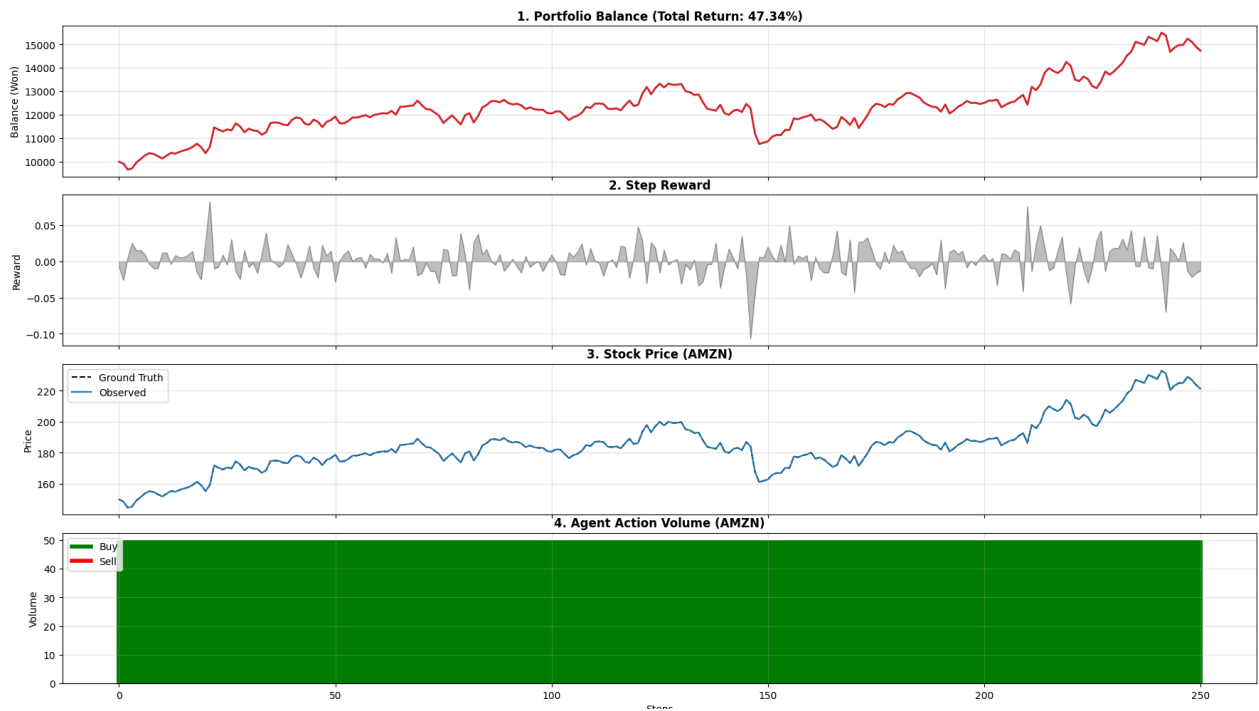
4. Experiments - Algorithm Comparison

Setup

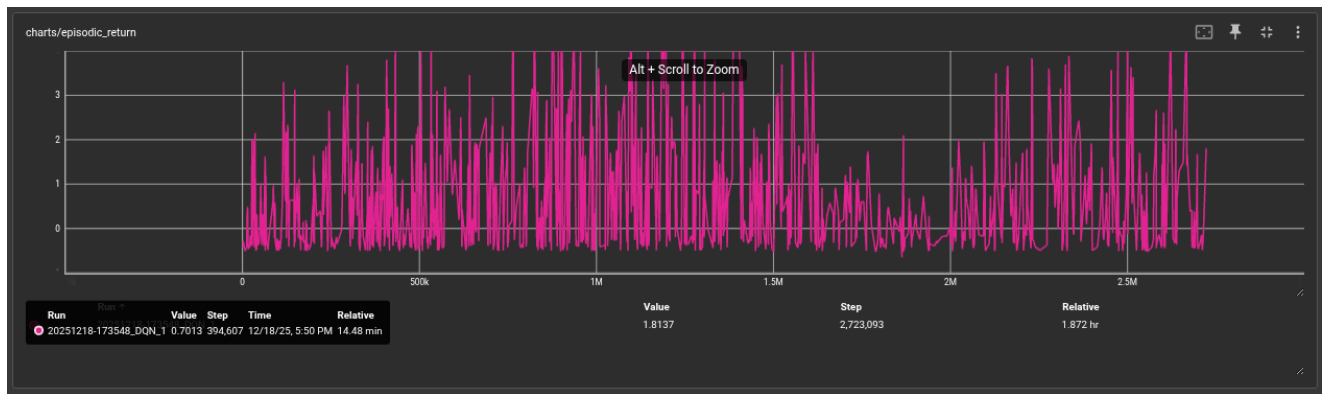
- Same as baseline

PPO Results

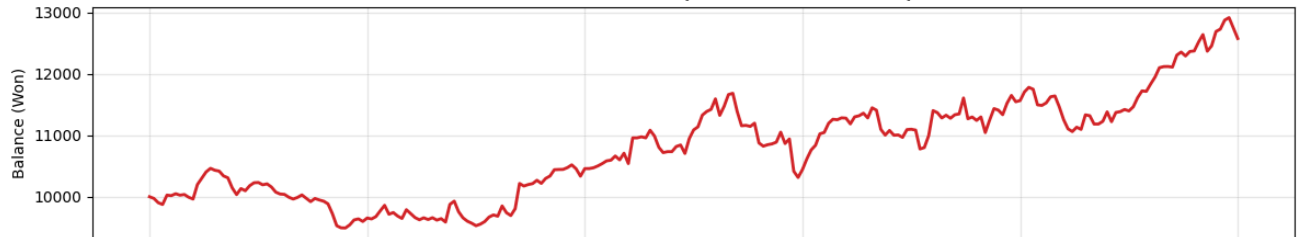




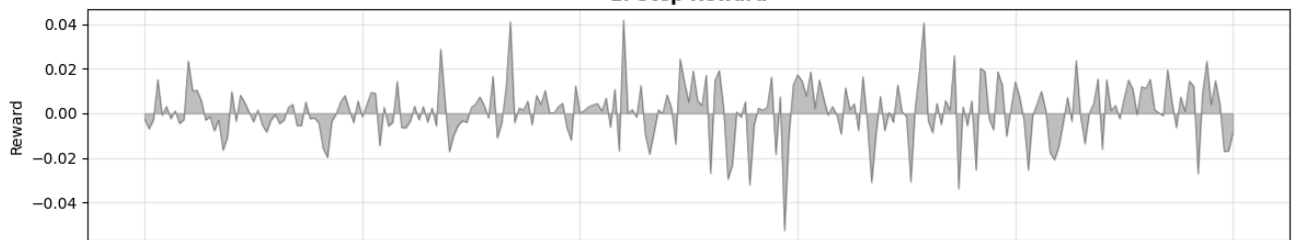
DQN Results



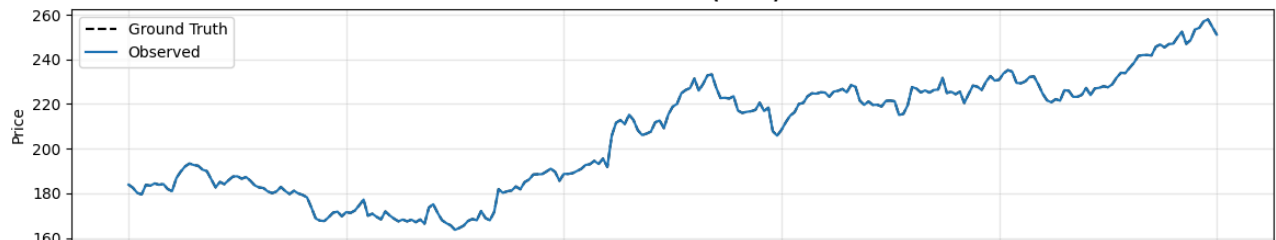
1. Portfolio Balance (Total Return: 25.75%)



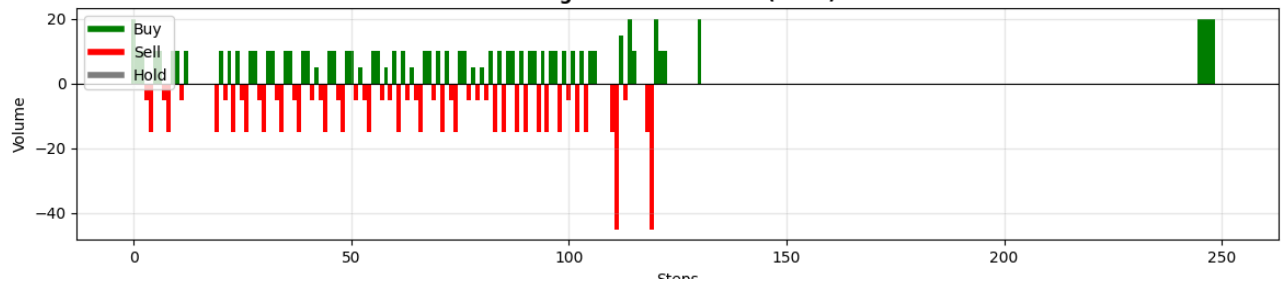
2. Step Reward



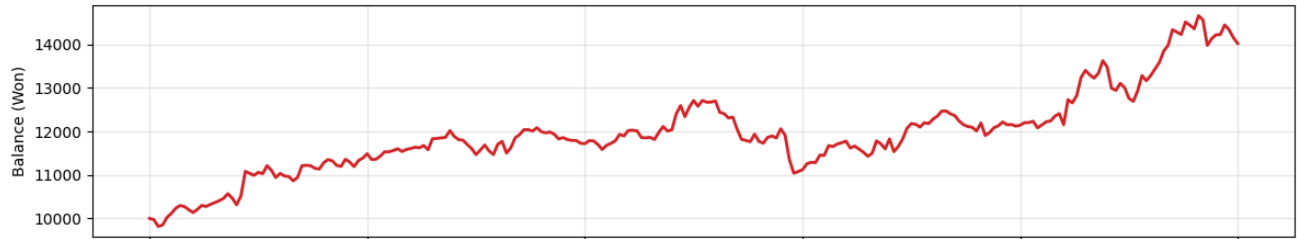
3. Stock Price (AAPL)



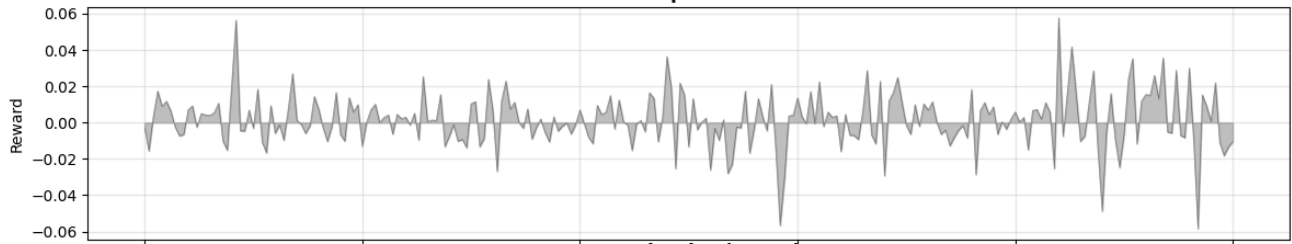
4. Agent Action Volume (AAPL)



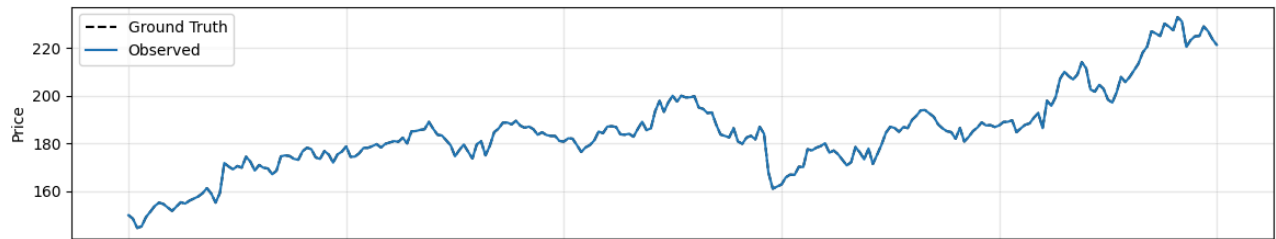
1. Portfolio Balance (Total Return: 40.23%)



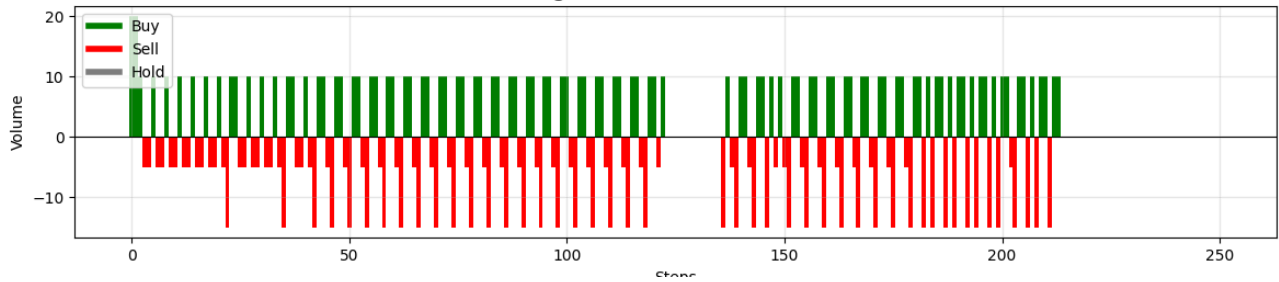
2. Step Reward

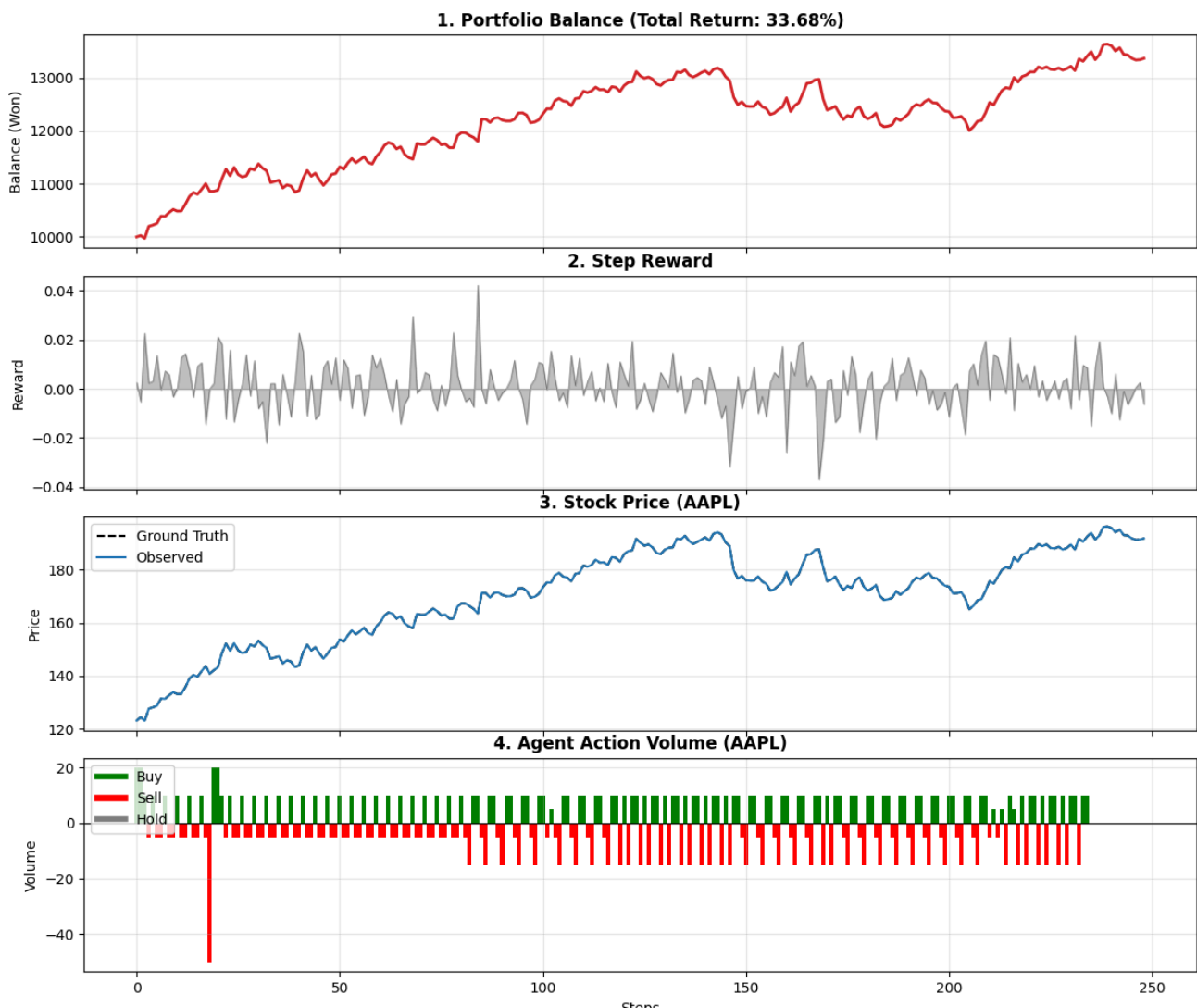


3. Stock Price (AMZN)



4. Agent Action Volume (AMZN)





Discussion

- To compare the learning graph tendency, as you can see in the PPO learning graph, it looks like it has similarly high variance from DQN's but the lower bound is slightly getting higher, as the learning proceeds. As PPO is using clipping loss, some radical updates can be prevented, so the variance plot of PPO can be seen as slightly more stable than DQN. However, for DQN's learning graph, it shows high picks even in the later phase of the learning.
- As we are using very limited data (at most around 1000 days), the performance of PPO and DQN actions are interesting. PPO just uses the data once and then gets rid of it. (Still there is a loop for reusing data multiple times; in my code, it's 4 times) However, DQN keeps all the data and reuses throughout the learning process. In addition to that, PPO's action space is continuous, while DQN's action space is discretized. Here, it's 21 options. Generally, using continuous actions will be more helpful for more sensitive investment, but given that there is very little data, it might not have enough data to explore. The action space might be too large.
- So the actions of PPO weren't able to learn different actions for different observation states, but just take action of full buying or full selling (exploitation). However, DQN has a small action space, with efficient data usage, which makes DQN to explore the action space enough, so it can output more diverse actions based on the observation state.

Future Works

- As I mentioned before, the main problem of all the trained policies is the dataset. It is too small and biased. It doesn't have enough data to make the policy explore enough. So I will try to get more unbiased datasets.
- Apart from dataset, I would like to take 2 different approach to improve the performance.
- First, company value. Specify all the informations related to the value of the company. By gathering those information and preprocessing and then using as a state, would help the agent to more accurately predict the value of the company.
- Second, tendency of the market. Not directly predicting the value of stock itself, I would like to take an approach interpreting the people's tendency. It may be possible by using analyst's informations as a state.