



Lo\_Gobang\_151

Updated versions:

- 1. Score optimization(frequent)
- 2. Double mode:

Two scoring policies. In the beginning or when chess pieces become dense, the agent should be more likely to defend than attack. Just like bosses in RPG games will change mode under certain situations.

3. Search and pruning(Fig. 2.):

The agent makes decision only according to the current state, which means it lacks ability to make temporary sacrifice to get higher score in the future. I tried to make a few simulations towards future to search for the better move.

Win rate are shown in the report.

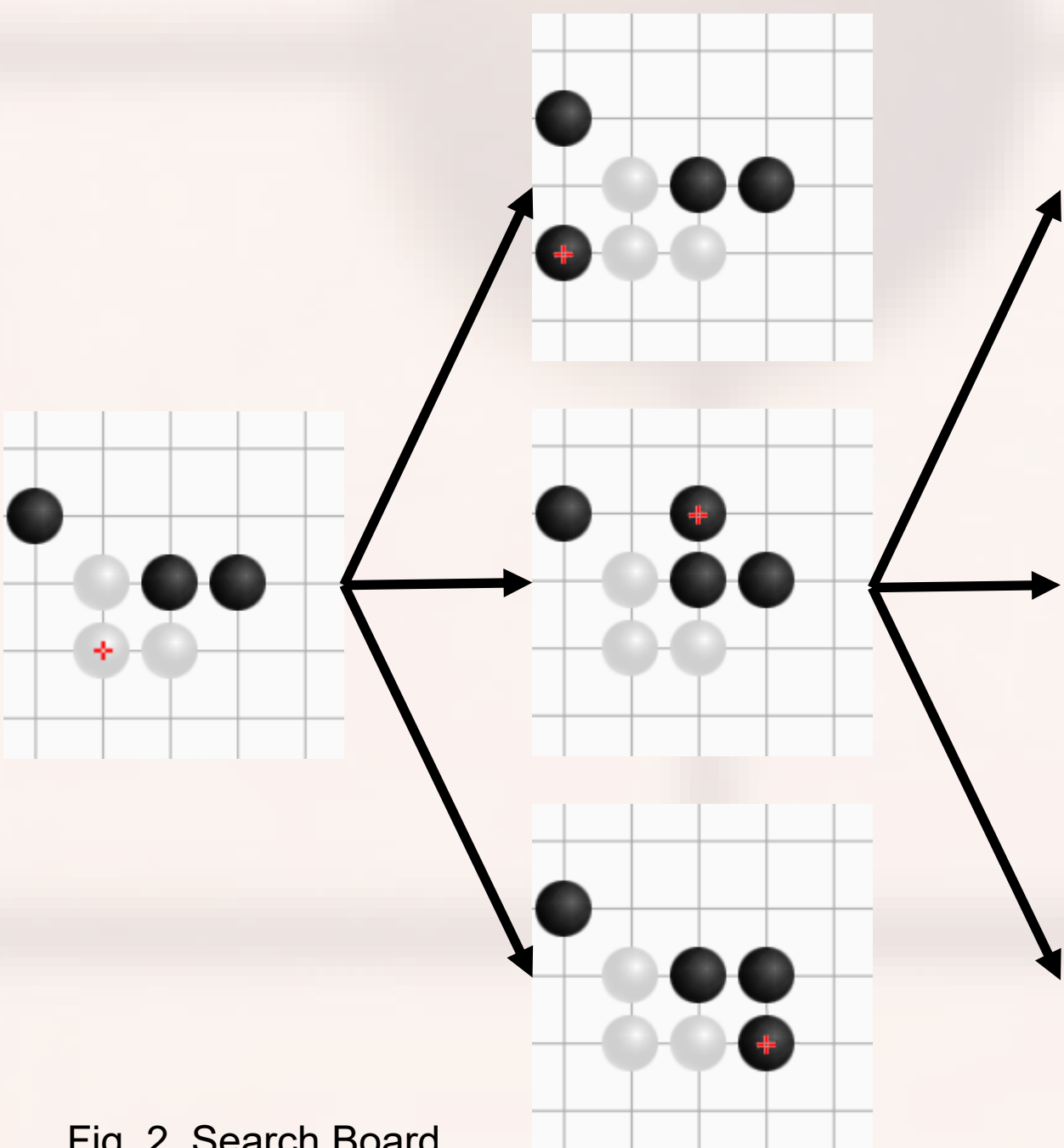


Fig. 2. Search Board

Feature:

First finished on 6.30.  
100% original code and scoring method.

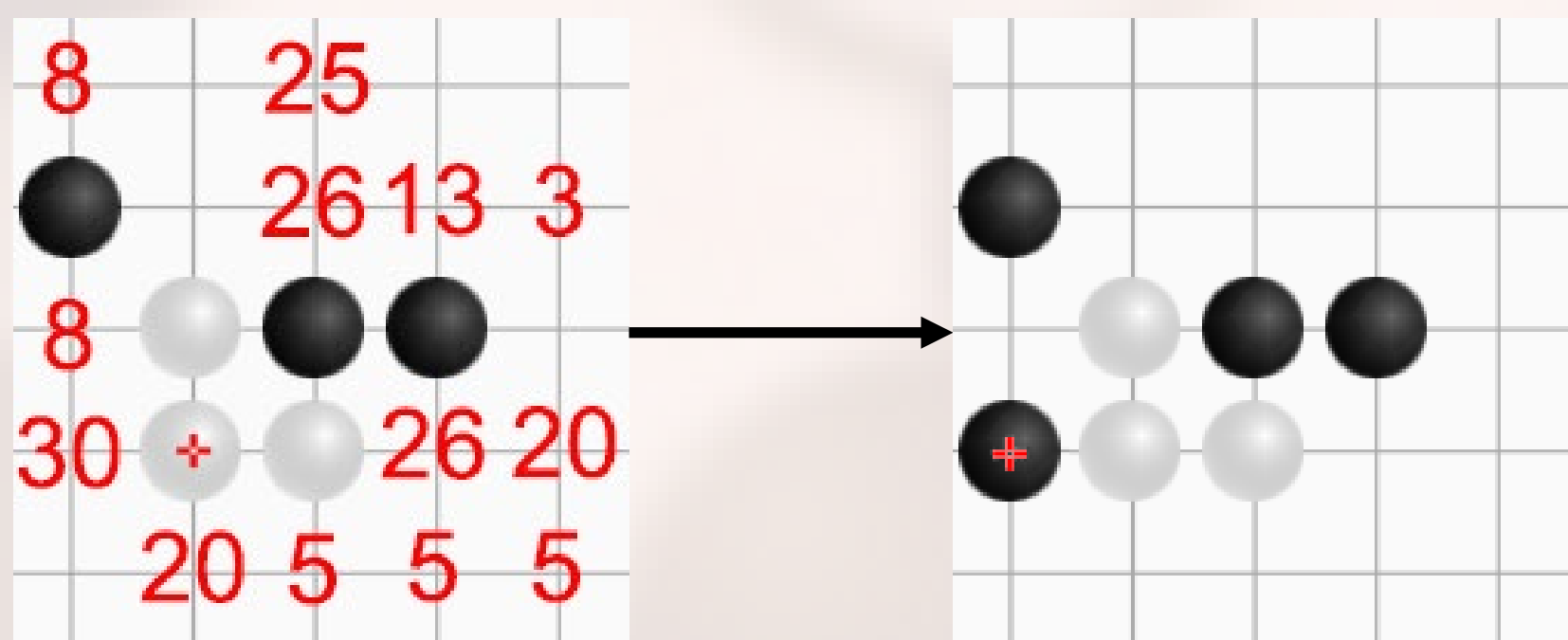


Fig. 1. Basic Logic

Basic logic(Fig. 1. ):

- 1. Scan the map in four directions
- 2. Add score to each position if certain situations are found.
- 3. Select position with max score.

$$\arg \max_p score(p)$$

Advantages:

- 1. Fast and easy.
- 2. Easy to explain.

Disadvantage:

- 1. Too easy, anyone can realize.
- 2. TOO MANY IF.
- 3. As for the course, this method use little python coding skills.
- 4. Once finished, its chess skill level is fixed.
- 5. Time limit exceed(for the search method only)

Lo\_Gobang\_Zero

Feature:

Reinforcement learning. Deep Q Network.  
First finished on 7.4 with pyTorch.  
Designed action reward function.  
Various training modes.  
100% original code.  
The **first** gobang AI based on TensorFlow 2.x on GitHub

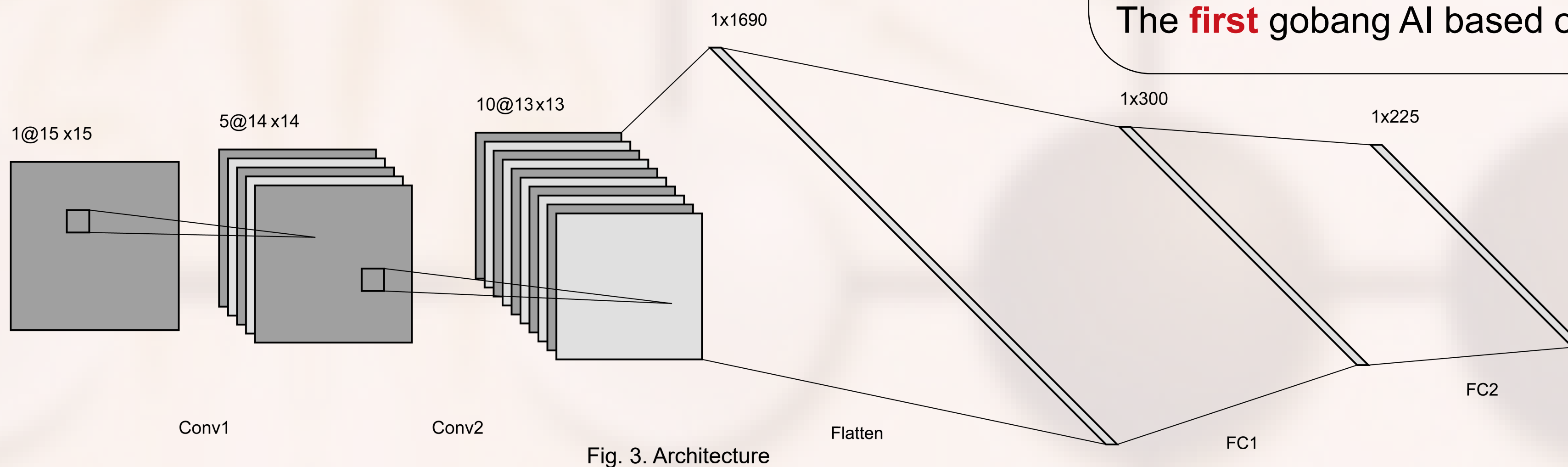


Fig. 3. Architecture

Chess game generation:

Use 2 Lo\_Gobang\_151s(151 below) to generate game experience of 6000 rounds, which is used to train a single agent. This agent is then trained with itself for 30000 rounds. The loss curve falls more steady(Fig. 4.).

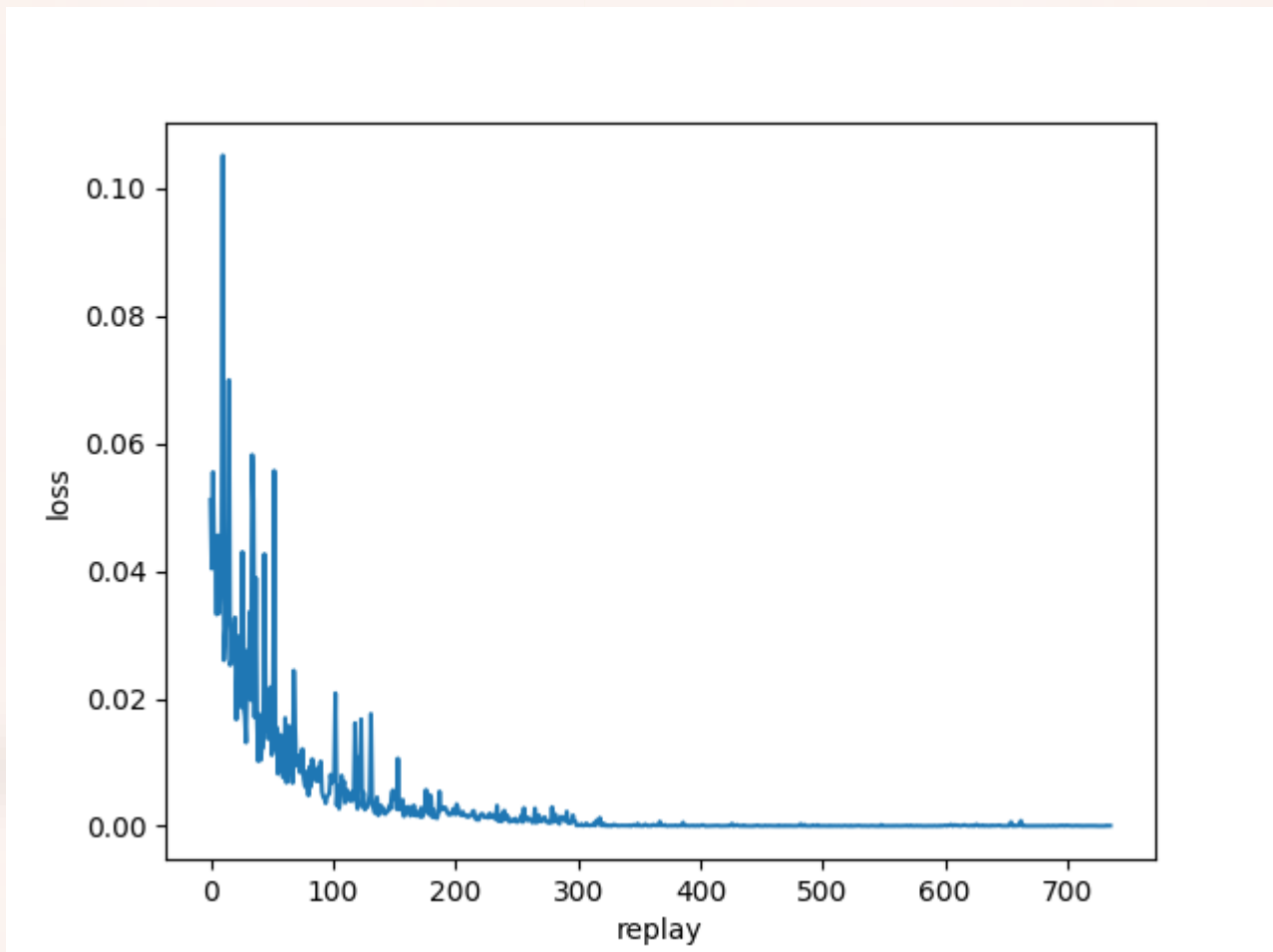


Fig. 4. Loss coach

About my DQN:

The network architecture is like Fig. 3. Conv is to attract feature of a chessboard block, and max pooling is avoided. In reinforcement learning, the agent needs to explore, and the explore step is given by a 151. The reward function is a scan-and-assess method, but it's different from the scoring method used in 151.

Key parameters:

Replay size: 50; Update frequency: 300  
Learning rate: 0.1; Q factor:0.5  
Training platform:  
Intel i9-11980HK, RTX3080laptop16gvr

Result of the model

trained with coach:  
Beats random 1000:0  
Can last over 20 steps playing with coach.

Self train:

In this case, two randomly initialized agents take turns to teach each other for 3000 rounds. The loss keeps low but is unstable. This model is also far better than random.

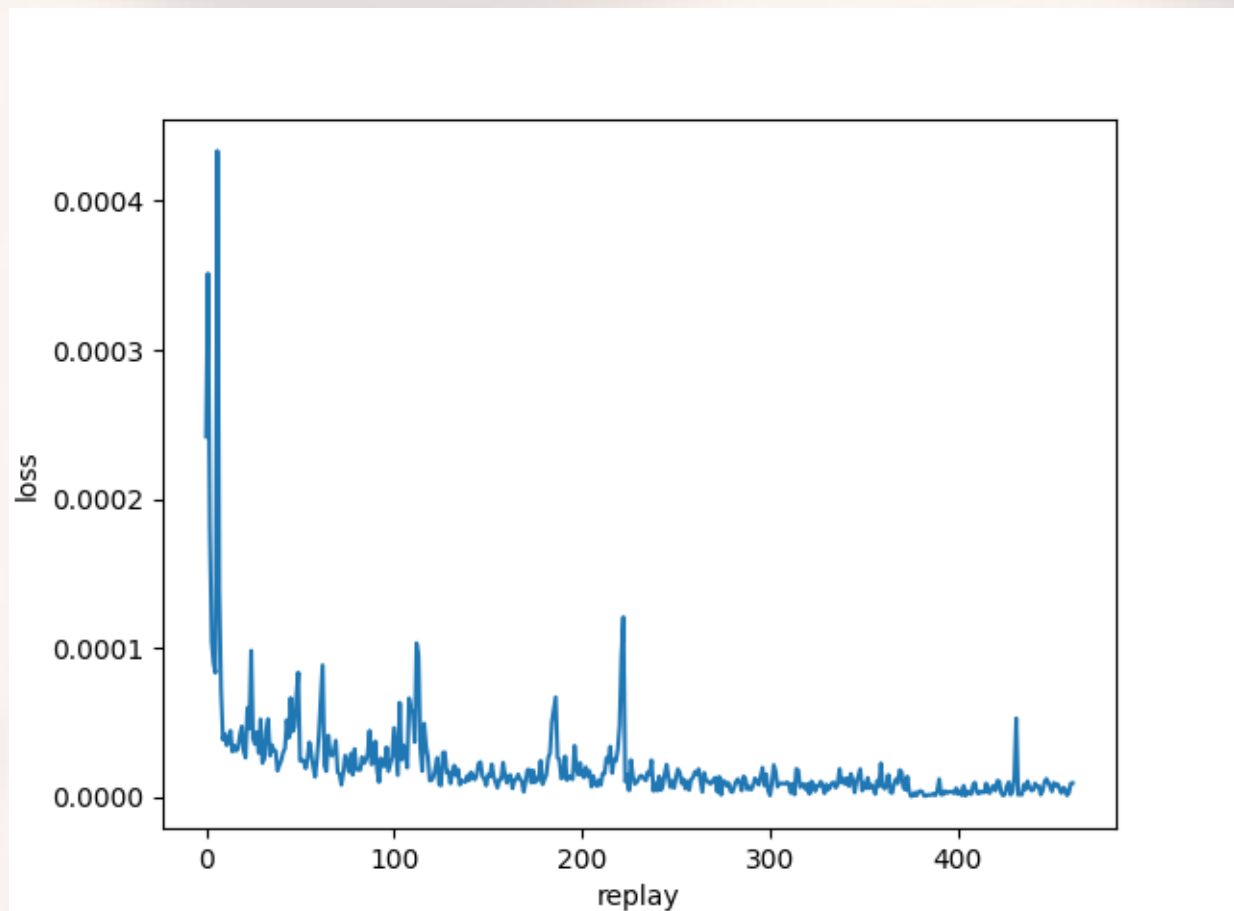


Fig. 5. Loss self