

,

\_\_\_\_\_ ( )

\_\_\_\_\_ ( )

\_\_\_\_\_ ( )

\_\_\_\_\_

\_\_\_\_\_

\_\_\_\_\_ ( )

:

\_\_\_\_\_ II \_\_\_\_\_ -20-1

\_\_\_\_\_ ( , )

123 « ' »

\_\_\_\_\_ ( )

\_\_\_\_\_ - \_\_\_\_\_ ( - - )

\_\_\_\_\_ ( )

: \_\_\_\_\_ ( , , )

\_\_\_\_\_ ( ) \_\_\_\_\_ ( , )

,

( )

123 « ' »

( )

-

( - - )

( )

:

“ ” 20 .

( , , )

1.

“ 5 ” 2021 . 1657  
13 2021 .

2.

3.

1)

2) Python

3) PyTorch

4) HandyRL

4.

1)

2)

« »

3)

Kaggle

5. \_\_\_\_\_ , \_\_\_\_\_ , \_\_\_\_\_ , \_\_\_\_\_ , \_\_\_\_\_  
 ( ) \_\_\_\_\_  
 - -13

---



---



---



---



---



---



---



---

6. \_\_\_\_\_ , \_\_\_\_\_ .1) ( \_\_\_\_\_ )

	( _____ , _____ , _____ , _____ )		

1		09.11.21-18.11.21	
2		19.11.21-22.11.21	
3		23.11.21-29.11.21	
4		29.11.21-30.11.21	
5		01.12.21-02.12.21	
6	-	03.12.21-04.12.21	

: 76 ., 25 ., 5 ., 28

REINFORCEMENT LEARNING, SNAKE GAME, DEEP Q-LEARNING,  
MONTE-CARLO SEARCH TREE, ARTIFICIAL NEURAL NETWORKS,  
PYTHON.

« »,

.

.

« »,

« »,

,

,

.

## ABSTRACT

Master's thesis: 76 pages, 20 figures, 5 tables, 28 sources.

REINFORCEMENT LEARNING, SNAKE, DEEP Q-LEARNING,  
MONTE-CARLO SEARCH TREE, NEURAL NETWORKS, PYTHON.

The aim of the work is to study the work of agents in the course of which it is necessary to create a game artificial intelligence for the game "Snake", built on adaptive game artificial intelligence, using training with reinforcement and depth of the neural network. The agent must learn from previous experience and demonstrate the result while playing with the user. As the number of players played with the agent grows, the difficulty of advancing to the next game should increase.

In the course of the qualification work, an adaptive game agent based on SNM for the game "Snake" was created, an environment for the agent to play the game "Snake" was created, the results of training and confirmation of the agent during the research were evaluated. extensive analysis of the work of the agent created interaction with competitors and the steps taken to improve the agent created in the qualification work.

	,	,	,	
	.....			8
	.....			9
1	.....			11
1.1	.....			11
1.2	.....			12
1.3	.....			15
1.4	.....			18
2	.....			20
2.1	CRISP-DM.....			20
2.2	Reinforcement learning.....			22
2.3	Q-learning.....			24
2.4	Deep Q-learning.....			27
2.5	PyTorch.....			28
2.6	Monte-Carlo tree search.....			30
2.7	Stacking.....			32
2.8	HandyRL.....			32
2.9	-.....			32
3	.....			34
3.1	.....			34
3.2	.....			39
3.3	.....			43
4	.....			50
4.1	.....			50
4.2	.....			57
4.2.1	.....			57
4.2.2	.....			57

.....	61
.....	62
.....	65
.....	73

, , ,

–

–

MC – ( ., Monte-Carlo)

MCST – - ( ., Monte-Carlo Search Tree)

RL – ( ., Reinforcement learning)

[1, 5].

(NPC),

1950-

1950-

[21-24].

Reinforcement learning (RL)

RL

[5].

RL,

RL

, RL,

,

, RL-

,

[5].

RL

Atari [4],

Go [19].

Deep RL

RL

Deep RL,

[25].

« » [26].

( ) Reinforcement learning

1

1.1

( )

[5].

« » [5].

« » –

( ),

[7].

« »

[7].

Kaggle [8]. Kaggle –

Google,

«Hungry Geese» [9, 11].

« »

( , ),

[13].

## 1.2

NPC,

( , , , , ) [5].

, .  
,  
, ,  
,  
.

:  
( ) [1].

,  
.  
. , « »  
,

, .  
, .

10

.  
,  
.  
,

, ,  
,  
.  
.

( ). RL

,

,

.

,

,

,

.

«

»

:

;

-

,

;

-

;

-

.

-

:

-

;

-

.

«

»

;

-

.

,

(

)

:

-

,

-

;

-

;

:

.

- ;  
 - ;  
 - ;

1.3

« ».

[7, 13, 26]:

- ;  
 - ;  
 - ;  
 - ;  
 - ;  
 - ;

« »

[26].

[26]: (greedy),  
 (risk averse greedy),

.  
 .  
 ,  
 ,

[26].  
(Flood fill) [26].

Monte-Carlo Search Tree [2].

Python

Monte-Carlo Search Tree

Python

++.

:

-

;

-

;

-

;

-

Flood fill

.

,

Kaggle ( 1.1).

1.1 –

« »

	.	
	,	
	,	.
	200	
	,	
	,	.

1.1

<p>a      Flood fill</p>	<p>.</p> <p>230</p> <p>,</p>	<p>.</p>
--------------------------	------------------------------	----------

«    ».

,

:

Flood

fill.

1.4

,

,

«    »

,

.

.

:

Deep Q-Learning

Monte-Carlo Search Tree.

[20].

« ».

5

Kaggle [8].

Deep Q-Learning

Monte-Carlo search Tree [17].

« ».

30





2.1 –

### CRISP-DM

CRISP-DM ( 2.1).



reinforcement learning –  
[6].

unsupervised [6].

supervise

[10].

supervised learning ( ),

[12].

« », AI

### 2.3 Q-learning

Q-learning – model-free . «Q» Q-learning

Q-learning value-based [14]. value-based ( policy-based )

Q-learning – off-policy . , on-policy learner

[15, 15, 22].

$$Q^*(s, a) = \sum_{s'} \sum_{a'} P(s', a' | s, a) [R(s, a) + \gamma V(s')]$$

### Q-learning Temporal Differences (TD)

$$Q^*(s, a)$$

Temporal Differences

Q-

$$Q(s, a) \leftarrow (1 - \alpha) Q(s, a) + \alpha (R_t + \gamma V(s_{t+1}) - Q(s, a))$$

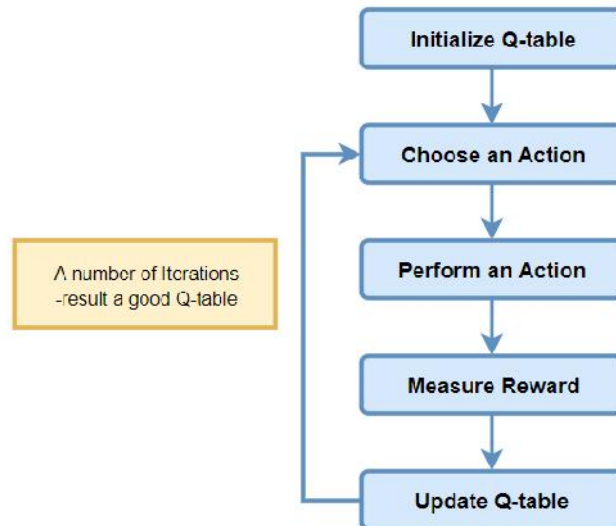
$$Q^\pi(s_t, a_t) = \underline{E}[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | s_t, a_t]$$

Q-Values for the state given a particular state
Expected discounted cumulative reward
Given the state and action

2.2 -

### Q-learning

$$(2.3)$$



2.3 -

### Q-Learning

Q-table,  $m$ ,  $m$  (2.4).

	↑	↓	→	←
Action				
Start	0	0	0	0
Idle	0	0	0	0
Correct Path	0	0	0	0
Wrong Path	0	0	0	0
End	0	0	0	0

2.4 – Q-table,

$Q[S, A]$ ,  $S$ ,  $A$ .  
 $Q[s, a]$   $Q^*(s,a)$ .  
 Q-Table –

Learning.

(a) (s).  
 0.

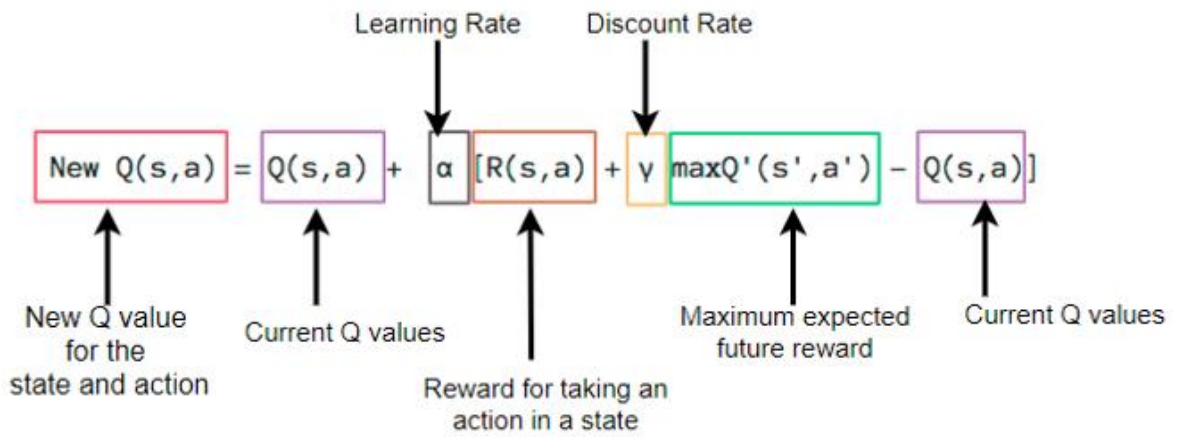
Epsilon.

Epsilon

Epsilon

[15, 16].

$$Q(s, a) = Q(s, a) + \alpha [R(s, a) + \gamma \max_{a'} Q(s', a') - Q(s, a)] \quad (2.5)$$



2.5 – Q-

### 2.4 Deep Q-learning

Deep Q-learning [14]

Q- , Q-

E

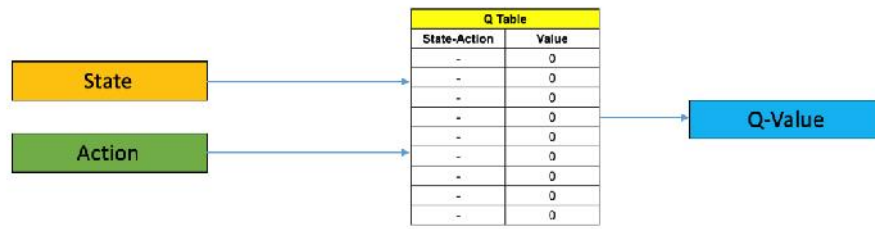
Q- (DQN)

Deep Q-learning

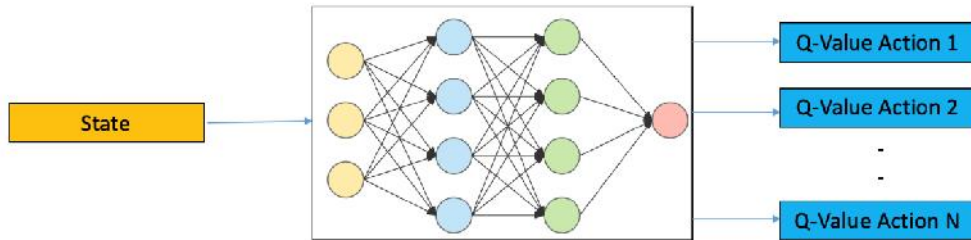
network.

Q- Q- – Q\*

Q-learning Deep Q-learning ( 2.6).



Q Learning



Deep Q Learning

2.6 – Q-Learning    Deep Q-Learning

Q, (2.1):

$$Q(S_t, A_t) = Q(S_t, A_t) + \alpha [R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t)] \quad (2.1)$$

, R, backprogration,

2.5 PyTorch

. PyTorch [8] –

, Pythonic .

« »,

,

. Pytorch

autograd

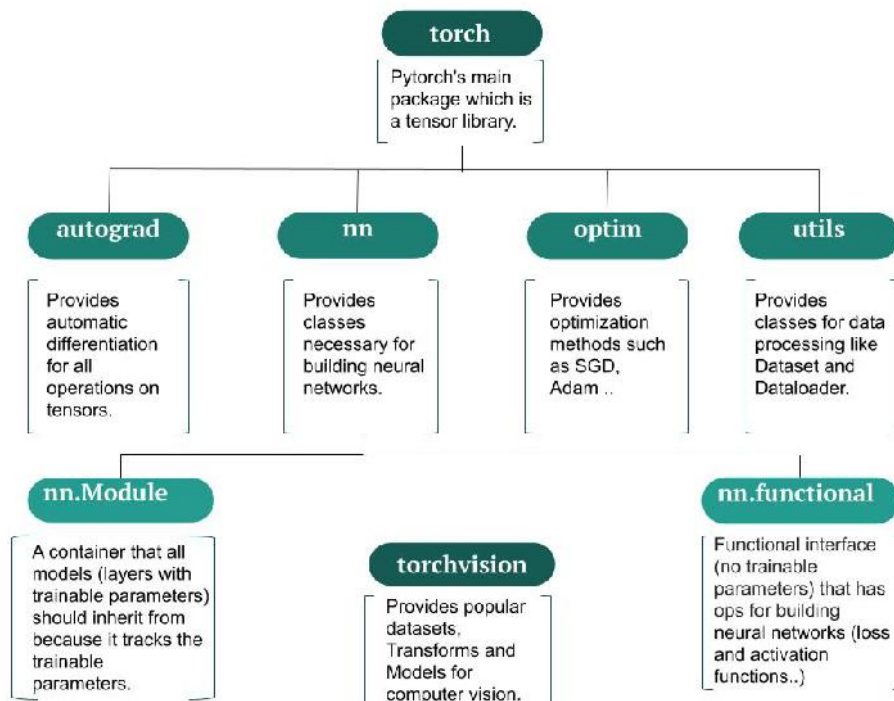
NumPy [11].

Matlab

Python

NumPy, SciPy Pandas.

( 2.7).



2.6 Monte-Carlo tree search

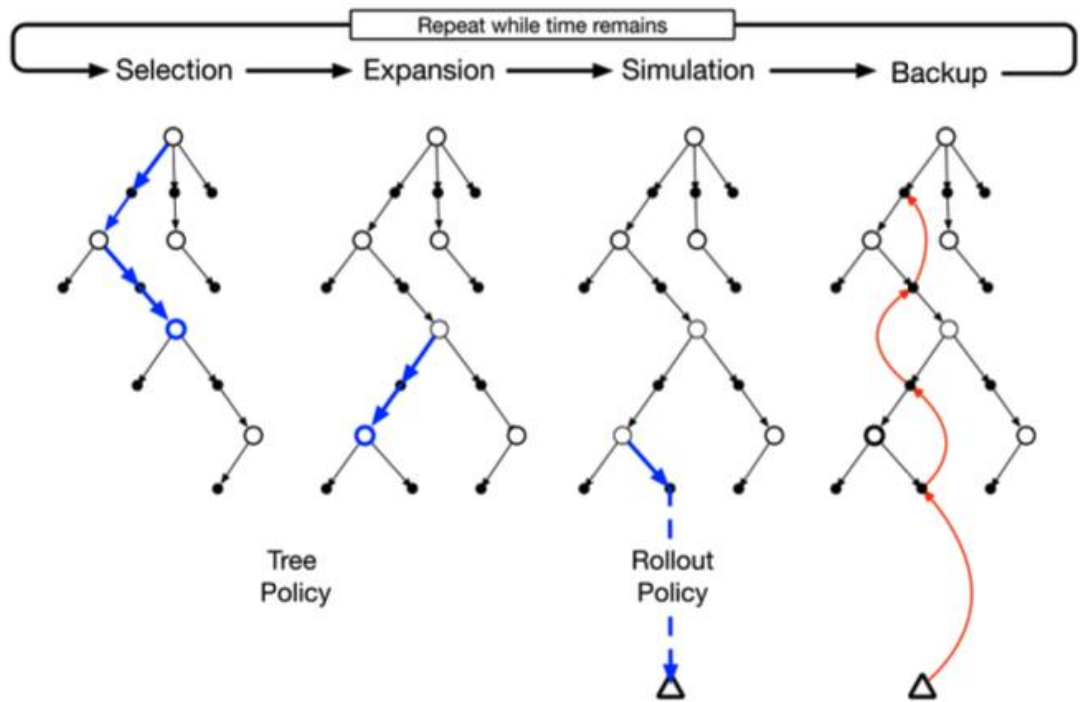
MCTS [18, 20]

. MCTS

[18].

tree policy

( 2.8).



2.8 –

Monte-Carlo tree search

Tree policy –

( )

(

).

(exploration vs

exploitation).

Reinforcement learning.

AlphaGo

Upper Confidence Bound (UCB) [18].

UCB,

UCB (2.2).

$$USB(node_i) = x_i + c \sqrt{\log N / n_i} \quad (2.2)$$

$x_i$  - ,  $n_i$  - ,  $N$  -

.

UCB, ,

UCB

.

.

.

.

.

,

,

.

+1, -

0 - 1.

.

.

.

,

.

.

## 2.7 Stacking

Stacking – ,

voting,

( , ).

Stacking

Kaggle [19].

## 2.8 HandyRL

HandyRL [3] –

, Python

PyTorch

,

HandyRL

HandyRL

«Policy gradient algorithm with

off-policy correction»

off-

policy

( - , TD( ))

(V-Trace, UPGO). HandyRL

learner-

worker, IMPALA [12]. –

,

( )

.

## 2.9 -

- (cross-validation, CV) –

,

[28].

:

( 2.2).

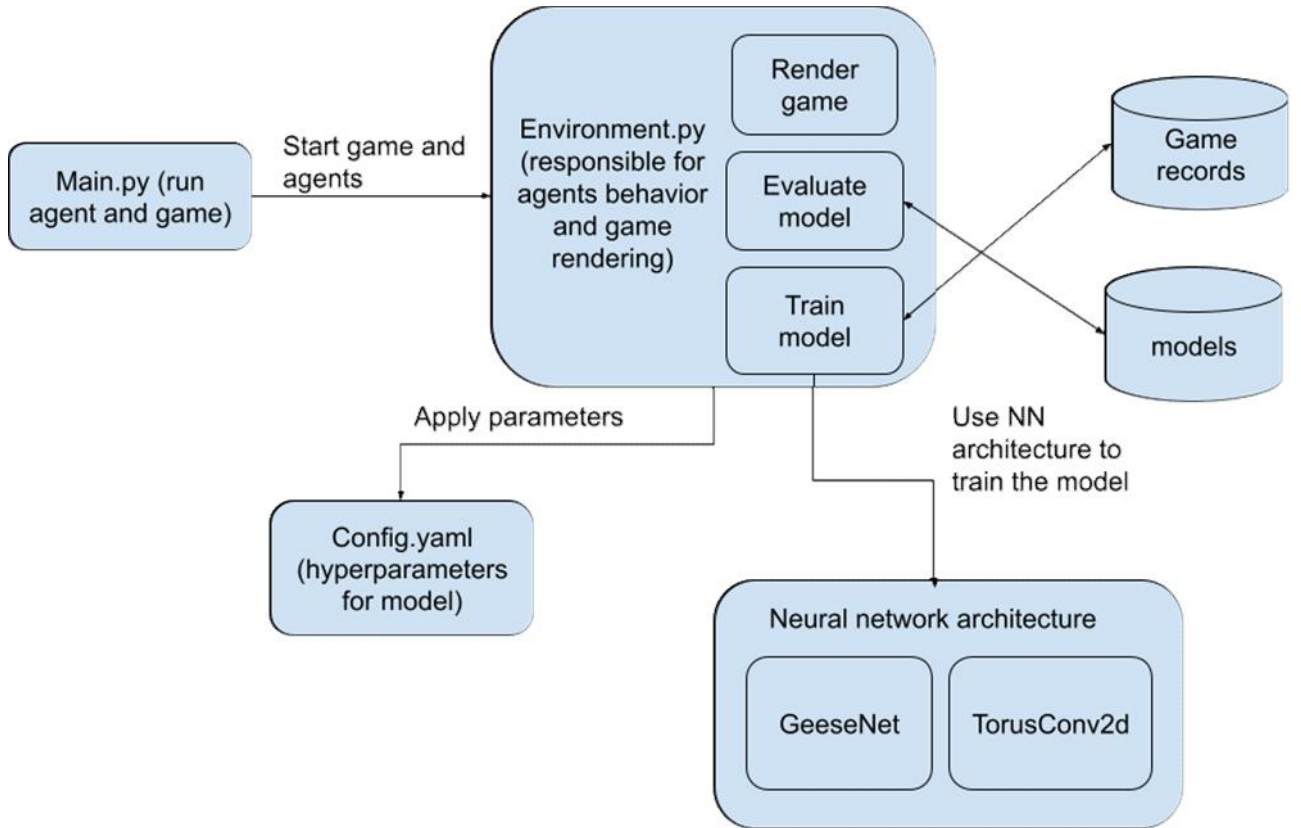
2.2 – c-

<p>Hold-out CV</p>	<p>-</p>
<p>Leave One Out (LOO)</p>	<p>Leave One Out ( LOO) –</p>
<p>K-Fold</p>	<p>Fold</p> <p>( )</p>



100 models,

game records.



3.1 –

kaggle\_environments,

kaggle\_environments

« »

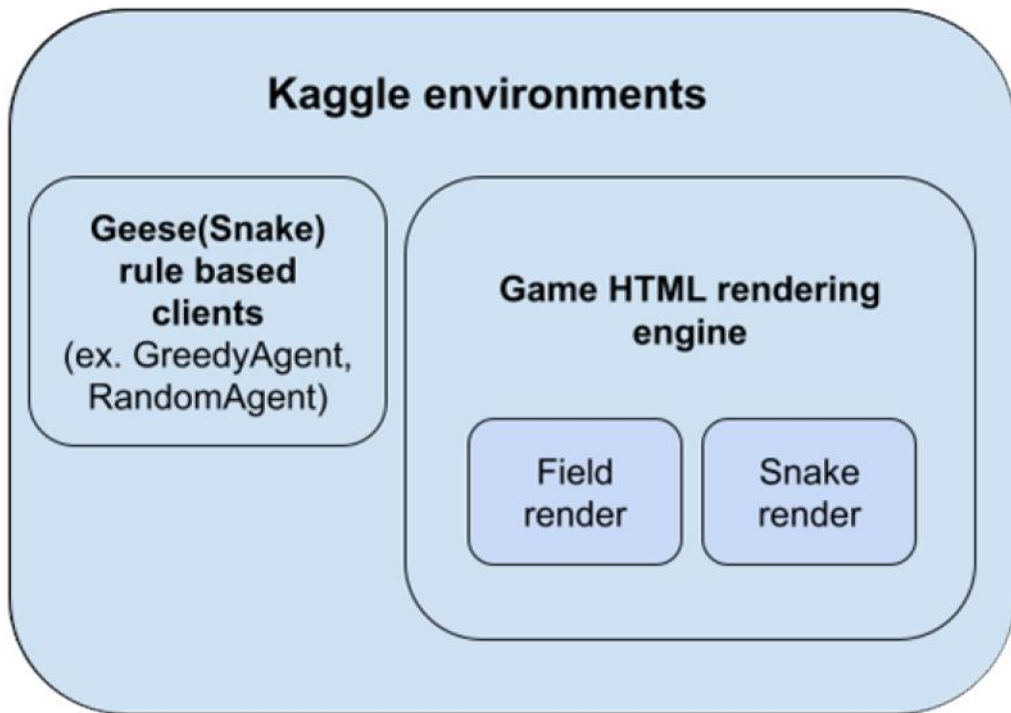
( 3.2). , GreedyAgent

API

( )

2

kaggle



3.2 –

kaggle\_enviroments

battles.

kaggle

( 3.3).



3.3 –

PyCharm

Python

:

- neural\_net.py;
- torus\_net.py;
- environment.py;
- main.py;
- config.yaml.

neural\_net.py

torus\_net.py

environment.py

HandyRL

main.py

( 3.1),

environment.py

config.yaml –

HandyRL.

MC.

game\_records, models battles,

game\_records

game\_records

.pickle

Python-

Pickle-

Models

battles

(

.ipynb),

## 3.2

Python,

Q-learning Deep Q Learning.

(policy) –

(value)

– ( 3.1).

3.1 –

```

class TorusConv2d(nn.Module):
    def __init__(self, input_dim, output_dim, kernel_size, bn):
        super().__init__()

        self.edge_size = (kernel_size[0] // 2, kernel_size[1] // 2)
        self.conv = nn.Conv2d(input_dim, output_dim, kernel_size=kernel_size)
        self.bn = nn.BatchNorm2d(output_dim) if bn else None

    def forward(self, x):
        h =
        torch.cat([x[:, :, :, -self.edge_size[1]:], x,
x[:, :, :, :self.edge_size[1]]], dim=3)
        h = torch.cat([h[:, :, -self.edge_size[0]:], h,
h[:, :, :self.edge_size[0]]], dim=2)
        h = self.conv(h)
        h = self.bn(h)
        if self.bn is not None
        else h
        return h

```

nn.Module,

PyTorch.

TorusConv2d

\_\_init\_\_

.

.

,

,

.

,

0,

–

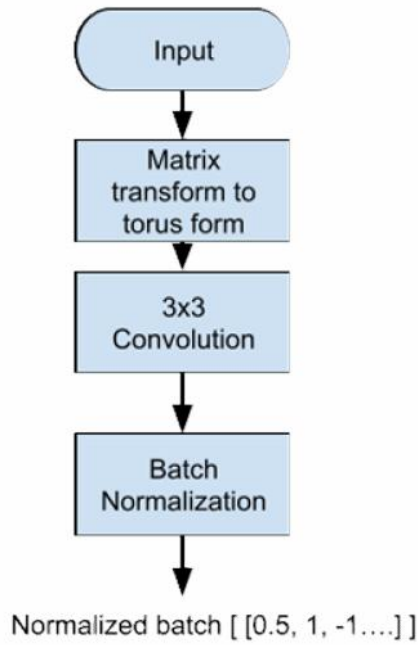
1.

,

\_\_forward\_\_

( 3.4).

**Residual block (Torus Conv)**



3.4 –

GeeseNet ( 3.2),  
 ,  
 3 3 12  
 3 3.

linear.

3.2 –

```
class GeeseNet(nn.Module):
    def __init__(self):
        super().__init__()
        layers, filters = 12, 32

        self.conv0 = TorusConv2d(17, filters, (3, 3), True)
        self.blocks = nn.ModuleList([TorusConv2d(filters, filters, (3, 3),
True) for _ in range(layers)])
        self.head_p = nn.Linear(filters, 4, bias=False)
        self.head_v = nn.Linear(filters * 2, 1, bias=False)

    def forward(self, x, _=None):
        h = F.relu_(self.conv0(x))
        for block in self.blocks:
            h = F.relu_(h + block(h))
        h_head = (h * x[:, :1]).view(h.size(0), h.size(1), -1).sum(-1)
        h_ v}
```

( 3.6)

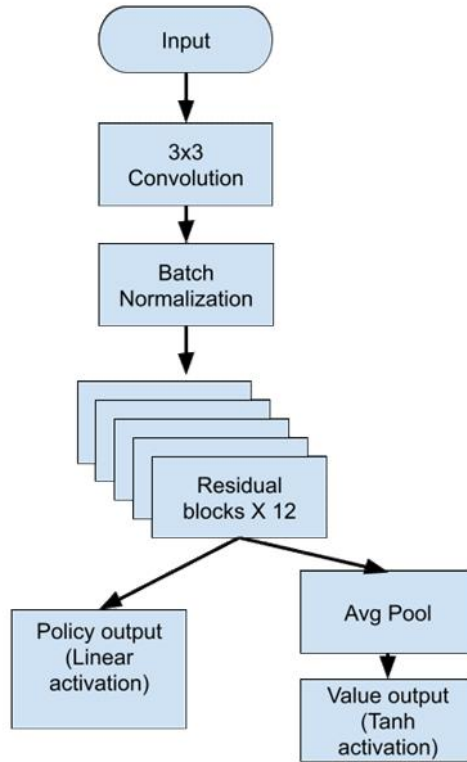
(residuals).

head extraction head expansion.

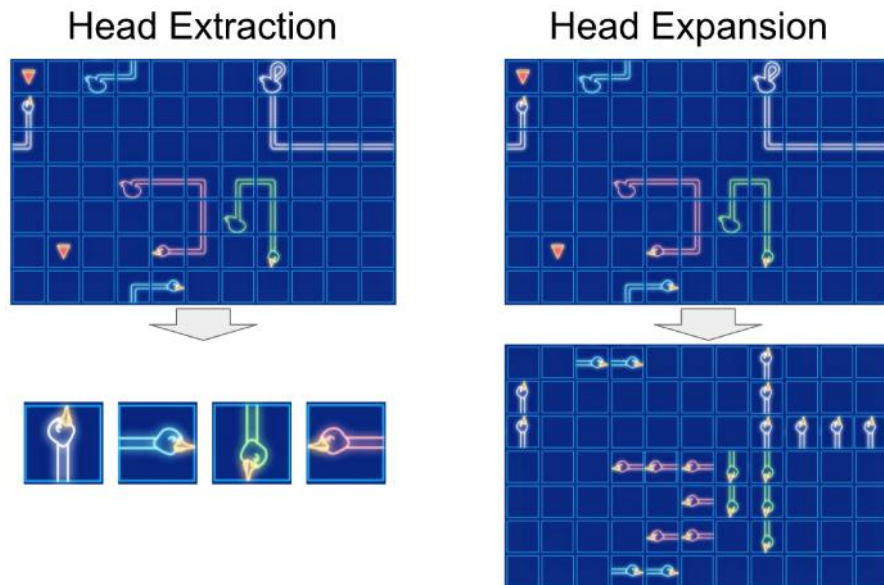
( )

( 3.7).

**Network architecture**



3.6 –



3.7 –

head extraction head expansion



## . BaseEnvironment

config.yaml.

```

BaseEnvironment :
-   ,           «   »
  kaggle_environments;
-   rule_based_action,
    «   »;
-   observation,
        ;
-   step,
    ;
-   turns,
        ;
-   terminal,           ,
        ;
-   outcomes,
.
        (   3.4),
,
.
ACTION ,
.
DIRECTION - ,
, -
, , «SOUTH»,
.
NUM_AGENTS ,
.

```

3.4 –

environment.py

```
ACTION = ['NORTH', 'SOUTH', 'WEST', 'EAST']
DIRECTION = [[-1, 0], [1, 0], [0, -1], [0, 1]]
NUM_AGENTS = 4
```

kaggle\_environments ( 3.5).

« »

«hungry\_geese»

make.

3.5 –

```
def __init__(self, args={}):
    super().__init__()
    self.env = make("hungry_geese")
    self.reset()
```

rule\_based\_action ( 3.6),

GreedyAgent

3.6 –

« »

```
def rule_based_action(self, player):
    from kaggle_environments.envs.hungry_geese.hungry_geese import
    Observation, Configuration, Action, GreedyAgent
    action_map = {'N': Action.NORTH, 'S': Action.SOUTH, 'W': Action.WEST,
    'E': Action.EAST}
    agent = GreedyAgent(Configuration({'rows': 7, 'columns': 11}))
    agent.last_action =
    action_map[self.ACTION[self.last_actions[player]][0]] if player in
    self.last_actions else None
    obs = {**self.obs_list[-1][0]['observation'], **self.obs_list[-
    1][player]['observation']}
    action = agent(Observation(obs))
    return self.ACTION.index(action)
```

observation ( 3.7)

net.

GeeseNet

HandyRL

observation

3.7 –

« »

```
def observation(self, player=None):
    if player is None:
        player = 0
    b = np.zeros((self.NUM_AGENTS * 4 + 1, 7 * 11), dtype=np.float32)
    obs = self.obs_list[-1][0]['observation']
    for p, geese in enumerate(obs['geese']):
        # head position
        for pos in geese[:1]:
            b[0 + (p - player) % self.NUM_AGENTS, pos] = 1
        # tip position
        for pos in geese[-1:]:
            b[4 + (p - player) % self.NUM_AGENTS, pos] = 1
        # whole position
        for pos in geese:
            b[8 + (p - player) % self.NUM_AGENTS, pos] = 1

    # previous head position
    if len(self.obs_list) > 1:
        obs_prev = self.obs_list[-2][0]['observation']
        for p, geese in enumerate(obs_prev['geese']):
            for pos in geese[:1]:
                b[12 + (p - player) % self.NUM_AGENTS, pos] = 1

    # food
    for pos in obs['food']:
        b[16, pos] = 1

    return b.reshape(-1, 7, 11)
def net(self):
    return GeeseNet
```

turns ( 3.8)

terminal.



```

outcomes = {p: 0 for p in self.players()}
for p, r in rewards.items():
    for pp, rr in rewards.items():
        if p != pp:
            if r > rr:
                outcomes[p] += 1 / (self.NUM_AGENTS - 1)
            elif r < rr:
                outcomes[p] -= 1 / (self.NUM_AGENTS - 1)
return outcomes

```

main.py ( 3.11).

« »

3.11 –

```

if __name__ == '__main__':
    e = Environment()
    for _ in range(200):
        e.reset()
        while not e.terminal():
            print(e)
            actions = {p: e.legal_actions(p) for p in e.turns()}
            print([[e.action2str(a, p) for a in alist] for p, alist in
actions.items()])
            e.step({p: random.choice(alist) for p, alist in actions.items()})
            print(e)
            print(e.outcome())

```

FloodfullResult.

Floodfill

,

( 3.12).

3.12 –

```

@dataclass
class FloodfillResult:
    field_dist: np.ndarray
    frontiers: List[List[Tuple[int, int]]]

def flood_fill(is_occupied: np.ndarray, seeds: List[Pos]) -> FloodfillResult:
    """
    Flood will start with distance 0 at seeds and only flow where
    is_occupied[x,y]==0
    """
    size_x, size_y = is_occupied.shape

```

```

field_dist = np.full(fill_value=-1, shape=(size_x, size_y))
frontier = [(s.x, s.y) for s in seeds]
frontiers = [frontier]
for seed in seeds:
    field_dist[seed] = 0
dist = 1

while frontier:
    new_frontier: List[Tuple[int, int]] = []
    for x, y in frontier:
        for dx, dy in [(0, 1), (0, -1), (1, 0), (-1, 0)]:
            new_x = (x + dx) % size_x
            new_y = (y + dy) % size_y
            if is_occupied[new_x, new_y] == 0 and field_dist[new_x,
new_y] == -1:
                field_dist[new_x, new_y] = dist
                new_frontier.append((new_x, new_y))
    frontier = new_frontier
    frontiers.append(frontier)
    dist += 1
return FloodfillResult(field_dist=field_dist, frontiers=frontiers)

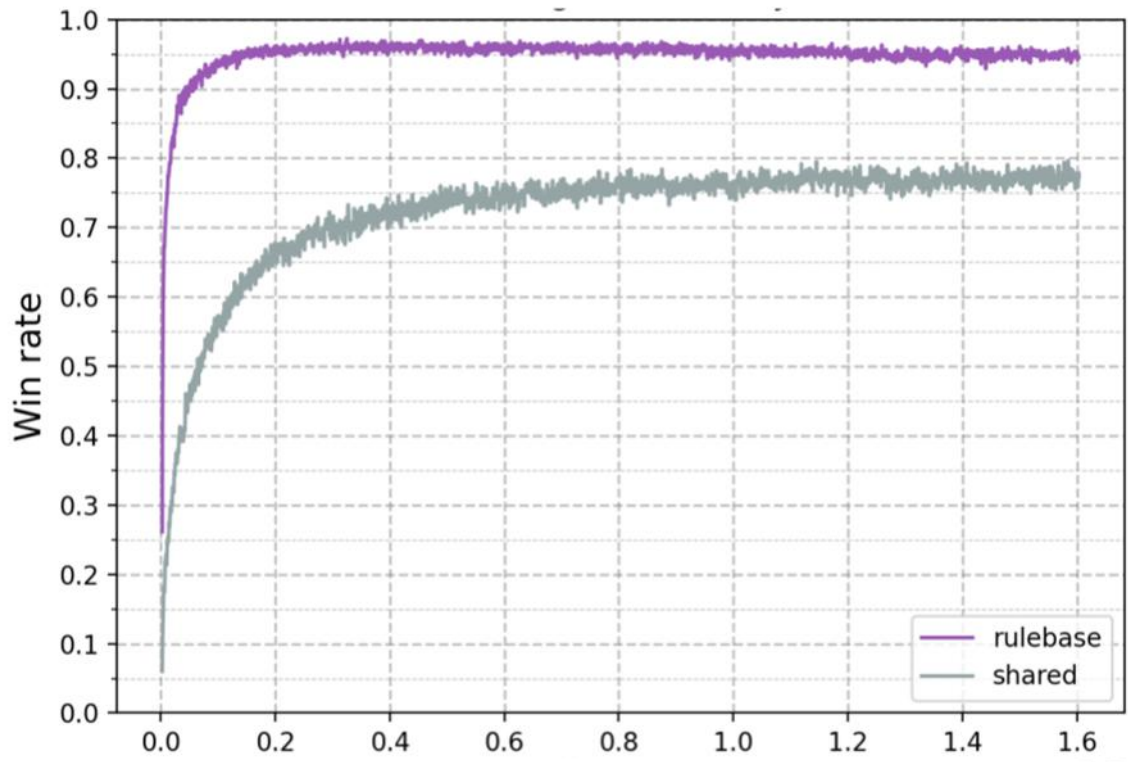
def get_dist(
    floodfill_result: FloodfillResult, test_func: Callable[[Tuple[int, int]],
bool]
) -> Optional[int]:
    for dist, frontier in enumerate(floodfill_result.frontiers):
        for pos in frontier:
            if test_func(pos):
                return dist
    return None

```

4

4.1

,  
 .  
 .  
 ( 4.1),  
 . Rulebase  
 .  
 .  
 20 000  
 ,  
 ( 100%).  
 ,  
 ( 4.1). Shared  
 .  
 ,  
 ,  
 ,  
 70 40000  
 ,  
 40000  
 ,  
 .  
 ,  
 .  
 0.7 0.8 , 4 .  
 ,



4.1 –

Rulebase-

Shared-

Kaggle

200

4.2.

10, 25, 100

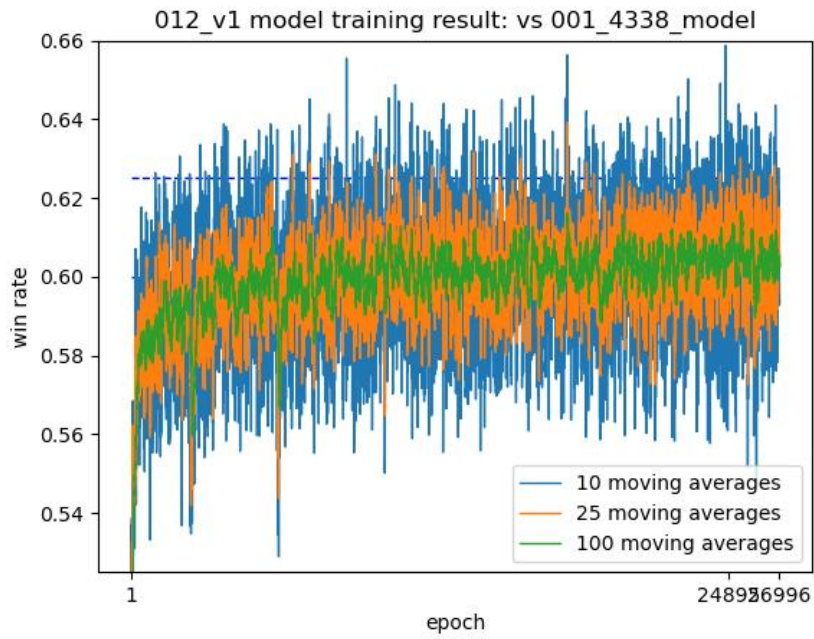
60

65

55

60

1000



4.2 –

2

2

:

(Smart Geese)

Kaggle.

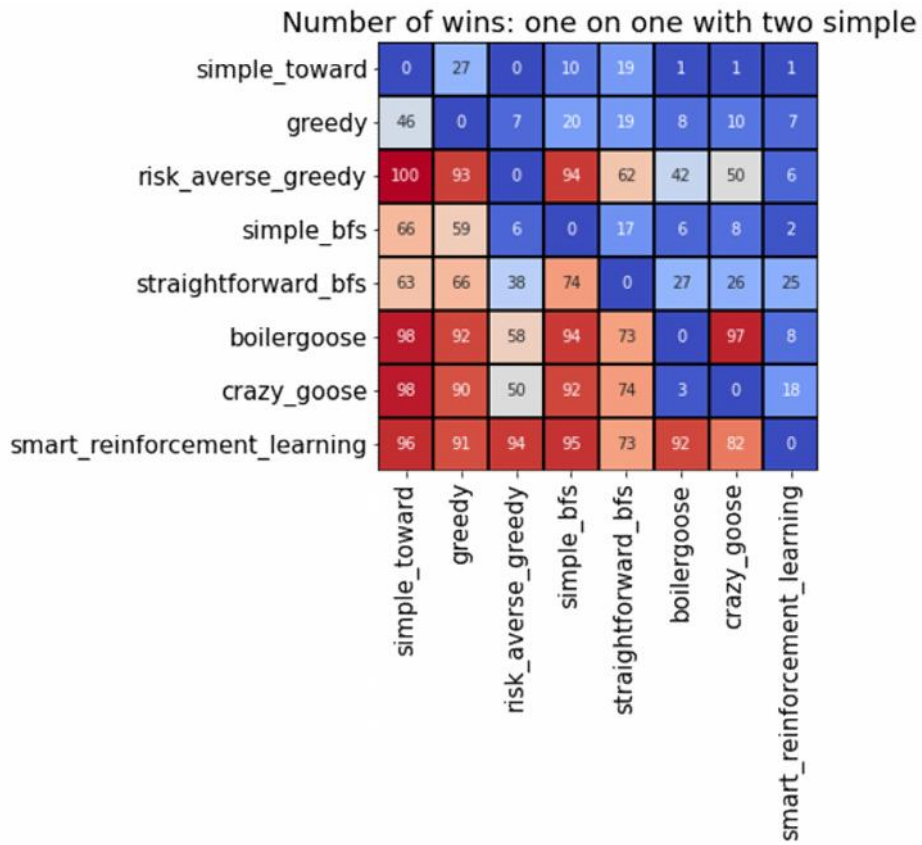
( 4.1).

## 4.1 –

		Smart Geese
simple_toward	,	96%
greedy		91%
risk_averse_greedy	, ( ,  , )	94%
simple_bfs		95%
crazy_goose	, ( ,  , )	82%

heatmap (

4.3):



### 4.3 – Heatmap

, 3  
,  
( 4.2).

### 4.2 –

« »

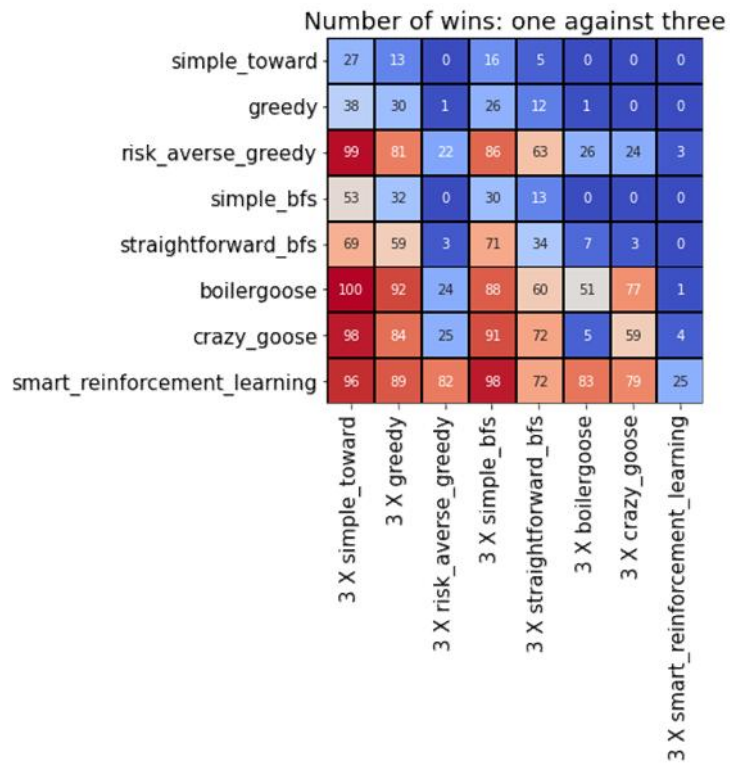
		Smart Geese
simple_toward		96%
greedy		89%

4.2

risk_averse_greedy	( , , )	82%
simple_bfs		98%
crazy_goose	( , , )	79%

heatmap (

4.4).



4.4 – Heatmap

Kaggle,

( 4.3).

4.3 –

« »

		Smart Geese
genetic_agent		67%
inclined_risk_agent		56%
risk_averse_greedy		60%



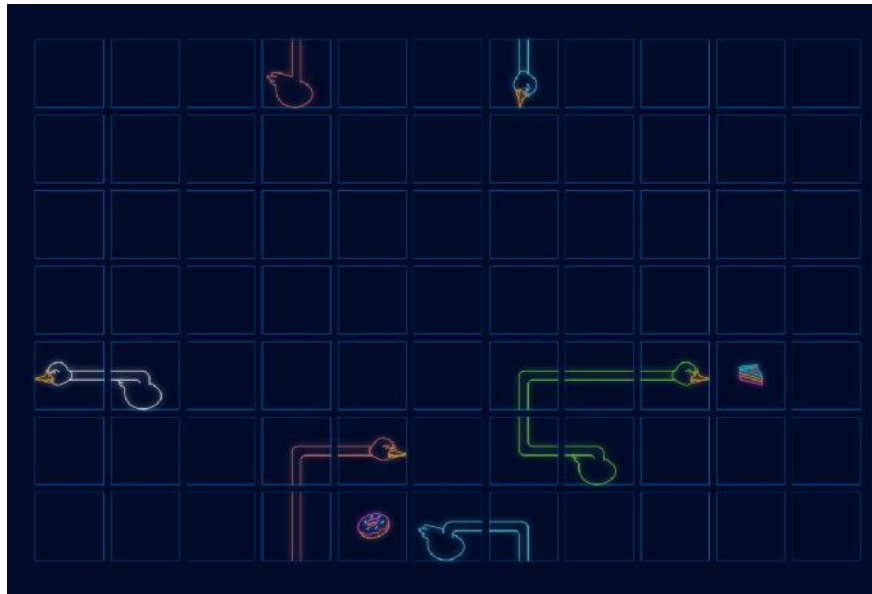
4 , , ,  
 , :  
 - ;  
 - ;  
 - .  
 1 , 2  
 , 3 ,  
 , ,  
 .



4.5 –

, Kaggle  
 ,  
 . « » (

4.6).



4.6 –

30

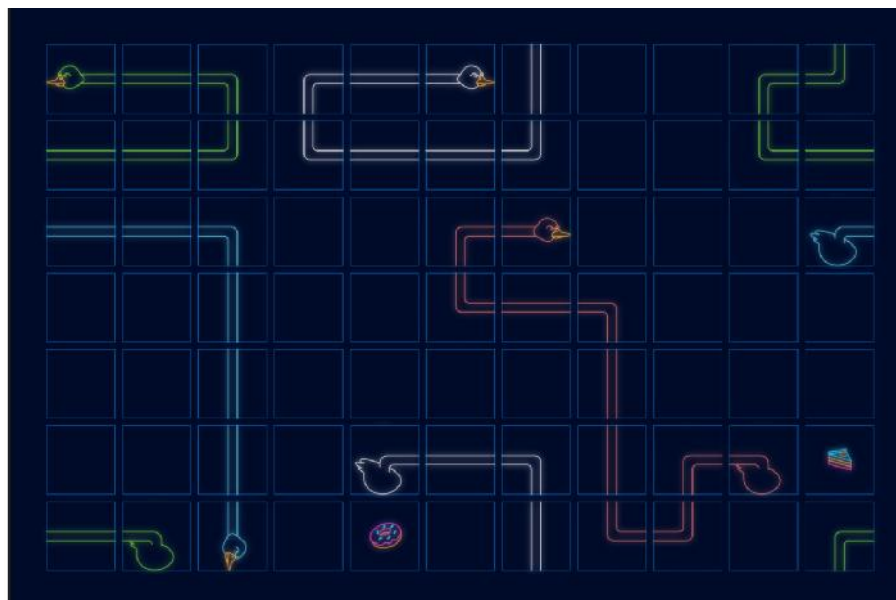
30

70

2

2

( 4.7).



4.7 –

70

90

-

«

»

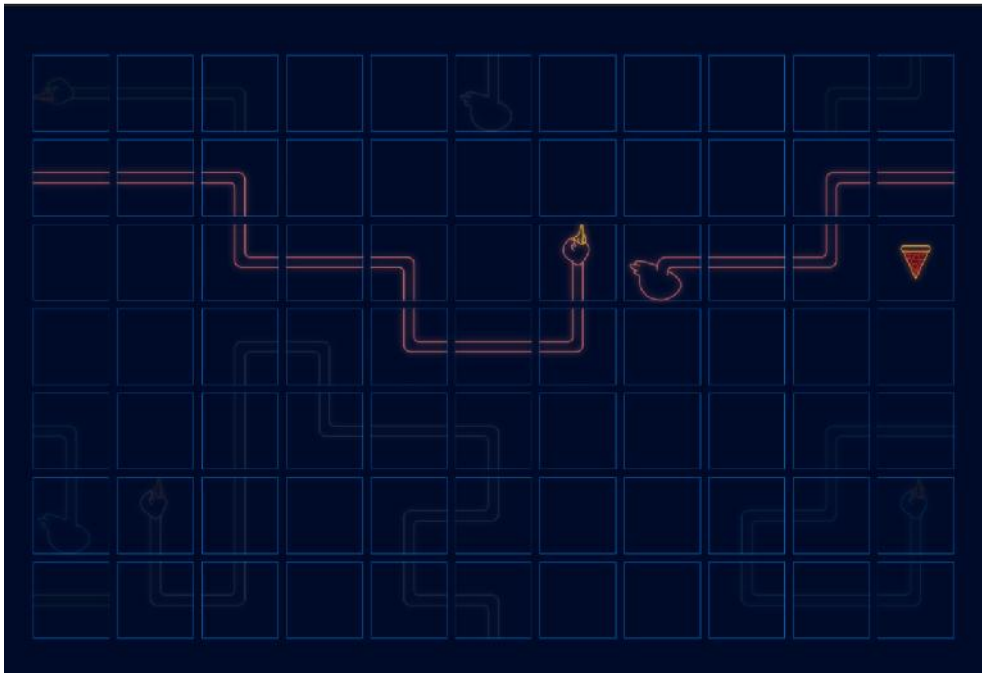
.

90

103

1

( 4.8).



4.8-

103

Kaggle,

4

« »  
(Reinforcement learning).

« »  
Deep Q-Learning neural network Monte-Carlo tree search,

, ,

60  
«Kaggle».

«Hungry Geese» «Kaggle».

.  
,  
.

.



Library [ ] / A. Paszke, S. Gross // NeurIPS. – 2019. – P. 1–12.

11. Python [ ] – : www/ URL: <https://www.python.org>. – 11.11.2021 . – . .

12. Scalable Distributed Deep-RL with Importance Weighted Actor-Learner Architectures [ ] – : www/ URL: [https://github.com/deepmind/scalable\\_agent/tree/6c0c8a701990fab9053fb338ede9c915c18fa2b1](https://github.com/deepmind/scalable_agent/tree/6c0c8a701990fab9053fb338ede9c915c18fa2b1). – 11.11.2021 . – . .

13. Bing Z. Energy-Efficient Slithering Gait Exploration for a Snake-like Robot based on Reinforcement Learning [ ] / Z. Bing, C. Lemke // IJCAI. – 2019. – P. 1–5.

14. Fitzek D. Deep Q-learning decoder for depolarizing noise on the toric code [ ] / D. Fitzek, M. Eliasson. – P. 1–17.

15. Swaminathan A. Off-Policy Policy Gradient with State Distribution Correction [ ] / A. Swaminathan, I. Liu. – 2019. – P. 1–17.

16. Islam R. Off-Policy Policy Gradient Algorithms by Constraining the State Distribution Shift [ ] / R. Islam, K. Teru. – 2019. – P. 1–13.

17. Chaslot G. Monte-Carlo Tree Search: A New Framework for Game AI [ ] / G. Chaslot, S. Bakkes // AIIDE. – 2008. – P. 1–2.

18. Gelly S. Monte-Carlo tree search and rapid action value estimation in computer Go [ ] / S. Gelly, D. Silver. – 2011. – P. 1–5.

19. Silver D. A general reinforcement learning algorithm that masters chess, shogi, and Go through self-play [ ] / D. Silver, T. Hubert. – 2018. – P. 1–3.

20. Hendrycks D. Gaussian error linear units (GELUS) [ ] / D. Hendrycks, K. Gimpel. – 2016. – P. 1–9.

21. Frey B. Adaptive dropout for training deep neural networks [ ] / B. Frey, L. Ba // NIPS. – 2013. – P. 1–2.

22. N Tsitsiklis J. analysis of temporal difference learning with function approximation. [ ] / J. N Tsitsiklis, V. Roy // NIPS. – 1996. – P. 1–2.

23. V J. analysis of temporal difference learning with function approximation [ ] / J. V, B. Van Ro // NIPS. – 1996. – P. 1–7.

24. Sallans B. Reinforcement Learning with Factored States and Actions [ ] / B. Sallans, G. E. Hinton // J. Mach. Learn. Res. – 2004. – P. 1–2.
25. Jen L. An application of SARSA temporal difference learning to Super Mario [ ] / Lucas Jen – 2004. – P. 1–2.
26. Hungry Geese – Agents Comparison [ ] – : www/ URL: <https://www.kaggle.com/ihelon/hungry-geese-agents-comparison>. – 11.11.2021 . – . .
27. Wirth R. CRISP-DM: Towards a Standard Process Model for Data Mining [ ] / R. Wirth, J. Hipp – 2005. – P. 1–5.
28. Cross-validation: evaluating estimator performance [ ] – : www/ URL: [https://scikit-learn.org/stable/modules/cross\\_validation.html](https://scikit-learn.org/stable/modules/cross_validation.html). – 11.11.2021 . – . .