```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from typing import NamedTuple
4 from google.colab import output
```

```
1 SEED = 0
2
3 BOARD_COL = 3
4 BOARD_ROW = 3
5 BOARD_SIZE = BOARD_COL * BOARD_ROW
6
7 """
8 Game board and actions are: {q, w, e, a, s, d, z, x, c}
9
10 q | w | e
11 --|---|--
12 a | s | d
13 --|---|--
14 z | x | c
15 """
16 ACTIONS_KEY_MAP = {'q': 0, 'w': 1, 'e': 2,
17                    'a': 3, 's': 4, 'd': 5,
18                    'z': 6, 'x': 7, 'c': 8}
```

```
1 np.random.seed(SEED)
```

## State Definition

Tic-Tac-Toe game state with methods for updating the board, checking the game's status, and determining valid actions

```
1 #Prints the current board state, Uses symbols 'x' for player 1, 'o' for player 2, ' ' for empty
2 def print_state(board, clear_output=False):
3   if clear_output:
4     output.clear()
5   for i in range(BOARD_ROW):
6     print('-------------')
7     out = '| '
8     for j in range(BOARD_COL):
9       if board[i, j] == 1:
10           token = 'x'
11       elif board[i, j] == -1:
12           token = 'o'
13       else:
14           token = ' '  # empty position
15       out += token + ' | '
16     print(out)
17   print('-------------')
```

```
1
2 # State Class : Represents the game state
3 # Stores:
4 # board → 3×3 NumPy array (1 = Player 1, -1 = Player 2, 0 = Empty)
5 # symbol → Current player's turn
6 # winner → Stores the winner (1, 2, or 0 for draw)
7 # end → Boolean flag indicating if the game is over
8 class State:
9   def __init__(self, symbol):
10     # the board is represented by an n * n array,
11     #  1 represents the player who moves first,
12     # -1 represents another player
13     #  0 represents an empty position
14     self.board = np.zeros((BOARD_ROW, BOARD_COL))
15     self.symbol = symbol
16     self.winner = 0
17     self.end = None
18
19   #to generates a unique hash for board states using a base-3 formula.
20   @property
21   def hash_value(self):
22     hash = 0
23     for x in np.nditer(self.board):
24       hash = 3*hash + x + 1   # unique hash
25     return hash
26
27   #to takes an action key, Converts it to board coordinates, Calls next_by_pos(i, j).
28   def next(self, action: str):
```

```
29    id = ACTIONS_KEY_MAP[action]
30    i, j = id // BOARD_COL, id % BOARD_COL
31    return self.next_by_pos(i, j)
32
33  #to generates a new game state after a move, asserts that the move is valid,swaps turns (self.symbol → -self.symbol).
34  def next_by_pos(self, i: int, j: int):
35    assert self.board[i, j] == 0
36    new_state = State(-self.symbol)      # another player turn
37    new_state.board = np.copy(self.board)
38    new_state.board[i, j] = self.symbol  # current player choose to play at (i, j) pos
39    return new_state
40
41  #to returns available moves based on empty board positions
42  @property
43  def possible_actions(self):
44    rev_action_map = {id: key for key, id in ACTIONS_KEY_MAP.items()}
45    actions = []
46    for i in range(BOARD_ROW):
47      for j in range(BOARD_COL):
48        if self.board[i, j] == 0:
49          actions.append(rev_action_map[BOARD_COL*i+j])
50    return actions
51
52  #to check rows and columns for a sum of 3 or -3 (win condition), diagonals for a winning condition, if moves are availa
53  def is_end(self):
54    if self.end is not None:
55      return self.end
56
57    check = []
58    # check row
59    for i in range(BOARD_ROW):
60      check.append(sum(self.board[i, :]))
61
62    # check col
63    for i in range(BOARD_COL):
64      check.append(sum(self.board[:, i]))
65
66    # check diagonal
67    diagonal = 0; reverse_diagonal = 0
68    for i in range(BOARD_ROW):
69      diagonal += self.board[i, i]
70      reverse_diagonal += self.board[BOARD_ROW-i-1, i]
71    check.append(diagonal)
72    check.append(reverse_diagonal)
73
74    for x in check:
75      if x == 3:
76        self.end = True
77        self.winner = 1   # player 1 wins
78        return self.end
79      elif x == -3:
80        self.end = True
81        self.winner = 2   # player 2 wins
82        return self.end
83
84    for x in np.nditer(self.board):
85      if x == 0:          # play available
86        self.end = False
87        return self.end
88
89    self.winner = 0       # draw
90    self.end = True
91    return self.end
92
93
```

## Environment

```
1 class Env:
2   #to initializes the environment, calls get_all_states() to precompute all possible game states, sets curr_state to the
3   def __init__(self):
4     self.all_states = self.get_all_states()
5     self.curr_state = State(symbol=1)
6
7   #to precomputes all valid board states and stores them in all_states, uses a recursive function (explore_all_substates(
8   # the dictionary all_states: Key is Unique "hash_value" of a board state and Value is Corresponding "State object".
9   def get_all_states(self):
10    all_states = {}  # is a dict with key as state_hash_value and value as State object.
11    def explore_all_substates(state):
```

```
12        for i in range(BOARD_ROW):
13          for j in range(BOARD_COL):
14            if state.board[i, j] == 0:
15              next_state = state.next_by_pos(i, j)
16              if next_state.hash_value not in all_states:
17                all_states[next_state.hash_value] = next_state
18                if not next_state.is_end():
19                  explore_all_substates(next_state)
20      curr_state = State(symbol=1)
21      all_states[curr_state.hash_value] = curr_state
22      explore_all_substates(curr_state)
23      return all_states
24
25    # to resets the game to the initial state and returns the new starting State
26    def reset(self):
27      self.curr_state = State(symbol=1)
28      return self.curr_state
29
30    #to moves to the next state based on a given action, Checks if the move is valid (must be in possible_actions), Retriev
31    #updates curr_state and returns the new state after the move and a reward (currently 0, can be modified for RL).
32    # def step(self, action):
33    #   assert action in self.curr_state.possible_actions, f"Invalid {action} for the current state \n{self.curr_state.prir
34    #   next_state_hash = self.curr_state.next(action).hash_value
35    #   next_state = self.all_states[next_state_hash]
36    #   self.curr_state = next_state
37    #   reward = 0
38    #   return self.curr_state, reward
39    def step(self, action):
40      assert action in self.curr_state.possible_actions, f"Invalid move {action}"
41      next_state_hash = self.curr_state.next(action).hash_value
42      next_state = self.all_states[next_state_hash]
43      self.curr_state = next_state
44
45      if self.curr_state.is_end():
46          if self.curr_state.winner == 1:
47              reward = 1   # Player 1 wins
48          elif self.curr_state.winner == 2:
49              reward = -1  # Player 2 wins
50          else:
51              reward = 0   # Draw
52      else:
53          reward = 0  # No immediate reward for non-terminal moves
54
55      return self.curr_state, reward
56
57    #to checks if the game has ended
58    def is_end(self):
59      return self.curr_state.is_end()
60
61    #to returns 'player1', 'player2', or 'draw' based on curr_state.winner.
62    @property
63    def winner(self):
64      result_id = self.curr_state.winner
65      result = 'draw'
66      if result_id == 1:
67        result = 'player1'
68      elif result_id == 2:
69        result = 'player2'
70      return result
71
```

Right now, reward = 0 (neutral), Update the rewards +1 if player wins. -1 if player loses.

## ⌄ Policy

```
1 class BasePolicy:
2   def reset(self):
3     pass
4
5   def update_values(self, *args):
6     pass
7
8   def select_action(self, state):
9     raise Exception('Not Implemented Error')
```

```
1 # human player to interact with the Tic-Tac-Toe environment by manually choosing actions
2 class HumanPolicy(BasePolicy):
3   def __init__(self, symbol):
```

```python
   4     self.symbol = symbol
   5
   6   def select_action(self, state):
   7     assert state.symbol == self.symbol, f"Its not {self.symbol} symbol's turn"
   8     print_state(state.board, clear_output=True)
   9     key = input("Input your position: ")
  10     return key
```

```python
   1 #randomly selects a valid action for the given player
   2 class RandomPolicy(BasePolicy):
   3   def __init__(self, symbol):
   4     self.symbol = symbol
   5
   6   def select_action(self, state):
   7     assert state.symbol == self.symbol, f"Its not {self.symbol} symbol's turn"
   8     return np.random.choice(state.possible_actions)
```

```python
   1 class ActionPlayed(NamedTuple):
   2   hash_value: str
   3   action: str
   4   # reward: float  #reward received after action
   5
```

```python
   1 class MenacePolicy(BasePolicy):
   2   def __init__(self, all_states, symbol, tau=5.0):
   3     self.all_states = all_states
   4     self.symbol = symbol
   5     self.tau = tau #tau is fixed, a decaying tau (annealing) can improve learning,this makes MENACE explore early, exploi
   6
   7     # It store the number of stones for each action for each state
   8     self.state_action_value = self.initialize()
   9     # variable to store the history for updating the number of stones
  10     self.history = []
  11
  12   #to initialize state-action values i.e. number of stones/beads
  13   def initialize(self):
  14     state_action_value = {}
  15     for hash_value, state in self.all_states.items():
  16       # initially all actions have 0 stones
  17       state_action_value[hash_value] = {action: 0 for action in state.possible_actions}
  18     return state_action_value
  19
  20   #resets all state-action values to zero
  21   def reset(self):
  22     for action_value in self.state_action_value.values():
  23       for action in action_value.keys():
  24         action_value[action] = 0
  25
  26   def print_updates(self, reward):
  27     print(f'Player with symbol {self.symbol} updates the following history with {reward} stone')
  28     for item in self.history:
  29       board = np.copy(self.all_states[item.hash_value].board)
  30       id = ACTIONS_KEY_MAP[item.action]
  31       i, j = id//BOARD_COL, id%BOARD_COL
  32       board[i, j] = self.symbol
  33       print_state(board)
  34
  35
  36   #to  adjusts stone counts based on game results
  37   def update_values(self, reward, show_update=False):
  38     # reward: if wins receive reward of 1 stone for the chosen action
  39     #         else -1 stone.
  40     # reward is either 1 or -1 depending upon if the player has won or lost the game.
  41
  42     if show_update:
  43       self.print_updates(reward)
  44     for item in self.history:
  45
  46         # ensure minimum value is 1 (MENACE cannot have negative beads)
  47         self.state_action_value[item.hash_value][item.action] = max(1, self.state_action_value[item.hash_value][item.acti
  48
  49     # clear history after updating values
  50     self.history = []
  51
  52   def select_action(self, state):
  53     assert state.symbol == self.symbol, f"Its not {self.symbol} symbol's turn"
  54     action_value = self.state_action_value[state.hash_value]
  55     max_value = action_value[max(action_value, key=action_value.get)]
  56     exp_values = {action: np.exp((v-max_value) / self.tau) for action, v in action_value.items()}
  57     normalizer = np.sum([v for v in exp_values.values()])
```

```
58        prob = {action: v/normalizer for action, v in exp_values.items()}
59        action = np.random.choice(list(prob.keys()), p=list(prob.values()))
60        self.history.append(ActionPlayed(state.hash_value, action))
61        return action
```

A MENACE (Matchbox Educable Noughts and Crosses Engine)-style learning agent for Tic-Tac-Toe:

It uses reinforcement learning principles by updating the number of "beads" (probability weights) in each matchbox (state-action pair) based on the game outcome. This policy learns by adjusting state-action values using a softmax-based selection.

## ⌄ Game Board

```
1 class Game:
2   def __init__(self, env, player1, player2):
3     self.env = env
4     self.player1 = player1
5     self.player2 = player2
6     self.show_updates = False
7
8   #to yields players in an alternating order (Player-1 to Player-2 to Repeat)
9   def alternate(self):
10    while True:
11      yield self.player1
12      yield self.player2
13
14  def train(self, epochs=1_00_000):
15    game_results = {'player1': 0, 'player2': 0, 'draw': 0}
16    player1_reward_map = {'player1': 1, 'player2': -1, 'draw': 0}
17
18    for _ in range(epochs):
19      result = self.play()
20      game_results[result] += 1
21
22      # if player1 wins add 1 stone for the action chosen
23      player1_reward = player1_reward_map[result]
24      player2_reward = -player1_reward   # if player2 wins add 1 stone
25
26      self.player1.update_values(player1_reward)
27      self.player2.update_values(player2_reward)
28
29    print(f"Training complete! Results over {epochs} games:")
30    print(game_results)
31
32  #to runs a single game until the end state is reached
33  def play(self):
34    alternate = self.alternate()
35    state = self.env.reset()
36    while not self.env.is_end():
37      player = next(alternate)
38      action = player.select_action(state)
39      state, _ = self.env.step(action)
40    result = self.env.winner
41    return result
```

## ⌄ Experiment

```
1 env = Env()
2
3 player1 = MenacePolicy(env.all_states, symbol=1)
4 player2 = MenacePolicy(env.all_states, symbol=-1)
5 # player2 = RandomPolicy(symbol=-1)
6
7 #training MENACE (player1) through self-play, it learns by playing against another MENACE agent
8 game = Game(env, player1, player2)
9 game.train(epochs=1_00_000)
10
```

```
⊋ Training complete! Results over 100000 games:
   {'player1': 95766, 'player2': 2113, 'draw': 2121}
```

There is an issue with the MENACE training process due to an overflow in the exponential calculation. The values used in the softmax function are growing too large, causing invalid probability distributions (NaN values).

I will adjust the update mechanism to keep the values stable.

```
1  #testing MENACE against a human after training, we set up a game where player1 (trained MENACE) faces a human (HumanPolic
2
3  game_with_human_player = Game(env, player1, HumanPolicy(symbol=-1))
4  game_with_human_player.play()
5
6  result = env.winner
7  print(f"winner: {result}")
8
9  player1_reward_map = {'player1': 1, 'player2': -1, 'draw': 0}
10 player1.update_values(player1_reward_map[result], show_update=True)
11
12 # """
13 # Game board and actions are: {q, w, e, a, s, d, z, x, c}
14
15 # q | w | e
16 # --|---|--
17 # a | s | d
18 # --|---|--
19 # z | x | c
20 # """
21
22
```

```
 ⮧   -------------
     | o |   | x |
     -------------
     |   | o |   |
     -------------
     | x |   | x |
     -------------
     Input your position: d
     winner: player1
     Player with symbol 1 updates the following history with 1 stone
     -------------
     |   |   | x |
     -------------
     |   |   |   |
     -------------
     |   |   |   |
     -------------
     -------------
     |   |   | x |
     -------------
     |   | o |   |
     -------------
     | x |   |   |
     -------------
     -------------
     | o |   | x |
     -------------
     |   | o |   |
     -------------
     | x |   | x |
     -------------
     -------------
     | o |   | x |
     -------------
     |   | o | o |
     -------------
     | x | x | x |
     -------------
```

```
1  #testing MENACE against a human after training, we set up a game where player1 (trained MENACE) faces a human (HumanPolic
2  game_with_human_player = Game(env, player1, HumanPolicy(symbol=-1))
3  game_with_human_player.play()
4  result = env.winner
5  print(f"winner: {result}")
6  player1_reward_map = {'player1': 1, 'player2': -1, 'draw': 0}
7  player1.update_values(player1_reward_map[result], show_update=True)
8
9  # """
10 # Game board and actions are: {q, w, e, a, s, d, z, x, c}
11
12 # q | w | e
13 # --|---|--
14 # a | s | d
15 # --|---|--
16 # z | x | c
17 # """
18
19
```

```
 ⮧   -------------
     | o | x | x |
```

```
           -------------
          |   | o |   |
           -------------
          | x | o | x |
           -------------
Input your position: d
winner: draw
Player with symbol 1 updates the following history with 0 stone
           -------------
          |   |   | x |
           -------------
          |   |   |   |
           -------------
          |   |   |   |
           -------------
           -------------
          |   |   | x |
           -------------
          |   | o |   |
           -------------
          | x |   |   |
           -------------
           -------------
          |   | x | x |
           -------------
          |   | o |   |
           -------------
          | x | o |   |
           -------------
           -------------
          | o | x | x |
           -------------
          |   | o |   |
           -------------
          | x | o | x |
           -------------
           -------------
          | o | x | x |
           -------------
          | x | o | o |
           -------------
          | x | o | x |
           -------------
```

```python
1  #testing MENACE against a human after training, we set up a game where player1 (trained MENACE) faces a human (HumanPolic
2  game_with_human_player = Game(env, player1, HumanPolicy(symbol=-1))
3  game_with_human_player.play()
4  result = env.winner
5  print(f"winner: {result}")
6  player1_reward_map = {'player1': 1, 'player2': -1, 'draw': 0}
7  player1.update_values(player1_reward_map[result], show_update=True)
8
9  # """
10 # Game board and actions are: {q, w, e, a, s, d, z, x, c}
11
12 # q | w | e
13 # --|---|--
14 # a | s | d
15 # --|---|--
16 # z | x | c
17 # """
18
19
```

```
           -------------
          | o |   | x |
           -------------
          |   | o |   |
           -------------
          | x |   | x |
           -------------
Input your position: d
winner: player1
Player with symbol 1 updates the following history with 1 stone
           -------------
          |   |   | x |
           -------------
          |   |   |   |
           -------------
          |   |   |   |
           -------------
           -------------
          |   |   | x |
           -------------
          |   | o |   |
           -------------
          | x |   |   |
```

```
 ------------
 ------------
 | o |   | x |
 ------------
 |   | o |   |
 ------------
 | x |   | x |
 ------------
 ------------
 | o |   | x |
 ------------
 |   | o | o |
 ------------
 | x | x | x |
 ------------
```

## modify the MENACE model to use k = 3

```python
1 class MenacePolicy(BasePolicy):
2   def __init__(self, all_states, symbol, tau=5.0):
3     self.all_states = all_states
4     self.symbol = symbol
5     self.tau = tau #tau is fixed, a decaying tau (annealing) can improve learning,this makes MENACE explore early, exploi
6
7     # It store the number of stones for each action for each state
8     self.state_action_value = self.initialize()
9     # variable to store the history for updating the number of stones
10    self.history = []
11
12   #to initialize state-action values i.e. number of stones/beads
13   def initialize(self):
14       state_action_value = {}
15       for hash_value, state in self.all_states.items():
16           # initially, each action starts with k = 3 stones (instead of 0)
17           state_action_value[hash_value] = {action: 3 for action in state.possible_actions}
18       return state_action_value
19
20
21   #resets all state-action values to zero
22   def reset(self):
23     for action_value in self.state_action_value.values():
24       for action in action_value.keys():
25         action_value[action] = 0
26
27   def print_updates(self, reward):
28     print(f'Player with symbol {self.symbol} updates the following history with {reward} stone')
29     for item in self.history:
30       board = np.copy(self.all_states[item.hash_value].board)
31       id = ACTIONS_KEY_MAP[item.action]
32       i, j = id//BOARD_COL, id%BOARD_COL
33       board[i, j] = self.symbol
34       print_state(board)
35
36
37   #to  adjusts stone counts based on game results
38   def update_values(self, reward, show_update=False):
39     # reward: if wins receive reward of 1 stone for the chosen action
40     #         else -1 stone.
41     # reward is either 1 or -1 depending upon if the player has won or lost the game.
42
43     if show_update:
44       self.print_updates(reward)
45     for item in self.history:
46
47         # ensure minimum value is 1 (MENACE cannot have negative beads)
48         self.state_action_value[item.hash_value][item.action] = max(1, self.state_action_value[item.hash_value][item.acti
49
50     # clear history after updating values
51     self.history = []
52
53   def select_action(self, state):
54     assert state.symbol == self.symbol, f"Its not {self.symbol} symbol's turn"
55     action_value = self.state_action_value[state.hash_value]
56     max_value = action_value[max(action_value, key=action_value.get)]
57     exp_values = {action: np.exp((v-max_value) / self.tau) for action, v in action_value.items()}
58     normalizer = np.sum([v for v in exp_values.values()])
59     prob = {action: v/normalizer for action, v in exp_values.items()}
60     action = np.random.choice(list(prob.keys()), p=list(prob.values()))
61     self.history.append(ActionPlayed(state.hash_value, action))
62     return action
```

```
1 env = Env()
2
3 player1 = MenacePolicy(env.all_states, symbol=1)
4 player2 = MenacePolicy(env.all_states, symbol=-1)
5
6
7 #training MENACE (player1) through self-play, it learns by playing against another MENACE agent
8 game = Game(env, player1, player2)
9 game.train(epochs=1_00_000)
```

⮑ Training complete! Results over 100000 games:
    {'player1': 95962, 'player2': 1358, 'draw': 2680}

```
1 #testing MENACE against a human after training, we set up a game where player1 (trained MENACE) faces a human (HumanPolic
2 game_with_human_player = Game(env, player1, HumanPolicy(symbol=-1))
3 game_with_human_player.play()
4
5 result = env.winner
6 print(f"winner: {result}")
7
8 player1_reward_map = {'player1': 1, 'player2': -1, 'draw': 0}
9 player1.update_values(player1_reward_map[result], show_update=True)
10
11 # """
12 # Game board and actions are: {q, w, e, a, s, d, z, x, c}
13
14 # q | w | e
15 # --|---|--
16 # a | s | d
17 # --|---|--
18 # z | x | c
19 # """
```

⮑      -------------
       |   |   | x |
       -------------
       | x | o | o |
       -------------
       |   |   | x |
       -------------
       Input your position: w
       -------------
       |   | o | x |
       -------------
       | x | o | o |
       -------------
       |   | x | x |
       -------------
       Input your position: z
       winner: draw
       Player with symbol 1 updates the following history with 0 stone
       -------------
       |   |   |   |
       -------------
       |   |   |   |
       -------------
       |   |   | x |
       -------------
       -------------
       |   |   | x |
       -------------
       |   | o |   |
       -------------
       |   |   | x |
       -------------
       -------------
       |   |   | x |
       -------------
       | x | o | o |
       -------------
       |   |   | x |
       -------------
       -------------
       |   | o | x |
       -------------
       | x | o | o |
       -------------
       |   | x | x |
       -------------
       -------------
       | x | o | x |
       -------------
       | x | o | o |
       -------------
       | o | x | x |
       -------------
```

Currently, MENACE always follows tau=5.0, if we change tau **it** can encourage early exploration, later exploitation

```
1 def train(self, epochs=100_000):
2     for epoch in range(epochs):
3         self.player1.tau = max(0.1, 5.0 * np.exp(-0.0001 * epoch))  # reduce tau
4         self.player2.tau = self.player1.tau
5         result = self.play()
6
```

```
1 env = Env()
2
3 player1 = MenacePolicy(env.all_states, symbol=1)
4 player2 = MenacePolicy(env.all_states, symbol=-1)
5
6
7 #training MENACE (player1) through self-play, it learns by playing against another MENACE agent
8 game = Game(env, player1, player2)
9 game.train(epochs=1_00_000)
```

⊋ Training complete! Results over 100000 games:
    {'player1': 94874, 'player2': 1568, 'draw': 3558}

```
1 game_with_human_player = Game(env, player1, HumanPolicy(symbol=-1))
2
3 game_with_human_player.play()
4
5 result = env.winner
6 print(f"winner: {result}")
7
8 player1_reward_map = {'player1': 1, 'player2': -1, 'draw': 0}
9 player1.update_values(player1_reward_map[result], show_update=True)
10
11 # """
12 # Game board and actions are: {q, w, e, a, s, d, z, x, c}
13
14 # q | w | e
15 # --|---|--
16 # a | s | d
17 # --|---|--
18 # z | x | c
19 # """
20
21
```

⊋
```
-------------
|   | x |   |
-------------
| x | o | o |
-------------
| x | o | x |
-------------
Input your position: q
winner: draw
Player with symbol 1 updates the following history with 0 stone
-------------
|   |   |   |
-------------
|   |   |   |
-------------
| x |   |   |
-------------

-------------
|   |   |   |
-------------
|   | o |   |
-------------
| x |   | x |
-------------

-------------
|   | x |   |
-------------
|   | o |   |
-------------
| x | o | x |
-------------

-------------
|   | x |   |
-------------
| x | o | o |
-------------
| x | o | x |
-------------

-------------
```

```
| o | x | x |
-------------
| x | o | o |
-------------
| x | o | x |
-------------
```

Initial Number of Stones, At the beginning, every possible action for each state starts with 0 stones. Updating Stones After a Game

When the policy wins,Each action taken during the game is rewarded with +1 stone.

When the policy loses,Each action taken during the game is penalized with -1 stone, but the minimum number of stones is 1 (i.e., MENACE cannot have negative beads). When the game ends in a draw,No changes are made to the stone count.

Could not connect to the reCAPTCHA service. Please check your internet connection and reload to get a reCAPTCHA challenge.

```
| o | x | x |
-------------
| x | o | o |
-------------
| x | o | x |
-------------
```