

Lexicographic Reinforcement Learning Benchmark on MuJoCo environment

-Appendix-

Symbolic and Evolutionary Artificial Intelligence a.a. 2022/2023

Michele Lisi

September 2023

Appendix Guide

In this appendix, all the data regarding testing and training are showed. All the plots are obtained using `tf_reader` and can be better visualized (and interacted with) simply by running the command `python3 main.py ./trainings`, assuming the trainings folder is in the same as `tf_reader`.

To gain a better understanding, here follows an explanation of the meaning of the model tags used:

- **[20S]** - a network with sample size of 20
- **[HHH,BBB]** - a network with layer size of HHH and Experience Replay Buffer's batch size of BBB
- **[NOB]** - the network does not use bias in its layers
- **[LT-X]** - the loss threshold used in the lexicographic network is 0.X
- **[SL-X]** - the slack used in the lexicographic network is 0.X

Tags not listed are not relevant to the purpose of the report.

The reward tags, instead, are defined as follows:

- **[HLT] Healthy Reward:** a fixed reward (1.0) given every timestamp that the ant is considered healthy (which is defined below).
- **[FWD] Forward Reward:** reward calculated as the difference between the x position of the Ant before and after the action, divided by `dt`, which is the time between frames. The reward is positive if the movement is in the positive direction of the x axis.
- **[CST] Control Cost:** a negative reward to penalize actions too large. It is calculated as the sum of the squared actions taken, multiplied by a weight factor
- **[ORG] Distance from Origin:** represents the distance from the spawn point and is calculated similarly to the Forward Reward, but using the xy position instead of the x position

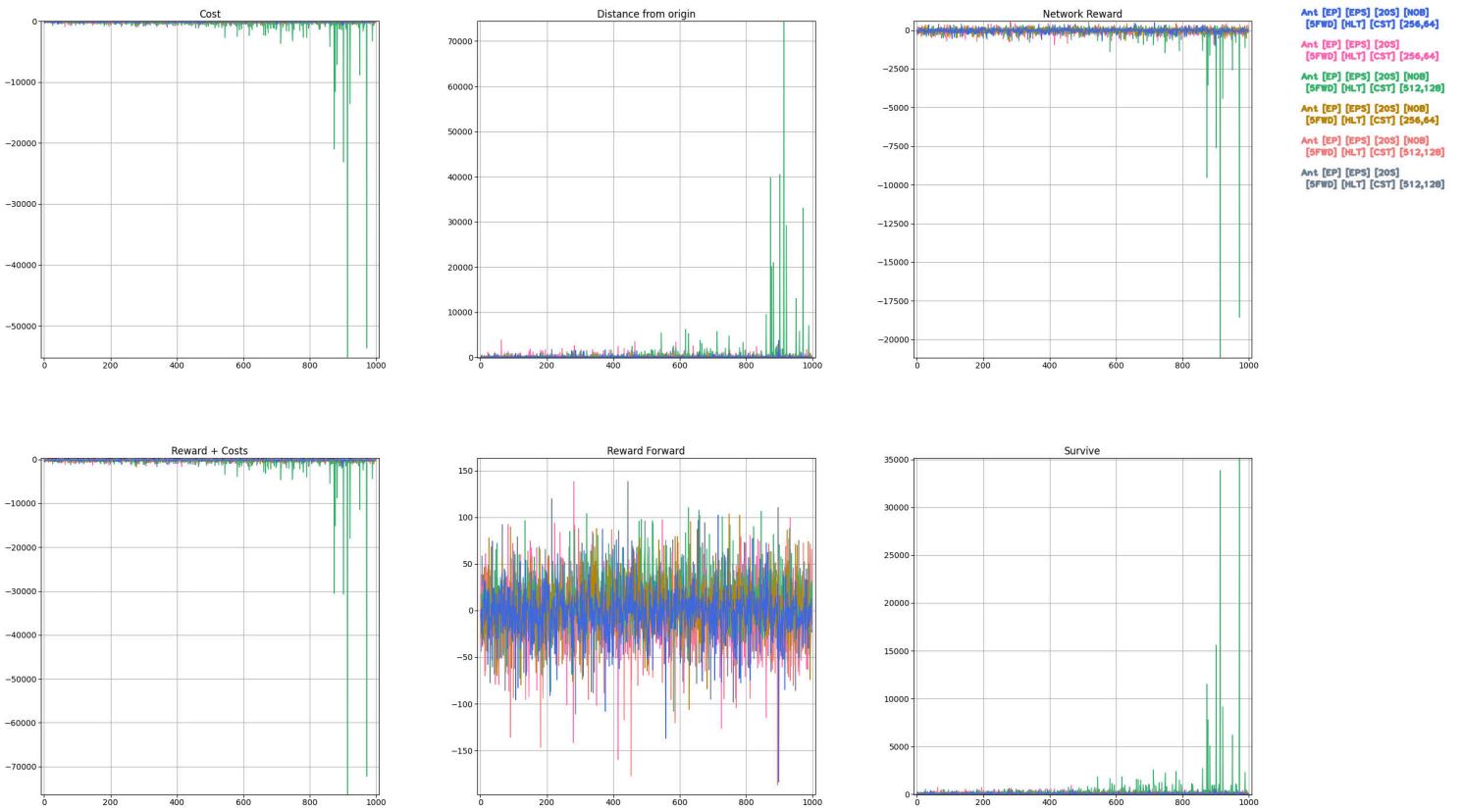
Finally, as for the name of the plots, the prefix dictates which metric they belong to:

- **No prefix** - cumulative reward obtained during training in every episode

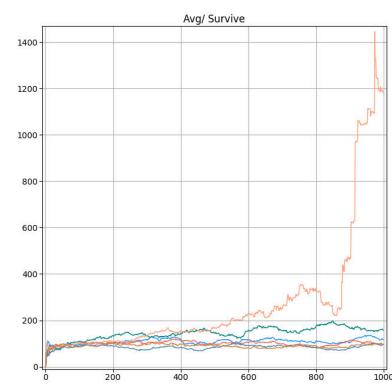
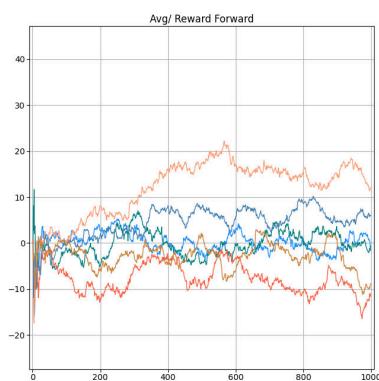
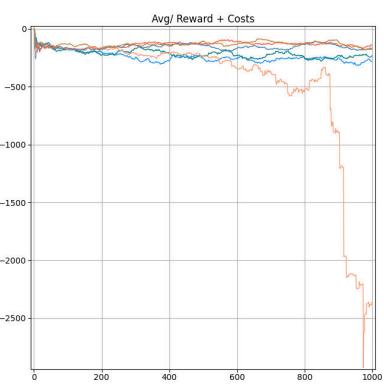
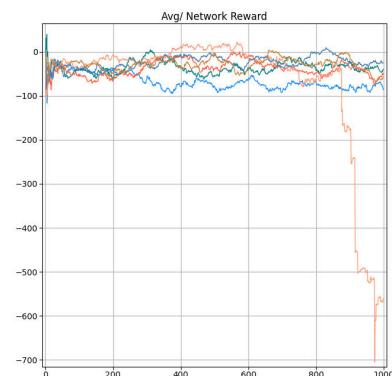
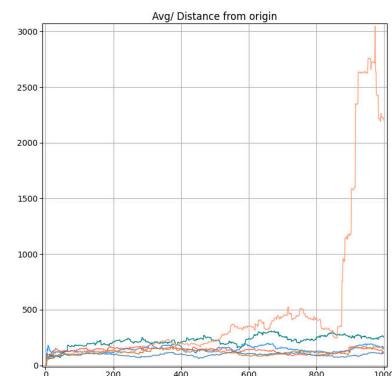
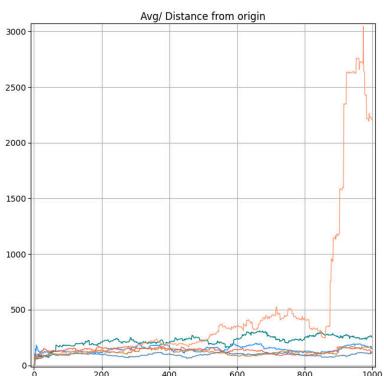
- **Avg/** - average cumulative reward obtained during training over 100 episodes
- **Test/** - cumulative reward obtained during test in every episode
- **Test Avg/** - average cumulative reward obtained during test over 100 episodes

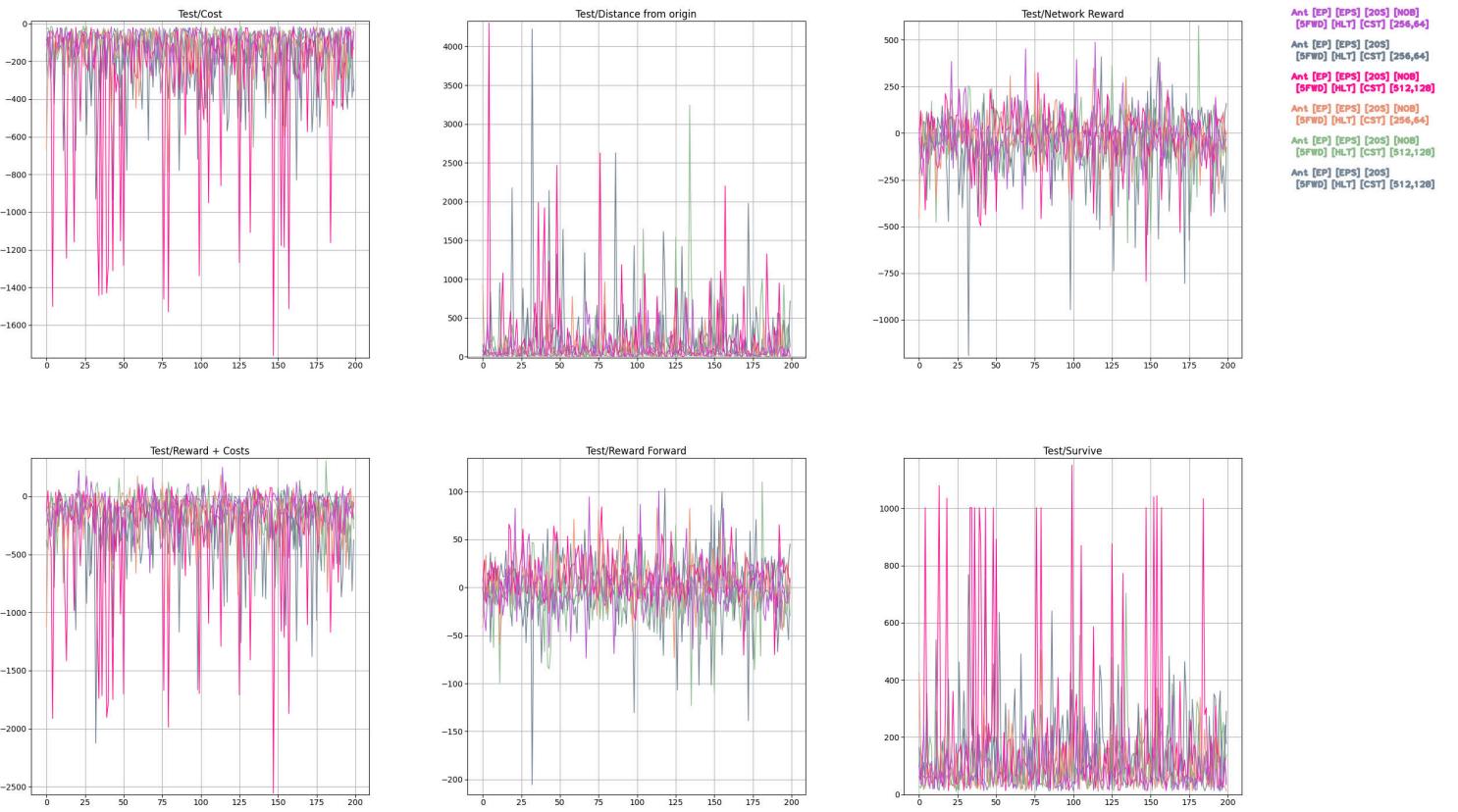
Appendix

Continuous DQN - [5FWD] [HLT] [CST]



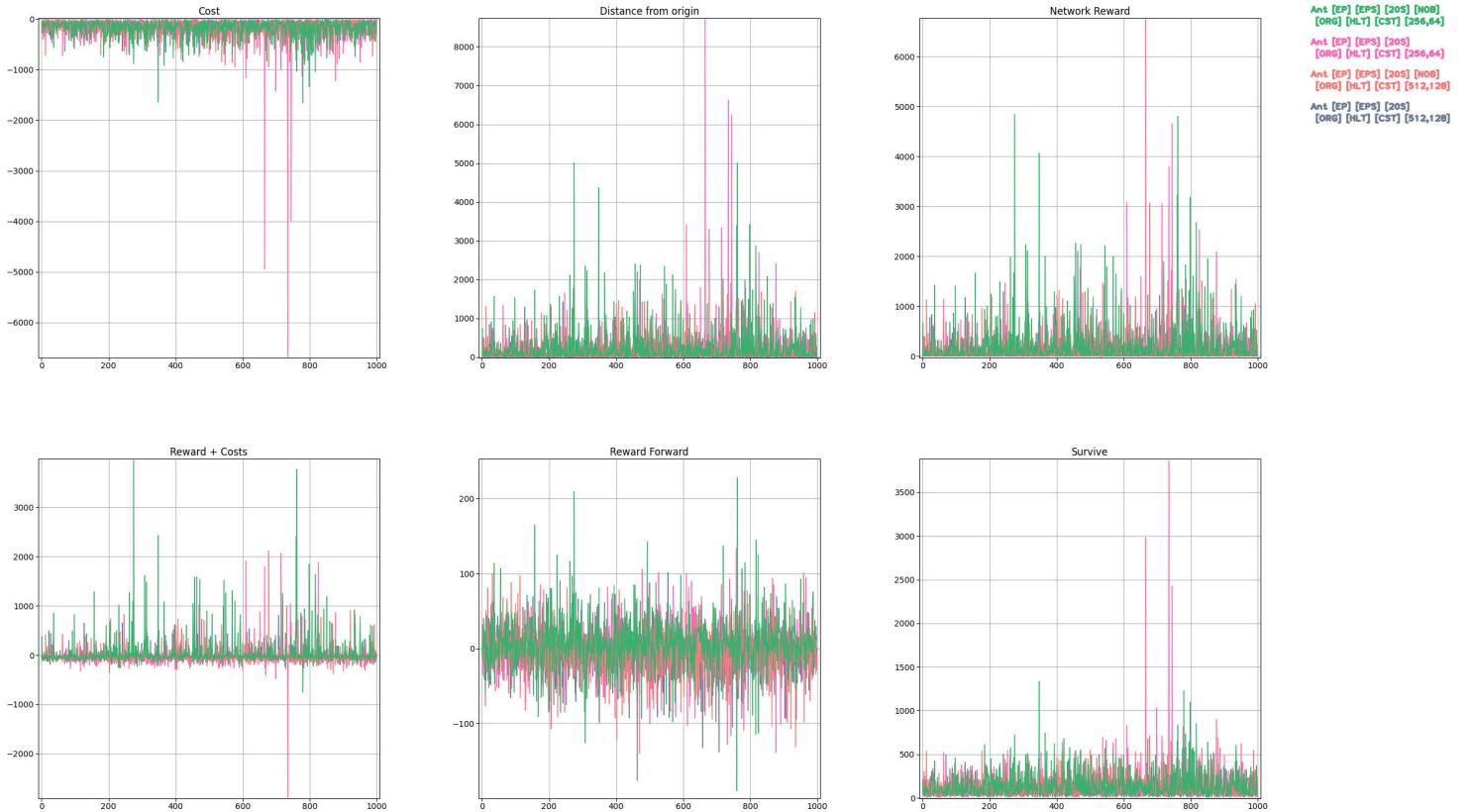
Ant [EP] [EPS] [205] [NOB]
 [SFWD] [MLT] [CST] [256,64]
 Ant [EP] [EPS] [205]
 [SFWD] [MLT] [CST] [256,64]
 Ant [EP] [EPS] [205] [NOB]
 [SFWD] [MLT] [CST] [512,128]
 Ant [EP] [EPS] [205] [NOB]
 [SFWD] [MLT] [CST] [256,64]
 Ant [EP] [EPS] [205] [NOB]
 [SFWD] [MLT] [CST] [512,128]

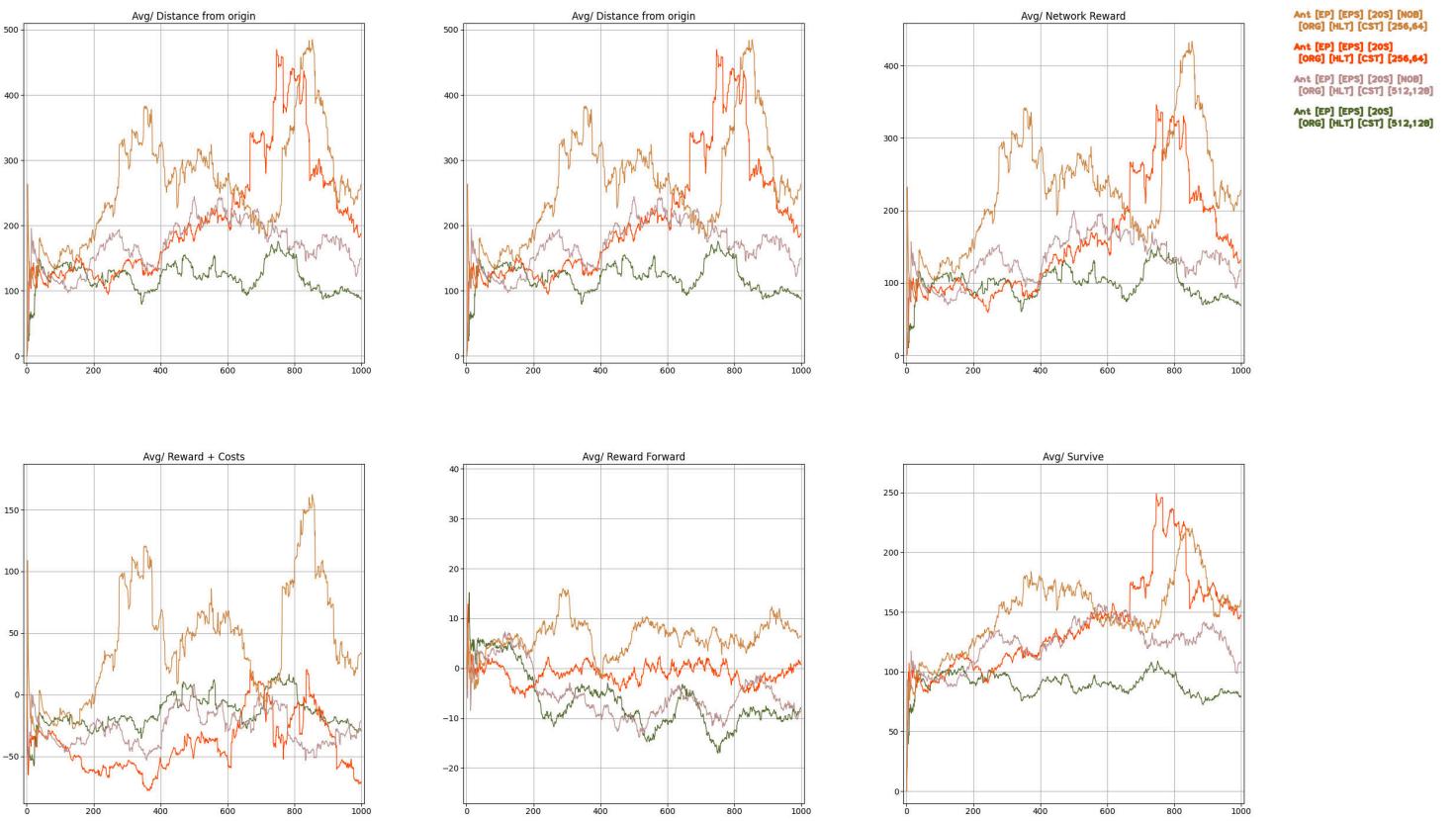


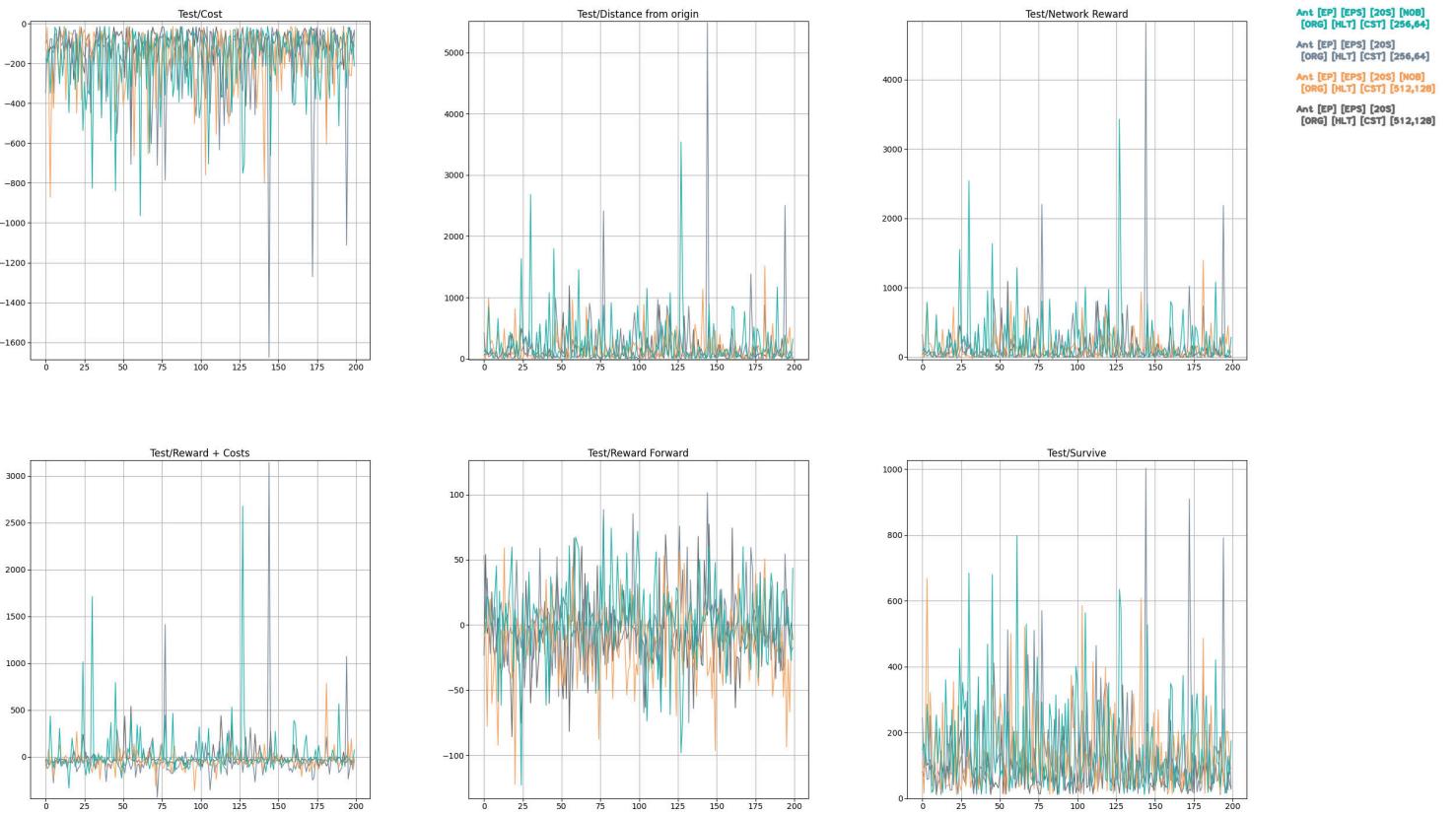




Continuous DQN - [ORG] [HLT] [CST]

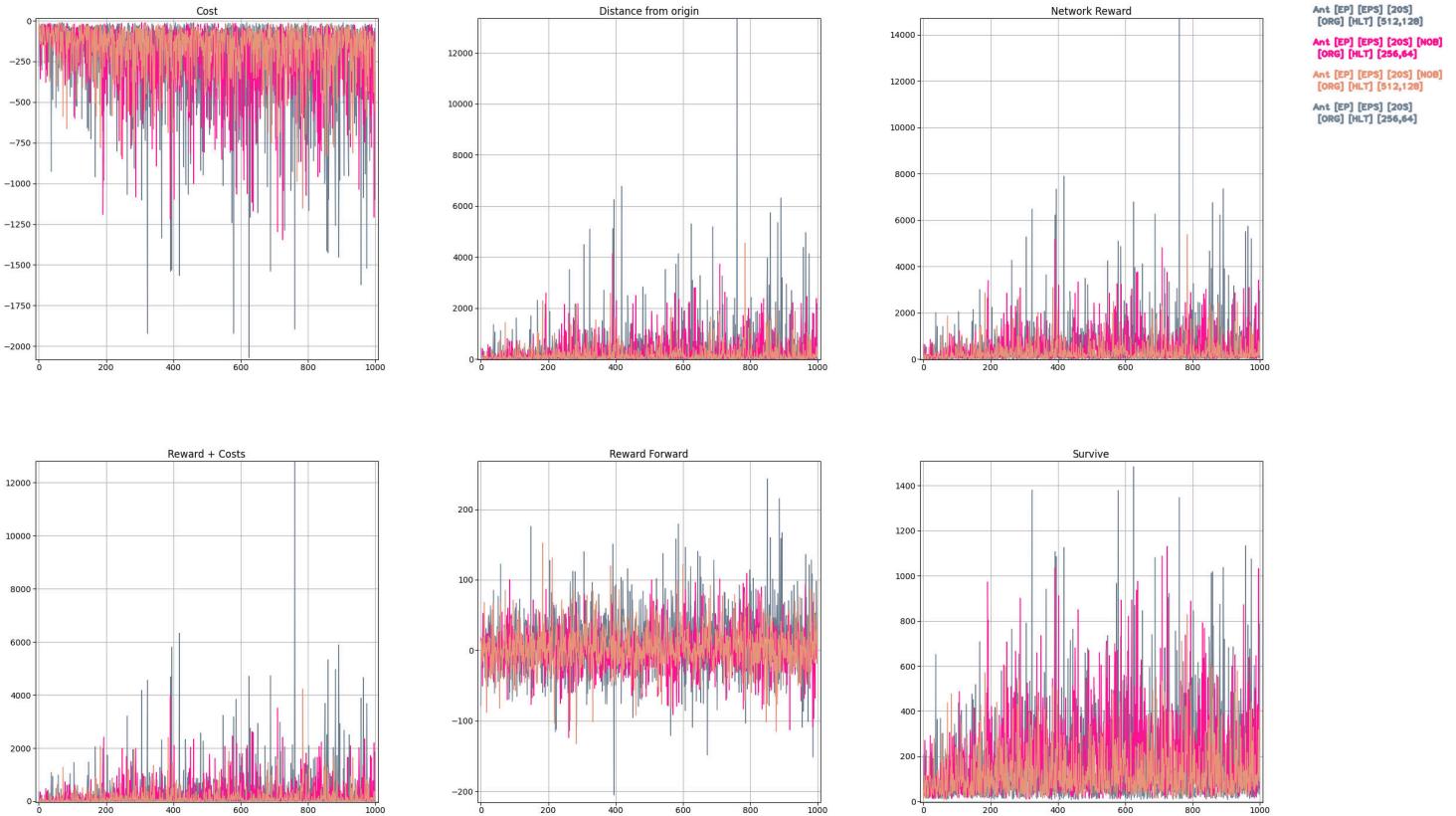




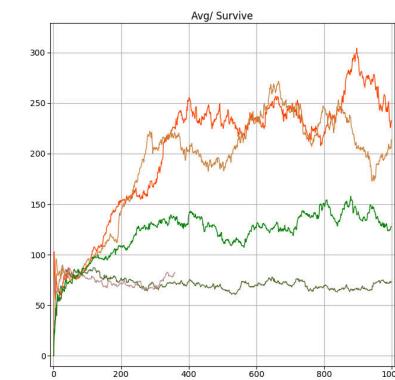
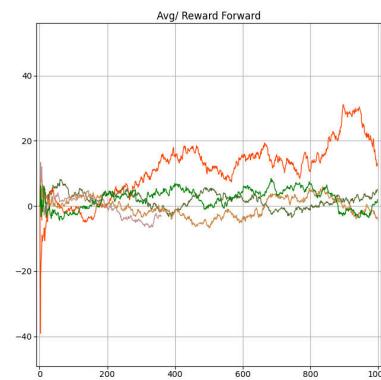
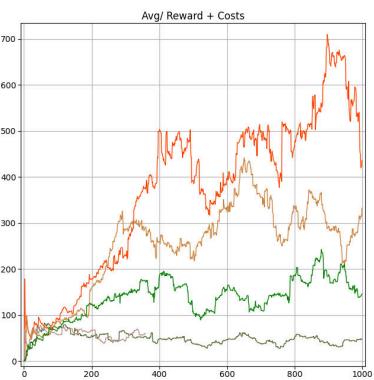
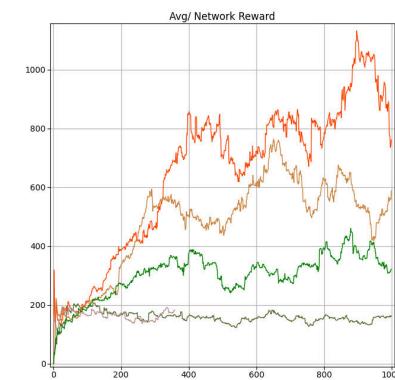
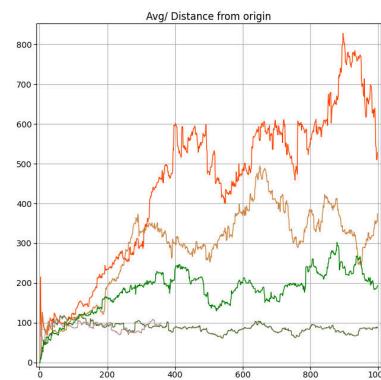
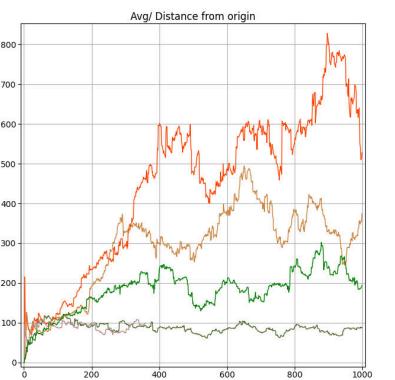


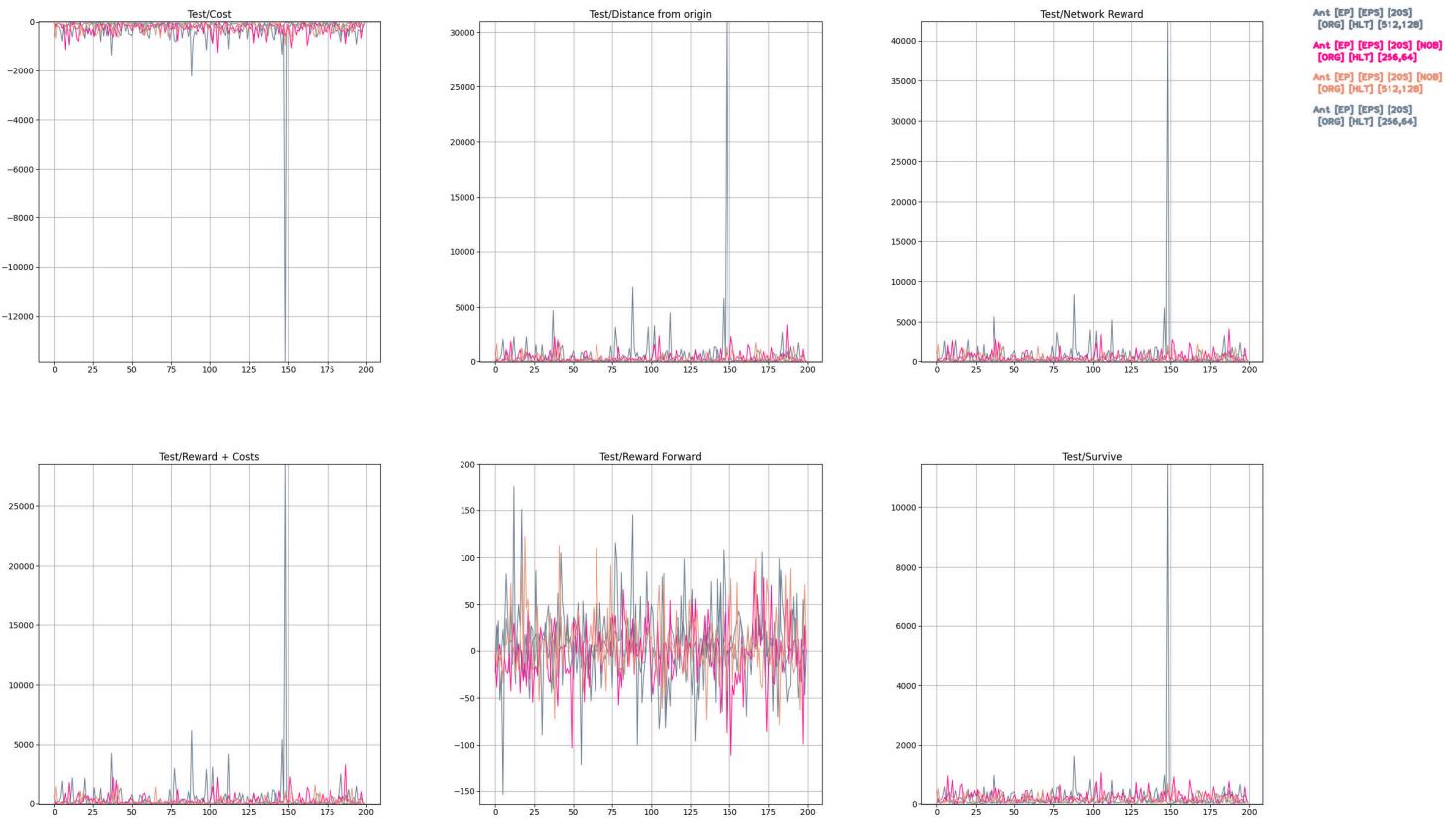


Continuous DQN - [ORG] [HLT]

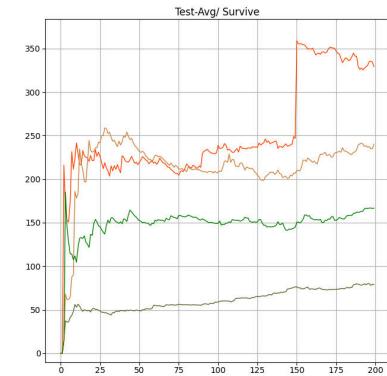
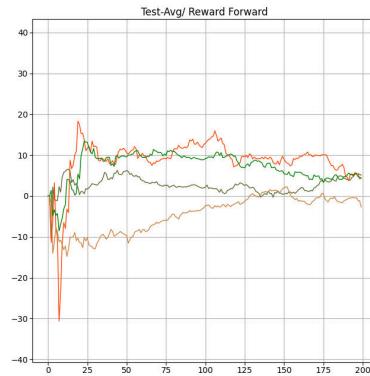
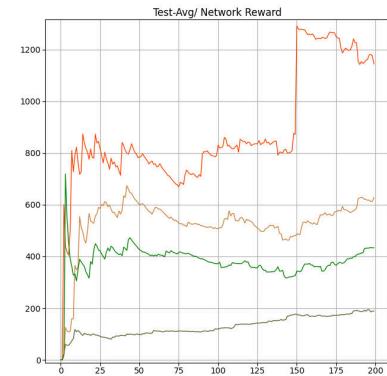
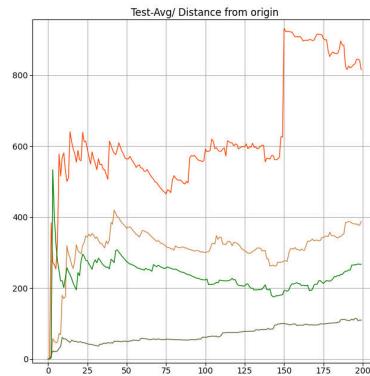
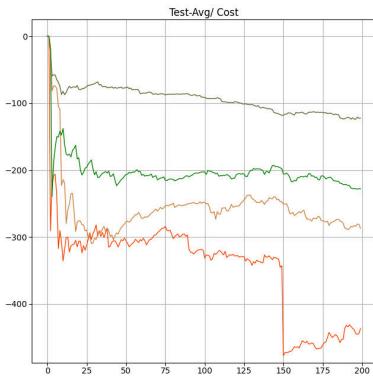


Ant [EP] [EPS] [208]
 [ORG] [HLT] [512,128]
 Ant [EP] [EPS] [208] [NOB]
 [ORG] [HLT] [256,64]
 Ant [EP] [EPS] [208] [NOB]
 [ORG] [HLT] [512,128]
 Ant [EP] [EPS] [208]
 [ORG] [HLT] [256,64]
 Ant [EP] [EPS] [208] [NOB]
 [ORG] [HLT] [256,64]

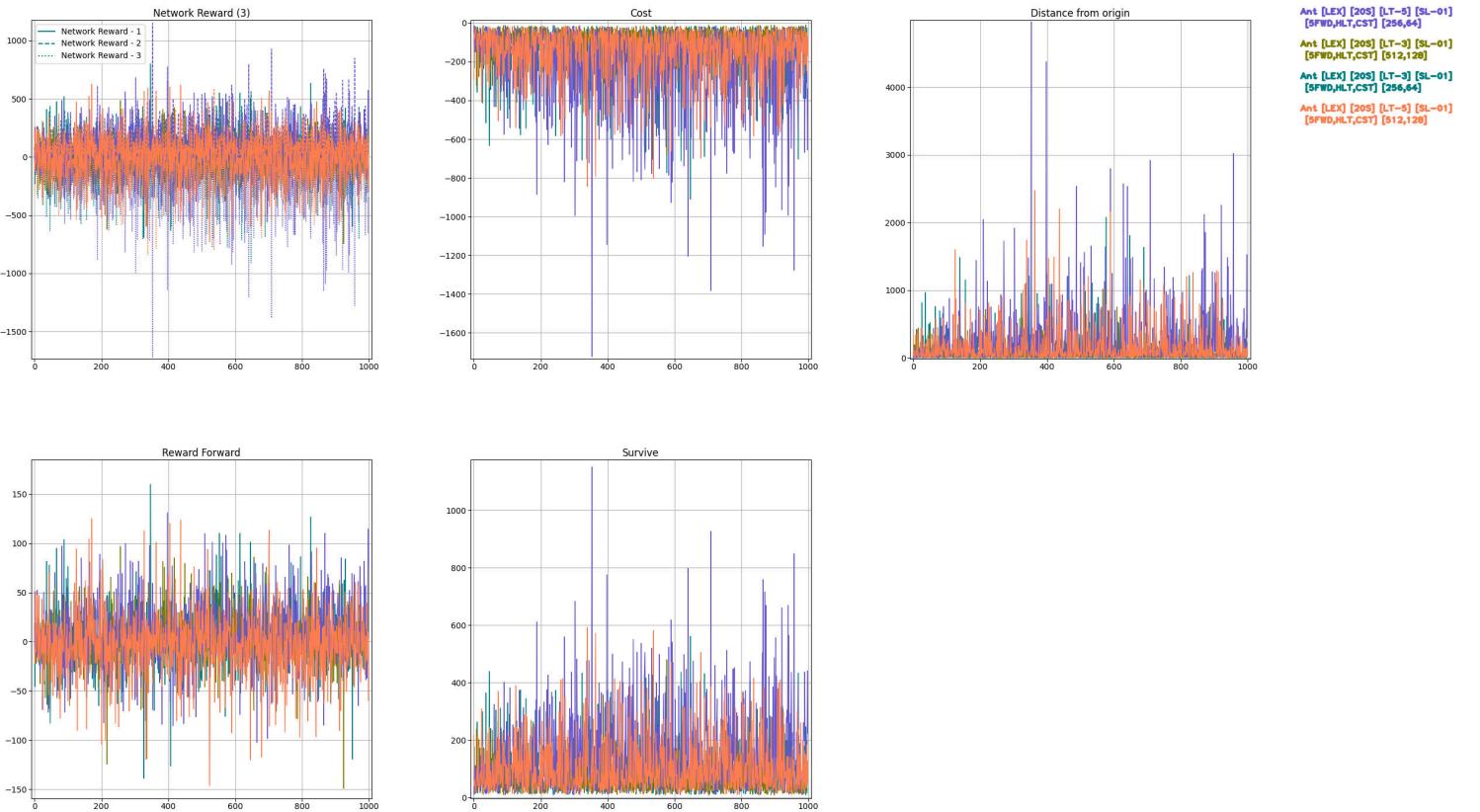


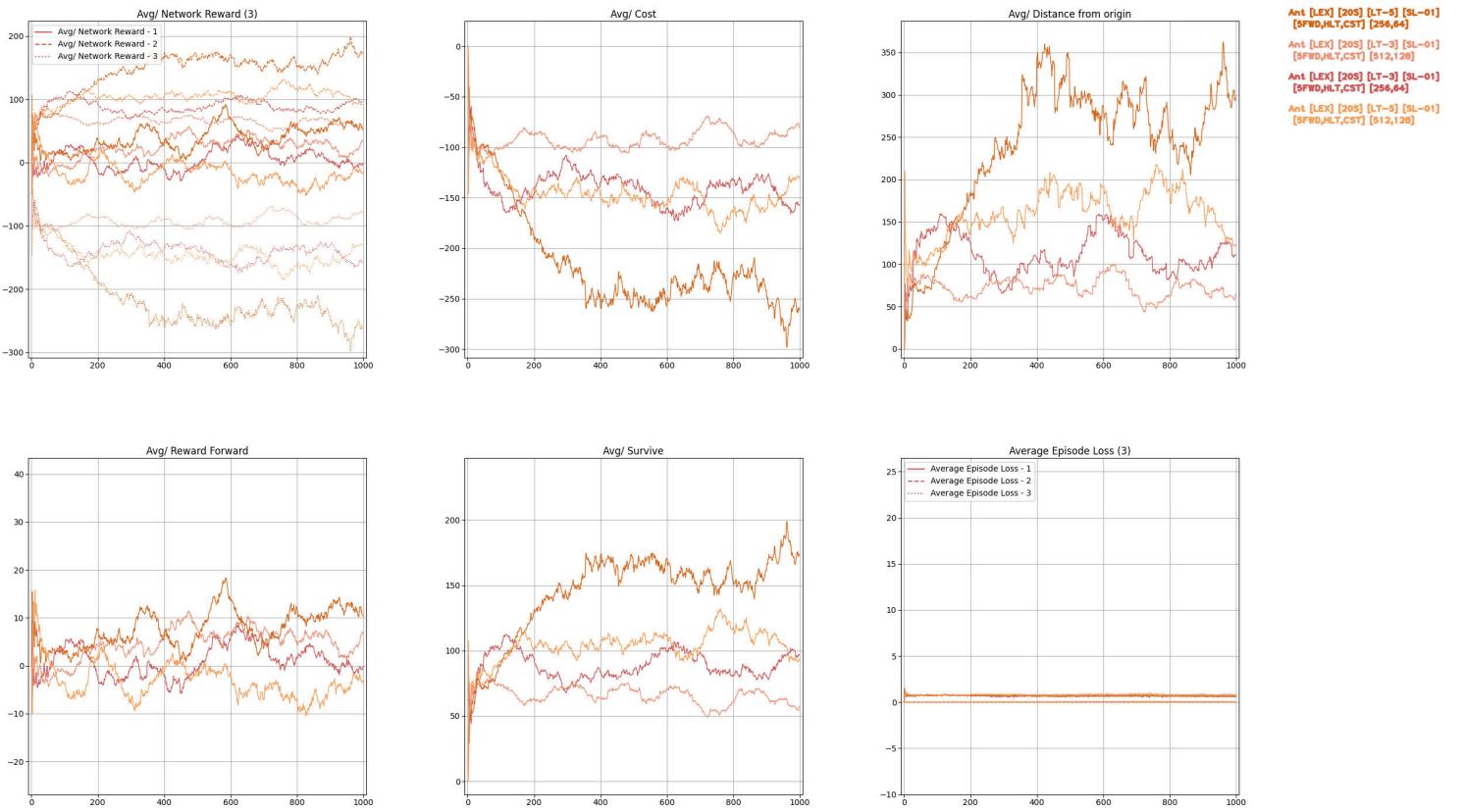


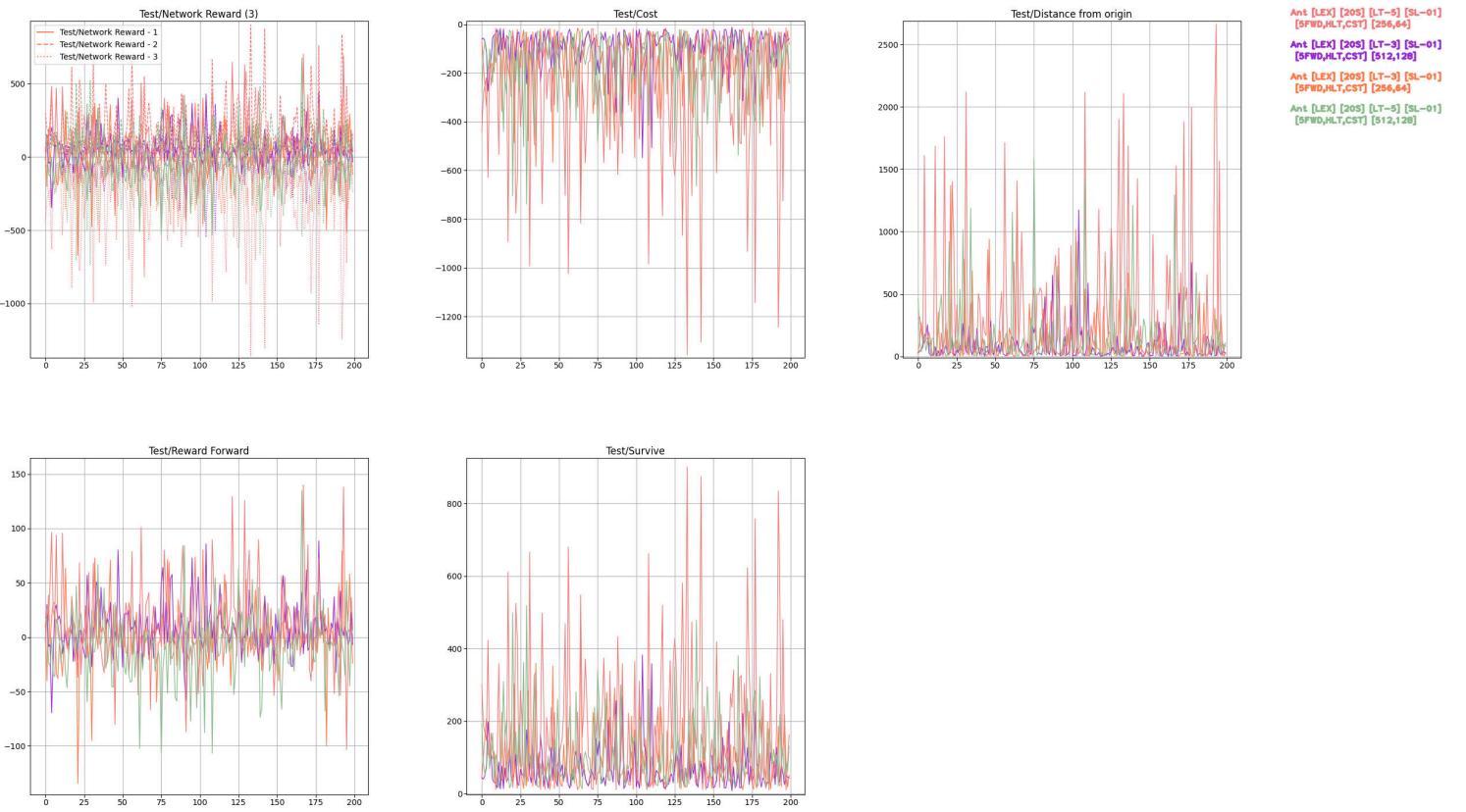
Ant [EP] [EPS] [208]
 [ORG] [HLT] [512, 128]
 Ant [EP] [EPS] [208] [NOB]
 [ORG] [HLT] [256, 64]
 Ant [EP] [EPS] [208] [NOB]
 [ORG] [HLT] [512, 128]
 Ant [EP] [EPS] [208] [NOB]
 [ORG] [HLT] [256, 64]
 Ant [EP] [EPS] [208] [NOB]
 [ORG] [HLT] [256, 64]

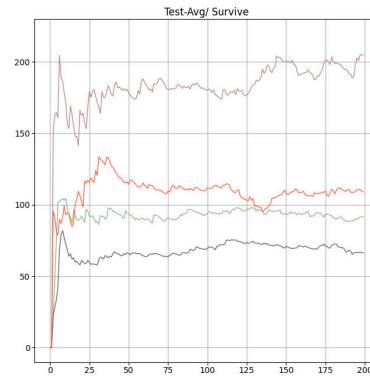
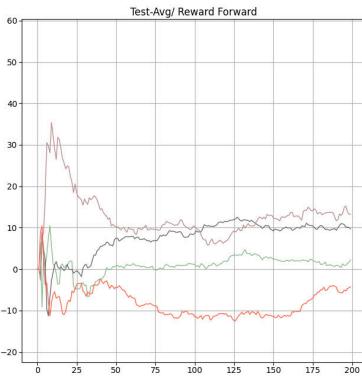
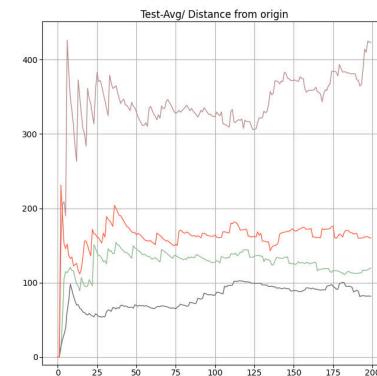
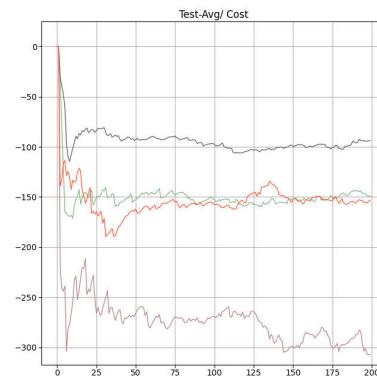
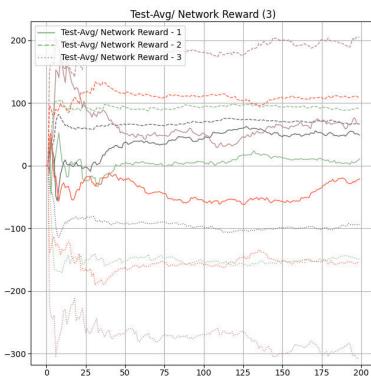


Lexicographic CDQN - [5FWD,HLT,CST]



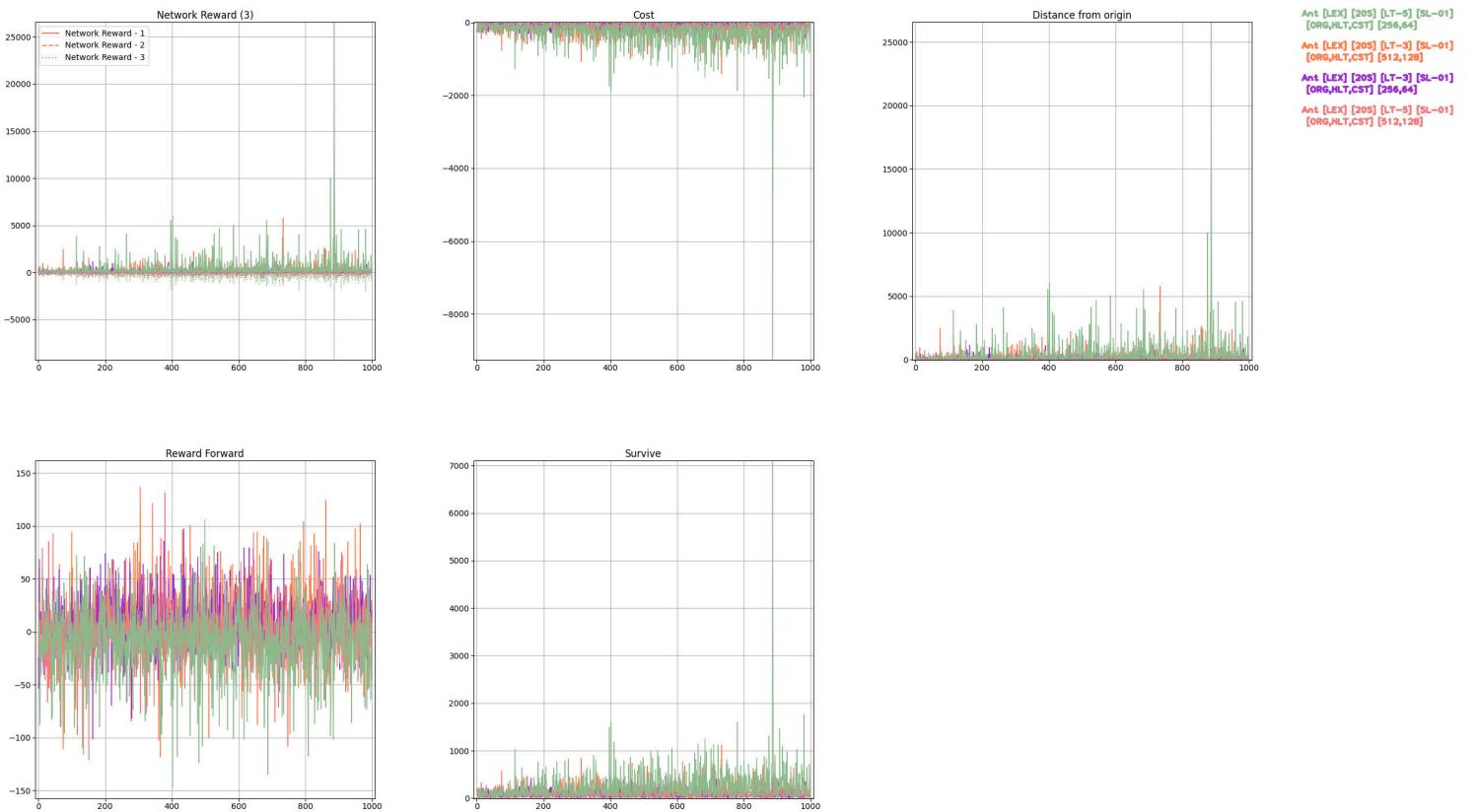


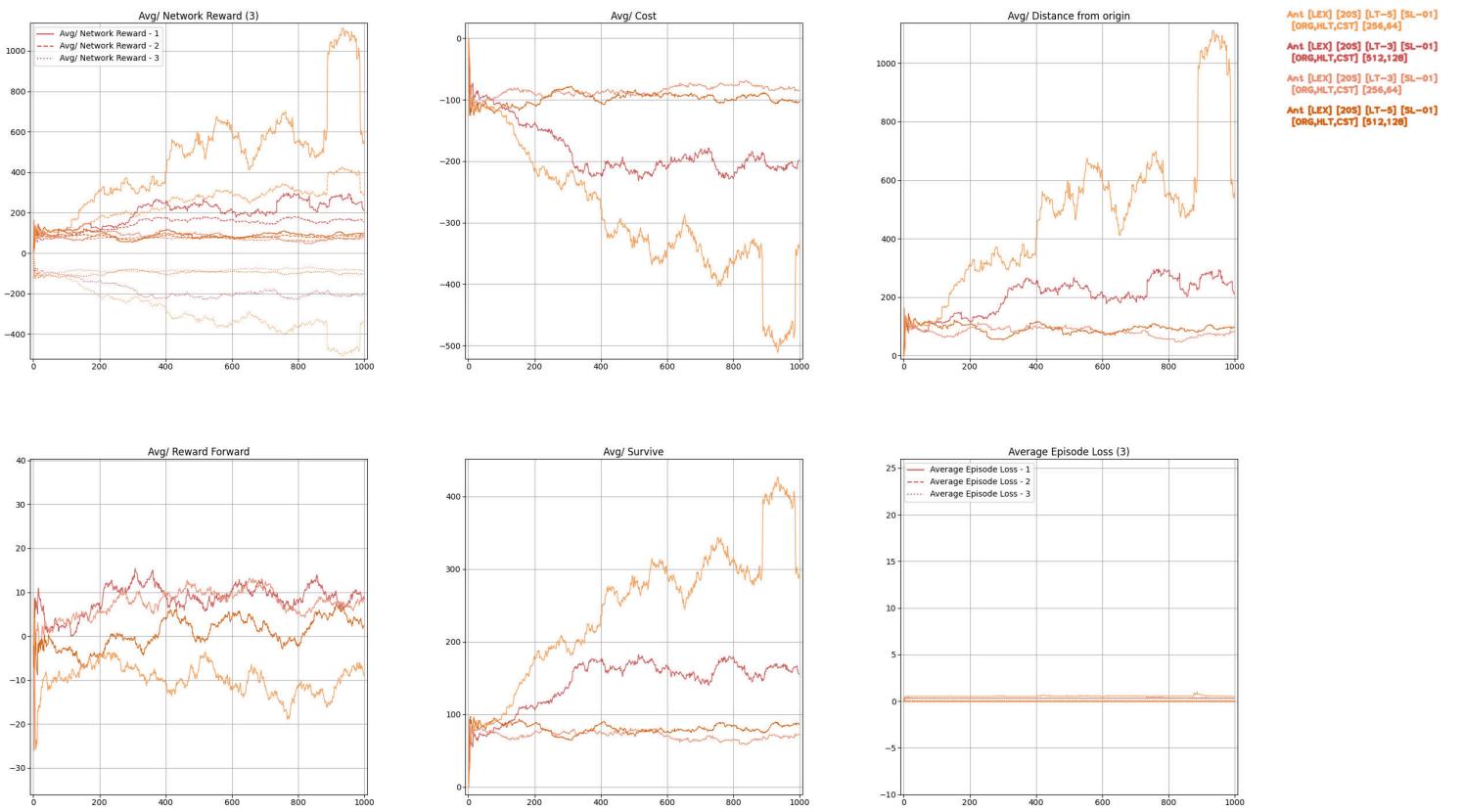


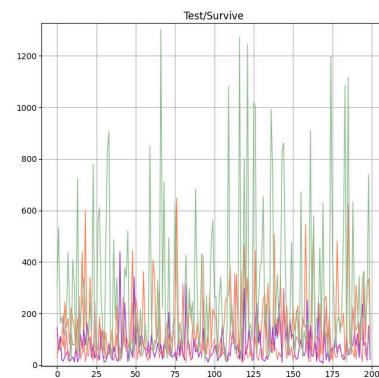
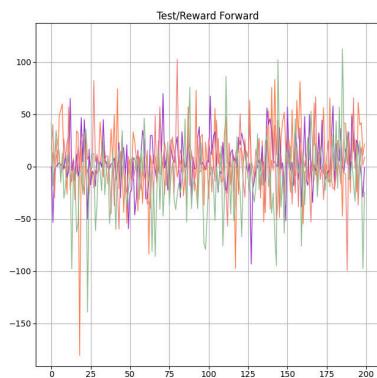
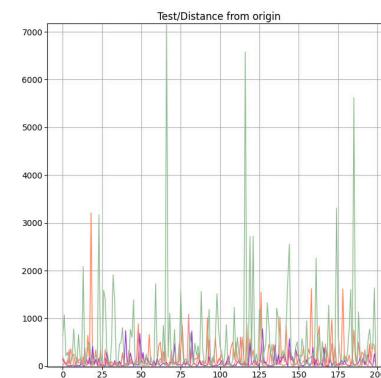
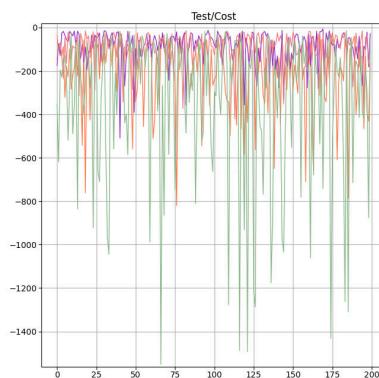
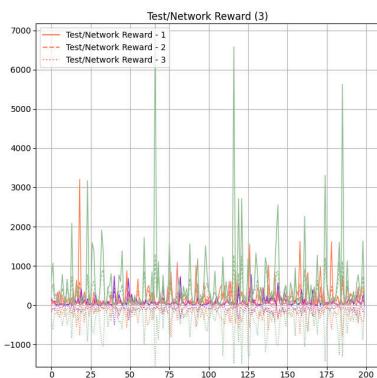


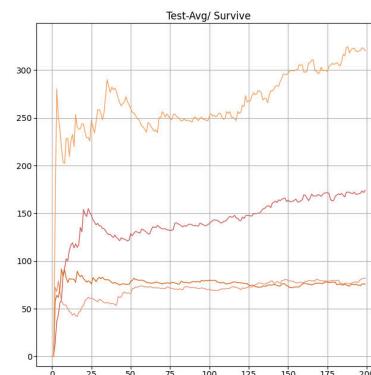
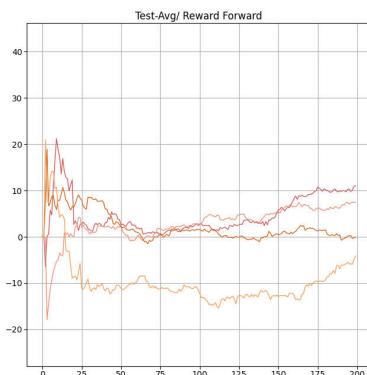
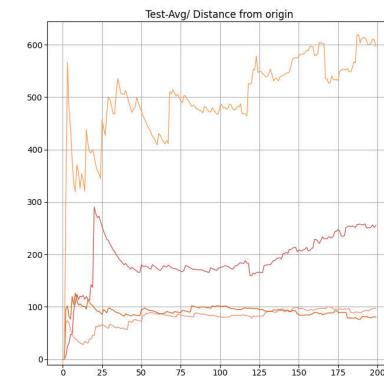
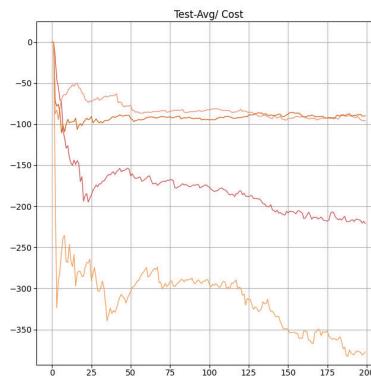
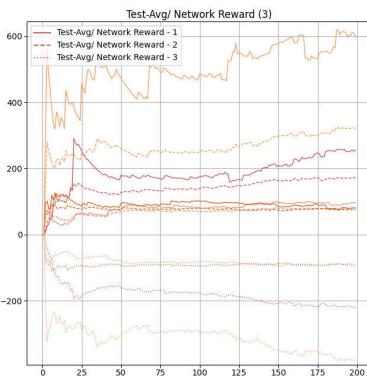
Ant [LEX] [20S] [LT-5] [SL-01]
[SFWD,HLT,CST] [256,64]
Ant [LEX] [20S] [LT-3] [SL-01]
[SFWD,HLT,CST] [512,128]
Ant [LEX] [20S] [LT-3] [SL-01]
[SFWD,HLT,CST] [256,64]
Ant [LEX] [20S] [LT-5] [SL-01]
[SFWD,HLT,CST] [512,128]

Lexicographic CDQN - [ORG,HLT,CST]

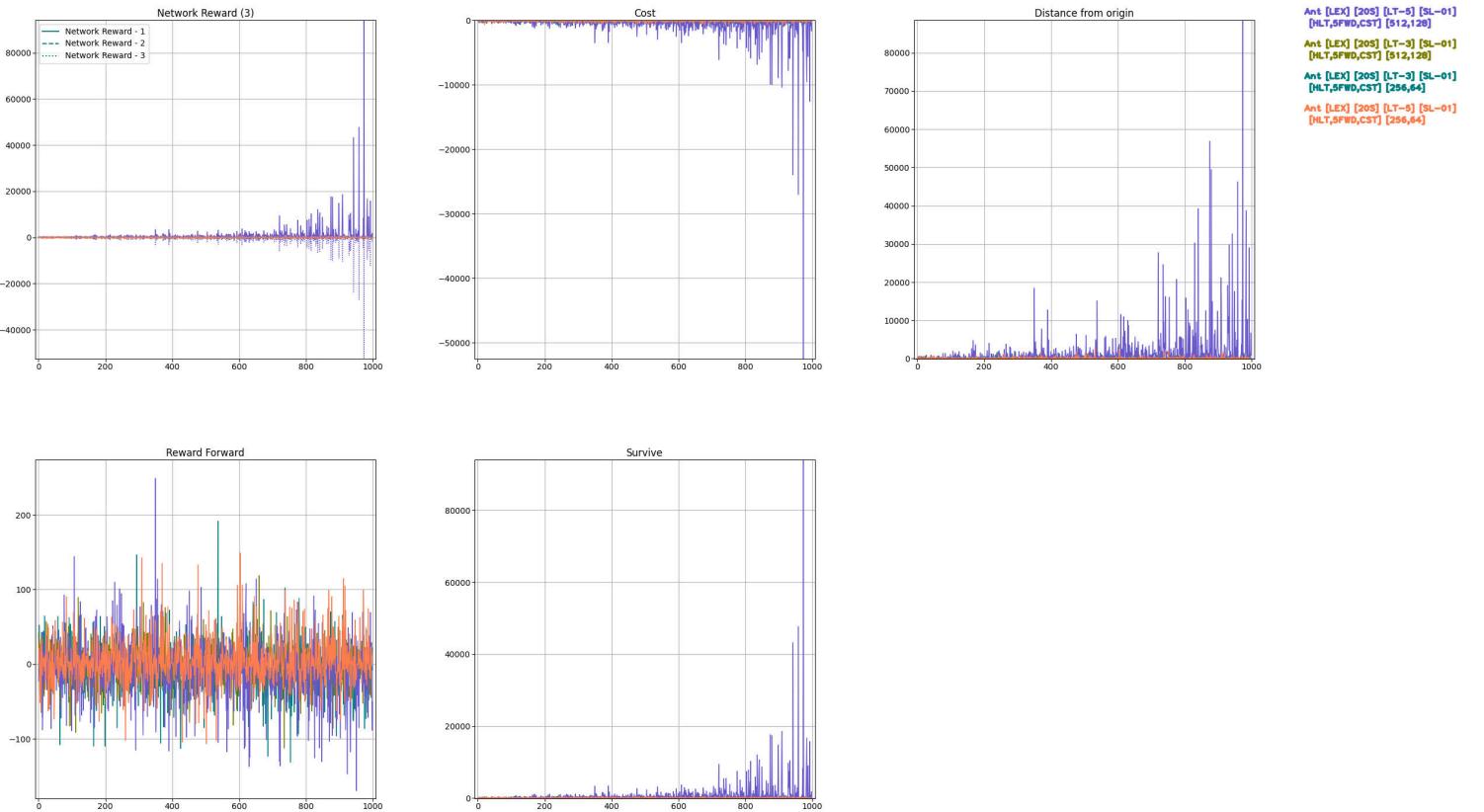


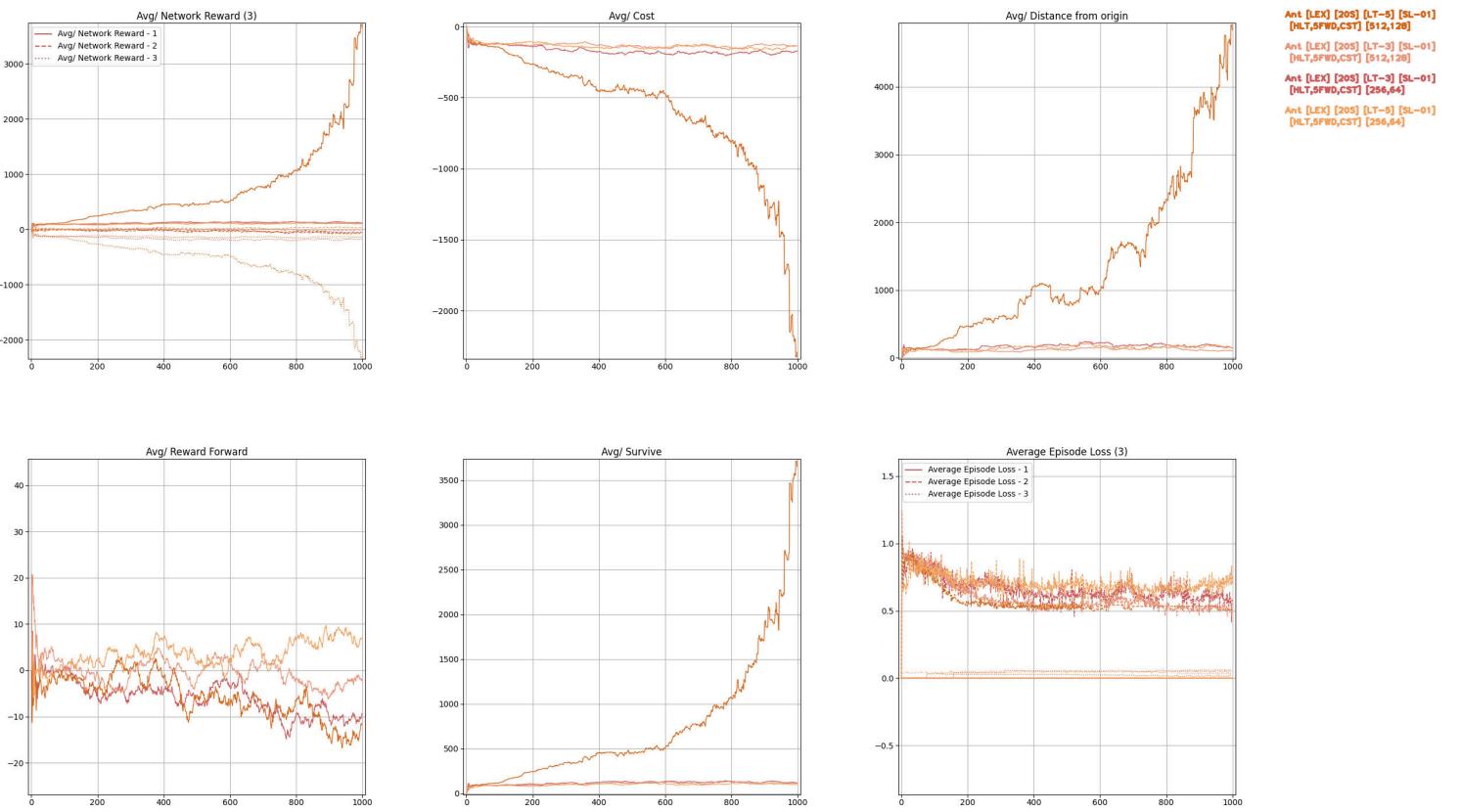


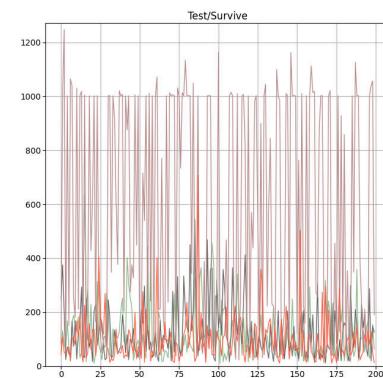
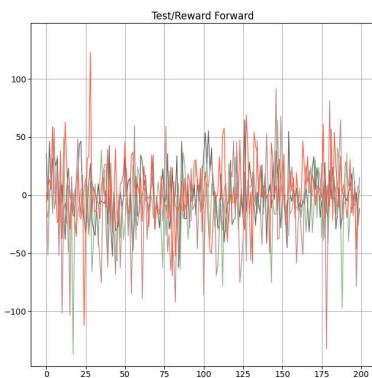
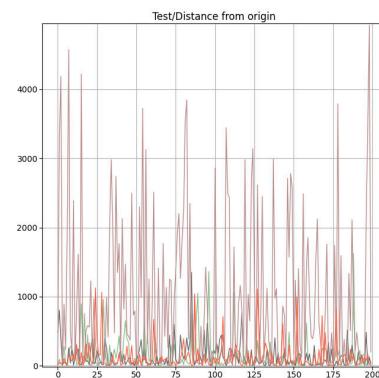
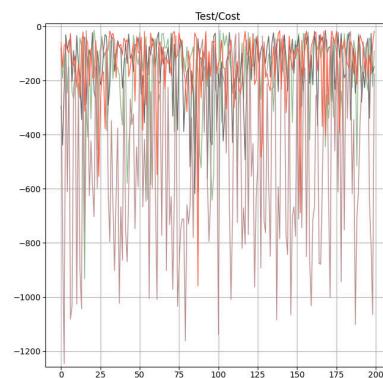
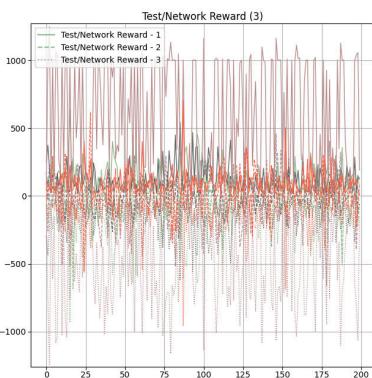


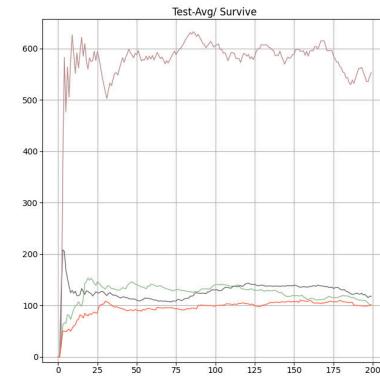
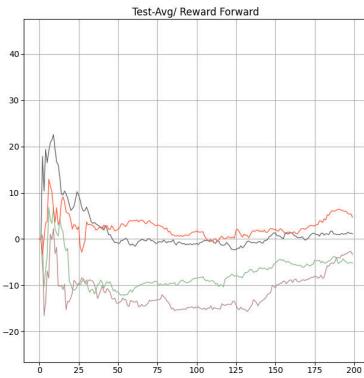
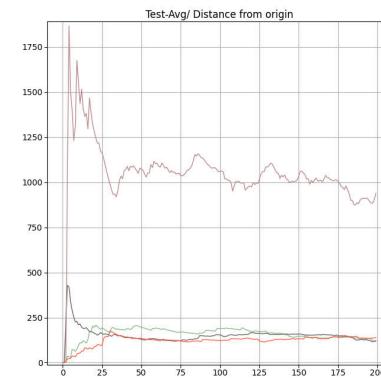
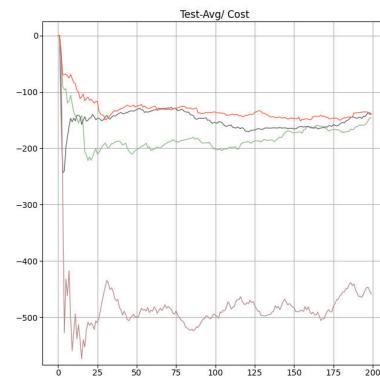
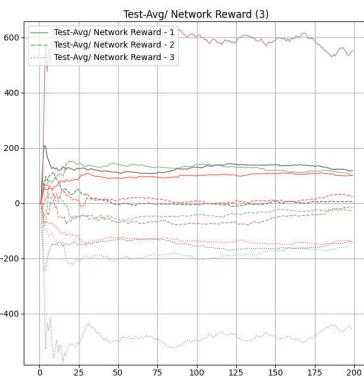


Lexicographic CDQN - [HLT,5FWD,CST]

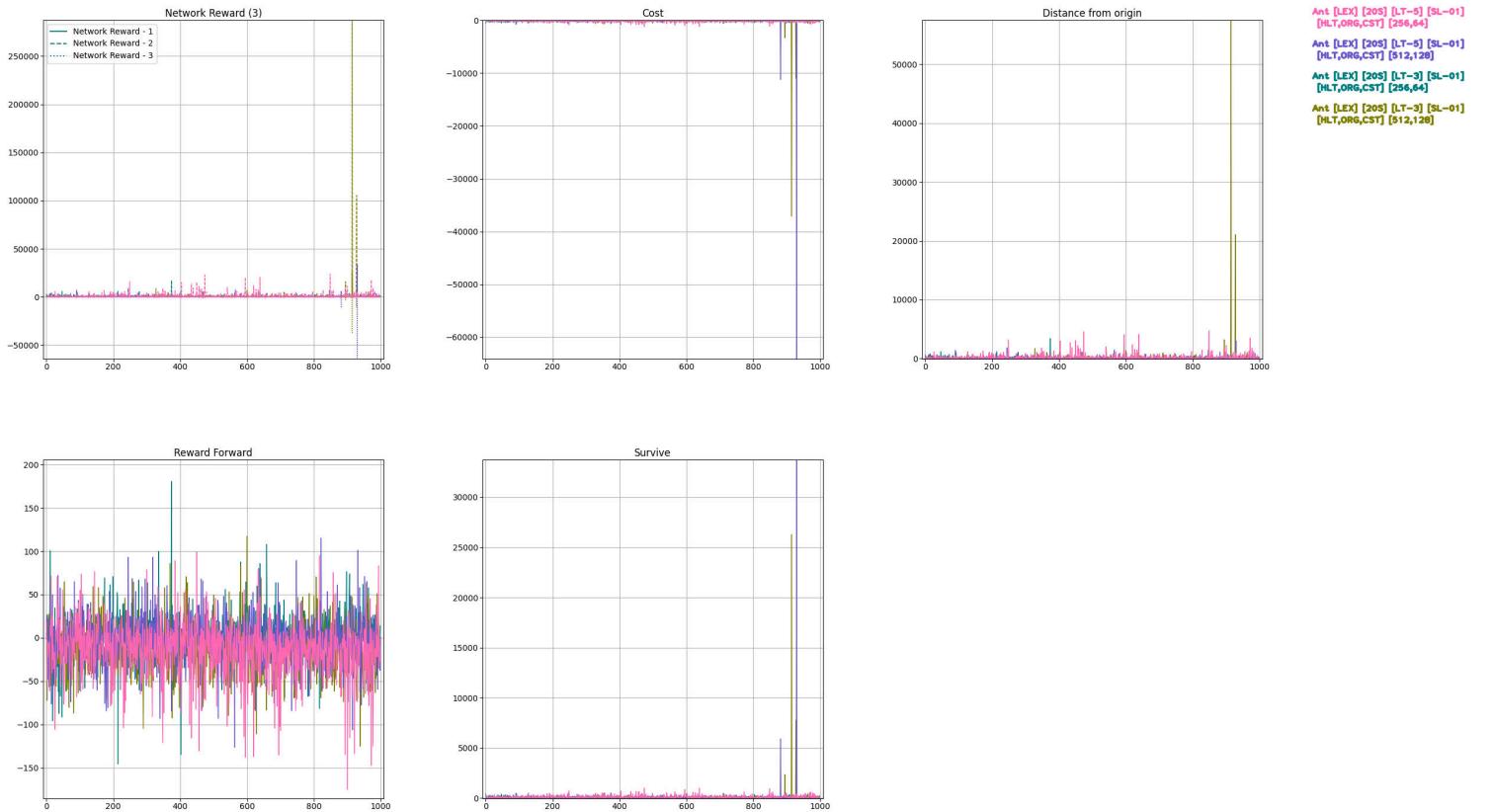


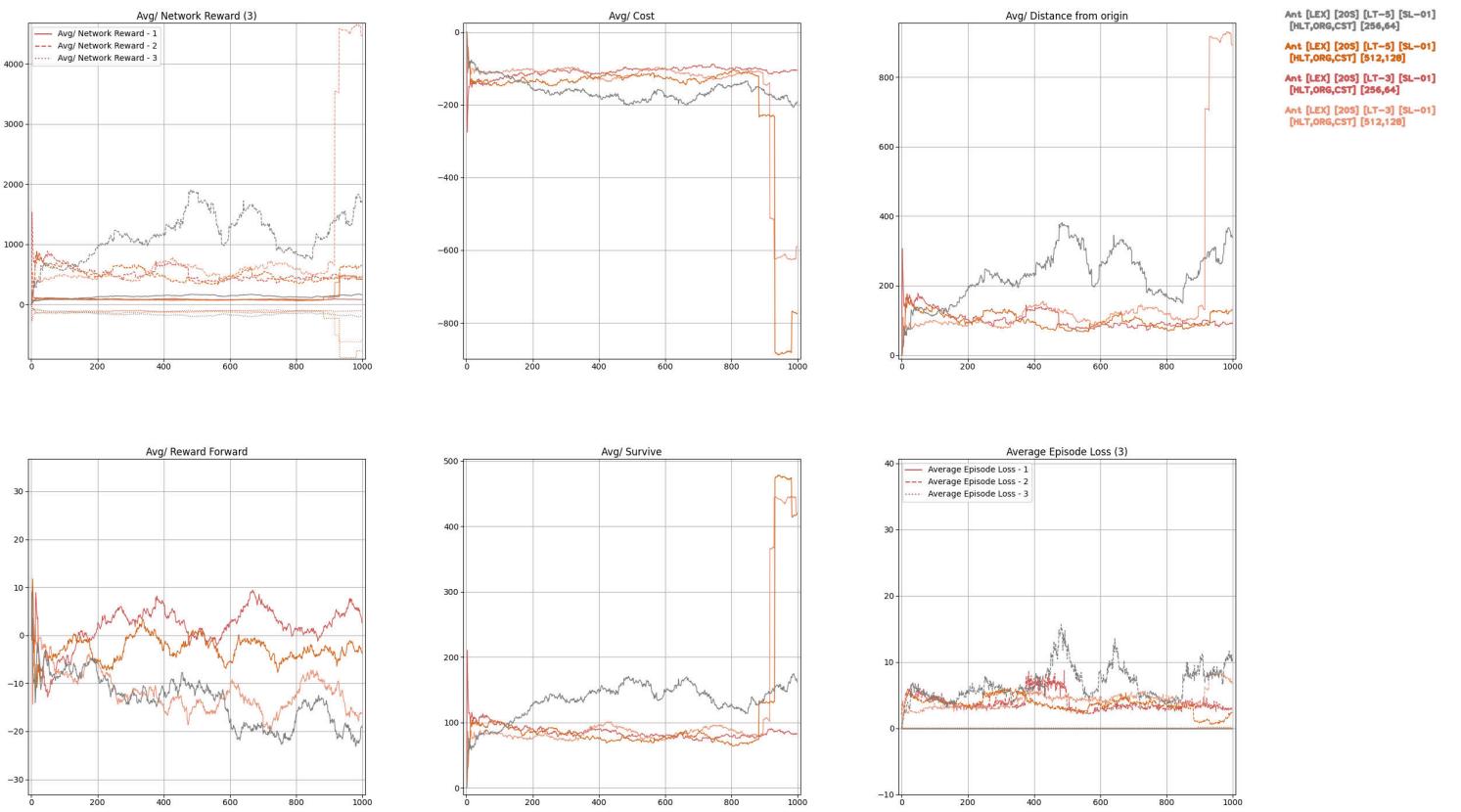


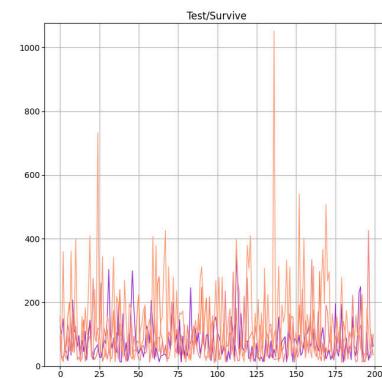
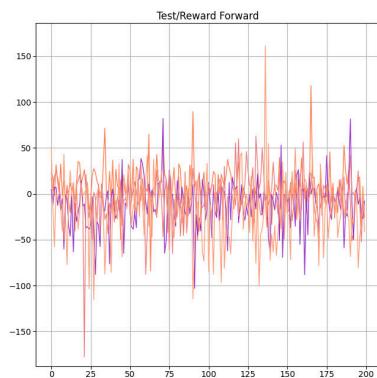
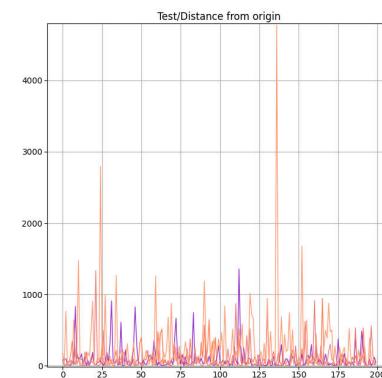
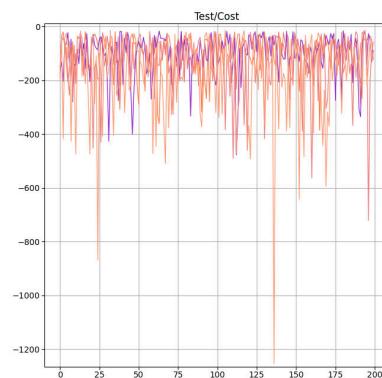
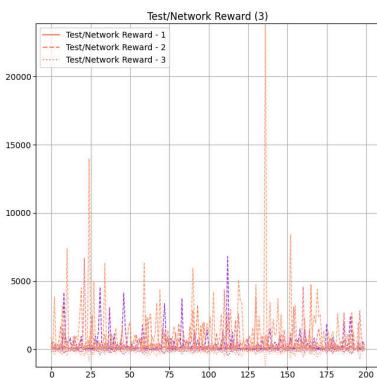


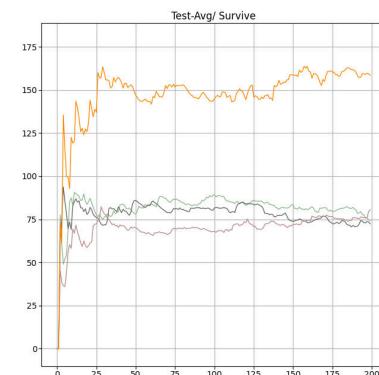
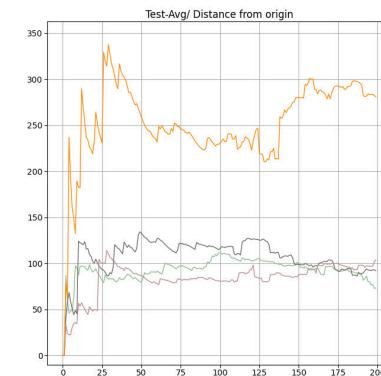
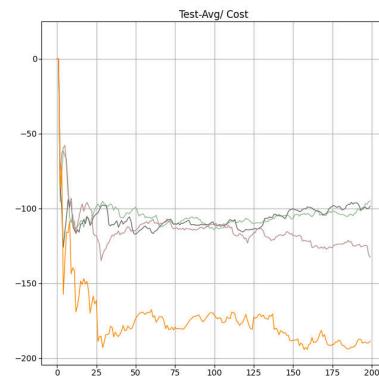
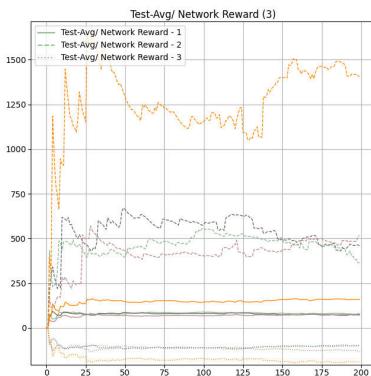


Lexicographic CDQN - [HLT,ORG,CST]

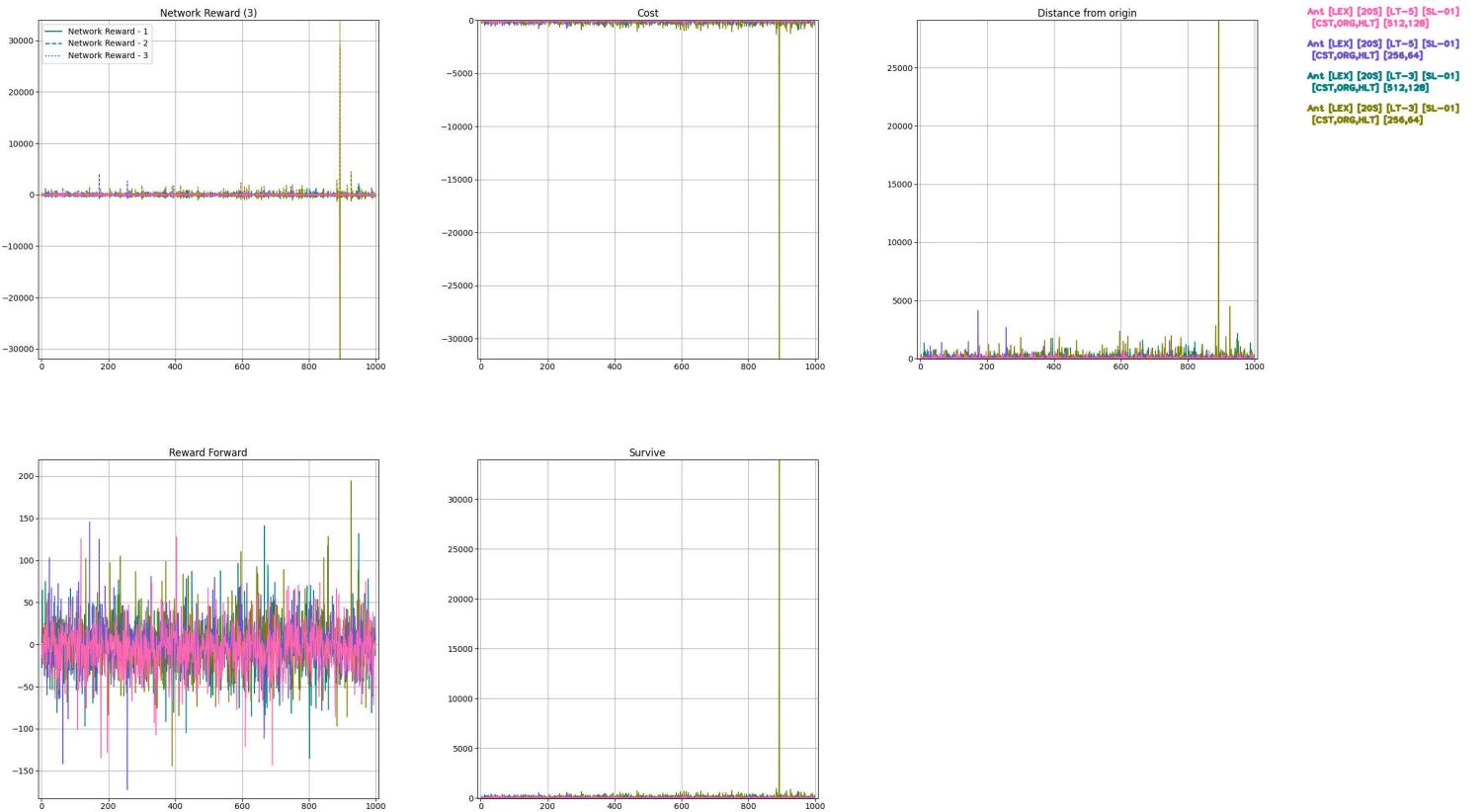


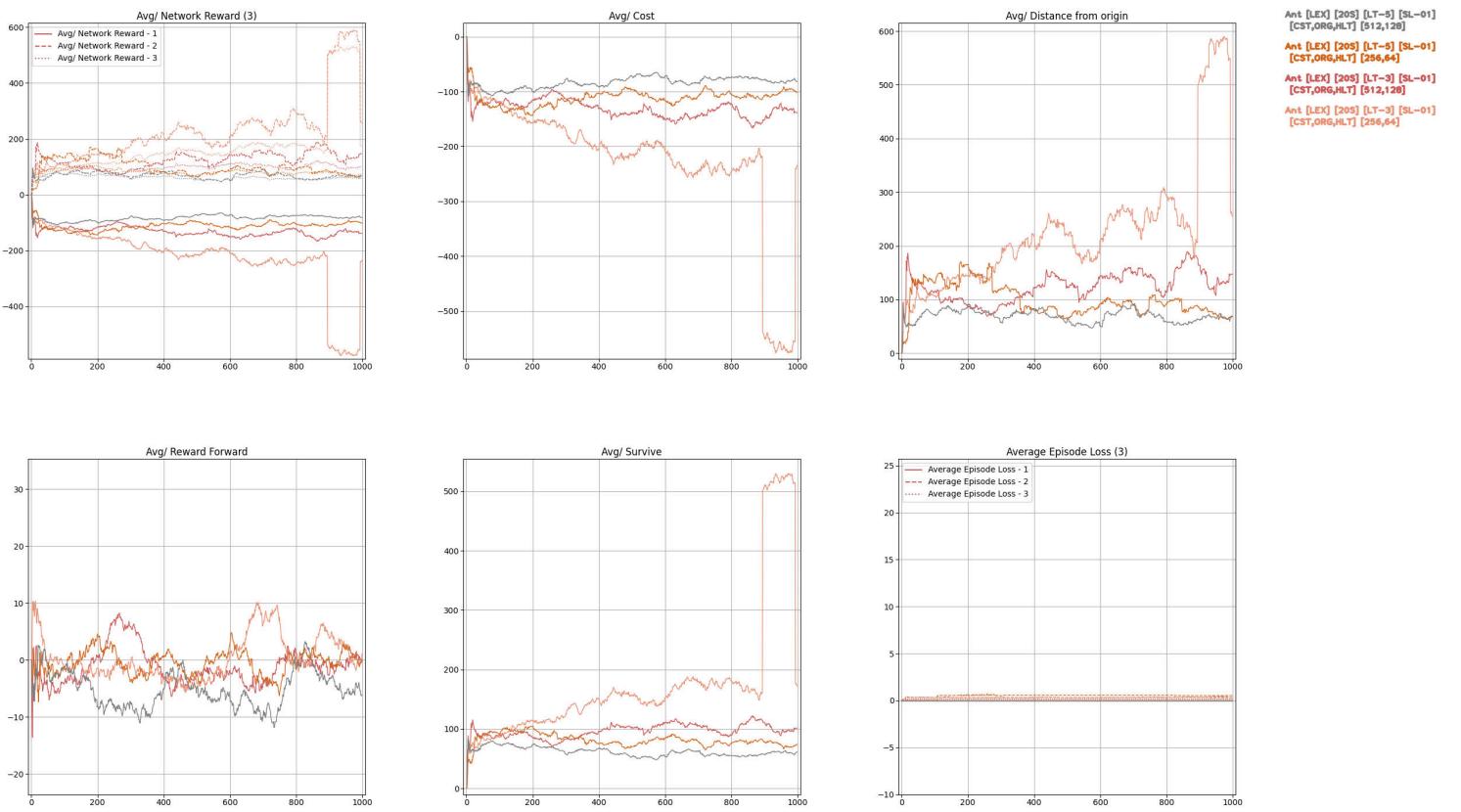


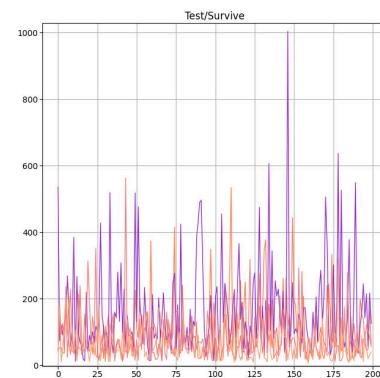
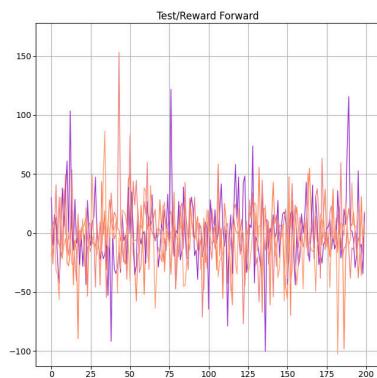
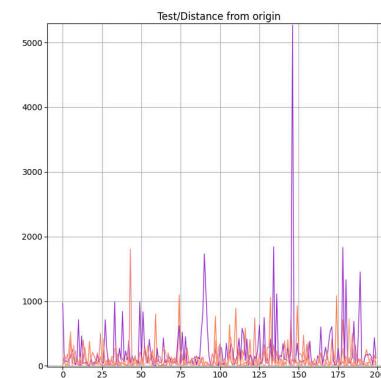
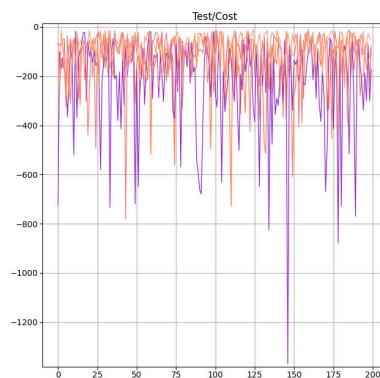
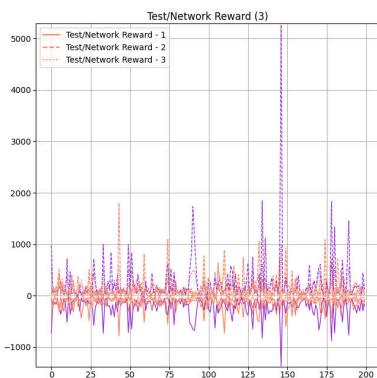




Lexicographic CDQN - [CST,ORG,HLT]







Ant [LEX] [20S] [LT-5] [SL-01]
[CST,ORG,HLT] [512,128]
Ant [LEX] [20S] [LT-5] [SL-01]
[CST,ORG,HLT] [256,64]
Ant [LEX] [20S] [LT-3] [SL-01]
[CST,ORG,HLT] [512,128]
Ant [LEX] [20S] [LT-3] [SL-01]
[CST,ORG,HLT] [256,64]

