

## 1 A Anonymized Code

2 The anonymized code of FOLLOWER is available at Learn-to-Follow-13223.

3 This repository contains everything needed to reproduce the results of FOLLOWER (train and eval  
4 scripts, pre-trained weights, dataset of maps, etc.).

5 Further in this Appendix we will refer to specific maps from our dataset by the names used in the  
6 repo.

## 7 B Hyperparameters

8 Table 1 presents the hyperparameters of FOLLOWER. The hyperparameters for which the tuning  
9 range is given (e.g. learning rate, LSTM hidden size,  $K$  (distance to the sub-goal), etc.) were  
10 optimized using Bayesian search. The observation radius was set to  $11 \times 11$  as it is commonly used in  
11 similar learning-based methods (with whom we compare). The parameters for the number of *rollout*  
12 *workers*, *environments per worker*, and *training steps* were empirically determined to decrease the  
13 overall learning time of the algorithm. For the remaining parameters (value loss coefficient,  $GAE_\lambda$ ,  
14 activation function, network initialization) we used the default values provided in the SampleFactory  
15 framework<sup>1</sup>.

16 We perform a hyperparameter sweep consisting of approximately 150 runs, totaling around 120 GPU  
17 hours. The final model was trained using a single TITAN RTX GPU in approximately 1 hour.

Table 1: The hyperparameters of FOLLOWER. The *tune range* column indicates the range for parameters adjusted through a hyperparameter optimization procedure.

Hyperparameter	Value	Tune range
Adam learning rate	0.000123	0.0001 – 0.0002
$\gamma$ (discount factor)	0.962983	0.95 – 0.99
Rollout	8	[4, 8, 16, 32]
Clip ratio	0.076785	0.05 – 0.2
Batch size	1024	[1024, 2048, 4096]
Optimization epochs	1	[1, 5, 10]
Entropy coefficient	0.014733	0.01 – 0.02
Value loss coefficient	0.5	-
$GAE_\lambda$	0.95	-
ResNet residual blocks	4	[2, 4, 6, 8]
ResNet number of filters	64	[32, 64, 128]
LSTM hidden size	512	[128, 256, 512]
Activation function	ReLU	-
Network Initialization	orthogonal	-
Number of agents	[16, 32, 64, 128]	-
Rollout workers	8	-
Environments per worker	2	-
Training steps	60000000	-
$K$ (sub-goal dist.)	2	[1, 2, 3, 4, 5, 6]
$H$ (sub-goal recalc.)	10	[2, 4, 8, 10, 12]
$C$ (add. transition cost)	0.4	[0.0, 0.2, 0.4, 0.6, 0.8, 1.0]
Observation radius	$11 \times 11$	-

## 18 C Weighted Transition Cost

19 As explained in the main body of the paper, each agent utilizes the heatmap of frequently used grid  
20 cells for individual pathfinding. The number of times the other agents were seen in a certain cell is  
21 multiplied by the user-defined parameter  $C$  and added to the transition cost to that cell. This helps  
22 each agent to avoid areas that are often used by everyone and thus to pro-actively avoid congestion.

<sup>1</sup>[github.com/alex-petrenko/sample-factory](https://github.com/alex-petrenko/sample-factory)

23 The effect of incorporating additional transition cost is visualized in Fig. 1. Here the cells with the  
 24 higher intensity of red indicate the areas visited (by the agents) more frequently during the episode.  
 25 These heatmaps were constructed from solving a single instance with 128 agents with FOLLOWER as  
 26 well as with Randomized A\* (i.e. trimmed FOLLOWER lacking a learnable policy).

27 Clearly, when the additional transition costs are not employed, i.e.  $C = 0.0$ , the agents tend to use  
 28 the central part of the map often. This makes it hard for the agents to avoid each other. On the other  
 29 hand, when additional transition costs are applied, the agents get evenly distributed across the map,  
 30 which prevents congestion and increase the performance (as confirmed by the experiments, reported  
 31 in the main body of the paper).



Figure 1: Heatmaps representations of how often the agents visited certain cells of the grid when solving a particular LMAPF instance containing 128 agents.

## 32 D Tuning RHCR Parameters

33 As it was mentioned in the main text, we have tuned the parameters of RHCR before conducting the  
 34 empirical comparison to our method. We varied the values of such RHCR parameters as planning  
 35 horizon (2, 5, 10, 20), re-planning rate (1,5) and time limit for each re-planning attempt (1s, 10s).  
 36 Planning horizon parameter controls for how many time steps the resultant plans will be collision-  
 37 free. E.g. when it equals 10 it is guaranteed that for the next 10 time steps the agents following  
 38 the constructed plans will not collide. Re-planning rate determines how frequently (in time steps)  
 39 reconstruction of the plans (for all agents) occurs. Time limit parameter restricts the amount of time  
 40 (in seconds) which is allotted for each re-planning attempt.

41 Fig. 2 demonstrates the results of different versions of RHCR (note that planning horizon cannot be  
 42 lower than re-planning rate). The best average throughput was achieved by RHCR with re-planning  
 43 rate 5 and planning horizon 20 (denoted as ( $w = 5, h = 20$ ) in the figure). The same values of these  
 44 parameters were also used for the experimental evaluation of RHCR on the warehouse map in the  
 45 original paper. Thus, the results of this version were included into the main part of our paper.

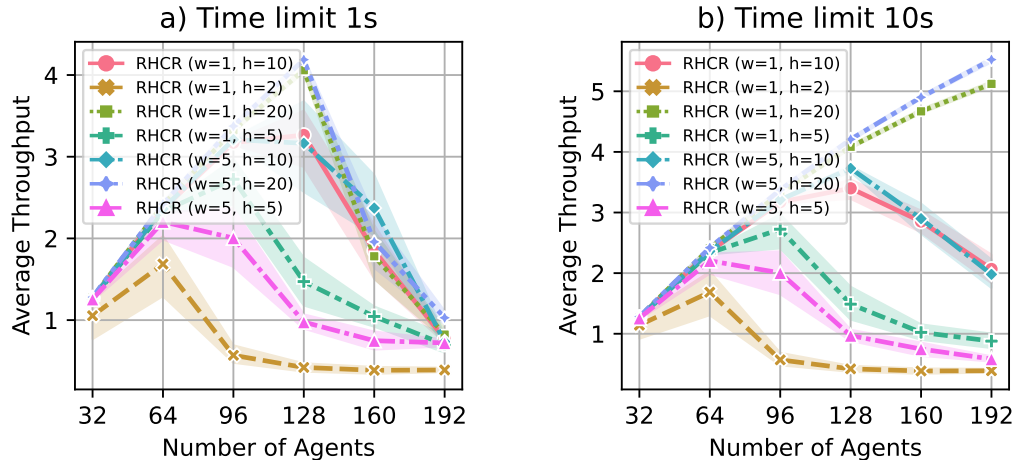


Figure 2: Evaluation of RHCR with varied parameter settings:  $w$  – re-planning rate,  $h$  – planning horizon.

## 46 E Impact Of the Episode Length

47 We set the episode length to 512 in all experiments.  
 48 Additionally, we examined the impact of episode  
 49 length on the throughput of FOLLOWER by running  
 50 additional experiments on the maze maps with  
 51 256 agents. The results are shown in Fig. 3. The  
 52 throughput increases first, starting from 1.43, but  
 53 then plateaus at 1.65. We attribute the initial increase  
 54 to the accumulation of knowledge regarding the transition  
 55 cost. We believe our choice of the episode  
 56 length (512) is reasonable.

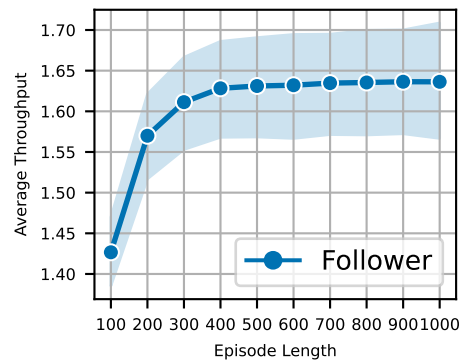
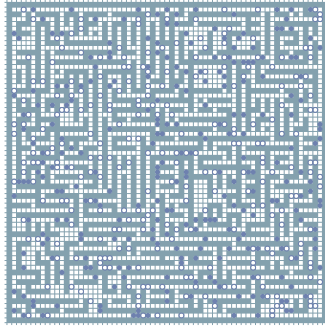


Figure 3: Impact of episode length on throughput of FOLLOWER with 256 agents.

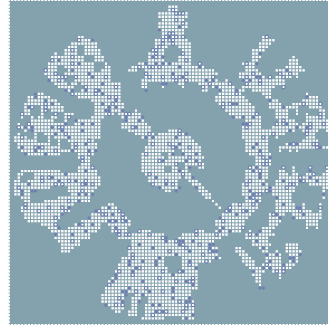
## 57 F Maps Visualizations

58 Fig. 4 illustrates examples of the maps used for testing.  
 59 The names of the maze-like maps and Pico-maps  
 60 are the same as in the repository.

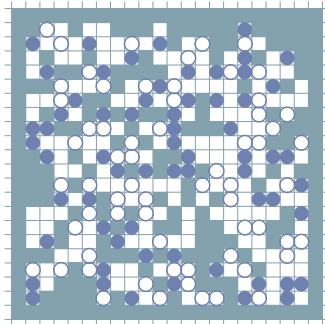
61 Initial positions of the agents are represented by the filled circles, while their (initial) goals are represented by the empty ones. Each agent is assigned a unique goal initially. Subsequent LMAPF goals are randomly generated, ensuring a feasible path from the agent’s current location to the goal exists. The goals for each agent are generated independently using a fixed seed, ensuring consistency and enabling fair testing of the algorithms (i.e. each algorithm gets the same start/goals locations).



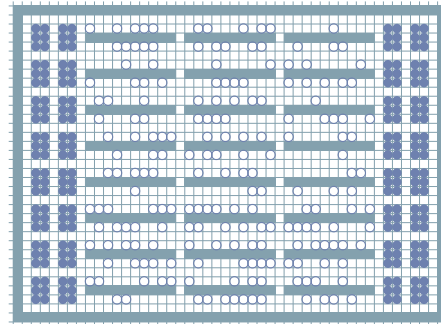
(a) Mazes-wc3-od70  
65 × 65



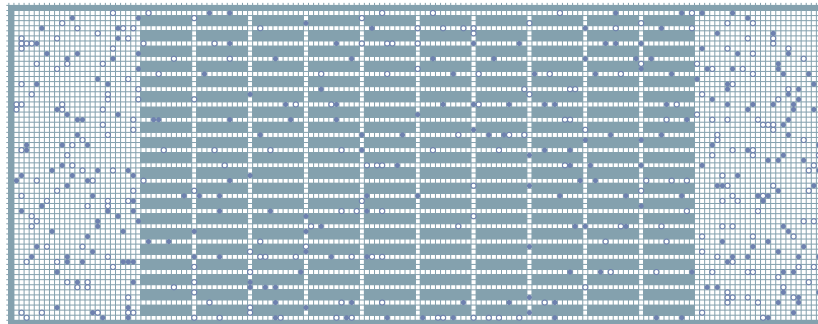
(b) Lak303d  
95 × 95



(c) Pico-s21-od30  
20 × 20



(d) Fulfillment warehouse map  
46 × 33



(e) Moving-ai-warehouse-10-20-10-2-1  
159 × 60

Figure 4: Visualizations of all the maps with a maximum number of agents used during the experimental evaluation.

## 66 G Detailed Comparison with PRIMAL2

67 The detailed results of a comparison between FOLLOWER and PRIMAL2 on the entire test set of  
 68 maze maps are illustrated in Fig. 5. FOLLOWER demonstrates superior performance on all maps. It is  
 69 important to note that the PRIMAL2 algorithm utilizes various heuristics to take advantage of the  
 70 topological characteristics of the maps (i.e. the presence of corridors).

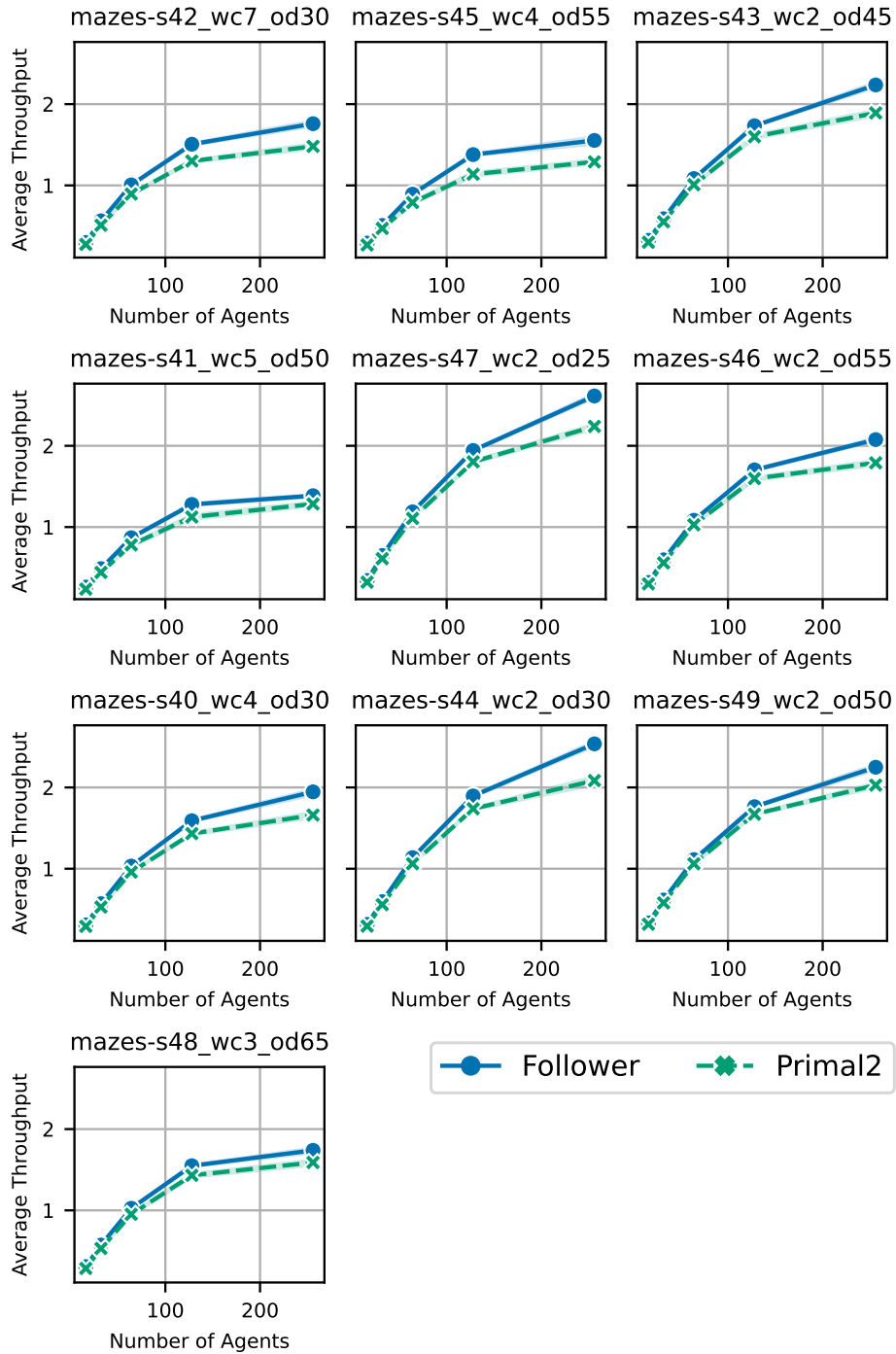


Figure 5: Average throughput on entire test set of the maze-like environments. The shaded area indicates 95% confidence intervals.