Multi-Agent Reinforcement Learning for Schedule-Constrained Automation

Kareem Eissa, Rayal Prasad, Ankur Kapoor

Siemens Healthineers

{kareem.abdelrahman, rayal.prasad, ankur.kapoor}@siemens-healthineers.com

Abstract

In modern automation settings, jobs are processed across numerous machines and are characterized by strong inter-task dependencies while adhering to limited equipment availability. When accounting for transportation of jobs between machines, this gives rise to a complex multi-agent routing problem with intricate operational limitations. Existing Multi-Agent Path Finding (MAPF) algorithms used for routing jobs already consider some aspects such as robustness, uncertainty, and plan execution. In this paper, we propose MAPF-SC - a lifelong variant of MAPF that incorporates scheduling constraints for a continuous stream of tasks. We propose to solve 14 MAPF-SC utilizing a Multi-Agent Reinforcement Learning (MARL) formulation with temporal and team reward. We investigate the effects of temporal and topological variations of various automation 18 scenarios on the performance of our method.

20 1 Introduction

21 State-of-the-art automation settings involve multiple interact-22 ing components that require coordinated actions. These sce-23 narios can be conceptualized as multiple jobs simultaneously 24 being transported to successions of machines for task pro-25 cessing. Whether it is logistics, manufacturing, or laboratory 26 and healthcare automation, the Multi-Agent Path Finding 27 (MAPF) formulation can be used to address important as-28 pects of these target domains such as collision avoidance and 45 decontamination. Finally, a job would have a specific se-29 travel time minimization to ensure optimal and safe motion 46 quence of operations to undergo, leading to strong prece-30 between machines. Recent works have shown significant pro- 47 dence constraints. Under the classical MAPF formulation, 32 al. 2021; Chudý et al. 2021].

One key assumption prevalent across MAPF variants is in- 50 cies of the automation pipeline. 34 stantaneous task processing which is atypical in many real-35 world situations. Even in the popular warehouse setting, a ro-36 bot agent delivering a payload would certainly require a non-37 instantaneous amount of time at the processing station. This 54 acterized by earliest arrival times, deadlines, non-instantane-38 discrepancy is amplified in domains such factory automation, 55 ous task processing times, and inter-task precedence con-39 where tasks demand specific non-trivial processing as a part 56 straints. The schedule is enforced by the automation pipeline 40 of the automation pipeline. The automation workflow may 57 capabilities i.e., machine operating characteristics. To handle 41 also be subject to restricted space due to costs or physical 58 this, we model Multi-Agent Path Finding with Scheduling



Figure 1: Example MAPF-SC instance. The bottom plane encodes the spatial information. Agents (circles) need to go to goals (black squares) according to schedules (colored rectangles). The elevation of the rectangles represents the earliest arrival time, and its height represents the duration of the time-window.

42 constraints such as ventilation, power outlets, structural con-43 straints, etc. Machines may further require additional time be-44 tween tasks for operations such as component switching or gress in different variants of MAPF [Li et al. 2021; Shahar et 48 agent congestions will arise as they fail to capture the prece-49 dence constraints and queuing effects induced by the intrica-

> 51 In this work, we formulate the MAPF-SC problem with ex-52 ecution schedules and precedence constraints. Such con-53 straints address the challenges of automation domains char

59 Constraints (MAPF-SC) as a sequential decision-making

60 problem that builds upon recent works [Sartoretti et al. 2019; 61 Damani et al. 2021; van Knippenberg et al. 2021] showing

62 promising results in adapting to uncertain and dynamical en-

63 vironments in a scalable decentralized manner.

Related Work 64 **2**

65 Classical MAPF is concerned with finding collision-free 113 cessing. $P_t(i,j)$ represents the task processing time during $_{66}$ paths for a group of agents from their start locations to their $_{114}$ which the agent is not allowed to move from the goal posi-67 goal locations. There are many lines of work on solving 115 tion. Figure 2 demonstrates relationship between these tem-68 MAPF [Surynek et al. 2019; Hoenig et al. 2016]. The related 116 poral components. Scheduling constraints may also involve 69 literature we focus on involves MAPF variants that factor in 117 precedence among tasks assigned to the same machine where 70 time, as well as Reinforcement Learning approaches to the 118 the order of arriving agents needs to be preserved. As in typ-71 problem.

72 2.1 MAPF with Time

73 Several works have studied temporal aspects of MAPF. 74 Wang et al. [Wang and Chen 2022] focus on minimizing the 123 the same. These situations arise when machines require some 75 violation of task specific due times. More recently, Gao et al. 124 postprocessing (e.g., replenish resource) after completing a 76 [Gao et al. 2023] focus on the maximizing the average cus- 125 task, leading to a lag in the execution of the schedule. Finally, 77 tomer satisfaction proportional to the degree of deadline vio- 126 agents are always present on the graph from the beginning 78 lation. Another line work [Ma et al. 2018; Zhang et al. 2022] 127 and do not disappear after completing their goals - the persis-79 focuses on the problem of maximizing the number of agents 128 tence of all agents on a space-constrained track is a key chal-80 that reach their goals within a global deadline or satisfying 129 lenge we aim to address. precedence-constraints among goal sequences. The Flatland 82 environment [Mohanty et al. 2020] simulates the railway set- 130 4 Method 83 ting as a multi-agent problem with the additional motion and 84 schedule constraints. However, trains are allowed to park off 85 the map before their departure and after their arrival, limiting 86 the complexity of the planning problem.

87 2.2 Reinforcement Learning for MAPF

88 Early works [Sartoretti et al. 2019; Damani et al. 2021] model 89 MAPF as a MARL problem and utilize expert training to train 90 decentralized policies with a mixture of imitation learning 91 and reinforcement learning. Another line of work focus on 92 solving MAPF with arrival and deadline constraints using 139 The observation space consists of two main components. The 93 Reinforcement Learning [Knippenberg et al. 2021] but are 140 first component is the agent-centered Local Field of View 94 limited to instantaneous task execution and may remove 141 (FOV) which builds upon earlier related works [Sartoretti et 95 agents outside their scheduled times similar to Flatland.

Problem Formulation 96 **3**

97 The MAPF-SC problem is defined by an undirected graph 146 experiments, we use a 7x7x5 FOV. 98 G = (V, E), a set of m agents $\{a_1, \dots, a_m\}$, and a schedule S. An 147 99 example instance is illustrated in figure 1. These agents rep-148 The spatial context of the task is represented through a three-100 resent the jobs within the schedule, each having a sequence 149 dimensional vector, G_s , as the difference in x-coordinates, 101 of associated tasks. The set of vertices V corresponds to loca- 150 difference in y-coordinates, and the magnitude of the vector constraints for each vertex. At each step, an agent can either $_{152}$ captured in another three-dimensional vector, τ_t , that pro-104 move to one of the unoccupied neighboring vertices or wait 153 vides the remaining time until the earliest arrival time A_t , the goals corresponding to vertices $\{v_i^1, v_i^2 \dots\}$ according to the 155 the deadline D_t has passed. The dimensions of this... schedule. The schedule represents time-windows of goal *j* for 108 agent *i* consisting of Earliest Arrival Time $A_t(i, j)$, Deadline 156 **4.2** Rewards 109 $D_t(i,j)$, and Goal Processing Time $P_t(i,j)$. An agent may ar-157 Contrary to classical MAPF, we do not assign a motion pen-110 rive at its goal prior to $A_t(i,j)$, but the task will not be pro-158 alty to minimize the travel distance. Our intuition is that an 111 cessed until the scheduled time $A_t(i,j)$. $D_t(i,j)$ indicates the



Figure 2: Different aspects of scheduling constraints.

112 latest time that an agent is allowed to reach its goal for pro-119 ical job shop automation workflows, the precedence con-120 straints have a higher priority than the time-window con-121 straints i.e., even if some agents are delayed past their dead-122 lines, the order in which the machine processes them remains

131 We model agent interactions with the environment as a Par-132 tially Observable Markov Decision Process (POMDP) 133 (S, O, A, P, R, γ) where S is the set of environment states, O 134 is the set of partial observations, A is the set of actions, P is 135 the transition probabilities, R is the reward function, and γ 136 is the discount factor. We aim to learn a decentralized policy 137 that can be deployed to different numbers of agents.

38 4.1 Observations

142 al. 2019; Damani et al. 2021]. The FOV consists of channels 143 denoting obstacles, the presence of other agents, and other 144 agents' target locations. An additional channel provides the 145 distance from each cell in the FOV to the agent's goal. In our

The second component encodes the agent's task as vectors. tions on the graph and the set of edges E represents motion 151 from the agent's current location to its goal. The context is at its current vertex. Each agent *i* is assigned a sequence of 154 duration until the deadline D_t , and a binary delay indicator if

159 agent may have surplus time until its earliest arrival time A_t . 212 • Redundancy: Lanes are connected to remove dead ends 160 Such an agent may potentially take a longer path to allow 213 and provide continuous coverage of the entire layout. 161 other agents to pass through. Taking this into consideration, 214 • Size: The number of machines and corresponding lanes are 162 minimizing only the travel distance may be detrimental to the 215 varied to create small (3) and large (6). throughput due to induced congestion. In our setup, each agent receives a positive reward when it reaches its goal within the scheduled time-window and a negative reward for each time step it occupies its goal before its earliest arrival 167 time A_t . Similarly, agents receive a negative reward for each 219 two degrees of freedom are padded with additional vertices. 168 time step it has not reached its goal past the deadline D_t to 169 incentivize arriving on time. However, agents may fail to 170 meet their schedules on-time especially early on in training. 221 To replicate the temporal intricacies found in automation do-171 Rather than terminating the episode, we permit soft deadline 222 mains, we introduce variations in schedule distributions 172 violations to allow the schedule's continued execution. All 223 across two dimensions, similar to the approach in the Flatland 173 precedence constraints of the schedule remain enforced 224 environment [Mohanty et al. 2020; Laurent et al. 2021]: 174 through the environment, which results in a queueing system 225 • Single-agent shortest distance (A* factor) which is the 175 for agents that have overlapping goal locations.

176 4.3 Architecture

177 We process the local FOV with a 3-layer convolutional neural 178 network, each layer followed by a max pooling operation (ex-179 cept for the last one which is followed by a global average 231 fluence the earliest arrival time for each agent providing min-180 pooling) to produce a FOV representation vector. The task 181 vectors (location and time) are then subsequently fed through fully connected layers to produce a combined task representation vector. Both the FOV and task vectors are then passed ²³⁵ fore completion. We investigate two distributions of sched-184 through a recurrent module to incorporate information about ²³⁶ ules we use: 185 past states as a mitigation strategy for partial observability 237 • Tight: A* factor is sampled uniformly from [1, 2] and task 186 problems. The model is optimized with Proximal Policy Op-187 timization which has shown great results on cooperative 239 • Relaxed: A* factor is sampled uniformly from [1.5, 3] and 188 multi-agent settings [Yu et al. 2022].

Experiments 189 5

190 We provide a comprehensive overview of the layout design's 191 attributes and its influence on input spatial characteristics.

192 Simultaneously, we adjust the time-windows' parameters to

193 induce variations in the resulting schedule slack.

194 5.1 Layouts

196 fixed layout is an important aspect of MAPF, we argue that it 249 dows to reach their goals. Table 1 summarizes the results for 197 does not hold the same value in the time-constrained scenario 250 different spatial and temporal configurations. 198 where the throughput of the entire system is limited by the 251 schedule capacity itself. In most automation settings, the 252 the agents learn to maintain a state of steady flow by circling 200 tasks are finite and are rate-limited by the machines that pro- 253 the track. Upon nearing their earliest arrival times, we ob-201 cess them. Adding more agents will simply lead to them wait 254 serve deviation from this behavior as agents move directly 202 idle and have sufficient time to navigate the track. This means 255 towards their goals. This continuous flow behavior globally 203 that as the number of agents increases in MAPF-SC, the bot- 256 reduces congestion. In the padded layouts with dead end 204 tleneck becomes then the schedule itself rather than the path 257 lanes (table1, rows 5-6), agents learn maneuvering behaviors 205 finding algorithm.

207 layout designs as shown in the Layout column in table 1. 260 "parking" to wait for their scheduled goals without blocking 208 These capture variations in restricted tracks connected by sin- 261 other agents.

209 gle lane corridors to capture situations with limited spatial re- 262

216 The agent count is adjusted proportionally and set to the num-217 ber of machines + 2.

218 • Padding: Decision regions where agents have more than

220 **5.2** Schedules

226 minimum distance traveled by an agent assuming there are no 227 other agents in the environment.

228 • Multi-agent congestion factor where we add the average of 229 all single agent shortest distances as a congestion estimate.

These dimensions represent scheduling constraints that in-232 imal necessary conditions for the schedule to be feasible. Ad-233 ditionally, we vary the size of the time-windows from earliest 234 arrival time to deadline, and the task processing runtime be-

238 durations are sampled uniformly from [4, 11].

240 task durations are sampled uniformly from [5, 25].

241 We generate schedules specific to each layout using OR-242 Tools [Perron and Furnon 2023] where the solver optimizes 243 for a randomly sampled list of agent tasks.

244 6 **Results and Discussion**

245 Through the experiments, we note that the model tends to 246 complete more goals on-time on a relaxed schedule versus a 247 tight one. Intuitively, this reflects a less constrained optimi-195 While scaling to a larger number of agents within the same 248 zation problem since agents have less restricted time-win-

In the closed-loop layouts (table 1, rows 1-4), we find that 258 through the padded sections surrounding the decision re-To investigate these effects, we focus on more restrictive 259 gions. Agents also tend to use these padding vertices as

In general, we find the most restrictive instances (table 1, 210 sources e.g., repurposing an existing automation plant. We 263 rows 7-8) to be less stable in training and the RL policy often

211 vary the layouts across three design dimensions:

| Layout | Schedule | On-Time % | | Late % | |
|----------|----------|------------------|----------|--------|-------|
| | | Mean | Std | Mean | Std |
| | Layout | ts with red | undancy | | |
| (padded) | Relaxed | 97.72 | 3.44 | 3.13 | 4.14 |
| | Tight | 96.22 | 3.40 | 4.71 | 3.82 |
| | Relaxed | 96.69 | 9.01 | 3.38 | 5.80 |
| (padded) | Tight | 98.01 | 2.63 | 2.58 | 3.04 |
| | Relaxed | 95.41 | 6.63 | 5.79 | 6.30 |
| | Tight | 93.61 | 8.72 | 7.00 | 7.35 |
| | Relaxed | 96.66 | 6.27 | 4.27 | 5.85 |
| | Tight | 88.98 | 11.41 | 12.22 | 10.83 |
| mm | Relaxed | 94.80 | 8.87 | 6.11 | 7.80 |
| (padded) | Tight | 89.28 | 8.41 | 11.72 | 6.46 |
| 111 | Relaxed | 97.06 | 7.17 | 3.72 | 7.24 |
| (padded) | Tight | 97.17 | 7.60 | 2.93 | 4.98 |
| | Layouts | without re | dundancy | | |
| | Relaxed | 46.00 | 22.46 | 31.98 | 18.28 |
| | Tight | 31.37 | 14.40 | 60.29 | 12.28 |
| \Box | Relaxed | 11.33 | 11.30 | 50.59 | 20.60 |
| | Tight | 14.85 | 10.05 | 67.78 | 10.49 |

Table 1: Results for layouts with different degrees of redundancy under two schedule distributions.

264 collapses without converging to acceptable performance lev-265 els. More restrictive optimization problems induce harder underlying MDPs which destabilizes the RL training [Hafner et 267 al. 2023]. The less restrictive scenarios (rows 1-6) almost 268 converge near 10k episodes while the most restrictive scenar-269 ios progress slower by an order of magnitude. We hypothesize that the performance degradation in the more restrictive scenarios has to do with the increasing planning complexity 272 arising from the restrictive maneuverability of the layouts. 304 [Hafner et al. 2023] Hafner, D.; Pasukonis, J.; Ba, J.; Lillic-274 stricted layouts yield configurations that require long-horizon 306 els. ArXiv preprint arXiv:2301.04104. 275 swapping maneuvers with higher degrees of agent coordina-276 tion, and we look to further explore...



Figure 3: Example failure cases. Agents are represented by squares with lines pointing to their scheduled goal locations.

Conclusions 77 7

78 In this paper, we proposed a variant of the MAPF problem 79 (MAPF-SC) to address scheduling constraints using shared 80 rewards and spatiotemporal observations. We evaluated our algorithm on a variety of automation layouts and schedule 82 distributions, and observed how the spatial constraints pre-83 sent a harder challenge than the temporal constraints. As next 84 steps, we aim to improve agent coordination to improve per-85 formance in the more challenging layouts.

86 Disclaimer

87 The concepts and information presented in this paper are 88 based on research results that are not commercially available. 89 Future commercial availability cannot be guaranteed.

90 References

91 [Chudý et al. 2021] Chudý, J.; Surynek, P. 2021. ESO-92 MAPF: Bridging Discrete Planning and Continuous Execu-93 tion in Multi-Agent Pathfinding. Proceedings of the AAAI Conference on Artificial Intelligence 35: 16014–16016.

295 [Damani et al. 2021] Damani, M.; Luo, Z.; Wenzel, E.; Sar-296 toretti, G. 2021. PRIMAL₂: Pathfinding Via Reinforcement 297 and Imitation Multi-Agent Learning - Lifelong. IEEE Robot-298 ics and Automation Letters 6: 2666–2673.

299 [Gao et al. 2023] Gao, J.; Liu, Q.; Chen, S.; Yan, K.; Li, X.; 300 Li, Y. 2023. Multi-agent path finding with time windows: 301 Preliminary results. Proceedings of the 2023 international 302 conference on autonomous agents and multiagent systems: 303 2586-2588.

Examples of failed instances are shown in figure 3. These re- 305 rap, T. 2023. Mastering diverse domains through world mod-

307 [Hoenig et al. 2016] Hoenig, W.; Kumar, T.K.; Cohen, L.; 308 Ma, H.; Xu, H.; Ayanian, N.; et al. 2016. Multi-Agent Path 309 Finding with Kinematic Constraints. Proceedings of the International Conference on Automated Planning and Scheduling 26: 477-485.

312 [van Knippenberg et al. 2021] van Knippenberg, M.; Holenderski, M.; Menkovski, V. 2021. Time-constrained 314 multi-agent path finding in non-lattice graphs with deep reinforcement learning. Proceedings of the 13th Asian confer-316 ence on machine learning 157: 1317-1332.

317 [Laurent et al. 2020] Laurent, F.; Schneider, M.; Scheller, C.; 318 Watson, J.; Li, J.; Chen, Z.; et al. 2021. Flatland competition 319 2020: MAPF and MARL for efficient train coordination on a 320 grid world. Proceedings of the NeurIPS 2020 competition and 321 demonstration track 133: 275–301.

322 [Li et al. 2021] Li, J.; Tinka, A.; Kiesel, S.; Durham, J.W.; 323 Kumar, T.K.S.; Koenig, S. 2021b. Lifelong Multi-Agent Path

324 Finding in Large-Scale Warehouses. Proceedings of the

- 325 AAAI Conference on Artificial Intelligence 35: 11272–326 11281.
- 327 [Mohanty et al. 2020] Mohanty, S.; Nygren, E.; Laurent, F.;
- 328 Schneider, M.; Scheller, C.; Bhattacharya, N.; et al. 2020.
- 329 Flatland-RL: Multi-Agent Reinforcement Learning on 330 Trains.
- [Perron and Furnon 2023] Perron, L.; Furnon, V. 2023. OR-332 Tools.
- 333 [Sartoretti et al. 2019]Sartoretti, G.; Kerr, J.; Shi, Y.; Wag-
- 334 ner, G.; Kumar, T.K.S.; Koenig, S.; et al. 2019. PRIMAL:
- 335 Pathfinding via Reinforcement and Imitation Multi-Agent 336 Learning. IEEE Robotics and Automation Letters 4: 2378–
- 337 2385.
- 338 [Shahar et al. 2021] Shahar, T.; Shekhar, S.; Atzmon, D.; Saf-
- 339 fidine, A.; Juba, B.; Stern, R. 2021. Safe Multi-Agent Path-
- 340 finding with Time Uncertainty. Journal of Artificial Intelli-341 gence Research 70.
- 342 [Surynek et al. 2019] Surynek, P.; Kumar, T.K.S.; Koenig, S.
- 343 2019. Multi-agent Path Finding with Capacity Constraints.
- 344 In: Alviano, M.; Greco, G.; Scarcello, F. (Eds.), AI*IA 2019 345 – Advances in Artificial Intelligence, Vol. 11946, Springer
- 346 International Publishing, Cham, p.235–249.
- 347 [Wang et al. 2022] Wang, H.; Chen, W. 2022. Multi-Robot
- 348 Path Planning With Due Times. IEEE Robotics and Automa-349 tion Letters 7: 4829–4836.
- 350 [Yu et al. 2022] Yu, C.; Velu, A.; Vinitsky, E.; Gao, J.; Wang,
- 351 Y.; Bayen, A.; et al. 2022. The surprising effectiveness of
- 352 PPO in cooperative multi-agent games. Advances in neural 353 information processing systems 35: 24611–24624.
- 354 [Zhang et al. 2022] Zhang, H.; Chen, J.; Li, J.; Williams,
- 355 B.C.; Koenig, S. 2022. Multi-agent path finding for prece-
- 356 dence-constrained goal sequences. Proceedings of the 21st 357 international conference on autonomous agents and multia-
- 358 gent systems: 1464–1472.