



Figure 1: Time-constrained evaluation in the grab a chair domain with 33 (left), 65 (middle) and 129 (right) agents.

1 We thank all reviewers for the encouraging and constructive feedback and the relevant pointers.

2 **Reproducibility (@R2, R3):** Our codebase is available at <https://www.dropbox.com/sh/sg2qyfpdwfmkfqd/AACq5R5BQS-jpSS60MFS3FrLa?dl=0> with source code, pretrained models and instructions to reproduce every figure
 3 in the submission. We will make the codebase available on Github after the review phase.
 4

5 **Comparison to exogenous variables (@R4):** Thanks for the pointers to the related work which indeed are similar in
 6 spirit. We will add the discussion of these works to our related work section. Nevertheless, we would like to point out
 7 that exogenous variables are fundamentally different from the non-local variables: by definition, exogenous variables
 8 refer to those variables that are **beyond the control of the agent** in the sense that **the values of which are not affected,**
 9 **directly or indirectly, by the agent's actions** (Boutilier et al., 1999; Zhang et al., 2017; Chitnis et al., 2019). However,
 10 the non-local variables in the influence-based abstraction (IBA) (Oliehoek et al., 2012) refer to those variables that **do**
 11 **not directly affect the agent's observation and reward.** Therefore, the exogenous variables and non-local variables
 12 are in general two different sets of variables that can be exploited to reduce the state space size. For instance, in
 13 the traffic problem of Figure 4a, there are *no* exogenous variables as our action can directly or indirectly effect the
 14 transitions at other intersections (by taking or sending vehicles from/to them). In other words, **our approach allows us**
 15 **to reduce the state space of this problem beyond the exogenous variables.**

16 **Clarifications on data collection and influence learning (@R3, R4):** In Section 3.1, we formalize influence prediction
 17 as a classification problem where inputs are trajectories of local states and actions, and outputs are influence source
 18 variables. In this work, the training trajectories are sampled **from the global simulator** (cf. line 162) with an exploratory
 19 policy and therefore no human annotation is required.

20 In all our experiments, we collected 1000 episodes (line 199) and trained the influence predictor to convergence, which
 21 did not take more than an hour. While such offline computation could be costly (as noted by R4), we believe that in
 22 many settings the real-time constraints during the online phase are the bottleneck for application and these can be
 23 significantly brought down with our approach.

24 **Clarifications on experimental design and results (@R1, R2, R3, R4):** As R1 and R3 noted, we only investigate
 25 the performance of our approach in those domains for which it was designed: problems with clear factorized structure.
 26 For these problems, we think our benchmarking domains are general enough so that the results can be easily translated.
 27 We certainly would not want to claim anything about a broader class of problems. We do stress, however, that there is a
 28 large literature on factored (multiagent) decision making that has been motivated by many different applications, such
 29 as smart grids, warehouse commissioning, etc, and that Bayesian networks (the tool we use to represent the factored
 30 structure) are one of the most applied AI techniques in history (even if perhaps overtaken by NNs by now).

31 When sufficient time is available for online planning, the benefits in the traffic domain are indeed moderate (R1).
 32 However, in many real-world decision problems, the online decision time is limited, and precisely for those cases our
 33 approach significantly outperforms the baselines (e.g., when <10s per action allowed in Fig. 4d). To further test this,
 34 Figure 1 shows real-time planning results in the grab a chair domain, as suggested by R3. These further demonstrate
 35 that the advantage of our approach in time-constrained settings, as the global model of the problem gets more complex.

36 **Further clarifications: R3:** learning models that matter with value, goal and task awareness is indeed relevant to the
 37 idea of capturing the impact from the rest of the environment by predicting only the influence source variables. However,
 38 these approaches do not perform explicit abstraction as ours does and therefore are complimentary. **R2:** we chose to
 39 use POMCP in unmodified form in our experiments to make it easier to interpret the benefits of our approach, but
 40 indeed our approach can be combined generally with other planners and more advanced particle filters, and therefore
 41 all known improvements from SIR can be applied. **R4:** it is not possible to define a baseline local simulator with "no
 42 influence" because by definition the influence destination variables have dependency on the source variables. We agree
 43 that attention mechanisms are a promising avenue for learning influence representations (not the focus of this paper)
 44 and thus are a promising direction of future work. As for the grab a chair domain (R4), we performed online planning
 45 for a finite horizon of 10 (line 198-199), which explains the choice for undiscounted reward $\gamma = 1$. At every time step,
 46 each agent can only decide to grab the chair on their left or right and therefore it is not possible for an agent to obtain
 47 more than one chair at a time step (line 33, 388-389).