# 1. Effect of $\alpha$ on Performance

Fig. 1 shows the effect of varying $\alpha$ in the indicator function. We can see if $\alpha$ is too low, performance degrades considerably because UTRAN degenerates to QMIX (the indicator function mistakes all targets for stochastic samples, and thus no target is transformed). In contrast, UTRAN replaces all targets with the best per-agent value if $\alpha$ is too high. Its poor performance in Fig. 1-a and c suggest difficulties in solving the stochastic state-action pairs.

Therefore, $\alpha$ is a task-specific hyperparameter related to the reward and the state scale. Our recommendation is to normalize the state, observation, and reward when we create an environment.
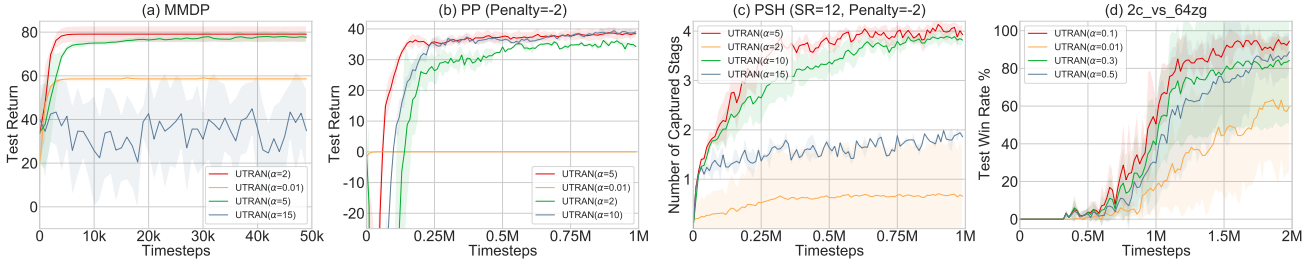


*Figure 1.* The results of UTRAN with different $\alpha$ on MMDP, PP, PSH, and SMAC.