1  We gratefully appreciate the efforts made by all the reviewers. Thanks to **Reviewers #1** and **#2** for bringing up
2  the missing references. Hughes et al. [2018] extend the inequity aversion model and define a shaped reward $r_i -$
3  $\frac{\alpha}{N-1} \sum \max(r_j - r_i, 0) - \frac{\beta}{N-1} \sum \max(r_i - r_j, 0)$. Wang et al. [2018] design a reward network generating intrinsic
4  reward, evolved based on the group's collective reward. Peysakhovich et al. [2018] propose a shaped reward
5  $r_i = \alpha r_i + (1 - \alpha) r_j$ for two-player Stag Hunts. The baseline `Avg` can be seen as its multi-player version as the authors
6  claim. These works aim to improve cooperation but cannot guarantee fairness. We compare against Hughes et al. [2018],
7  `Inequity Aversion`, in job scheduling. Table 1 shows the CV of `Inequity Aversion` is better than `Independent`
8  but still much worse than `FEN`, and the resource utilization is much lower. That shows `Inequity Aversion` cannot
9  solve job scheduling fairly and efficiently. More details will be included in the final version. We will also include the
10 review of the missing references in the final version.

<table>
<tr><td colspan="3">Table 1: Job scheduling</td></tr>
<tr><td></td><td>resource utilization</td><td>CV</td></tr>
<tr><td>Independent</td><td>96% ±11%</td><td>1.57 ±0.26</td></tr>
<tr><td>FEN</td><td>**90%** ±5%</td><td>**0.17** ±0.05</td></tr>
<tr><td>Inequity Aversion</td><td>72% ±9%</td><td>0.69 ±0.17</td></tr>
</table>

| Table 2: Hierarchy | | |
|---|---|---|
|  | resource utilization | CV |
| Min w/ Hierarchy | 62% ±9% | 0.31 ±0.11 |
| Avg w/ Hierarchy | 84% ±6% | 0.61 ±0.14 |
| Min+$\alpha$Avg w/ Hierarchy | 71% ±8% | 0.28 ±0.09 |
| FEN | **90%** ±5% | **0.17** ±0.05 |

12 **Reviewer #1** To verify the effectiveness of the hierarchy, we use the hierarchy with other baselines in job scheduling.
13 Table 2 shows their performance has a certain degree of improvement, especially the resource utilizations of `Min` and
14 `Min+`$\alpha$`Avg` raise greatly and the CV of `Min+`$\alpha$`Avg` reduces significantly. That demonstrates the effect of the hierarchy.
15 However, these baselines with the hierarchy are still worse than `FEN` in both resource utilization and CV, verifying the
16 effectiveness of the fair-efficient reward.

17 The intuition of the fair-efficient reward is to maximize the resource utilization while punish the agent's utility deviation
18 from the average, taking both fairness and efficiency into consideration. Also, the fair-efficient reward is suitable for
19 decentralized training, which can be easily coordinated by $\bar{u}$. We design the fair-efficient reward, prove it satisfies the
20 criteria of Propositions 1 and 2, and empirically verify it really works well.

21 The main hyperparameters are contained in the Appendix, we will make a further supplement in the final version. All
22 the results are obtained by five runs with different random seeds and presented with standard deviation (line 234), and
23 we will make it clearer.

24 **Reviewer #2** It is really a constructive suggestion to analyze the behavior of the controller. Figure 1 visualizes the
25 probability of selecting sub-policy $\phi_1$ and other sub-policies in terms of the utility deviation from average, $(u_i - \bar{u})/\bar{u}$.
26 It shows when the agent's utility is below average, the controller is more likely to select $\phi_1$ to occupy the resources, and
27 when the agent's utility is above average, the controller tends to select other sub-policies to improve fairness. Thus, it
28 can be seen that the balance of these two kinds of sub-policies depends on the current fair-efficient reward.

29 Sorry for the confusion induced by "decentralized training." By that we mean the
30 training of each agent requires only limited information exchanging with neighboring
31 agents. We will use more precise expression to replace that in the final version. At test
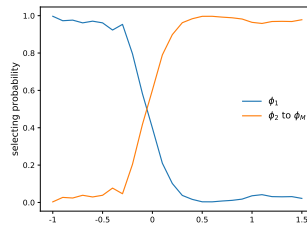32 time, it is the same as training where the controller chooses one sub-policy every $T$
33 timesteps.



Figure 1: Selecting probability over $(u_i - \bar{u})/\bar{u}$.

34 **Reviewer #3** We explain the necessity of the hierarchy from three aspects. First, the
35 hierarchy reduces the difficulty of learning both efficiency and fairness. Since the
36 problem is a multi-objective optimization, the learning difficulty for a single neural
37 network cannot be neglected. In the hierarchy, each sub-policy focuses on its own easy
38 objective and there is no conflict; the controller focuses on the fair-efficient reward by
39 selecting the sub-policies, without directly interacting with the environment. Second, the fair-efficient reward changes
40 slowly since it is slightly affected by the sub-policy's action in one timestep. Thus, the controller can plan over a
41 long-time horizon to optimize both objectives. Third, in all the three experiments, FEN with hierarchy outperforms
42 the version without hierarchy, verifying the hierarchy helps the learning greatly. We also use the hierarchy with other
43 baselines and their performance improves as shown in Table 2, which also shows the effectiveness of the hierarchy.

44 Although the training is decentralized, the agent can obtain the average utility $\bar{u}$ by Gossip, and hence each agent knows
45 the utility deviation from the average. To make its own fair-efficient reward higher, the agent with lower $u_i$ must occupy
46 more resources and the agent with higher $u_i$ than $\bar{u}$ will choose to not occupy resources. That is why the policies can be
47 coordinated. Each agent only focuses on its own fair-efficient reward and the fairness could be achieved.

48 This paper focuses on achieving both fairness and efficiency in the context of *common resources* which is one of
49 important fields in MARL. Propositions 1 and 2 have proved that the resources will be *fully* occupied and *equally*
50 allocated *in infinite-horizon sequential decision-making*, thus the CV is also minimized according to its definition.