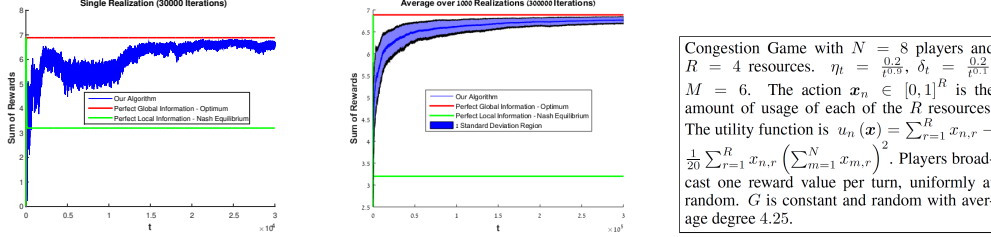


1 We thank the reviewers for their detailed and constructive comments, especially during these challenging times.

2 **Numerical Validation (R1,R3,R4):** We will add the figures below, displaying the typical behavior of a realization and
 3 the statistical behavior over time. This also shows that the convergence time is reasonable for this challenging problem
 4 with limited information. For specific applications, convergence can be accelerated by replacing our generic gradient
 5 estimator with an estimator that exploits the structure of the specific problem (e.g., parameterized regression).



6
 7 **Variance (R1):** R1 is correct that the variance increases like $\frac{1}{\delta_t^2}$, but the step size η_t multiplying the (estimated)
 8 gradient overcomes that. This is why $p - q > 0.5$, for which, $\sum_t \eta_t^2 / \delta_t^2 = \sum_t t^{-2(p-q)} < \infty$ is required.

9 **The Perturbations (R1):** are indeed being independently generated across players. The reason is that the perturbation
 10 of player n , z_n , is only meant to estimate the gradient of the sum of rewards wrt x_n . The gradient approximation
 11 argument is used N times independently, and no two perturbations ever serve the same purpose. Hence, for player
 12 $m \neq n$, z_n serves as undesired (but necessary) noise in their communicated reward value (so no expectation is 0).

13 **Lipschitz Factor (R1):** A computation very similar to that suggested by R1 is done in (45) to bound the bias from
 14 the M steps delay. L appears in almost all the variance and bias terms in our proofs, affecting the convergence time.
 15 However, the theorem guarantees convergence for all L , and η_t (or δ_t) do not have to be tuned with respect to L .

16 **Ergodic Markov Graph (R2,R4):** S_G is a set of states of the process $G(t)$, so it is a set of graphs. We now denote
 17 this set \mathcal{G} to avoid confusion. The union is between graphs i.e., $\bigcup_{G \in S_G} G$ as opposed to $\bigcup_{G \in S_G} \{G\}$ (using standard
 18 notation), which we now explicitly mention. R2 is correct that $G(t)$ transitions into $G(t+1)$ with no restrictions
 19 except that it constitutes an ergodic Markov chain. It is also correct that the algorithm does not know the graph. The
 20 purpose of this Markovian model is to capture a general time-varying process. This is the reality with physical agents,
 21 that move in space or that objects in space move around them, changing the configuration of the network. We agree
 22 that this is not the most important part of the paper. A useful and simple special case of this model (R4) is indeed
 23 when G is constant in time and connected or strongly connected (if directed). While not many algorithms can work
 24 under this general and challenging model, we are not the first to consider this well established model (see [34]). We
 25 hope that our explanation and the new notation make clear that the formulation is indeed rigorous and precise.

26 **M versus N (R2):** We discuss the requirement on M in Section 7.1 in the appendix. For a constant graph G , M needs
 27 to be larger than the diameter of G , and then R2 is correct that $M \geq N$ is always enough (and does not hurt much).
 28 However, for some Markov Chain (MC) $G(t)$, $M = N$ might not be enough if, for example, player 1 is isolated
 29 most of the turns and has to wait to get connection (edges). Hence, rigorously, a large enough M exists if the MC is
 30 ergodic, which might be an unknown function of the stationary distribution of $G(t)$. We now clarify this point in the
 31 theorem statement, and provide a simple example to illustrate this. We do not claim that M **needs** to be smaller than
 32 N . However, typically $diameter(G) \ll N$ ($O(\log N)$ for a random Erdős-Rényi G), so the memory complexity is
 33 low. We also claimed that $\mathcal{U}_n^M(t)$ can be much smaller than \mathcal{N} , meaning that at turn t , player n knows much less than
 34 the rewards of all players. This can happen even with $M = N$, if players broadcast only one reward value.

35 **Known Actions (R4):** Two practical examples where actions are never revealed: (1) Multi-agent Deep RL: like in
 36 GANs, the actions are the weights of a DNN that determine the “physical actions” of the agent (e.g. an autonomous
 37 vehicle), and the reward is a designed performance score over period of time. (2) Wireless networks: the action is
 38 the transmission power P_m , and device n measures the interference $\sum_m g_{m,n} P_m$ and cannot measure an individual
 39 action P_m . Nevertheless, our algorithm is appealing even with observable neighbors’ actions, since it’s unclear how
 40 one can efficiently use this information. First, each player needs to collect the entire action profile x_t from all over the
 41 network, for each t . A random subset will not do since $\sum_n u_n(x_t)$ is non-linear in x_t (but it is in $u_n(x_t)$, which is
 42 why our algorithms works). x_t are of higher dimensionality than the rewards so the communication overhead is larger.

43 **Convergence Time Analysis (R2,R3):** is highly complicated and requires novel techniques and proofs. Hence, we
 44 believe it should be pursued in a different paper. Our numerical experiments give some indication that the convergence
 45 time is reasonable, and behaves like SGD (the complicated part is the effect of G and the communication protocol).

46 **Minor Comments (R1,R2):** We resolved all minor comments (e.g., typos, using “optimal” more carefully).