

Noise-Tolerant Community Enforcement and the Strength of Small Stakes

Drew Fudenberg* Alexander Wolitzky, †

April 4, 2024

Abstract

We study community enforcement in a large population with noisy monitoring. We focus on equilibria in the prisoner’s dilemma that are *coordination-proof*, meaning that matched partners never play a Pareto-dominated Nash equilibrium in the one-shot game induced by the equilibrium continuation payoffs at their current histories. We show that a noise-tolerant version of contagion strategies is optimal among all coordination-proof equilibria. Welfare under tolerant contagion strategies decreases in the noise level and the gain from defection faster than welfare in a fixed partnership does. Thus, community enforcement has a comparative advantage in supporting “low-stakes” relationships.

*Department of Economics, MIT, drew.fudenberg@gmail.com

†Department of Economics, MIT, alexander.wolitzky@gmail.com

1 Introduction

Repeated game models of decentralized cooperation in large societies—“community enforcement”—have been used to explain cooperation in settings such as merchant coalitions (Milgrom, North, and Weingast, 1990, Greif, 1993), credit and risk-sharing (Klein, 1992, Karlan et al., 2009, Bhaskar and Thomas, 2019), cooperation in village economies (Jackson, Rodriguez-Barraquer, and Tan, 2012) and online markets (Friedman and Resnick, 2001, Tadelis, 2016). In all of these settings, in reality a partnership has a significant chance of failing even when both partners act in good faith. However, this feature—which we simply call *noise*—is largely absent from canonical community enforcement models: existing models are often robust to introducing a small amount of noise, but they are typically ill-suited to studying cooperation when noise is substantial and causes welfare to fall short of the first best. Consequently, existing models cannot assess how welfare under community enforcement compares to that under other social or institutional arrangements when noise is present, or how this comparison depends on the noise level and other parameters. This is an important shortcoming, because these comparisons influence which kinds of economic transactions are more likely to be mediated by community enforcement, rather than alternatives such as repeated interaction in a fixed partnership or small group.

This paper develops a simple model of community enforcement under noise. We consider the prisoner’s dilemma with random matching and perfectly complementary actions, where matched partners who cooperate obtain a success with probability $p < 1$, while success is impossible if either partner defects. We adapt the continuum-player model of Clark, Fudenberg, and Wolitzky (2021) (CFW) by specifying that each player observes their partner’s history of successes and failures, and that the population distribution of histories is in a steady state. Thus, while each partnership is subject to noise, the noise washes out in aggregate. We also follow CFW in focusing on equilibria that are “coordination-proof,” meaning that matched partners never play a Pareto-dominated Nash equilibrium in the one-shot game induced by the equilibrium continuation payoffs at their current histories.

Our first result is that a noise-tolerant version of the contagion strategies introduced by Kandori (1992) is optimal among all coordination-proof equilibria. (Our strategies differ from Kandori’s in that failure leads to punishment only probabilistically, by conditioning on the outcome of a randomization device.) When players’ time horizons are sufficiently long, welfare under these strategies is given by a simple formula, and is decreasing in both the noise level and the gain from defection. This gives a simple theory of how welfare under community enforcement depends on noise and the defection gain.

We then compare welfare under community enforcement with welfare when players interact in fixed partnerships (without rematching). Our second result is that, as either noise or the defection gain increases, welfare under community enforcement falls faster than welfare in fixed partnerships. Thus, community enforcement is less robust to noise (or to increases in the defection gain) than is cooperation in fixed partnerships. Intuitively, noise inevitably causes some players to switch from cooperation to defection—in either community interactions or fixed partnerships—but in community interactions contagion additionally causes defection to spread to some innocent players.

Our results speak to the classic question of how productive activities should be divided between small-scale groups, such as fixed partnerships, and larger communities or markets. The key tradeoff between these modes of production is thought to be that trust is easier to sustain in a fixed partnership, while wider interactions allow greater specialization and productive efficiency. Thus, a typical conclusion is that if agents are sufficiently forward-looking to sustain trust in community-wide interactions, then these interactions are more efficient than interactions in fixed partnerships; while if agents are more myopic then it is more efficient to retreat into fixed partnerships where trust is easier to sustain.¹ Our results instead imply that if community interactions have an exogenous productivity advantage over fixed partnerships (e.g., due to specialization), then overall social welfare is higher under community interactions if noise is sufficiently low and the defection gain is sufficiently small, and is higher under fixed partnerships otherwise. In particular, if we adopted the “variable stakes” framework of Ghosh and Ray (1996), Kranton (1996), or Ali and Miller (2016), where the

¹Arguments along these lines have been made by many scholars, including Putnam (1993), Greif (1994), Dixit (2003), Karlan et al. (2009), and Seabright (2010).

defection gain is relatively larger in higher-stakes relationships, then community interactions would have a comparative advantage in supporting low-stakes relationships.² This finding echoes the classic intuitions of Granovetter (1973) and Putnam (2000) that “weak ties” or “bridging social capital”—i.e., low-stakes but nonetheless valuable interactions, such as advice or job recommendation networks—are key benefits of community interactions.³

Related literature. We contribute to the community enforcement literature by developing a tractable model where maximum equilibrium welfare depends on the amount of noise. Classic community enforcement models like Kandori (1992) and Okuno-Fujiwara and Postlewaite (1995) exclude noise altogether, as do many subsequent papers including Dixit (2003), Karlan et al. (2009), Jackson, Rodriguez-Barraquer, and Tan (2012), and Wolitzky (2013). Ellison (1994), Takahashi (2010), Deb, Sugaya, and Wolitzky (2018), Heller and Mohlin (2018), and CFW establish folk theorems in the limit of vanishingly little noise.⁴ Bhaskar and Thomas (2019) consider a model with one-sided moral hazard, where there is no scope for contagion, and efficiency under community enforcement can be as high as in a fixed partnership. Finally, Clark, Fudenberg, and Wolitzky (2020) analyze the performance of a class of tolerant trigger strategies that are similar to the tolerant contagion strategies we consider. However, that paper has a different information structure (players only observe their current partner’s past actions), under which tolerant trigger strategies are suboptimal and efficiency is determined by stage-game strategic complementarity (which is irrelevant in the current paper), rather than noise.

²We do not formally introduce variable stakes or a productivity difference between community and partnership interactions in our model, as this point is straightforward given our results for the standard prisoner’s dilemma.

³Earlier models such as Dixit (2003), Karlan et al. (2009), Ali and Miller (2013), or Wolitzky (2013) capture related intuitions, but the logic of these models is very different because they do not involve noise.

⁴These papers make different population structure and observability assumptions. Our doubly-infinite time model with long player histories is most similar to Heller and Mohlin (2018) and CFW, as well as earlier papers on learning in games such as Fudenberg and Levine (1993) and Fudenberg and He (2018).

2 Cooperation in a Large Community

This section develops our model of cooperation in a large community and characterizes maximum welfare in this setting. Section 3 will then compare this welfare level with the maximum attainable welfare in a fixed partnership.

2.1 Matching and Pairwise Interactions

There is a unit mass of players, each of whom has a geometrically distributed lifespan with continuation probability $\gamma \in (0, 1)$, with exits balanced by an inflow of new entrants of size $1 - \gamma$. The time horizon is doubly infinite, so there is no fixed start date.

Each period, the players randomly match in pairs to play the prisoner's dilemma stage game, with expected payoffs

$$\begin{array}{cc}
 & C & D \\
 C & 1, 1 & -\ell, 1 + g \\
 D & 1 + g, -\ell & 0, 0
 \end{array} \tag{1}$$

where $g, \ell > 0$. The *public outcome* of each bilateral interaction is either S (a *success*) or F (a *failure*). We assume that the probability of a success when the partners take actions $a \in \{C, D\}^2$ is given by

$$\Pr(S|a) = \begin{cases} p & \text{if } a = (C, C), \\ 0 & \text{otherwise,} \end{cases} \tag{2}$$

where $0 < p < 1$. Thus, the public outcome can only be a success if both partners cooperate, but success is never assured. As we explain in Section 4, our results extend to the case where $\Pr(S|a) = q$ for all $a \neq (C, C)$, for sufficiently small $q > 0$. In addition to generating a public outcome, each bilateral interaction also generates an independent Uniform $[0, 1]$ random variable z , which, like the public outcome, is publicly observed at the end of the period. These additional random variables, which we call *pairwise randomizations*, amount to having a separate public randomizing device for each matched pair. We discuss the interpretation and role of these randomizations in Section 2.5.

When two players meet, each observes the other's entire history of past outcomes and random draws (and no further information), which is an element of $H = \emptyset \cup \bigcup_{t=1}^{\infty} (\{S, F\} \times [0, 1])^t$,

where a new entrant has the null history \emptyset , and an experienced player’s history is the record of their past outcomes along with the corresponding values of z . Players also recall their own histories. Thus, a *pure strategy* is a function $\sigma : H \times H \rightarrow \{C, D\}$, with the convention that the first coordinate is a player’s own history and the second is the opponent’s history.⁵ We restrict attention to symmetric pure strategy profiles, where each player uses the same pure strategy.⁶ We henceforth omit the qualifiers “symmetric” and “pure” without further comment, and use the same notation σ for an individual strategy and a (symmetric) strategy profile.

2.2 Aggregate Behavior and Equilibrium

The *state* $\mu \in \Delta(H)$ of the community describes the share of players with each possible history.⁷ Since there is a continuum of players, under any strategy profile σ the state evolves according to a deterministic transition function $f_\sigma : \Delta(H) \rightarrow \Delta(H)$. A *steady state* under σ is a state μ such that $f_\sigma(\mu) = \mu$.⁸ One can then define the expected continuation payoff $(1 - \gamma) \sum_t \gamma^t \mathbb{E}[u_t]$ of a player with any history $h \in H$ at a steady state μ , when all other players in the population follow σ and the player under consideration plays an arbitrary strategy σ' . An *equilibrium* is a pair (σ, μ) such that μ is a steady state under σ and, for any history h , the expected continuation payoff of a player with history h is maximized by taking $\sigma' = \sigma$. Finally, *welfare* at an equilibrium (σ, μ) is defined as the average payoff in the population at state μ under strategy σ . Note that welfare is equal to the expected lifetime payoff of a new entrant, because the steady state distribution μ also describes the fraction of periods in which an entrant expects to have each history.

⁵In principle, a player could condition on their own past actions in addition to their history of successes and failures, but as we allow public randomization there is no benefit to doing so. Since H is a continuum due to the random variables z , we also formally require that strategies are measurable functions on $H \times H$, where H is endowed with the weak* topology.

⁶Restricting to pure equilibria simplifies the analysis and also captures a type of robustness. We could alternatively further restrict attention to strict equilibria. This would yield almost the same analysis except for some technicalities resulting from the fact that the set of strict equilibria is not closed. In contrast, mixed strategies would allow (C, D) and (D, C) to be played on the equilibrium path, which as discussed below would greatly complicate the analysis.

⁷Formally, since H is a continuum, the state describes the measure of players with each measurable set of histories.

⁸Section 2.5 discusses the details involved in constructing the transition function and the issue of the existence of a steady state.

A preliminary observation is that in any (pure) equilibrium, only (C, C) and (D, D) are played along the equilibrium path. To see this, observe that equation (2) implies that when the opponent defects, the probability of success is independent of a player’s own action. Since future opponents will observe only whether the current outcome is a success or a failure (as well as the uniform variable z), and D is dominant in the stage game, this implies that D is the unique best response of a player who anticipates that their opponent will play D . Thus, only (C, C) and (D, D) can be played along the equilibrium path.

The fact that only (C, C) and (D, D) are played on path implies that equilibrium incentives can be provided only through surplus creation and destruction—switching from (D, D) to (C, C) or vice versa—rather than surplus transfers between matched partners, which would occasionally require (C, D) or (D, C) to be played. This feature greatly simplifies the analysis, as well as ensuring that first-best efficiency is unattainable regardless of the players’ expected lifespans.⁹

2.3 Coordination-Proof Equilibria

We say that an equilibrium is *coordination-proof* if matched partners never play a Pareto-dominated Nash equilibrium in the one-shot game induced by the equilibrium continuation payoffs at their current histories.¹⁰

We restrict attention to coordination-proof strategies throughout our analysis. The motivation for imposing this refinement is that an equilibrium that is not coordination-proof would break down if a pair of matched partners could manage to coordinate on the efficient equilibrium in their interaction (taking behavior in the rest of the population as given). The following lemma describes the key implication of coordination-proofness for our analysis (indeed, the only implication of coordination-proofness that we will use).¹¹

Lemma 1. *In any coordination-proof equilibrium (σ, μ) , the set of all histories H can be partitioned into two sets, H^C and H^D , such that:*

⁹Our results would be similar if $\Pr(S|(C, D)) \neq \Pr(S|(D, D))$ but ℓ is sufficiently large, as then too only (C, C) and (D, D) are played on path.

¹⁰See Definition 4 in CFW for the formal definition of coordination-proofness.

¹¹The lemma relies on our assumption that actions are perfect complements: $\Pr(S|(C, D)) = \Pr(S|(D, D))$. Without this assumption, more equilibria can be coordination-proof, including ones where (C, D) is played on the equilibrium path.

1. *Players with histories in H^C cooperate against opponents with histories in H^C : $\sigma(h, h') = C$ for all $h, h' \in H^C$.*
2. *Players with histories in H^C defect against opponents with histories in H^D : $\sigma(h, h') = D$ for all $h \in H^C, h' \in H^D$.*
3. *Players with histories in H^D defect against all opponents: $\sigma(h, h') = D$ for all $h \in H^D, h' \in H$.*

Proof. Fix a coordination-proof equilibrium (σ, μ) , and define $H^C = \{h \in H : \exists h' \in H \text{ s.t. } \sigma(h, h') = C\}$ and $H^D = H \setminus H^C$. By definition, $\sigma(h, h') = D$ for all $h \in H^D, h' \in H$. Moreover, since C is never a best response against an opponent who plays D (as $\Pr(S|(C, D)) = \Pr(S|(D, D))$ and D is dominant in the stage game), we have $\sigma(h, h') = D$ for all $h \in H^C, h' \in H^D$.

It remains to show that $\sigma(h, h') = C$ for all $h, h' \in H^C$. Fix any $h, h' \in H^C$. Since C is never a best response against an opponent who plays D , the fact that $h, h' \in H^C$ implies that for a player with history h or h' , C is a best response against an opponent who plays C . Hence, when players with histories h and h' meet each other, both (C, C) and (D, D) are equilibria in the induced one-shot game. Next, since C is sometimes a best response for these players even though D is dominant in the stage game, their expected continuation payoffs must each be higher when the outcome of their match is S rather than F . This implies that the (C, C) equilibrium yields higher expected continuation payoffs for both players (since the probability of S is higher) as well as higher stage game payoffs for both players, relative to the (D, D) equilibrium. Therefore, coordination-proofness requires that $\sigma(h, h') = \sigma(h', h) = C$. ■

Given Lemma 1, a coordination-proof equilibrium is entirely described by the partition $\{H^C, H^D\}$. Henceforth, given a coordination-proof equilibrium, we simply refer to players with histories in H^C as *cooperators*, and to players with histories in H^D as *defectors*.

Our first main result is a bound on the payoff of any coordination-proof equilibrium.¹² We will subsequently show that when γ is sufficiently large this bound is tight, and is attained by a version of contagion strategies. Hence, the bound \bar{W} derived in the proposition will be key for comparing welfare in communities and partnerships in Section 3.

¹²The proof of this results builds on the proof of Lemma 11 of Clark, Fudenberg, and Wolitzky (2020).

Proposition 1. *For any continuation probability γ , welfare in any coordination-proof equilibrium is bounded above by $\bar{W} = \bar{\mu}^2$, where*

$$\bar{\mu} = \begin{cases} \frac{1+g}{2} + \sqrt{\left(\frac{1+g}{2}\right)^2 - \frac{g}{p}} & \text{if } p \geq \frac{4g}{(1+g)^2}, \\ 0 & \text{otherwise.} \end{cases}$$

Note that the condition $p \geq 4g/(1+g)^2$ can only be satisfied if $g \leq 1$, regardless of p . Thus, no cooperation is possible in a coordination-proof community equilibrium if $g > 1$. Moreover, as noise increases, cooperation becomes impossible even for smaller values of g : for example, if $p = 8/9$ then cooperation is impossible whenever $g \geq 1/2$.

Proof. Fix a coordination-proof equilibrium with partition $\{H^C, H^D\}$. Let $\mu^C = \int_{H^C} d\mu$ denote the share of cooperators. Suppose that $\mu^C > 0$.

For any $h \in H$, let $V(h)$ denote the expected continuation payoff of a player with history h . Note that $V(h) \geq 0$ for all $h \in H$, since a player's minmax payoff is 0. Next, let $\bar{V} = \sup_{h \in H} V(h)$. Note that $\bar{V} > 0$, since $V(h) \geq 0$ for all h and $\mu^C > 0$. We claim that $\bar{V} = \sup_{h \in H^C} V(h)$. Otherwise, there would exist $h \in H^D$ such that $V(h) > \sup_{h' \in H^C} V(h')$, but since all defectors obtain a stage game payoff of 0, we have $V(h) \leq \gamma \sup_{h' \in H^C} V(h')$.

Now consider a cooperator with history $h \in H^C$. Let $V(h, S) = \mathbb{E}_z[V(h, S, z)]$ and $V(h, F) = \mathbb{E}_z[V(h, F, z)]$ denote this player's expected continuation payoff when their current-period outcome is a success or a failure, respectively. Since this player cooperates against opponents in H^C and defects against opponents in H^D , we have

$$V(h) = (1 - \gamma)\mu^C + \gamma(p\mu^C V(h, S) + (1 - p\mu^C)V(h, F)).$$

At the same time, since the player prefers to play C against an opponent who plays C , we have

$$\gamma p(V(h, S) - V(h, F)) \geq (1 - \gamma)g.$$

Combining these inequalities, we have

$$\frac{p}{1 - p\mu^C} \left(\mu^C - V(h) + \frac{\gamma}{1 - \gamma} (V(h, S) - V(h)) \right) \geq g.$$

This inequality holds for all $h \in H^C$ and $\bar{V} = \sup_{h \in H^C} V(h) \geq \sup_{h \in H^C} V(h, S)$, so

$$\frac{p}{1 - p\mu^C} (\mu^C - \bar{V}) \geq g.$$

Moreover, the expected lifetime payoff of a new entrant equals $(\mu^C)^2$ (the share of matches that cooperate), so $\bar{V} \geq (\mu^C)^2$. Hence, we have

$$\begin{aligned} \frac{p\mu^C(1 - \mu^C)}{1 - p\mu^C} &\geq g && \iff \\ (\mu^C)^2 - (1 + g)\mu^C + \frac{g}{p} &\leq 0. \end{aligned}$$

This implies that $\mu^C \leq \bar{\mu}$. Since welfare equals $(\mu^C)^2$, we conclude that welfare is bounded above by $\bar{\mu}^2$. ■

2.4 Tolerant Contagion Strategies

We will show that the following class of coordination-proof strategies attains the welfare bound \bar{W} .

Definition 1. *In a tolerant contagion strategy profile, there is a parameter $\phi \in (0, 1)$ such that:*

1. *New entrants are cooperators: $\emptyset \in H^C$.*
2. *If the outcome of a cooperator's interaction is (S, z) for any z , they remain a cooperator: if $h \in H^C$ then $h \times (S, z) \in H^C$.*
3. *If the outcome of a cooperator's interaction is (F, z) , they remain a cooperator if $z \geq \phi$, and become a defector if $z < \phi$: if $h \in H^C$ then $h \times (F, z) \in H^C$ for $z \geq \phi$, and $h \times (F, z) \in H^D$ for $z < \phi$.*
4. *A defector remains a defector forever.*

Note that these strategies utilize our assumption that a player's past values of z are observed by all of their partners. We refer to ϕ as the *transition probability*.

The next lemma characterizes the equilibrium conditions and the steady state share of cooperators under tolerant contagion strategies.

Lemma 2. *Under a tolerant contagion strategy profile with transition probability ϕ , an equilibrium with cooperator share μ^C and cooperator payoff $V > 0$ exists if and only if*

$$V = (\mu^C)^2, \quad (\text{PK})$$

$$\gamma p \phi V \geq (1 - \gamma) g, \quad \text{and} \quad (\text{IC})$$

$$\mu^C = 1 - \gamma + \gamma \mu^C (1 - \phi + p \mu^C \phi). \quad (\text{SS})$$

Moreover, the steady state share of cooperators is unique.

Proof. (PK) (“promise keeping”) says that the cooperator payoff equals $(\mu^C)^2$. This is necessary because entrants are cooperators, and an entrant’s payoff equals social welfare, $(\mu^C)^2$. (IC) is the incentive constraint, which is necessary because deviating to D against a cooperator yields a gain of $(1 - \gamma) g$, but increases the probability of a failure by p , which implies an expected future loss of $\gamma p \phi V$. (SS) is the steady state condition, which is necessary because in a steady state, the share of cooperators μ^C must equal the share of new entrants $1 - \gamma$ (who are all cooperators) plus the share of surviving cooperators $\gamma \mu^C$ who remain cooperators (as all surviving defectors remain defectors), and this latter share is equal to $1 - \phi$ (the share of surviving cooperators who obtain outcomes with $z > \phi$), plus $p \mu^C \phi$ (the share of surviving cooperators who obtain outcomes (S, z) with $z < \phi$). Conversely, when all three conditions are satisfied, (σ, μ) is an equilibrium. Finally, the steady-state equation (SS) is quadratic in μ^C , and only the smaller solution is in the required $[0, 1]$ range.¹³ ■

Now we show that tolerant contagion strategies attain the maximum welfare level \bar{W} , whenever γ is sufficiently large.

Proposition 2. *For any $\gamma \geq \bar{\gamma}$, there exists a tolerant contagion equilibrium that yields welfare \bar{W} , where $\bar{\gamma} = (1 + p\bar{\mu}^2/g)^{-1} \in (0, 1)$.*

¹³In particular, the unique steady state share of cooperators is

$$\frac{1 - \gamma + \gamma \phi - \sqrt{(1 - \gamma - \gamma \phi)^2 - 4\gamma(1 - \gamma)p\phi}}{2\gamma p \phi}.$$

Proof. Consider tolerant contagion strategies with transition probability

$$\phi = \frac{(1 - \gamma)g}{\gamma p \bar{\mu}^2}.$$

Note that $\phi \leq 1$ (so the strategy profile is well-defined) iff $\gamma \geq \bar{\gamma}$. Substituting this value of ϕ into (SS) gives

$$\begin{aligned} \mu^C &= 1 - \gamma + \gamma \mu^C \left(1 - \frac{(1 - \gamma)g}{\gamma p \bar{\mu}^2} + \frac{(1 - \gamma)g}{\gamma p \bar{\mu}^2} p \mu^C \right) && \iff \\ \mu^C \bar{\mu} - \bar{\mu} - g \frac{(\mu^C)^2}{\bar{\mu}} + \frac{g \mu^C}{p \bar{\mu}} &= 0. \end{aligned}$$

Observe that $\mu^C = \bar{\mu}$ solves this equation. Thus, μ^C is the steady state share of cooperators under tolerant contagion strategies with transition probability ϕ .¹⁴ Moreover, steady state welfare equals $\bar{W} = \bar{\mu}^2$ by (PK), and (IC) holds with equality by construction of ϕ . Hence, the steady state corresponds to an equilibrium with the desired properties. ■

We henceforth assume that $\gamma \geq \bar{\gamma}$, so that maximum welfare under coordination-proof community enforcement is \bar{W} .

2.5 Interpretation and Technical Details

Here we provide an interpretation of pairwise randomizations, and discuss some technical details that we deferred above.

Interpretation of pairwise randomization. The role of pairwise randomizations is to introduce some noise tolerance into contagion strategies, by reducing the probability that failure causes players to switch to defection. A possible interpretation of tolerant contagion strategies is that these randomizations determine whether a failure leads to a “dispute,” where society remembers which players have been involved in disputes (but not necessarily the precise values of the associated random z ’s). It is also possible to dispense with randomizing devices altogether, at the cost of some additional complexity. For example, we could consider strategies with several cooperative states, where players become defectors only after

¹⁴The steady state is unique by Lemma 2, although the current proof does not require this fact.

experiencing some number $K > 1$ of failures. We analyzed such “GrimK” strategies in Clark, Fudenberg, and Wolitzky (2020). In the current setting, we conjecture that GrimK strategies are approximately as efficient as tolerant contagion strategies when γ is sufficiently large. However, they cannot exactly attain efficiency for any fixed γ due to the constraint that K must be an integer, and they are harder to analyze because we have to keep track of the share of players with each number of failures. Allowing pairwise randomizations thus considerably simplifies the analysis, and we believe it does not substantially affect the results.

Definition of the update map; existence of a steady state. CFW give the equation for the update map and establish existence of a steady state in a model without pairwise randomizations where players observe their partner’s “record,” which can include additional information such as the past actions and records of the current partner’s past partners. The update map and existence proof generalize to pairwise randomizations, but the notation is somewhat complicated and the existence of a steady state involves some measurability issues. However, in the current paper the only behaviorally relevant aspect of the update map is the update rule for the share of cooperators, and the only relevant aspect of a steady state is the steady state share of cooperators. These are both given by the steady-state equation (SS).

Observability of own payoffs. The payoffs given by equation (1) cannot be written as an expectation over the outcomes S and F with probabilities given by (2) of a utility function that depends only a player’s own action and the outcome. Thus, to interpret the model as one where players observe their own payoffs, we must also let each player privately observe their own payoff. Adding this information does not affect the analysis in the current section, because players can never gain by conditioning on it. Adding such information to the fixed-partnership repeated game considered in the next section could expand the set of all Nash equilibria, but not the set of *perfect public equilibria* (Fudenberg, Levine, and Maskin, 1994), where players condition only on public signals. Thus, when players observe their own payoffs, the analysis in the next section remains valid for PPE.¹⁵

¹⁵In addition, if we restrict attention to strict equilibria and assume that players’ additional private signals are independent conditional on actions and the public signal, then focusing on PPE is without loss

3 Comparison of Community and Partnership

We now compare \bar{W} with the maximum welfare that can be attained in a fixed partnership. Suppose the prisoner's dilemma stage game (1) is played repeatedly by a fixed pair of players with a fixed start date $t = 1$, with outcome distribution as in (2) and public randomization, with discount factor $\delta \in (0, 1)$. As in our analysis of cooperation in a large community, we restrict attention to pure-strategy equilibria. In place of tolerant contagion strategies, we now show that the upper bound on welfare can be achieved using *tolerant grim trigger strategies*.

Definition 2. *In a tolerant grim trigger strategy profile, there is a parameter $\phi \in (0, 1)$ such that:*

1. *The players cooperate in period 1.*
2. *If the players cooperate in period t and the outcome is (S, z) for any z , they continue cooperating in period $t + 1$.*
3. *If the players cooperate in period t and the outcome is (F, z) , they continue cooperating in period $t + 1$ if $z \geq \phi$, and switch to defecting if $z < \phi$.*
4. *If the players defect in period t , they continue defecting forever.*

The next result (which is standard) characterizes maximum welfare in a fixed partnership, and shows that it is attained by tolerant grim trigger strategies.

Proposition 3. *In a fixed partnership, for any discount factor δ , welfare in any pure-strategy Nash equilibrium is bounded above by*

$$\hat{W} = \max \left\{ 1 - \frac{1-p}{p}g, 0 \right\}.$$

Moreover, whenever $\hat{W} > 0$, for any $\delta \geq \hat{\delta}$, there exists a tolerant grim trigger equilibrium that yields welfare \hat{W} , where $\hat{\delta} = g / ((1 + g)p) \in (0, 1)$.

of generality, because players never have a strict incentive to condition on their private signals. Our analysis is unchanged under a restriction to strict equilibria, except that strict tolerant contagion equilibria can only approximate the maximum welfare level \bar{W} rather than exactly attaining it.

Proof. As in the community enforcement game considered above, only (C, C) and (D, D) are played on-path in any pure Nash equilibrium, so every pure Nash equilibrium is *strongly symmetric*. By standard arguments, in the optimal strongly symmetric equilibrium, continuation payoffs at every history are either some $V \geq 0$ or 0. When public randomizations are available, equilibria of this form are precisely tolerant grim trigger equilibria. It therefore remains to characterize the optimal tolerant grim trigger equilibria. There is a tolerant grim trigger equilibrium with transition probability $\phi \in [0, 1]$ and payoff V iff

$$V = 1 - \delta + \delta(1 - \phi + p\phi)V \quad \text{and} \quad (\text{PK}')$$

$$\delta p\phi V \geq (1 - \delta)g. \quad (\text{IC}')$$

Taking ϕ to satisfy (IC') with equality and substituting into (PK'), we have

$$V = 1 - \frac{1-p}{p}g.$$

Substituting for V in (IC') shows that the required value of ϕ is less than 1 (so the strategy profile is well-defined) iff $(1-p)(1+g) \leq 1$ (i.e., $\hat{W} > 0$) and $\delta \geq \hat{\delta}$. ■

With Propositions 1 and 3 in hand, we can now compare maximum welfare under community interactions and fixed partnerships.

Proposition 4. *Assume that $p > 4g/(1+g)^2$, so that $\bar{W} > 0$. Then we have:*

1. $\bar{W} < \hat{W}$.
2. $\partial\bar{W}/\partial p > \partial\hat{W}/\partial p$.
3. $\partial\bar{W}/\partial g < \partial\hat{W}/\partial g$.

Thus, welfare is lower in community interactions than in fixed partnerships, and also decreases faster as noise increases (i.e., p decreases) and the defection gain g increases. The intuition for why $\bar{W} < \hat{W}$ is that, since a matched cooperator and defector take (D, D) in community interactions, and some players inevitably become defectors when their interactions are “hit by noise” (i.e., when a partnership fails despite mutual effort), there is a certain

unavoidable amount of contagion, as even players whose own interactions are never hit by noise become defectors as a result of matching with other players who were hit by noise. In contrast, in a fixed partnership, the partners only switch to defection when they themselves are hit by noise. A similar intuition explains why $\partial\bar{W}/\partial p > \partial\hat{W}/\partial p$ and $\partial\bar{W}/\partial g < \partial\hat{W}/\partial g$. Noise is more harmful in community interactions, because it affects both players whom it directly hits and players who match with players whom it hits. Finally, the defection gain g determines how much future cooperation must be lost when noise hits to preserve incentives, so increasing g has a similar effect as increasing noise.

Proof. Observe that

$$\frac{\partial\bar{W}}{\partial p} = 2\bar{\mu}\frac{\partial\bar{\mu}}{\partial p} = \left(\frac{1+g}{2\sqrt{\left(\frac{1+g}{2}\right)^2 - \frac{g}{p}} + 1} \right) \frac{g}{p^2} > \frac{g}{p^2} = \frac{\partial\hat{W}}{\partial p},$$

where the inequality uses $p > 4g/(1+g)^2$. Moreover, if $p = 1$ (contrary to our assumptions) then $\bar{\mu} = \bar{W} = \hat{W} = 1$. Since $p \in (0, 1)$, the first two parts of the proposition follow.

For the third part of the proposition, observe that

$$\begin{aligned} \frac{\partial\bar{W}}{\partial g} &= 2\bar{\mu}\frac{\partial\bar{\mu}}{\partial g} = \bar{\mu} \left(1 + \frac{\frac{1+g}{2} - \frac{1}{p}}{\sqrt{\left(\frac{1+g}{2}\right)^2 - \frac{g}{p}}} \right) = \frac{\bar{\mu} \left(\bar{\mu} - \frac{1}{p} \right)}{\bar{\mu} - \frac{1+g}{2}}, \quad \text{and} \\ \frac{\partial\hat{W}}{\partial g} &= -\frac{1-p}{p}. \end{aligned}$$

Hence, $\partial\bar{W}/\partial g < \partial\hat{W}/\partial g$ iff

$$\begin{aligned} \bar{\mu} \left(\frac{1}{p} - \bar{\mu} \right) &> \frac{1-p}{p} \left(\bar{\mu} - \frac{1+g}{2} \right) && \iff \\ p\bar{\mu}(1-\bar{\mu}) &> -(1-p)\frac{1+g}{2}, \end{aligned}$$

which holds as the LHS is positive and the RHS is negative. ■

Proposition 4 implies that if community interactions have an exogenous productivity advantage over interactions in fixed partnerships, then society should allocate activities involving low noise and low defection gains to community interactions, while allocating activities

involving high noise and high defection gains to fixed partnerships. (Absent a productivity edge for community interactions, all production should take place in fixed partnerships.) Intuitively, community interactions have a comparative advantage in activities with low noise and low defection gains, because the failures that inevitably occur in the presence of noise are more harmful in community interactions than in fixed partnerships, as they necessarily trigger some degree of community-wide contagion. As a consequence, if we adopt the standard assumption that the defection gain is relatively higher in high-stakes activities, we can conclude that community interactions have a comparative advantage in supporting low-stakes activities.¹⁶

4 Discussion

We conclude by discussing some possible variants and extensions.

Non-Coordination-Proof Equilibria. Our comparison of welfare in communities and partnerships relies on focusing on coordination-proof equilibria in community interactions. We have argued that this equilibrium refinement is reasonable, since equilibria built on within-match miscoordination are arguably fragile. We have also shown that tolerant contagion strategies are optimal under this refinement, so the refinement also supports focusing on this class of strategies, which generalize the usual contagion strategies that are a centerpiece of the community enforcement literature.

Nonetheless, it is worth noting that, at least for some parameters p and g , the maximum partnership welfare level of \hat{W} can be attained in community interactions using non-coordination-proof strategies. In particular, consider the following variant of tolerant contagion strategies.

1. The set of histories H is partitioned into cooperators and defectors. Entrants are cooperators.

¹⁶We are not aware of direct empirical evidence on whether the defection gain rises faster than the cooperative payoff in higher-stakes activities, but this is a standard assumption: see, e.g., Ghosh and Ray (1996), Kranton (1996), Ali and Miller (2013, 2016). Another similar setup is Dixit (2003), which considers a random matching model where information about high-value interactions spreads more slowly, so that incentive constraints bind more in higher-value interactions.

2. Matched partners always cooperate, unless they are both defectors, in which case they both defect. (This is a key difference from coordination-proof strategies, where a matched cooperator and defector both defect.)
3. If the outcome of a cooperator's interaction is (S, z) for any z , they remain a cooperator. If the outcome of a cooperator's interaction is (F, z) , they remain a cooperator if $z \geq 1 - \phi$, and become a defector if $z < \phi$.
4. If the outcome of a defector's interaction is (F, z) for any z , they remain a defector. If the outcome of a defector's interaction is (S, z) , they remain a defector if $z \geq 1 - \phi$, and become a cooperator if $z < \phi$.

These strategies are not coordination-proof, because (C, C) is a Pareto-dominant Nash equilibrium in an interaction between two defectors, but the strategies prescribe (D, D) in these interactions. Nonetheless, it can be shown that if g is sufficiently small, there exists $\tilde{\gamma}$ such that, for all $\gamma > \tilde{\gamma}$, there exists a value for ϕ such that these strategies form an equilibrium that delivers welfare \hat{W} .¹⁷ This shows that coordination-proofness is an essential part of our theory. Intuitively, coordination-proofness implies that a matched cooperator and defector must take (D, D) , as if they took (C, C) then matched defectors would also take (C, C) (as (C, C) would be a Pareto-dominant Nash equilibrium in their interaction), which in turn would destroy incentives. Matched cooperators and defectors taking (D, D) is the source of contagion discussed following Proposition 4, which accounts for the gap between welfare in communities and partnerships. If instead matched cooperators and defectors take (C, C) as in the above strategies, there is no contagion, and welfare in communities and partnerships is equal.

Two-Sided Noise. Our assumption that success requires cooperation from both partners is restrictive. A natural generalization is to assume that

$$\Pr(S|a) = \begin{cases} p & \text{if } a = (C, C), \\ q & \text{otherwise,} \end{cases} \quad (3)$$

¹⁷It can also be shown that, as in a fixed partnership, welfare in any pure Nash equilibrium in the community enforcement game cannot exceed \hat{W} .

where $0 < q < p < 1$. Under this assumption, the partners' actions remain perfect complements in delivering a success, in that $\Pr(S|(C, D)) = \Pr(S|(D, D))$, but it is possible to obtain a success even if one or both partners defect. Under (3), the same argument as in the proof of Proposition 1 gives an upper bound for welfare that converges to $\bar{\mu}^2$ as $q \rightarrow 0$. Similarly, the same argument as in the proof of Proposition 2 implies that, for sufficiently high γ , there exist a sequence of tolerant contagion equilibria whose welfare converges to $\bar{\mu}^2$ as $q \rightarrow 0$, but now there is a gap between the upper bound and the welfare level that can be attained by tolerant contagion strategies.¹⁸ This gap can be shown to be second-order in q , so tolerant contagion strategies are quite robust to a small probability that the players obtain a success even when one or both of them defects.

The intuition for why tolerant contagion strategies are exactly optimal when $q = 0$, but not when $q > 0$, is as follows. When $q = 0$, obtaining a success proves that both partners cooperated. It is therefore optimal to assign each partner the highest possible continuation payoff following a success, which is what tolerant contagion strategies do. Instead, when $q > 0$, a success may be due to luck. It may therefore be better to assign a player the highest possible continuation payoff only after a series of successes. However, the advantage of such strategies over tolerant contagion strategies is small when q is small.

More General Outcome Distributions. It is more challenging to generalize the outcome structure to allow $\Pr(S|(C, D)) \neq \Pr(S|(D, D))$, so that the partners' actions are not perfect complements. If $\Pr(S|(C, D)) \neq \Pr(S|(D, D))$ and the parameter ℓ is sufficiently small, then we conjecture that the asymmetric action profiles (C, D) and (D, C) can be supported in equilibrium. This enables society to provide incentives through continuation payoff transfers (as in Fudenberg, Levine, and Maskin, 1994), even though with only two outcomes pairwise full rank cannot hold. Moreover, because players face a new partner each period, there can be a positive mass of players whose continuation payoffs are above the

¹⁸The upper bound and the maximum welfare level attainable by tolerant contagion strategies are, respectively,

$$\frac{1+g}{2} + \sqrt{\left(\frac{1+g}{2}\right)^2 - \frac{1-q}{p-q}g} \quad \text{and} \quad \frac{1 + \frac{p}{p-q}g}{2} + \sqrt{\left(\frac{1 + \frac{p}{p-q}g}{2}\right)^2 - \frac{1}{p-q}g}.$$

highest feasible and individually rational payoff in a fixed partnership. This in turn raises the question of what general properties of the outcome structure suffice for the folk theorem in our steady-state community interaction model, either with or without a restriction to coordination proof equilibrium. This is an interesting question for future work.¹⁹

Avoidance. A modification of the prisoner’s dilemma stage game that seems realistic in some contexts is that in every interaction each player has third action, *Avoid*, which corresponds to refusing to interact with the current partner. Suppose that *Avoid* gives both players a payoff of zero regardless of their partner’s action, and that the outcome of an interaction where either partner takes *Avoid* is recorded as *Avoided*, or alternatively is not recorded at all. It can be shown that in this modified game, when γ is sufficiently high, maximum welfare in (coordination-proof) community enforcement increases to \hat{W}^2 . To see the intuition, note that with avoidance, we can modify tolerant contagion strategies so that players with histories in H^D play *Avoid* rather than D , while players still transition from H^C to H^D only after a failure. This modification eliminates contagion: now, players only transition from H^C to H^D when they themselves are hit by noise, because when they match with partners in H^D the outcome is now *Avoided* rather than F . Consequently, the transition probability from H^C to H^D can be set at the minimum level required to provide incentives, which implies that the steady-state share of players with histories in H^C can be as high as \hat{W} , exactly as in a fixed partnership. Nonetheless, the maximum welfare level of \hat{W}^2 is still less than the welfare level of \hat{W} that is attainable in a fixed partnership, because, since noise and matching are independent of each other in community interactions, the share of matches where both partners have histories in H^C is only \hat{W}^2 .

More General Records. We have assumed that players observe no information other than their partner’s history of past outcomes and random draws. It is also interesting to consider settings where a player gets some information about their partner’s past partners, such as their identities or past outcomes. In CFW, we showed that such “interdependent

¹⁹Theorem 2 of CFW provides a partial folk theorem under almost-perfect monitoring.

records” provide little advantage beyond observing the partner’s history of outcomes, when outcomes are observed with very little noise. However, with non-trivial noise, interdependent records might help reduce contagion and improve efficiency. This seems an interesting direction for future research.

Richer Social Structure. Finally, we have only compared two extreme social structures: random matching in a large population, and a fixed partnership. Reality lies in between these extremes: we interact with friends and colleagues more than with strangers, but we do sometimes meet strangers. Incorporating richer social structures—e.g., weighted random matching or networked interactions—into our analysis is another promising direction for future work.

References

- Ali, S. N. and D. A. Miller (2013). “Enforcing Cooperation in Networked Societies.” *Working Paper*.
- (Aug. 2016). “Ostracism and Forgiveness.” *American Economic Review*, 106, 2329–48.
- Bhaskar, V. and C. Thomas (2019). “Community Enforcement of Trust with Bounded Memory.” *Review of Economic Studies*, 86, 1010–1032.
- Clark, D., D. Fudenberg, and A. Wolitzky (2020). “Indirect reciprocity with simple records.” *Proceedings of the National Academy of Sciences*, 117, 11344–11349.
- (2021). “Record-Keeping and Cooperation in Large Societies.” *The Review of Economic Studies*, 88, 2179–2209.
- Deb, J., T. Sugaya, and A. Wolitzky (2018). “The Folk Theorem in Repeated Games with Anonymous Random Matching.” *mimeo*.
- Dixit, A. (2003). “Trade Expansion and Contract Enforcement.” *Journal of Political Economy*, 111, 1293–1317. (Visited on 10/12/2023).
- Ellison, G. (1994). “Cooperation in the prisoner’s dilemma with anonymous random matching.” *The Review of Economic Studies*, 61, 567–588.
- Friedman, E. J. and P. Resnick (2001). “The Social Cost of Cheap Pseudonyms.” *Journal of Economics & Management Strategy*, 10, 173–199.
- Fudenberg, D. and K. He (2018). “Learning and Type Compatibility in Signaling Games.” *Econometrica*, 86, 1215–1255.
- Fudenberg, D., D. Levine, and E. Maskin (1994). “The Folk Theorem with Imperfect Public Information.” *Econometrica*, 62, 997–1039.

- Fudenberg, D. and D. K. Levine (1993). “Steady state learning and Nash equilibrium.” *Econometrica*, 547–573.
- Ghosh, P. and D. Ray (1996). “Cooperation in Community Interaction Without Information Flows.” *Review of Economic Studies*, 63, 491–519.
- Granovetter, M. (1973). “The Strength of Weak Ties.” *American Journal of Sociology*, 78, 1360–1380.
- Greif, A. (1993). “Contract Enforceability and Economic Institutions in Early Trade: The Maghribi Traders’ Coalition.” *The American Economic Review*, 83, 525–548.
- (1994). “Cultural Beliefs and the Organization of Society: A Historical and Theoretical Reflection on Collectivist and Individualist Societies.” *Journal of Political Economy*, 83, 912–950.
- Heller, Y. and E. Mohlin (2018). “Social learning and the shadow of the past.” *Journal of Economic Theory*, 177, 426–460. URL: <http://www.sciencedirect.com/science/article/pii/S0022053118303533>.
- Jackson, M., T. Rodriguez-Barraquer, and X. Tan (2012). “Social Capital and Social Quilts: Network Patterns of Favor Exchange.” *American Economic Review*, 102, 1857–1897.
- Kandori, M. (1992). “Social Norms and Community Enforcement.” *The Review of Economic Studies*, 59, 63–80.
- Karlan, D. et al. (2009). “Trust and Social Collateral.” *Quarterly Journal of Economics*, 124, 1307–1361.
- Klein, D. B. (1992). “Promise Keeping in the Great Society: A Model of Credit Information Sharing.” *Economics & Politics*, 4, 117–136.
- Kranton, R. (1996). “The Formation of Cooperative Relationships.” *Journal of Law, Economics, and Organization*, 12, 314–233.
- Milgrom, P. R., D. C. North, and B. R. Weingast (1990). “The Role of Institutions in the Revival of Trade: The Law Merchant, Private Judges, and the Champagne Fairs.” *Economics & Politics*, 2, 1–23.
- Okuno-Fujiwara, M. and A. Postlewaite (1995). “Social Norms and Random Matching Games.” *Games and Economic Behavior*, 9, 79–109.
- Putnam, R. (1993). *Making Democracy Work: Civic Traditions in Modern Italy*. Princeton University Press.
- (2000). *Bowling Alone*. Simon & Schuster.
- Seabright, P. (2010). *The Company of Strangers: A Natural History of Economic Life*. Princeton University Press.
- Tadelis, S. (2016). “Reputation and Feedback Systems in Online Platform Markets.” *Annual Review of Economics*, 8, 21–340.
- Takahashi, S. (2010). “Community Enforcement when Players Observe Partners’ Past Play.” *Journal of Economic Theory*, 145, 42–62.
- Wolitzky, A. (2013). “Cooperation with Network Monitoring.” *Review of Economic Studies*, 80, 395–427.