

SCIENTIFIC REPORTS



OPEN

The Evolution of Covert Signaling

Paul E. Smaldino¹, Thomas J. Flamson² & Richard McElreath³

Human sociality depends upon the benefits of mutual aid and extensive communication. However, diverse norms and preferences complicate mutual aid, and ambiguity in meaning hinders communication. Here we demonstrate that these two problems can work together to enhance cooperation through the strategic use of deliberately ambiguous signals: covert signaling. Covert signaling is the transmission of information that is accurately received by its intended audience but obscured when perceived by others. Such signals may allow coordination and enhanced cooperation while also avoiding the alienation or hostile reactions of individuals with different preferences. Although the empirical literature has identified potential mechanisms of covert signaling, such as encryption in humor, there is to date no formal theory of its dynamics. We introduce a novel mathematical model to assess when a covert signaling strategy will evolve, as well as how receiver attitudes coevolve with covert signals. Covert signaling plausibly serves an important function in facilitating within-group cooperative assortment by allowing individuals to pair up with similar group members when possible and to get along with dissimilar ones when necessary. This mechanism has broad implications for theories of signaling and cooperation, humor, social identity, political psychology, and the evolution of human cultural complexity.

Much of the research on human cooperation has focused on the free-rider problem: how to maintain cooperation when individuals' interests are opposed to those of the group. However, individual interests are often *aligned* with those of the group, and these mutualistic scenarios may be equally important in understanding human social evolution^{1–4}. Even without incentives for individuals to defect, mutualism provides a different dilemma. When individuals differ in preferences or norms, it is harder to efficiently coordinate. The formation of reliable expectations of partner behavior that make coordination possible is therefore essential for the evolution of mutualism⁵.

Consider, for example, a couple planning their Saturday. Chris wants to go to the opera; Pat wants to go to the monster truck rally. Each would rather do something with their partner than go it alone, but each has a different preference⁶. If such mismatches are sufficiently frequent, Chris and Pat might be better off finding new partners with better-aligned interests. Successful cooperation requires resolution of this clash of preferences. Human societies are replete with dilemmas of this kind⁷, and the need to efficiently coordinate extends to many forms of collective action⁸. Institutions like punishment convert other social dilemmas into coordination dilemmas, expanding their importance. If individuals could assort on preferences and norms, cooperative payoffs may be increased. But often these traits are impossible to directly observe. When preferences are consciously held, individuals can merely signal them. But often individuals are not conscious of their preferences or realize their relevance too late to signal them.

One solution is the evolution of ethnic marking. Anthropologists have long argued that ethnic markers may signal group membership and improve cooperative outcomes⁹. An extensive formal literature has developed exploring how arbitrary signals can facilitate assortment on unconscious norms and preferences^{10–15}. Language can also serve as a marker for social coordination¹⁵, as can visible purchasing and fashion choices^{16,17}.

Communication is implicated in all these solutions. However, much communication is ambiguous. Is this ambiguity merely the result of constraints on the accuracy of communication? It may naïvely appear that communication should have clarity as its goal. However, purposeful ambiguity may allow signalers flexibility and plausible deniability^{18–20}. Previous work has illustrated how leaders may use ambiguous language to rally diverse followers¹⁸, politicians may use vague platforms to avoid committing to specific policies²¹, and suitors may mask their flirtations to be viewed innocuously (or at least to provide plausible deniability) if their affections are unreciprocated^{19,22}.

We propose more broadly that ambiguity may enable coordination and thereby enhance cooperation. Consider again the use of signals for assortment on cooperative norms. Overt signals like ethnic markers are useful in some

¹Cognitive and Information Sciences, University of California, Merced, Merced, CA, USA. ²Unaffiliated, Sacramento, CA, USA. ³Department of Human Behavior, Ecology and Culture, Max Planck Institute for Evolutionary Anthropology, Leipzig, Germany. Correspondence and requests for materials should be addressed to P.E.S. (email: paul.smaldino@gmail.com)

Received: 1 March 2018

Accepted: 2 March 2018

Published online: 20 March 2018

contexts, where the adaptive problem is to delimit a set of partners who subscribe to the same broad behavioral norms and to categorically avoid interaction with those who do not. However, the “all-or-nothing” character of such signals can also foreclose valuable partnerships in different contexts. Signals that communicate similarity can also communicate difference, and this can be damaging for within-group cooperation. Individuals may benefit from not foreclosing relationships with less similar group members, so as to successfully cooperate with them in other contexts. Although any two individuals within a group can cooperate when it is mutually beneficial, pairs who are more similar can cooperate more effectively, generating larger benefits^{23–27}. To be clear, we focus on those aspects of individual variation for which similarity enhances cooperation—these include norms, values, and identity^{28,29}. Similarity on these dimensions matter even when the cooperative benefit is maximal for diverse rather than homogenous groups³⁰. Scenarios in which individuals are unable to effectively assort on norms or attitudes are common, especially in complex societies (for example in business or education settings), but also in smaller societies. A signaling system that enables group members to communicate relative similarity only when similarity is high while retaining a shroud of ambiguity when similarity is low is likely to have been advantageous for much of human history.

Covert signaling is the transmission of information that is accurately received by its intended audience but obscured when perceived by others. A common example is “dog-whistling,” in which statements have one meaning for the public at large and a more specialized meaning for others³¹. Such language attempts to transmit a coded message while alienating the fewest listeners possible. A possibly much more common form of covert signaling is humor. According to the encryption model of humor^{32,33}, a necessary component of humorous production is the presence of multiple, divergent understandings of speaker meaning, some of which are dependent on access to implicit information. Only listeners who share access to this information can “decrypt” the implicit understandings and understand the joke. Because the successful production of a joke requires access to that implicit information, humor behaves in manner similar to “digital signatures” in computer cryptography, verifying the speaker’s access to that information without explicitly stating it. Laughter may also serve as an honest signal of a shared sense of humor³⁴, and thus of similarity on a variety of traits. While not all humor has this form, a substantial amount of spontaneous, natural humor does^{32,33}. We allow that other types of identity signals¹⁶ may be also be covert.

We propose that covert signaling serves an important function in facilitating effective cooperation within groups by allowing individuals to assort on norms when possible while avoiding conflict with dissimilar individuals when their assistance is necessary. In the remainder of this paper, we define the logic of covert signaling. We analyze the conditions for selection to favor covert signaling over overt signaling, in which information about an individual’s traits is more transparent. Using a formal model, we show that covert signaling can be favored. It sacrifices transparency for the sake of maintaining working relationships with dissimilar individuals. Although covert signals are less accurate than overt signals, we show that the increased ambiguity can in some cases be advantageous. This characterization of covert signaling in terms of cooperative assortment may therefore help to explain forms of communication and coordination such as coded speech and humor, as well as for the flexibility of human sociality more generally. Nevertheless, covert signaling is not always advantageous. For example, if it is possible to freely choose cooperative partners from a large pool, overt signaling may be more advantageous, as individuals will be better able avoid dissimilar partners. Our model yields specific predictions about default attitudes toward strangers in the absence of clear signals, with implications for understanding differences between contemporary political affiliations.

Model Description

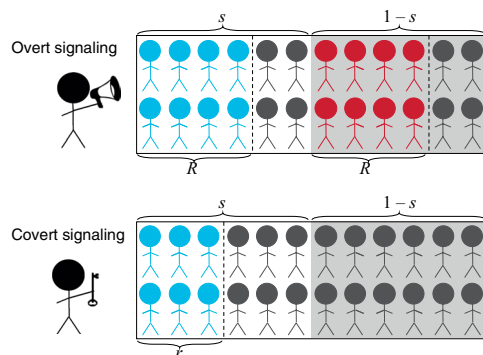
We consider a large population of individuals who have already solved the first-order cooperation problem of suppressing free riders, and can instead focus on maximizing the benefit generated by cooperation. Although individuals all belong to the same group, they also vary along many trait dimensions and thus share more in common with some individuals than others. Pairs of individuals whose trait profiles overlap to some threshold degree are deemed *similar*. Otherwise they are deemed *dissimilar*. Pairs of similar individuals can more effectively coordinate and obtain higher payoffs. The probability that two randomly selected individuals have similar trait profiles is given by s .

Our model proceeds in discrete generations, with each generation subdivided into two stages. In the first stage, individuals signal information about their trait profiles to the other members of the group, and those who receive signals form attitudes about their senders. In the second stage, individuals interact in one of two ways and receive payoffs conditional upon similarity as well as attitudes formed in the first stage.

Stage 1: Signaling. Individuals produce either an overt or covert signal of their underlying traits. We study a family of continuous signaling strategies in which covert signals are produced a fraction p of the time. Overt signals are received by a fraction R of the population and explicitly signal similarity or dissimilarity. Covert signals in contrast are received by a fraction $r < R$ of the population and have content contingent upon the similarity of the sender and receiver. See Fig. 1. When the sender and receiver are similar, covert signals are received as signaling similarity. Otherwise, the receiver does not notice the signal and acts as if no signal was received at all.

Receivers have a default attitude towards all individuals in the population and update this attitude upon receiving a signal. Receiver strategy maps three signal states—similar, dissimilar, no signal—to an attitude. We consider three discrete attitudes: like, dislike, and neutral. These correspond to the hypothesis that covert signals help individuals to avoid being disliked while also achieving sufficient positive assortment by type. Two attitudes would be too few, because it would force agents to adopt either like or dislike as a default attitude, removing any incentive for covert signals. Three is the minimum required to model the hypothesis. We allow a continuous family of receiver strategies. Each strategy parameter a_{XY} indicates the probability of mapping signal $X \in \{\text{Similar, None, Dissimilar}\}$ to attitude $Y \in \{\text{Like, Neutral, Dislike}\}$. The total receiver strategy can be represented by a table (Table 1).

Signal transmission



Attitude formation

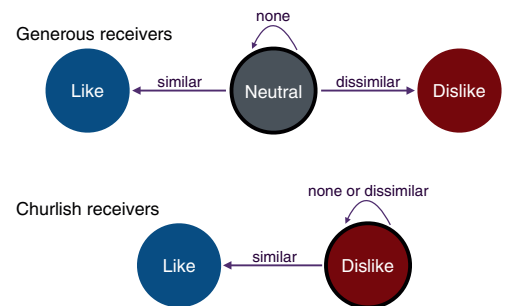


Figure 1. Illustration of the signal transmission and attitude formation components of the model. Left: A proportion s of the population is similar. Overt signals are received by a proportion R of the total population, while covert signals are received by a proportion $r < R$ of similar individuals only. Right: Generous receivers default to a neutral attitude in the absence of a signal, while churlish receivers dislike the signaler unless they know he or she is similar.

	Like	Neutral	Dislike
Similar	a_{SL}	a_{SN}	a_{SD}
None	a_{NL}	a_{NN}	a_{ND}
Dissimilar	a_{DL}	a_{DN}	a_{DD}

Table 1. Matrix of receiver strategies.

The three parameters in each row are constrained to sum to one. Although receiver strategies can vary continuously, we can think of the space of possible strategies as bounded by two extremes. *Generous* receivers default to a neutral attitude in the absence of a signal, while *churlish* receivers dislike the signaler unless they receive a signal of similarity. As we will see, these default receiver strategies are very important.

Stage 2: Interaction. After attitudes are established, pairs of individuals interact. There are two interaction contexts, each with its own mode of dyad formation. In a *free choice* scenario, dyads form conditional on the attitudes of both individuals. This is akin to those circumstances in which individuals can choose their partners based on pre-existing knowledge or established relationships. In contrast, in a *forced choice* scenario, an individual must seek help from whomever happens to be around and dyads are not conditional upon shared attitudes. Such scenarios are common in contemporary industrialized societies, as when two students are paired by their teacher to work on a jointly-graded project. Under these circumstances, it may be important not to have burned bridges, since this will limit the likelihood of effective coordination.

In the free choice context, dyads form based on joint attitudes, but attitudes do not directly influence payoffs. Instead, payoffs are determined by the underlying similarity of the dyad. Specifically, similar dyads receive an average payoff of 1, establishing a baseline measurement scale. Dissimilar dyads receive a payoff of zero. An individual in a dyad who *likes* the other individual increases the proportional odds of that dyad forming by a factor $w_L > 1$. For each individual who *dislikes* the other, the proportional odds of the dyad forming are reduced by a factor $w_D < 1$. This implies five possible kinds of dyads that might interact: LL, LN, NN, ND, and DD. The proportional odds of each, relative to random assortment, are: w_L^2 , w_L , 1, w_D , and w_D^2 . These parameters are fixed features of the social environment, not aspects of strategy. This prevents strategy dynamics from generating perfect assortment. Receiver strategies that assign attitudes in ways that make good use of signal information will achieve better assortment and receive higher payoffs, conditional on the assortment constraints determined by w_L and w_D .

In the forced choice context, dyads form at random with respect to attitudes, but attitudes directly influence payoffs. This context entails a baseline payoff of 1 for both individuals. However, attitudes adjust payoffs, because negative attitudes make it harder to interact. When one individual dislikes the other, he makes the interaction more difficult than it must be and thereby imposes a cost $-d$ on the other individual. When both individuals dislike one another, their difficulties act synergistically, inducing an additional cost $-\delta$ on each. This cost could result from spite or from uncontrollable inefficiency, a negative consequence of second-order common knowledge (*sensu* Chwe³⁵). As we show later, these synergistic costs are very important to the overall signaling dynamics.

Let q be the relative importance of the free choice context and $1-q$ the relative importance of the forced choice context. These two contexts are starkly different, presenting the clearest investigation of the hypothesis that covert signals trade worse performance in assortment contexts, in which norms influence payoffs, for better performance in forced contexts in which attitudes influence payoffs. Real contexts are some mix of these extremes, and the parameter q allows us to explore the range of mixes. We note also that payoffs depend on similarity and attitudes, which result in part due to sender and receiver strategies, but the payoffs are not directly derived from those

strategies. When dyads form for cooperation, both parties are senders *and* receivers, each having had the opportunity to signal to and receive from one another. While we can imagine cases of sender-receiver asymmetries in terms of cooperative payoffs, we have chosen not to model these in this particular analysis.

Payoff expression. With the assumptions above, we can define a general payoff expression for a rare individual with signal strategy p' and receiver strategy matrix \mathbf{a}' in a population with common-type strategy $\{p, \mathbf{a}\}$. The expected payoff to this individual is:

$$W(p', \mathbf{a}') = \Omega + q\text{Pr}(\text{similar}|p', \mathbf{a}') + (1 - q)(1 - \text{Pr}(\text{disliked}|p', \mathbf{a}')(\delta + \text{Pr}(\text{dislike}|p', \mathbf{a}')\delta))$$

where Ω is an expected payoff due to other activities. The work lies in defining the probabilities $\text{Pr}(\text{similar}|p', \mathbf{a}')$, $\text{Pr}(\text{dislike}|p', \mathbf{a}')$, and $\text{Pr}(\text{disliked}|p', \mathbf{a}')$. In the mathematical appendix, we show how to define these probabilities, using the assumptions above. The resulting general payoff expression is very complicated. In the following section, however, we are nevertheless able to analyze it by considering invasion and stability of relevant combinations of signaling and receiver strategies.

Analysis and Results

The motivating hypothesis is that covert signals can proliferate because they allow sufficient assortment in the free choice context and also reduce being disliked in the forced choice context. To evaluate the logic of this idea, we proceed by asking when covert signals can be stable, when they can invade, and which receiver strategies are necessary for their stability or invasion. The following conditions favor covert signals.

1. Covert signals require a sufficient proportion of receivers to default to neutral attitudes. If everyone defaults to disliking, then covert signals can produce no benefit. Defaulting to neutral is favored under a wide range of conditions, provided that covert signals are sufficiently hard to receive (r is not too large) and avoidance of disliked individuals is not too efficient (w_D is not too small).
2. Covert signals require that the cost of being disliked in the forced choice context be sufficiently high. This also means that baseline similarity in the population (s) must be sufficiently low, because this creates the risk of being disliked by dissimilar individuals.
3. Overt signalers cannot have too large an advantage in the free choice context. This requires that assortment with liked individuals not be too accurate. The accuracy of assortment is influenced by the reception probabilities of both signal types, R and r , as well as the proportional odds assortment factors, w_L and w_D .

In the remainder of this section, we derive these results and provide intuition for why they hold. First we derive simple evolutionary dynamics for these payoffs. This allows us to submit the model to invasion and stability analysis, asking both when covert signals can be stable and when they may invade a population of overt signals. Then we proceed by considering the dynamics within each interaction context—the forced choice context and the free choice context—separately. Then we summarize the joint dynamics of the full model with both contexts.

Evolutionary dynamics. We generate evolutionary dynamics for the strategy space by assuming that rare invading strategies increase in frequency when they achieve higher payoffs than a common-type strategy. We define selection gradients for both p and the attitude parameters. The gradient for signaling is defined by:

$$g(p) = \left. \frac{\partial W(p', \mathbf{a}')}{\partial p'} \right|_{p'=p, \mathbf{a}'=\mathbf{a}} \quad (1)$$

The gradient for each receiver parameter is defined similarly. A number of different mechanisms can generate such dynamics. For example, an individual could acquire its strategy from successful individuals through learning. Genetically coded strategies that influence biological fitness would also generate this dynamic. We remain agnostic about inheritance and transmission mechanism, because the point of our modeling exercise is to explore the design aspects of covert signals. This is best achieved by a form of analysis that abstracts away from transmission details, even though of course in any real system such details will turn out to influence which strategies are possible and how they evolve³⁶. We also note that if the mechanism of transmission is cultural, replicators are not strictly necessary but can serve as meaningful approximations of lower-fidelity transmission channels³⁷.

The potential space of receiver strategies is very large. However, the *relevant* space of strategies is fairly small. In the mathematical appendix, we show that payoff dynamics always favor mapping *similar* signals to *like* attitudes, implying $a_{SL} = 1$. The reason is that maximizing the probability of assortment for similar individuals also maximizes payoffs, and the *like* attitude maximizes the probability of interacting in the free choice context. On the other hand, payoff dynamics do not always favor mapping *dissimilar* signals to *dislike* attitudes. The reason is that the forced choice context disfavors disliking whenever $\delta > 0$. We therefore constrain further analysis to the relevant situations in which the penalty for mutual dislike, δ , is small enough that assortment incentives favor mapping *dissimilar* signals to *dislike* attitudes. We reemphasize this constraint in the discussion, because constraints of this sort help in producing predictions. Finally, with respect to default attitudes, formed when no signal is received, payoff dynamics never favor assigning *like*, because this erodes the value of assigning *like* to *similar* signals.

The remaining default receiver parameters are free to evolve. Therefore, for most of the analysis to follow, we assume that $a_{SL} = 1$, $a_{DD} = 1$, $a_{NL} = 0$, $a_{NN} = 1 - \alpha$, and $a_{ND} = \alpha$. This allows us to use the gradient on α , defined by:

$$g(\alpha) = \left. \frac{\partial W(p', \alpha')}{\partial \alpha'} \right|_{p'=p, \alpha'=\alpha} \tag{2}$$

to ask when evolution favors assigning *dislike* to no signal, $\alpha > 0$, effectively defaulting to disliking everyone. As noted earlier, it will be convenient to refer to $\alpha = 0$ as the *generous receiver* strategy and $\alpha = 1$ as the *churlish receiver* strategy.

Dynamics of the forced choice context. In this context, incentives favor covert signals, because such signals are better at avoiding being disliked. However, this advantage depends upon incentives favoring generous receiver strategies that do not dislike by default. That said, receiver incentives in this context *always* favor generous receiver strategies as long as there is any negative synergy, $\delta > 0$. Therefore the forced choice context favors generous receivers, $\alpha = 0$, which in turn favor covert signals, $p = 1$.

The gradients in this context are:

$$g(p)|_{q=0} = (\alpha rs + R(1 - \alpha - s))(d + \delta(\alpha(1 - prs) + R(1 - p)(1 - \alpha - s))) \tag{3}$$

$$g(\alpha)|_{q=0} = -\delta(1 - (1 - p)R - prs)(\alpha(1 - prs) + (1 - p)R(1 - \alpha - s)) \tag{4}$$

These expressions seem complex at first, but produce fairly simple dynamics. First let's ask when p increase. When $\alpha = 0$, the generous receiver strategy is common, and covert signals can increase when:

$$d + \delta(1 - s)(1 - p)R > 0 \tag{5}$$

This is satisfied for any allowable values of the parameters. Note also that it does not require both a direct cost of being disliked, d , and a synergistic cost of mutual dislike, δ . Either one is sufficient to favor covert signals, as long as α is small. Next consider when $\alpha = 1$, the churlish receiver strategy is common. Then covert signals can increase when:

$$d + \delta(1 - s(pr + (1 - p)R)) < 0 \tag{6}$$

And this is never satisfied, for any p . Therefore covert signals are favored when $1 - \alpha$, the amount of generous receiver behavior, is sufficiently high. The threshold value is found where $g(p) = 0$:

$$\hat{\alpha} = \frac{1 - s}{1 - \frac{r}{R}s} \tag{7}$$

When α is above this value, overt signals are favored. When it is below it, covert signals are favored. Why? When receivers are relatively generous, and there is sufficient dissimilarity in the population, covert signals reduce costs by avoiding being disliked. If generous receivers are relatively rare, however, then covert signals can actually do worse than overt signals, because they are received less often than overt signals, $r < R$. If r/R is sufficiently small, overt signals are favored for a wide range of values of α . If however $r = R$ and covert signals have no disadvantage in audience size, then overt signals are never favored in this context, no matter the amount of similarity s .

Now the crucial question is when α will fall below this threshold $\hat{\alpha}$. The condition for payoff dynamics to favor smaller values of α , in the forced choice context only, is just $\delta > 0$. Therefore the forced choice context always favors smaller values of α , provided there is any negative synergy between disliking and being disliked. Otherwise α is neutral and does not move at all, based on payoff dynamics. Why does this context always favor generous receivers? There is no advantage to be had in disliking people in this context, because assortment does not depend upon attitudes. Payoffs depend upon attitudes, however, and mutual dislike results in poor payoffs. Therefore, it pays to be generous in attitudes towards those one has no information about.

Dynamics of the free choice context. In the free choice context, attitudes influence assortment but do not directly influence payoffs. Instead, hidden norm similarity influences payoffs. The free choice context favors overt signals over covert signals, because overt signals increase assortment—such signals are easier to receive and more effectively discriminate similarity. The churlish receiver strategy, $\alpha = 1$, is favored by this context, because it also increases assortment with similar individuals. Therefore this context is hostile to covert signals and to the receiver strategy that favors them.

To support the above statements, we demonstrate that the gradient for p in this context is always negative and that the gradient for α in this context is always positive. The gradients in this context are:

$$g(p)|_{q=1} = (-1) \frac{s(1 - s)}{Z^2} (1 - (\alpha + (1 - \alpha)(1 - p)R)v_D) \tag{8}$$

$$((1 - \alpha v_D) + (R(1 - p) + pr)(\alpha v_D + v_L))$$

$$(R(1 - \alpha v_D)(v_D + v_L) - r(1 - v_D)(\alpha(1 - R) + R))(\alpha v_D + v_L)$$

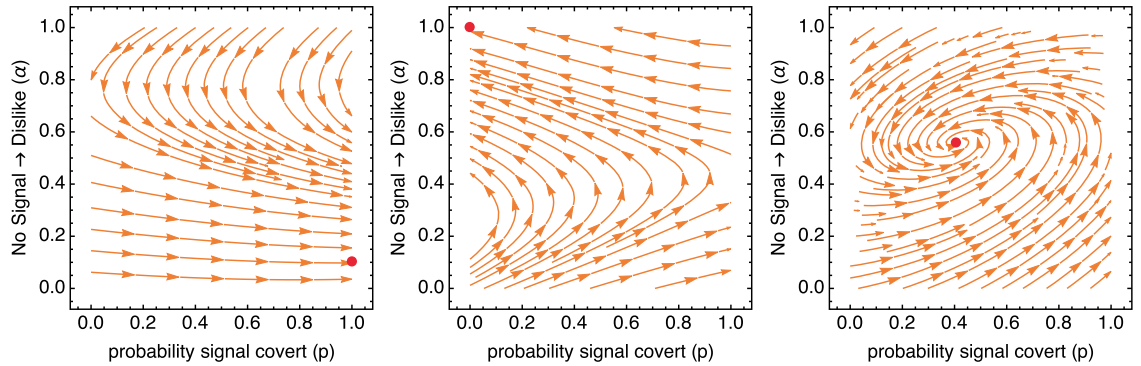


Figure 2. The three dynamic regimes that arise from the joint dynamics. In each plot, the paths show the evolutionary trajectories in each region of the phase space defined by the probability of covert signals (p , horizontal) and the probability of churlish receivers (α , vertical). The red points show equilibria. In all three plots: $d = 0.1$, $\delta = 0.01$, $q = 0.5$. Left: $s = 0.1$, $r/R = 0.25$, $w_L = 1.1$, $w_D = 0.9$. Middle: Same as left, but $w_D = 0.6$. Right: $s = 0.2$, $r/R = 0.75$, $w_L = 1.25$, $w_D = 0.8$.

$$g(\alpha)|_{q=1} = \frac{s(1-s)}{Z^2} (1-w_D)(1-(1-w_D)(\alpha+(1-\alpha)(1-p)R)) \\ ((1-p)R(w_L-w_D)(1-pr-(1-p)R)+prw_L) \\ ((1-(1-p)R-pr)(1-\alpha(1-w_D))+(R(1-p)+pr)w_L) \quad (9)$$

where Z is a normalizing term and the symbols $v_D = 1 - w_D$ and $v_L = w_L - 1$ are used for compactness of notation. By inspection, every term after the leading $(-)$ in $g(p)$ is positive, for all allowed values of the variables, and so the gradient is always negative. Similarly, every term in $g(\alpha)$ is positive, and so the gradient is always positive. Therefore payoff incentives in the free choice context never favor covert signals and always favor churlish receivers.

While this context always favors overt signalers and churlish receivers, the strength of the incentives may vary. First, both gradients are proportional to the variance in similarity, $s(1-s)$. This indicates that intermediate similarity more strongly favors overt signals, unlike the situation in the forced choice context, in which high similarity favored overt signals. The reason that the variance matters now is that payoffs depend directly upon similarity, not upon attitudes. The more variance in similarity in the population, the greater the advantage of efficient assortment.

Overt signals have the advantage in this context, because they are better at assortment. Therefore, any change in variables that reduces the accuracy of assortment overall will reduce overt signalers' advantage. The important variables are R , the probability an overt signal is received, and the assortment proportional odds w_D and w_L . Reducing R reduces overt signalers' advantage, because it makes signals less valuable overall. Making either w_D or w_L closer to 1 makes assortment, based on attitudes, less effective. This also reduces overt signalers' advantage.

Joint dynamics: When do covert signals evolve? When both contexts matter, the joint dynamics take on one of three characteristic regimes. First, covert signals both invade and are evolutionarily stable. Second, overt signals both invade and are evolutionarily stable. Third, a mixed equilibrium exists at which covert and overt signals coexist in the population. Figure 2 illustrates these three regimes. These three examples all weigh the forced choice and free choice contexts equally, $q = 0.5$. The other parameters then shift the strength of incentives in each context to influence overall dynamics. For other values of q , the strength of incentives would have to shift as well to overcome weight given to each context.

When covert signals are sufficiently noisy (r/R low), similarity is sufficiently rare (s low), and assortment (w_L, w_D) not too efficient, covert signals can both invade and are an ESS. This situation is shown in the lefthand plot. While dynamics do not favor covert signals when α is large, near the top of the phase space, dynamics in that region favor smaller values of α . Eventually, α becomes small enough to allow covert signals to invade and reach fixation. In many cases, a small amount of churlish receiver strategy, $\alpha > 0$, persists.

When the conditions outlined above are not met, incentives favor instead overt signals. The middle plot illustrates a case essentially the opposite of the one on the left. Here, $w_D = 0.6$, making assortment efficient. When assortment is efficient, it may pay to dislike by default. This sets up a dynamic that eventually favors overt signals. While covert signals are still favored when α is low, the fact that larger values of α are favored everywhere leads eventually to invasion and fixation of overt signals.

Finally, the plot on the right shows an intermediate case, in which conditions favor both signaling strategies. Here $s = 0.2$, $r/R = 0.75$, $w_L = 1.25$, and $w_D = 0.8$. In this regime, the conditions that favor covert signals also favor more churlish receivers. Similarly, the conditions that favor overt signals also favor fewer churlish receivers. In total, the population comes to rest with a mixture of signaling and receiving strategies.

A more general view of the dynamics is available by considering the boundary conditions that make covert signaling an ESS. Recall that q is the relative importance of free choice scenarios. Define a threshold \hat{q} as the larg-

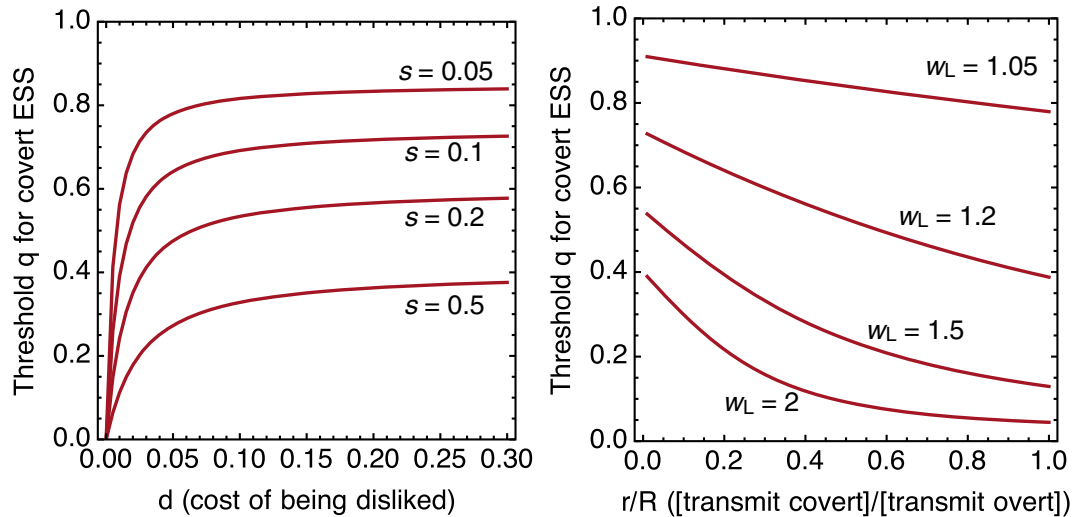


Figure 3. Plots of the largest value of q , \hat{q} , that allows covert signaling to be an ESS. Each curve represents a set of parameter values. Points below each curve make covert signals uninvadable by overt signals. Points above each curve allow overt signaling to invade. Left: The cost of being disliked, d , for four values of the baseline rate of similarity, s . $r/R = 0.5$, $w_L = 1.1$, $w_D = 0.9$, and $\delta = 0.01$. Right: The ratio of covert transmission rate, r , to the rate of overt transmission rate, R , for four values of $w_L = 1/w_D$, $s = 0.1$, $R = 0.5$, $d = 0.1$, and $\delta = 0.01$.

est value of q for which covert signals can resist invasion by overt signals. This is defined by values of $q = \hat{q}$ and $\alpha = \hat{\alpha}$ that satisfy $g(p)|_{p=1} = 0$ and $g(\alpha)|_{p=1} = 0$. These cannot in general be solved analytically. So we solve the system numerically, in order to illustrate the range of joint dynamics. Values of q less than \hat{q} make covert signals evolutionarily stable—overt signals cannot invade. Values of q greater than \hat{q} allow overt signals to invade, though we note that covert signals may nevertheless remain in the population, at an internal stable value p . Therefore \hat{q} provides a useful metric of how strongly a parameter configuration favors covert signals.

We use \hat{q} to summarize the tradeoffs in the signaling model. Recall that the cost of being disliked, d , is needed to favor covert signals. Therefore increasing d makes it easier for covert signals to be an ESS. However, the rate of similarity, s , favors overt signals. It is of value to note that d cannot compensate for s —if s is large, then steeply increasing costs d will not favor covert signals. We show this relationship in Fig. 3, lefthand plot. Each curve in this plot is a threshold \hat{q} , below which covert signaling is an ESS. For each level of s , the impact of increasing d diminishes rapidly. Therefore some cost d is necessary for covert signals to evolve and be stable, but these costs cannot easily compensate when similarity is sufficiently common.

Consider another important pair of dimensions: the ratio of transmission rates r/R in covert/overt signals and the efficiency of assortment, as measured by w_L and w_D . Figure 3, righthand side, shows \hat{q} curves for four values of $w_L = 1/w_D$, as functions of r/R . Covert signals are favored when r/R is small, as explained in the previous sections. But when assortment is very efficient, such as $w_L = 2$ near the bottom of the plot, it requires very low values of r/R to compensate in favor of overt signaling.

Discussion

The dynamics of cooperation are more complicated than implied by models in which maximal benefits accrue to those who can simply avoid free riders. Not all cooperators are equal. Individuals vary, making assortment among cooperators important. Circumstances also vary. When individuals must occasionally collaborate with those outside their circles of friends, it can be critical to avoid burning bridges with dissimilar members of one's group. Covert signaling makes this possible, and this may be why phenomena like humor are observed in all human societies, at both small and large scales^{38,39}.

We have shown that covert signaling is favored when forced choice scenarios are common, when similarity is low, when the cost of being disliked is high, and when covert signals are sufficiently noisy to make the meaning of a signal's absence ambiguous. We emphasize that covert signaling can be favored even though it is less effective than overt signaling at communicating similarity, because it simultaneously avoids communicating dissimilarity. Although we have focused our attention on the initial establishment of cooperative relationships via signaling, we also note that people can change over time; they may grow more similar to one another or further apart. Covert signals may be important for the continued maintenance of a relationship, or for its reestablishment after prolonged absence.

Our model points to interesting transitions from inter- to intra-group assortment dynamics. As noted, overt signaling systems are favored when the ability to assort on attitudes is high, and when being disliked by dissimilar individuals carries little risk. This is precisely the kind of situation that is assumed to obtain in inter-group assortment, where overt signals such as ethnic markers are used to discriminate between similar and dissimilar individuals. In these between-group contexts, the difference between similar and dissimilar individuals is so great that attempting to coordinate with dissimilar others is not worth the effort, and one can afford to burn bridges

with them in order to ensure that similar others are aware of their similarity¹³. In fact, it might be argued that burning bridges with dissimilar out-group members is as much a goal of overt signals like ethnic markers as is attracting similar in-group members.

Intra-group assortment, however, is not simply a matter of scaling down inter-group dynamics. Rather, we must already presume some baseline level of similarity resulting from inter-group assortment; for there to be a group within which to assort, some degree of similarity should already be in place that defines that group, such as the shared interaction norms, communication systems, etc. that ethnic markers are thought to ensure. The benefits of further assorting on the basis of more nuanced similarity are therefore likely to be marginal relative to random assortment within the group. When such benefits are small but the costs of being disliked are high, covert signaling is favored.

Relatedly, we emphasize that the probability of similarity, s , need not reflect some number of discrete types in the population, but can instead refer to a level of selectivity in how much a given pair of individuals must have in common to be considered “similar.” That is, s refers to the proportion of the population that could be considered similar to a focal individual in a given context, with higher values indicating a looser concept of “similar” than lower ones. Changes to s can have a significant impact on the overall dynamics of the system. When s is large, the focus is on avoiding rare dissimilar individuals, and overt signals will be favored. As s decreases, the criteria for considering a potential partner sufficiently similar to reap the benefits of enhanced coordination become stricter, and the utility of covert signals increases.

An interesting direction for future exploration is how these dynamics might respond to increased social complexity. In the larger and more complex societies associated with the development of agriculture, and particularly in the last few centuries, interactions with strangers are more frequent and occur across many contexts, necessitating strategies for temporary assortment^{28,29,40}. Consequently, expected similarity will be lower, signal fidelity will be noisier, and assortment on attitudes will be less efficient. These are precisely the conditions in our model associated with the evolution of covert signaling. In large, diverse populations, covert signaling may sustain social cohesion and prevent burning bridges between individuals or groups that must occasionally collaborate. That said, covert signals are not necessarily rare in small-scale societies. Our own experiences in the field and conversations with other researchers indicate that they occur with some regularity. Our model can help to identify contexts in which covert signaling should or should not be expected.

Identity signaling, whether overt social markers or more covert communication, can be used by individuals looking to find others similar to themselves and to avoid being mistaken for something they are not^{16,28,29}. If the need to cooperate with dissimilar individuals is unlikely or if similar individuals are common, then overt declarations of identity should be expected. On the other hand, if burning bridges is both costly and likely given an overt signaling strategy, we should expect identity to be signaled much more subtly. In reality, increasing levels of specificity may be signaled in increasingly covert ways, and without all received signals actively inducing a change in disposition toward the sender. A related signaling strategy, not covered by our model, might facilitate liking between similar individuals but only indifference otherwise. Casual, coarse-grain identity signaling may often take this form, as in cases of fashion adoption or pop culture allegiances. It would be interesting to investigate how common these “semi-covert” signals are in small-scale communities, as they seem pervasive in complex industrialized societies.

Our model additionally helps make sense of findings from political psychology suggesting that people in the industrialized West who identify as conservative or right-leaning tend to view ambiguous people as hostile, while those identifying as liberal or left-leaning tend to view ambiguous people as neutral^{41–43}. In our model, a default attitude to dislike was linked with overt signaling, which we in turn associate with the preservation of strong between-group boundaries. In contrast, a default attitude of neutral was associated with covert signaling, and with the avoidance of burned bridges to facilitate more widespread within-group cooperation. As a broad generalization, our analysis suggests that conservatives may be operating under the assumptions of stronger ingroup/outgroup boundaries, increased expectations of similarity toward those they signal, and lower costs to being disliked by dissimilar individuals. In contrast, liberals may be operating under the assumptions of a more broadly defined ingroup, limited expectations for similarity toward those they signal, and higher costs to being disliked by dissimilar individuals. Lending modest support to this idea is the finding that conservatives appear to have a stronger “need for cognitive closure” (reviewed in ref.⁴²), which is associated with, among other things, a distaste for uncertainty and ambiguity. The modeling framework we present in this paper may thus be useful in understanding patterns of differences between groups, including but not limited to political affiliation.

Ours is the first formal model of covert signaling. As such, it necessarily involves simplifying assumptions concerning the nature of signaling and cooperative assortment. For example, while we have allowed for covert signaling errors in the form of failed transmission to similar individuals, we have not included the converse form of error, where dissimilar individuals *are* able to detect the signal some of the time, and therefore update their disposition to disliking the covert signaler. Adding an additional parameter to account for this possibility does not qualitatively change our analysis. But it may create conditions where a non-signaling “quiet” strategy could invade. In addition, we ignore the possibility of strategic action on the part of the receiver to either improve coordination or to avoid partnering with dissimilar individuals entirely. We assumed that a pairing of dissimilar partners would simply lead to an unsuccessful collaboration, but such a pairing might instead lead each individual to pursue more individualistic interests. At the population level, we assumed that all individuals had an equal probability of encountering similar individuals, and that all similar and dissimilar individuals were equivalent. In reality, some individuals may be more or less likely to encounter similar individuals, perhaps related to differences in the tendency to be conformity- versus distinctiveness-seeking⁴⁴, or reflecting minority-majority dynamics^{45,46}. Exploration of this variation opens the door to evaluating signaling and assortment strategies in stratified groups. All of these limitations provide avenues for future research that will build upon the central findings reported here.

In a population where individuals vary and burning bridges is costly, overtly announcing precisely where one stands entails venturing into a zone of danger. Covert signaling, as in the case of humor or otherwise encrypted language, allows individuals to effectively assert when possible while avoiding burned bridges when the situation calls for partnerships of necessity.

References

1. Skyrms, B. *The Stag Hunt and the Evolution of Social Structure* (Cambridge University Press 2004).
2. Calcott, B. The other cooperation problem: Generating benefit. *Biology and Philosophy* **23**, 179–203, <https://doi.org/10.1007/s10539-007-9095-5> (2008).
3. Tomasello, M., Melis, A. P., Tennie, C., Wyman, E. & Herrmann, E. Two key steps in the evolution of human cooperation. *Current Anthropology* **53**, 673–692 (2012).
4. Smaldino, P. E. The cultural evolution of emergent group-level traits. *Behavioral and Brain Sciences* **37**, 243–295, <https://doi.org/10.1017/S0140525X13001544> (2014).
5. Schelling, T. C. *The Strategy of Conflict* (Harvard University Press 1960).
6. Luce, R. D. & Raiffa, H. *Games and Decisions* (Wiley 1957).
7. Boyd, R. & Richerson, P. J. The evolution of norms: An anthropological view. *Journal of Institutional and Theoretical Economics* **150**, 72–87, <http://www.jstor.org/stable/40753018> (1994).
8. Ostrom, E. Collective action and the evolution of social norms. *The Journal of Economic Perspectives* **14**, 137–158, <https://doi.org/10.1257/jep.14.3.137> (2000).
9. Barth, F. Introduction. In Barth, F. (ed.) *Ethnic Groups and Boundaries*, 9–38 (Little, Brown, New York 1969).
10. Castro, L. & Toro, M. A. Mutual benefit cooperation and ethnic cultural diversity. *Theoretical Population Biology* **71**, 392–399, <https://doi.org/10.1016/j.tpb.2006.10.003> (2007).
11. Efferon, C., Lalive, R. & Fehr, E. The coevolution of cultural groups and ingroup favoritism. *Science* **321**, 1844–1849, <https://doi.org/10.1126/science.1155805> (2008).
12. Mace, R. & Holden, C. J. A phylogenetic approach to cultural evolution. *Trends in Ecology and Evolution* **20**, 116–121 (2005).
13. McElreath, R., Boyd, R. & Richerson, P. J. Shared norms and the evolution of ethnic markers. *Current Anthropology* **44**, 122–130, <https://doi.org/10.1086/345689> (2003).
14. Moffett, M. W. Human identity and the evolution of societies. *Human Nature* **24**, 219–267, <https://doi.org/10.1007/s12110-013-9170-3> (2013).
15. Nettle, D. & Dunbar, R. Social markers and the evolution of reciprocal exchange. *Current Anthropology* **38**, 93–99, <http://www.jstor.org/stable/2744442> (1997).
16. Berger, J. & Heath, C. Who drives divergence? Identity signaling, outgroup dissimilarity, and the abandonment of cultural tastes. *Journal of Personality and Social Psychology* **95**, 593, <https://doi.org/10.1037/0022-3514.95.3.593> (2008).
17. Smaldino, P. E., Janssen, M. A., Hillis, V. & Bednar, J. Adoption as a social marker: Innovation diffusion with outgroup aversion. *Journal of Mathematical Sociology* **41**, 26–45, <https://doi.org/10.1080/0022250X.2016.1250083> (2017).
18. Eisenberg, E. M. Ambiguity as strategy in organizational communication. *Communication Monographs* **51**, 227–242 (1984).
19. Pinker, S., Nowak, M. A. & Lee, J. J. The logic of indirect speech. *Proceedings of the National Academy of Sciences* **105**, 833–838, <https://doi.org/10.1073/pnas.0707192105> (2008).
20. Santana, C. Ambiguity in cooperative signaling. *Philosophy of Science* **81**, 398–422 (2014).
21. Aragonès, E. & Neeman, Z. Strategic ambiguity in electoral competition. *Journal of Theoretical Politics* **12**, 183–204, <https://doi.org/10.1177/0951692800012002003> (2000).
22. Gersick, A. & Kurzban, R. Covert sexual signaling: Human flirtation and implications for other social species. *Evolutionary Psychology* **12**, 549–569, <https://doi.org/10.1177/147470491401200305> (2014).
23. Kaufman, H. Similarity and cooperation received as determinants of cooperation rendered. *Psychonomic Science* **9**, 73–74 (1967).
24. Wolosin, R. J. Cognitive similarity and group laughter. *Journal of Personality and Social Psychology* **32**, 503–509 (1975).
25. Fischer, I. Friend or foe: Subjective expected relative similarity as a determinant of cooperation. *Journal of Experimental Psychology: General* **138**, 341–350 (2009).
26. Hruschka, D. J. *Friendship: Development, ecology, and evolution of a relationship* (University of California Press, Berkeley 2010).
27. Toma, C., Corneille, O. & Yzerbyt, V. Holding a mirror up to the self: Egocentric similarity beliefs underlie social projection in cooperation. *Personality and Social Psychology Bulletin* **38**, 1259–1271 (2012).
28. Smaldino, P. E. The evolution of the social self: Multidimensionality of social identity solves the coordination problems of a society. In Love, A. C. & Wimsatt, W. C. (eds) *Beyond the Meme: Development and Structure in Cultural Evolution* (University Minnesota Press 2018).
29. Smaldino, P. E. Social identity and cooperation in cultural evolution. *Behavioural Processes*, <https://doi.org/10.1016/j.beproc.2017.11.015> (2018).
30. Hong, L. & Page, S. E. Groups of diverse problem solvers can outperform groups of high-ability problem solvers. *Proceedings of the National Academy of Sciences* **101**, 16385–16389, <https://doi.org/10.1073/pnas.0403723101> (2004).
31. López, I. H. *Dog Whistle Politics: How Coded Racial Appeals Have Reinvented Racism & Wrecked the Middle Class* (Oxford University Press 2014).
32. Flamson, T. & Barrett, H. The encryption theory of humor: A knowledge-based mechanism of honest signaling. *Journal of Evolutionary Psychology* **6**, 261–281, <https://doi.org/10.1556/JEP.6.2008.4.2> (2008).
33. Flamson, T. J. & Bryant, G. A. Signals of humor: Encryption and laughter in social interaction. In Dynel, M. (ed.) *Developments in Linguistic Humour Theory*, vol. 1, 49–73 (John Benjamins Publishing, Amsterdam 2013).
34. Bryant, G. A. & Aktipis, C. A. The animal nature of spontaneous human laughter. *Evolution and Human Behavior* **35**, 327–335, <https://doi.org/10.1016/j.evolhumbehav.2014.03.003> (2014).
35. Chwe, M. S.-Y. *Rational Ritual: Culture, Coordination, and Common Knowledge* (Princeton University Press 2001).
36. Grafen, A. Natural selection, kin selection and group selection. In Krebs, J. & Davies, N. B. (eds) *Behavioural ecology: An evolutionary approach*, 62–84 (Blackwell 1984).
37. Henrich, J. & Boyd, R. On modeling cognition and culture: Why cultural evolution does not require replication of representations. *Journal of Cognition and Culture* **2**, 87–111, <https://doi.org/10.1163/156853702320281836> (2002).
38. Apte, M. *Humor and laughter: An Anthropological Approach* (Cornell University Press, Ithaca, NY 1985).
39. Brown, D. *Human universals* (Temple University Press, Philadelphia 1991).
40. Johnson, A. W. & Earle, T. K. *The Evolution of Human Societies: From Foraging Group to Agrarian State* (Stanford University Press 2000).
41. Vigil, J. M. Political leanings vary with facial expression processing and psychosocial functioning. *Group Processes & Intergroup Relations* **13**, 547–558 (2010).
42. Hibbing, J. R., Smith, K. B. & Alford, J. R. Differences in negativity bias underlie variations in political ideology. *Behavioral and Brain Sciences* **37**, 297–350 (2014).

43. Holbrook, C., López-Rodríguez, L., Fessler, D. M. T., Vázquez, A. & Gómez, A. Gulliver's politics: Conservatives envision potential enemies as readily vanquished and physically small. *Social Psychological and Personality Science* **8**, 670–678, <https://doi.org/10.1177/1948550616679238> (2017).
44. Smaldino, P. E. & Epstein, J. M. Social conformity despite individual preferences for distinctiveness. *Royal Society Open Science* **2**, 140437, <https://doi.org/10.1098/rsos.140437> (2015).
45. Wimmer, A. *Ethnic Boundary Making: Institutions, Power, Networks* (Oxford University Press 2013).
46. Bunce, J. A. & McElreath, R. Sustainability of minority culture when inter-ethnic interaction is profitable. *Nature Human Behaviour* <https://doi.org/10.1038/s41562-018-0306-7> (2018).

Acknowledgements

Thanks to Patrick Barclay, William Baum, John Bunce, and several anonymous reviewers for helpful comments. This work was supported by National Science Foundation Grant BCS 1357240 (to T.F. and R.M.) and by the Division of Social Sciences Dean's Office at the University of California Davis.

Author Contributions

All authors conceived the project. P.S. and R.M. designed and analyzed the model. All authors wrote and reviewed the manuscript.

Additional Information

Supplementary information accompanies this paper at <https://doi.org/10.1038/s41598-018-22926-1>.

Competing Interests: The authors declare no competing interests.

Publisher's note: Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2018