

Research



Cite this article: Bohn M, Liebal K, Oña L, Tessler MH. 2022 Great ape communication as contextual social inference: a computational modelling perspective. *Phil. Trans. R. Soc. B* **377**: 20210096.

<https://doi.org/10.1098/rstb.2021.0096>

Received: 14 September 2021

Accepted: 4 April 2022

One contribution of 14 to a theme issue 'Revisiting the human 'interaction engine': comparative approaches to social action coordination'.

Subject Areas:

behaviour, evolution

Keywords:

communication, primates, social cognition, evolution, computational modelling

Author for correspondence:

Manuel Bohn

e-mail: manuel_bohn@eva.mpg.de

Correspondence concerning this article should be addressed to Manuel Bohn, Max Planck Institute for Evolutionary Anthropology, Deutscher Platz 6, 04103 Leipzig, Germany.

Great ape communication as contextual social inference: a computational modelling perspective

Manuel Bohn¹, Katja Liebal², Linda Oña³ and Michael Henry Tessler⁴

¹Department of Comparative Cultural Psychology, Max Planck Institute for Evolutionary Anthropology, 04103 Leipzig, Germany

²Institute of Biology, Leipzig University, 04103 Leipzig, Germany

³Naturalistic Social Cognition Group, Max Planck Institute for Human Development, 14195 Berlin, Germany

⁴Department of Brain and Cognitive Sciences, Massachusetts Institute of Technology, Cambridge, MA 02139-4307, USA

id MB, 0000-0001-6006-1348; KL, 0000-0003-2447-8327; MHT, 0000-0003-3945-0239

Human communication has been described as a contextual social inference process. Research into great ape communication has been inspired by this view to look for the evolutionary roots of the social, cognitive and interactional processes involved in human communication. This approach has been highly productive, yet it is partly compromised by the widespread focus on how great apes use and understand individual signals. This paper introduces a computational model that formalizes great ape communication as a multi-faceted social inference process that integrates (a) information contained in the signals that make up an utterance, (b) the relationship between communicative partners and (c) the social context. This model makes accurate qualitative and quantitative predictions about real-world communicative interactions between semi-wild-living chimpanzees. When enriched with a pragmatic reasoning process, the model explains repeatedly reported differences between humans and great apes in the interpretation of ambiguous signals (e.g. pointing or iconic gestures). This approach has direct implications for observational and experimental studies of great ape communication and provides a new tool for theorizing about the evolution of uniquely human communication.

This article is part of the theme issue 'Revisiting the human 'interaction engine': comparative approaches to social action coordination'.

1. Introduction

When discussing the origins of human communication, Levinson and colleagues [1,2] introduced the idea of a human *interaction engine*. This metaphorical engine is assembled from a range of social-interactional parts that, when put together, enable uniquely human forms of communication, including conventional language. Each part was assumed to have deep roots in our evolutionary history and might therefore—in one form or the other—also be found in other primates. Inspired by these ideas, this paper introduces a computational model that specifies the role that social-interactional processes play in great ape and human communication.

What are the parts that the human interaction engine is built from? First and foremost, human communication is seen as intentional. Senders produce signals to convey intentions and receivers use these signals to infer the sender's intentions [3–6]. As such, communication is deeply linked to reasoning about mental states. Signals, including conventional language, are used to express intentions but the link between signals and intentions is not rigid. There is always residual ambiguity that requires communicators to make additional (pragmatic) inferences—a second key feature of human communication. Such

inferences are licensed by a set of assumptions that humans hold about the nature of communication and social interaction more broadly. One such assumption is that communication occurs within some form of common ground—a shared body of knowledge and beliefs that builds up during social interaction and serves as the background against which signals are interpreted [7,8]. Another assumption is that communication is cooperative such that senders choose their signals so that the receiver is more likely to infer the underlying intention [9]. The receiver takes this into account when interpreting the signal.

The engine assembled from these—and many other—parts is independent of any particular modality. Multimodality is seen as the norm, not an exception in human communication. The system is also highly flexible. Sometimes a tiny hand gesture might be enough to get a message across; at other times, the same meaning might require a long, elaborate utterance comprised of multiple signals that are combined according to conventional rules (grammar). Or as Levinson & Holler [2] put it, ‘The system remains highly flexible, allowing us to shift the burden from words to gestures as required by the current communicative needs.’ Many roads lead to Rome in human communication and *what* works *when* depends on the social-interactional embedding. The system is also independent of the availability of conventional (or evolved) signals. Conventional language is assumed to rely on the engine in just the same way as non-conventional communication. New signals can be invented and understood on the spot and later even conventionalized into new languages [10–18].

The picture that emerges here provides an interesting starting point for an evolutionary research program because it decouples human communication from conventional language. The idea is that there is probably no direct link between the kinds of signals our ancestors used (which might be comparable to what we see in great apes) and human language. The link lies in *how* signals are used, that is, the social and cognitive underpinnings of communication. Once the interaction engine was in place, our ancestors started using and creating signals that, via intermediate proto-languages, evolved to become what we today see as conventional languages [19–23]. Thus, in addition to looking for structural features in animal communication that directly resemble aspects of conventional language (e.g. arbitrary sound-to-meaning mappings or combinatorial syntax [24–28]), comparative researchers can also ask which social-interactional processes underlie communication in other animals. In the next section, we will briefly summarize research in this tradition, with a focus on great ape communication.

2. A comparative approach to human language: the intentional nature of great ape communication

It is beyond the scope of this paper to give a comprehensive summary of existing research on primate communication. We will focus on two aspects that have received considerable attention in comparative research: signallers’ intentional signal production and receivers’ extraction of the intended meaning of a signal. We will show that research on these two aspects of great ape communication varies drastically

depending on whether the focus is on vocal, gestural, or facial signals. To make matters worse, there are also marked differences between research on the production versus the perception or comprehension of signals.

To identify acts of intentional communication in great apes and other non-human primates, Leavens *et al.* [29] suggested a set of criteria derived from research on pre-linguistic communication in human infants [30]. These include the sender’s sensitivity to the presence of other individuals, visual orienting behaviour and monitoring of the receiver, the adjustment of signal use to the receiver’s attentional state and the use of attention-getting behaviours if receivers are not visually attending. Finally, senders are expected to continue signaling and to elaborate signal use in case initial communicative attempts fail.

There is now ample evidence that great apes are intentional communicators in that sense, not only in the gestural modality [31,32]. For example, several species of great apes adjust their signal use to the attentional state of the receiver and only deploy visual gestures if the receiver is attending [29,33]. They also wait for a response and persist in their communicative attempts and might even elaborate their gesture use if the receiver does not react [29,34,35]. Sumatran orangutans use gestures and also some facial expressions flexibly to achieve a variety of social goals [36,37]. Furthermore, wild chimpanzees are more likely to produce alarm calls when other individuals are unaware of a potential threat [38,39].

However, which and how many of the criteria for intentional communication are applied does not only vary across studies but also across modalities [31]. While intentional use is an integral part of defining a gesture, until more recently, this aspect was not considered important in vocal and facial research [40], resulting in the common but unjustified dichotomy between intentional gestures and emotional vocalizations and facial expressions [6].

The different theoretical and methodological approaches in vocal, gestural and facial research have serious downstream consequences for research on primate communication more broadly. Gesture researchers focus on the behaviour of the sender because of the importance of intentional signal production, while vocal and to a lesser extent also facial researchers focus on signal perception and how receivers extract a signal’s meaning. Vocal researchers, for example, frequently use playback experiments to study receivers’ reactions to a very specific call to identify the meaning or function of this call [41]. As a consequence, vocal researchers are interested in context-specific signals, with very specific meanings, while gesture researchers investigate the flexible use of one signal across different contexts and argue that the information conveyed by a gesture might differ depending on the context in which it is used. Gesture researchers further largely ignore context-specific signals because this would not fulfil the criterion of flexible usage, which is often considered an additional marker of intentional use [31,36].

Meaning is also conceptualized very differently across modalities, depending on whether the focus is on the signaler’s or receiver’s behaviour [40]. While gesture researchers focus on the message the signaler intends to communicate, vocal (and partly also facial) researchers focus on the ‘meaning’ extracted by the receiver [42,43]. As a consequence, it is difficult—if not impossible—to compare findings across modalities with regard to how non-human primates’ communicative interactions are shaped by contextual information

and how they ‘make sense’ of others’ communicative attempts. Only more recently has there been some cross-fertilization in both vocal and gesture research. Vocal researchers report that some vocalizations are less context-specific than previously thought [44], while gesture researchers started to assign specific meanings to individual gestures [45,46].

Despite these recent developments, it is important to highlight that research on primate communication has almost exclusively used a uni-modal approach: the majority of research focused either on gestural, vocal or facial signals, and only very few studies investigated more than one signal modality simultaneously [47–51]. There are a number of different reasons why researchers artificially break up the communicative process into components and study each of them in isolation [52]. For example, researchers are trained in the theoretical approach and methods of their focal modality; methods used to study one modality (e.g. playback experiments) are not easily applicable to another modality.

There is, however, a deeper and more fundamental problem: we lack a theoretical account of how the different components integrate with one another. For human communication, Enfield [53], for example, proposed that composite utterances, incorporating multiple signals of multiple types, ‘[...] are interpreted through the recognition and bringing together of these multiple signs under a pragmatic unity heuristic or co-relevance principle, i.e. interpreter’s steadfast presumption of pragmatic unity despite semiotic complexity’. In other words, the recognition of each component’s (encoded) meaning is enriched by (the interpretation of) additional information, such as the meaning provided by the context in which this utterance is embedded. For primate communication, an equivalent theoretical account is still missing and many of the following questions remain unsolved. How do different signals relate to one another? That is, how does the combination of a gesture with another signal (e.g. gesture, facial expression or vocalization) change the meaning or usage of the initial gesture? What role does the social context play? Our goal for the rest of the paper is to sketch out such a theoretical account in the form of a computational model. As a first step, we will briefly introduce the Rational Speech Act (RSA) framework that formalizes some of the reasoning processes implied by the interaction engine and from which we took inspiration.

3. Computational models of inferential communication in humans

A core challenge for a multi-layered, multi-modal system is to specify how the different information sources—the aspects of the utterance and the context that relate to the message being communicated—flow together [53–56]. The RSA framework sees communication as a socially guided inference process [57,58]. A hypothetical receiver in the model is assumed to reason about the intention that underlies the sender’s production of an utterance in context.¹ Importantly, the receiver assumes that the sender is communicating in a cooperative way, choosing utterances that are maximally informative for the receiver given the context. This assumption allows the receiver to go beyond the literal meaning of the words that are used and to make pragmatic inferences.

The RSA framework has been successfully used to model a range of language understanding phenomena as pragmatic inferences including scalar and *ad hoc* implicatures, non-literal language, politeness and vagueness, among others [57,59–63]. More recently, it has been used to predict how adults and children integrate different information sources to make inferences about what a sender is referring to [64]. In one study, Bohn *et al.* [65] measured children’s developing sensitivity to different information sources, for example, their linguistic knowledge or their sensitivity to common ground. Then they used an RSA-type model to predict what should happen when children are confronted with multiple information sources at once. When they compared these predictions to new experimental data, they saw a very close alignment between the two, both qualitatively and quantitatively. To learn more about the integration process itself, they formalized a range of alternative models that varied in their assumptions about which information sources children used and how they integrate them. They found that children’s behaviour was best predicted by a model that assumed rational integration of all available information sources. Interestingly, the integration process was best described as stable across development. That is, even though children might change in how sensitive they are to different information sources, the way they integrate them seems not to change as they develop. These studies illustrate how computational models can be used as a tool to study multi-layered communication.

For the model we describe below, we take inspiration from the RSA framework. The connection is mainly conceptual: we see communication as a socially guided inference process that relies on multiple, context-dependent information sources. There is, however, little structural overlap in terms of the implied cognitive mechanisms. In §6, we explore how the social reasoning processes that are structural characteristics of RSA can be used to explain differences between great ape and human communication when it comes to interpreting novel and ambiguous signals.

4. Formal models of primate communication

Our main goal in this paper is to formulate a computational model of great ape communication. We focus on the in-the-moment comprehension of communicative acts. We ask how a receiver makes inferences about the intentions of a sender based on information contained in the signals that make up an utterance, the relationship between communicative partners, and the social context. The process of in-the-moment comprehension has received little attention in previous modelling work in primate communication. We briefly review some of the earlier literature before laying out our approach.

Most formal work in primate communication has focused on modelling the production of different primate calls [66,67]. Though relevant for answering questions about the evolution of speech, this work does not help us understand the social-interactive nature of primate or ape communication. In a very ambitious project, Stuart Altmann² [68] used stochastic models to predict the socio-communicative behaviour of rhesus monkeys (*Macaca mulatta*). He observed large groups of monkeys living on Cayo Santiago for two years with the goal to develop an ethogram of the species’ social behaviour.

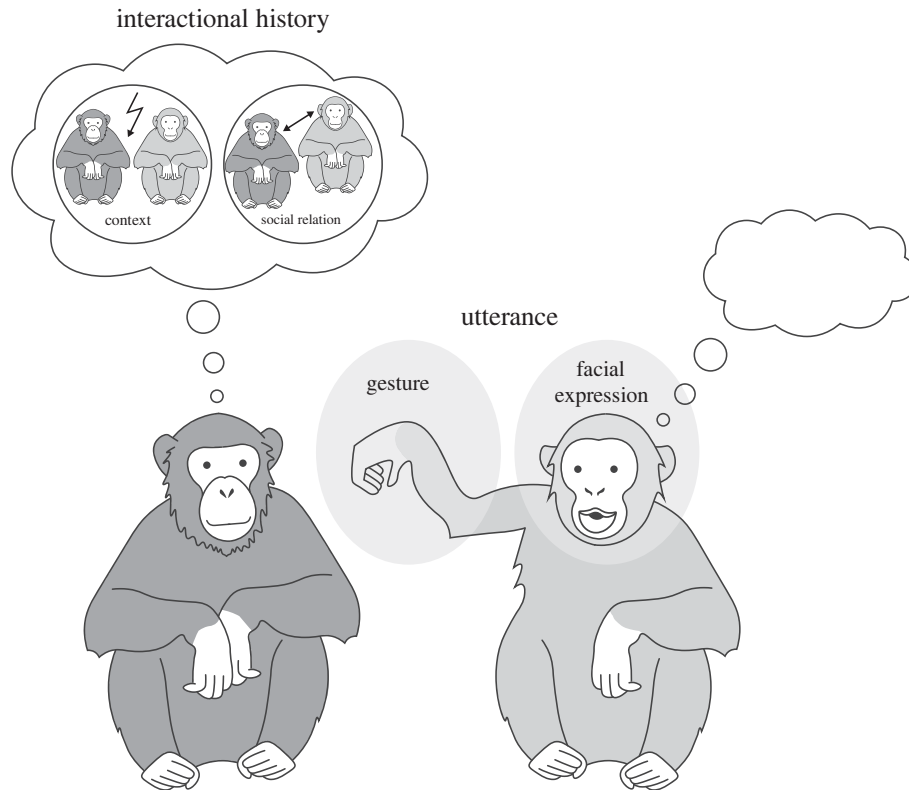


Figure 1. Schematic overview of the computational model. The sender (right) is producing an utterance and the receiver (left) tries to infer the intention of the sender based on the information sources available. The model takes in information provided by the utterance (gesture and facial expression) and the interactional history (immediate social context and dominance relation).

Next, he used his observations to define transitional probabilities between different behaviours. That is, he asked how well one can predict an individual's behaviour if the previous behaviour (by the same or another individual) is known. He did this for pairs of behaviours, but also for longer sequences. Perhaps unsurprisingly, he found that the behavioural stream is not a random sequence of events, but that behaviours cluster in a systematic way. In a very broad sense, we take this as an inspiration to look for a wider set of determinants when trying to predict in-the-moment comprehension and reactions.

Arbib and colleagues [69–72] focused specifically on gestural communication. Their main goal, however, was to model the ontogeny of gestures. Their model shows how behavioural patterns can evolve into communicative gestures during direct, physical interaction. Given their specific aim, the authors saw the gesture as the sole cause of changes in the receiver's behaviour. Comprehension is treated as an associative learning process during which the observation of a particular action becomes paired with a particular reaction (i.e. change in the receiver's goal state). The result is a linear mapping between observing a gesture and producing an outcome. In our model, we loosen this assumption and take into account that multiple information sources influence the response to a gesture.

5. A computational model of chimpanzee communication

In this section, we introduce a Bayesian computational model of great ape communication. In contrast to standard statistical procedures (e.g. linear regression) that describe a particular

dataset, our model describes the inference processes we assume to underlie great apes' interpretation of communicative signals in context. These inference processes are built into the model structure and the model provides an account of the process that generated the data. Such a generative model can be used to predict and explain datasets (see below), but its main purpose is to provide a theoretical account of the phenomenon in question. In what follows, we first present a very general formulation of our model and then further specify it to capture a particular type of communicative interaction. We then evaluate the model based on an existing dataset.

We see great ape communication as a contextualized social inference problem. That is, the sender produces an utterance that the receiver uses to make inferences about the sender's intention (figure 1). Utterances can be composed of different types of signals coming from different modalities (e.g. gestures, vocalizations, facial expressions etc.). Inferences are contextualized in that not just the utterance, but also the social context of the utterance as well as the relationship between the sender and receiver influence the receiver's interpretation. Thus, multiple information sources have to be integrated. We explore the hypothesis that this integration process occurs via a rational Bayesian procedure. This contrasts with the use of the term rational as describing a rule-based (i.e. logical) form of drawing conclusions. Here, we assume that the receiver's a posteriori belief is optimal given the receiver's prior beliefs and the constituent information sources they receive [73–75]. Given the simplicity of our model, we do not assume any limitations with respect to the cognitive resources that our communicative agents have at their disposal. However, our approach could easily be extended in this direction, for

example, with resource-rational considerations [76]. The model is formally defined as

$$P(i | u) \propto P(u | i)P(i), \quad (5.1)$$

with $P(i | u)$ being the probability that the sender has intention i given utterance u . This decomposes into the likelihood of producing an utterance given an intention $P(u | i)$ (e.g. raising one's arm when wanting to be groomed) and the prior probability of having an intention in the first place $P(i)$ (e.g. wanting to be groomed). This very general formulation can be used as a framework to evaluate different hypotheses about which social information sources contribute to the likelihood and the prior; that is, which information sources play an important role in great ape communication.

Next, we spell out one variant of the model, which was in part determined by the dataset that we had available for evaluation. As mentioned above, the general framework could be used with more, fewer, or different information sources. For the purpose of the current paper, the likelihood is defined by the semantics associated with a gesture, $\mathcal{L}(g, i)$, and a facial expression, $\mathcal{L}(f, i)$, which independently contribute to make up the utterance

$$P(u | i) = P(g, f | i) = \mathcal{L}(g, i | \theta_g)\mathcal{L}(f, i | \theta_f). \quad (5.2)$$

Signals have 'soft semantics', that is, in contrast to a truth-functional (Boolean) semantics, we assume a probabilistic mapping between a signal and an intention (defined by the parameters θ_g and θ_f [77]; where θ_g is the strength of association between the gesture and the intention and θ_f that of the facial expression and the intention). The utterance is contextualized by the prior probability of the intention, $P(i)$, which we take to be a function of the context, and the social relation between individuals, $P(i | c, s)$

$$P(i) = P(i | c, s) = \rho_c \rho_s. \quad (5.3)$$

The direction and strength of the context and social relation components are defined by the parameters ρ_c and ρ_s (where ρ_c denotes the association between the context and the intention and ρ_s that between the social relation and the intention). In the example below, we provide more information about the interpretation of these parameters.

To evaluate the model, we used it to predict the outcome of communicative interactions between semi-wild-living chimpanzees (*Pan troglodytes*). The data are taken from the study by Oña *et al.* [50] in which the authors observed two groups of chimpanzees (72 individuals) living in the Chimfunshi Wildlife Orphanage Trust in Zambia. They investigated if signal combinations were used in different contexts and/or elicited different responses compared to signals used alone. For every communicative interaction, they recorded the signals the sender produced, the context in which they were used and the reaction of the receiver. More specifically, they coded the type of manual gesture using a form-based coding scheme, differentiating between morphological configurations of the joints of the arm, hand and fingers. Using this procedure, they identified two frequently occurring gesture types: *stretched-arm*, consisting of an extended arm with both the arm and hand stretched, and *bent-arm*, with either hand or forearm bent and the back of the hand or arm directed at the receiver. Facial expressions were coded using a modified version of the human Facial Action Coding Scheme (FACS) [78] developed

to identify facial movements of chimpanzees (chimpFACS) [79]. The *bared-teeth* face, with the mouth either closed or slightly opened and the mouth corners laterally retracted and teeth fully exposed, was identified in addition to the *funneled-lip* face, consisting of an open, rounded mouth with protruded lips. When one of the gestures was combined with either of these facial expressions, this was considered a gesture-facial expression combination. When the gesture was used without a facial expression, the face was coded as neutral. Facial expressions produced in isolation, without an accompanying gesture, were not included. The social context of the interaction was coded as either positive (e.g. greeting, grooming, play) or negative (e.g. physical conflicts, harassment). The social relationship between the sender and receiver was considered by coding whether signals were directed towards a lower- or higher-ranking individual. Finally, the outcome of the interaction (i.e. the response of the receiver) was classified as either affiliative (receiver approaches the sender and shows behaviours such as embracing, grooming or play) or avoidant (receiver is avoiding or ignoring the sender, e.g. by turning away from, hitting or pushing the sender).

As noted above, in our model, the gesture and the facial expressions contribute to the utterance (the likelihood) and the social context and the relationship contribute to the prior. We assigned parameter values to each of the components of the communicative interactions. The goal was to show that by choosing intuitive parameter values, our model can give rise to the data we observed. These values range between 0 and 1 and represent the degree to which a component is indicative of a positive (affiliative; 0–0.5) or negative (avoidant; 0.5 = 1) interpretation. We assumed the stretched-arm gesture to be weakly negative ($\theta_{gs} = 0.53$) and the bent-arm gesture to be weakly positive ($\theta_{gb} = 0.47$). Neutral facial expressions were set to be neutral ($\theta_{fn} = 0.5$), bared-teeth expressions were set to be weakly negative ($\theta_{fb} = 0.6$), and funneled-lip expressions to be strongly negative ($\theta_{ff} = 0.9$). A negative context was set to be negative ($\rho_{cn} = 0.7$) and a positive to be positive ($\rho_{cp} = 0.3$). Finally, we assumed that a positive reaction was likely for a dominant sender ($\rho_{sd} = 0.25$) and a negative outcome likely for a subordinate sender ($\rho_{ss} = 0.75$).

We want to highlight that even though these parameter values are inspired by prior work and common sense, they are to some extent arbitrary and should not be taken to reflect a strong commitment to the role the individual components might play in a different context. Their main purpose is to capture the idea that different components of the communicative interaction are more or less associated with a particular response. Ideally—and hopefully in future work—these parameters would be directly estimated based on a training dataset and then used to predict a test dataset. Given the size of the dataset we had available, this approach was not possible here. The code that spells out the model architecture and the processing algorithms and that can be used to reproduce the results is available in the associated online repository: <https://github.com/manuelbohn/RSApes>.

Based on the model and the parameter settings, we generated predictions for all possible combinations of gestures, facial expression, dominance relationship and social context. We compared these predictions to the observations made by Oña *et al.* [50]. Our model makes predictions about the receiver's interpretation of the utterance in context. The

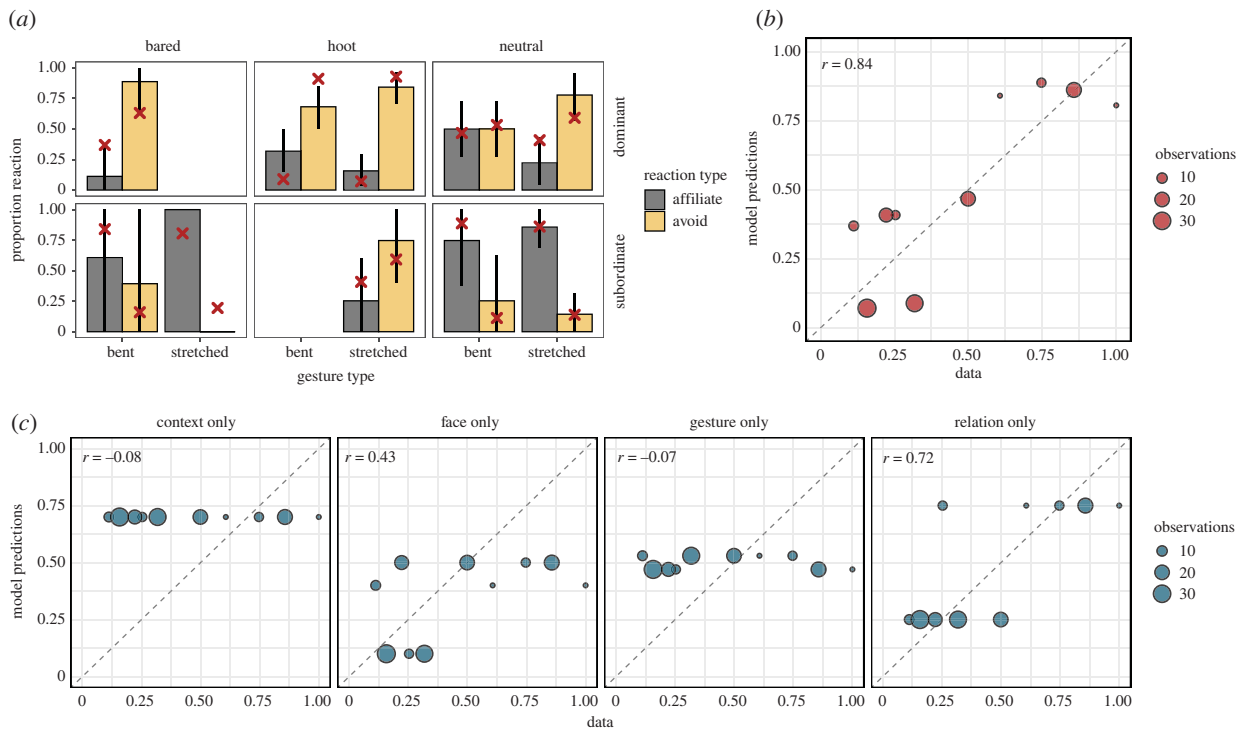


Figure 2. Model predictions compared to data from [50]. (a) The mean proportion (bars) of affiliative and avoidant reactions for combinations of gesture, facial expression, relationship and social context in the data. Only combinations with more than five observations are shown. Error bars are 95% confidence intervals based on a non-parametric bootstrap. Red crosses show model predictions. (b) Correlations between model prediction and data for avoidant reactions. The size of each point is proportional to the number of observations for a particular combination in the data. (c) Correlations for reduced models that focus only on a single component (with all other parameters set to 0.5). (Online version in colour.)

data, however, only recorded the receivers' reactions—as interpreted by the human coders. We assume that the receiver's reaction is guided by their interpretation of the utterance: when inferring a negative intention, the receiver shows an avoidant reaction and when inferring a positive intention, they show an affiliative reaction. Thus, for the purpose of the model comparison, we assume a one-to-one mapping between the interpretation of the sender's message and the receiver's reaction.

Observations in the data were not equally distributed across all possible combinations. To evaluate the model predictions, we focused on combinations that had at least five observations. All combinations that fulfilled this criterion were observed in a negative social context. When we compare the model predictions to the data, we therefore only visualize the negative context (figure 2). Note, however, that our model also generated predictions for the positive context.

In figure 2, we can see that the full model explains the data well, both quantitatively and qualitatively. The model predictions go in the same qualitative direction as the data, predicting more negative reactions when more were observed. Furthermore, many of the model predictions also align quantitatively with the data, resulting in a high correlation between the two (figure 2b). Let us take a closer look at some of these patterns. In most cases, the qualitative pattern in the data was the same for both gesture types. For example, in a negative context (figure 2 only includes the negative context), with a subordinate sender and a neutral facial expression, no matter if a bent or a stretched-arm gesture was used, there were more affiliative reactions. Our model predicts this pattern despite the fact that we took the stretched-arm gesture to be associated with a negative intention. The reason for this is that both gestures were assumed to have weak meanings.

As a consequence, they had very little predictive power when a different, stronger information source (the dominance relationship in this case) was also available.

Next, we used this modelling framework to illustrate the theoretical point made above, namely that a focus on a single aspect of great ape communication is likely to yield an incomplete picture of the interaction. We formulated four reduced models, which use the same parameter settings as above, but selectively focused only on one of the components (all other parameters set to 0.5). When comparing the predictions from these reduced models to the data, we saw that none of them captured the data equally well compared to the full model (figure 2c).³ For example, the models focusing only on the context or the gesture completely fail to capture any structure in the data. These results, however, should be taken with a grain of salt given the—rather arbitrary—way in which we chose the parameter values. Nevertheless, we think the results nicely illustrate how computational modelling can be used as a powerful tool to study great ape communication. In the next section, we explore ways in which we can use this tool to theorize about some potential differences between ape and human communication.

6. Pragmatics as an amplifier

In their description of the interaction engine, Levinson & Holler [2] point out that 'language is the tip of an iceberg riding on a deep infrastructure of communicational abilities'. Part of this deep infrastructure is pragmatics. As noted in §1, the central idea is that utterances are not interpreted at face value, but that receivers go beyond the literal and make inferences about why the sender produced a particular utterance in context. A cornerstone of this reasoning is the assumption

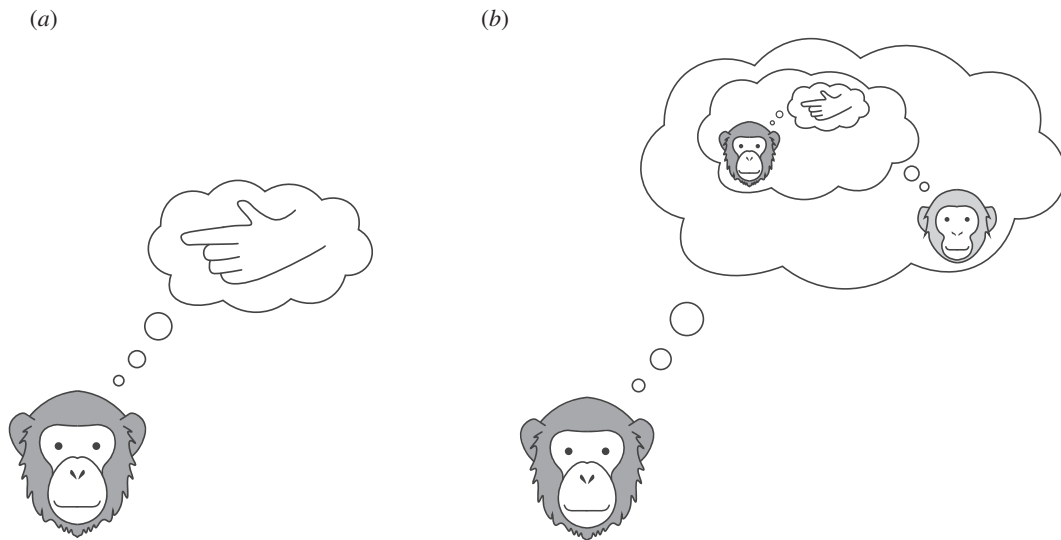


Figure 3. Schematic depiction of the added pragmatic reasoning component. The literal receiver (a) only reasons about the gesture whereas the pragmatic receiver (b) reasons about why the sender produced that particular gesture. The pragmatic receiver further expects the sender to produce the gesture with the goal of being informative.

that the sender is cooperative and informative; they produce utterances that help the receiver to infer their intention.

In the following, we enrich our model of great ape communication by pragmatics—i.e. cooperative social reasoning. From an evolutionary perspective, we may say that our great ape model stands in for the last common ancestor of great apes and humans. To recapitulate, we assume that this ancestor (and modern great apes) rationally integrated different information sources to make inferences about the sender's intentions. This includes information contained in the utterance as well as the social context and the relationship between communicators. The pragmatic abilities are built on top of this basic infrastructure to provide modern human communication.

To evaluate this pragmatically enriched model, we want to focus on some peculiar differences that have been reported for the communicative abilities of great apes and humans. Numerous studies have shown that great apes struggle to spontaneously understand ambiguous signals, for example, pointing or novel iconic gestures [10,80–88] (with some particular exceptions [89,90]). That is, when confronted with a novel gesture or a new context, great apes usually fail to spontaneously use the gesture. These findings are peculiar because these gestures are naturally meaningful in that they either index (pointing) or resemble (iconic gestures) the referent. What is more, human children understand them spontaneously already very early in life [91–93]. Apes also seem to be somewhat sensitive to the natural meaning of these gestures. In the case of pointing, they often look in the direction the experimenter is pointing [94]. And in one study, iconic gestures were learned faster compared to arbitrary ones [95].

Why do apes struggle with spontaneous comprehension of these gestures? The results of the model above can be taken to suggest that the social context and the relationship between sender and receiver play an important role in great ape communication. In the experimental set-ups of studies on pointing or iconic gesture comprehension, these components are controlled for and therefore offer no information about the sender's intention [10,83,86]. Great apes are left with only the gesture. If that gesture was initially only vaguely associated with one or the other outcome, it

would not provide sufficient information for apes to infer the sender's intention and thus to systematically select the referred-to object.

Why do humans spontaneously understand these gestures? We think that the notion of pragmatics as spelled out above can act as an amplifier of vague literal meanings. That is, a human receiver assumes that the sender produced a particular gesture in a cooperative and informative manner to inform them about their intention. The additional social reasoning singles out the gesture as a communicative act that was produced with the sole purpose to express a given intention (figure 3). This line of argument is of course reminiscent of the idea that humans—but not great apes—are sensitive to cooperative communicative intentions [6]. However, we assume that pragmatic inferences are just one information source that can be exploited and that they are graded—not all or nothing. Taken together, the degree to which pragmatic reasoning amplifies a meaning depends on (a) the presence of a social reasoning mechanisms and (b) expectations about how cooperative the sender is. Next, we substantiate these ideas via our modelling framework.

The RSA framework introduced above is built around the assumptions that (a) receivers reason about why senders produce certain utterances and (b) receivers assume that senders communicate in a cooperative and informative way. This social reasoning component is formalized by embedding the model of the (zero-order) literal receiver (short-hand notation: P_{R_0}) in a model of the sender, P_{S_1} . This *pragmatic* sender chooses utterances so that they are informative for the literal receiver, while the literal receiver simply interprets utterances in line with their literal semantics. This literal receiver behaves exactly like in the great ape model (figure 3). This illustrates the way in which our model of human communication is built around our model of great ape communication. At the highest level, we now have a pragmatic receiver, P_{R_1} . These additions change our model as follows:

$$P_{R_1}(i | u) \propto P_{S_1}(u | i)P(i), \quad (6.1)$$

$$P_{S_1}(u | i) \propto P_{R_0}(i | u)^\alpha \quad (6.2)$$

$$\text{and } P_{R_0}(i | u) \propto \mathcal{L}(u, i | \theta_u). \quad (6.3)$$

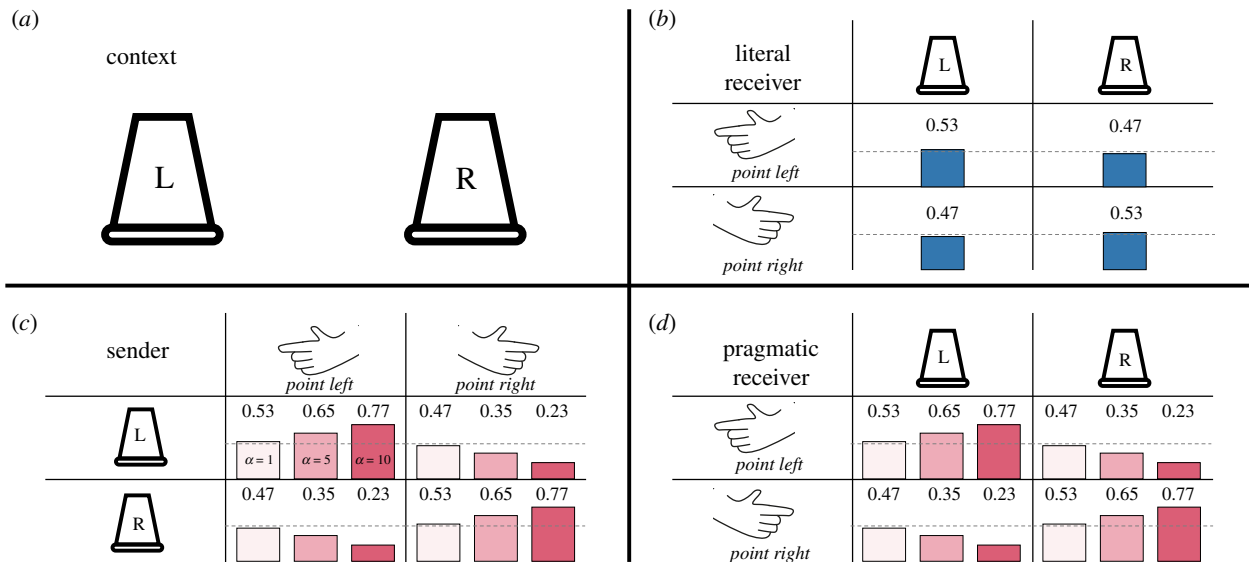


Figure 4. Application of the pragmatically enriched model to an object-choice task with pointing gestures. (a) The context with the two locations (L = left and R = right) that can be referred to. Panel (b) gives the interpretation probabilities of a literal receiver. (c) The production probabilities for the pragmatic sender for values of $\alpha = 1, 5$ and 10 . (d) The interpretation probabilities of the pragmatic receiver based on the production probabilities in (c). Coloured bars visualize the probabilities in reference to chance (grey dashed line). Different shades in (c,d) correspond to the magnitude of α . (Online version in colour.)

Equation (6.2) above shows that the degree to which the sender is assumed to be informative depends on the parameter α . The higher α , the more informative the sender is assumed to be. The effect of α , however, depends on the presence of the sender model, which represents the additional social reasoning component that we think is characteristic of human communication.

When we adapt such a model to a situation in which the receiver is faced with a vaguely meaningful gesture (e.g. a point or an iconic gesture; $\theta_u = 0.53$) without any additional contextual information, we see that the literal interpretation of the gesture simply reflects this vague meaning (figure 4b). We also see that pragmatic reasoning amplifies the initially vague meaning (figure 4d). As noted above, this is not due to the additional social reasoning component alone but critically depends on the receiver's expectation about cooperative communication (the parameter α , figure 4c). This highlights the graded relation between assumptions about cooperativeness and pragmatic inference. Once again, we would like to point out that the specific parameter values we picked here are arbitrary and do not reflect a strong commitment to how great apes or humans interpret pointing gestures. They simply serve to illustrate the point that pragmatics may amplify vague natural meanings.

7. Implications and future directions

With the modelling exercise presented above we had two overarching goals. The first was to show that great ape communication is best thought of (and studied) as a multi-faceted, multi-modal, social inference process. We saw that the outcome of a communicative interaction was best predicted when signals, as well as contextual components, were taken into account. We do not say that studying these components in isolation is fruitless, but we do emphasize that focusing exclusively on, for example, the gesture or vocalization produced makes it less likely that the unfolding interaction will be understood. From our perspective, the different components play complementary roles in an integrated inference process.

Our hope is that our model proves to be a useful tool—or at least an inspiration—for future research. The approach by Oña *et al.* [50], in which many different aspects of a communicative interaction are coded, seems to be especially promising. Such work could easily be done using already existing video recordings. Models like the one presented here could then be used to specify how the different components work together. In addition, our framework provides a new way to test competing hypotheses. Instead of relying on qualitative predictions, alternative hypotheses can be formalized as alternative models and then directly compared in a quantitative way. Across studies, it would be interesting to see if general patterns emerge. For example, models that emphasize social-contextual components could make better predictions compared to models emphasizing information provided by the utterance. Or models prioritizing facial expressions could be found to outcompete models that more strongly emphasize gestures. Or *vice versa* in both cases. Experimental studies could gradually vary the information provided by signals and the social context to examine how they trade-off with one another. Such an approach might reveal quantitative differences between humans and other primates where we currently assume qualitative ones. In all of this, we think that the study of great ape communication would benefit from an interdisciplinary approach in which computational modellers work together with primatologists and comparative psychologists. Hopefully, this will allow the field to move away from asking somewhat artificial questions about the importance of individual gestures, facial expressions or vocalizations and instead move towards more comprehensive theories of the actual processes that underlie communicative interactions.

We see our model as a first step that needs to be expanded in the future. The process that we capture in our model is in-the-moment comprehension, which is only a part of communicative interaction. An easy extension would be to look at the sender: we assume our model to be symmetric and so it could be easily used to generate predictions about what types of gestures, facial expressions and vocalizations the

sender should produce in different contexts given the intention they want to communicate. Furthermore, it would be interesting to extend our model to capture the temporal dynamics of communication—that is, to include mechanisms that are used to clarify or emphasize a message. Candidate behaviours in primates could be acts of persistence, repetition or elaboration that are often seen in naturalistic and experimental settings [29,35]. Including this aspect might have consequences for the cognitive architecture of the model. For example, van Arkel *et al.* [96] have suggested that a simple repair mechanism drastically changes the computational demands in human communication.

Our second goal was to demonstrate how pragmatic reasoning can act as a gradual amplifier for signals with vague meanings. This perspective might be helpful for theorizing about the gradual transition from animal to human communication. For example, Sterelny [22] has argued that the transition from animal to human communication involved shifting from code-based to ostensive inferential communication [22,97]. During this process, the tight signal–response coupling characteristic for code-based communication was loosened. This brought an increase in flexibility, allowing senders to use the same signal for different and potentially novel purposes. However, it also introduced ambiguity to the signal, which, according to Sterelny, was compensated by relying on social reasoning processes. This transition shifted the locus of selection from specific signal–response couplings to communicative behaviour more broadly, with downstream consequences for other forms of cooperative interaction [9]. Our model formalizes the trade-off between ambiguity in the signal—which is characteristic of human communication [21,98]—and social reasoning. As such, it could be used as a starting point to formalize the gradual evolution of human ostensive-inferential communication.

The gradual emergence of pragmatic social reasoning in the evolution of human communication might have had further downstream consequences for the emergence of conventional communication systems. Recently, Hawkins *et al.* [99] embedded an RSA model of pragmatic in-the-moment inferences in a model of convention formation and showed how signals with vague meanings can give rise to conventional communication systems. The meaning of a signal can get fixed (e.g. further amplified) when it is repeatedly used within dyadic communicative interactions. Conventions form when partner-specific communicative conventions are gradually transferred, via a hierarchical Bayesian model, to novel communicative partners. Work by Woensdregt *et al.* [100] suggests that the presence of conventional communication systems further facilitates in-the-moment inferences about communicative intentions, leading to a cascading coevolution of conventional communication systems and social reasoning.

Finally, our modelling approach informs discussions about the modality in which human language has evolved. For decades, there has been a strong divide between researchers arguing for a vocal or a gestural origin of language [20,47,52,101]. Recently, the idea that language origins were

multi-modal has gained traction [47,101]. Our model provides a way of thinking about multi-modal communication. The model does not make any principled distinction between different modalities: for every signal, it simply asks how indicative it is for different intentions the sender might have. This explains how different signals influence each other during in-the-moment comprehension and could also be used to investigate how the burden may have shifted between modalities during the course of evolution.

8. Conclusion

Inspired by work on the human interaction engine, we have described a computational approach for how to study great ape communication in context. Our model assumes that great apes rationally integrate different information sources to make inferences about the intention behind a sender's utterance in context. Using existing data, we have shown that our model makes accurate predictions about the outcome of multi-modal communicative interactions between chimpanzees in different social contexts. Based on the idea that pragmatic reasoning—social reasoning paired with assumptions about cooperative communication—acts as an amplifier for vague meanings, we suggested an explanation for some peculiar differences between the ways that great apes and humans interpret ambiguous signals. This approach illustrates some deep similarities between human and great ape communication, but also specifies in what way the human interaction engine might be equipped with some special parts.

Data accessibility. All data and model code are available in a public online repository: <https://github.com/manuelbohn/RSApes>.

Conflict of interest declaration. We declare we have no competing interests.

Authors' contributions. M.B.: conceptualization, data curation, formal analysis, investigation, methodology, project administration, software, visualization, writing—original draft, writing—review and editing; K.L.: data curation, methodology, writing—original draft, writing—review and editing; L.O.: methodology; M.H.T.: methodology, visualization, writing—original draft, writing—review and editing.

All authors gave final approval for publication and agreed to be held accountable for the work performed therein.

Funding. No funding has been received for this article.

Acknowledgements. We thank Matthias Allritz, Mark Dingemans and an anonymous reviewer for their helpful comments on previous versions of this paper.

Endnotes

¹The RSA framework usually uses *speaker* and *listener* to describe the agents involved. Here, we continue to use the terms *sender* and *receiver* instead to be more inclusive of non-human and human multi-modal communication.

²We are grateful to David Leavens for pointing us to Altmann's work.

³In the online repository, we also include a model in which the strength of the meaning of gestures and facial expressions was switched. That is, gestures were assumed to have a rather strong meaning and facial expressions a weak one. This model makes worse qualitative and quantitative predictions compared to one presented in the paper.

References

1. Levinson SC. 2006 On the human 'interactional engine'. In *Roots of human sociality: culture, cognition and interaction* (eds N Enfield, S Levinson), pp. 39–69. Oxford, UK: Berg.
2. Levinson SC, Holler J. 2014 The origin of human multi-modal communication. *Phil.*

- Trans. R. Soc. B* **369**, 20130302. (doi:10.1098/rstb.2013.0302)
3. Grice HP. 1991 *Studies in the way of words*. Cambridge, MA: Harvard University Press.
 4. Levinson SC. 2000 *Presumptive meanings: the theory of generalized conversational implicature*. Cambridge, MA: MIT Press.
 5. Sperber D, Wilson D. 2001 *Relevance: communication and cognition*, 2nd edn. Cambridge, MA: Blackwell Publishers.
 6. Tomasello M. 2008 *Origins of human communication*. Cambridge, MA: MIT Press.
 7. Bohn M, Köymen B. 2018 Common ground and development. *Child Dev. Perspect.* **12**, 104–108. (doi:10.1111/cdep.12269)
 8. Clark HH. 1996 *Using language*. Cambridge, UK: Cambridge University Press.
 9. Melis AP, Rossano F. 2022 When and how do non-human great apes communicate to support cooperation? *Phil. Trans. R. Soc. B* **377**, 20210109. (doi:10.1098/rstb.2021.0109)
 10. Bohn M, Call J, Tomasello M. 2019 Natural reference: a phylo- and ontogenetic perspective on the comprehension of iconic gestures and vocalizations. *Dev. Sci.* **22**, e12757. (doi:10.1111/desc.12757)
 11. Bohn M, Kachel G, Tomasello M. 2019 Young children spontaneously recreate core properties of language in a new modality. *Proc. Natl Acad. Sci. USA* **116**, 26 072–26 077. (doi:10.1073/pnas.1904871116)
 12. Brentari D, Goldin-Meadow S. 2017 Language emergence. *Annu. Rev. Linguist.* **3**, 363–388. (doi:10.1146/annurev-linguistics-011415-040743)
 13. Fay N, Lister CJ, Ellison TM, Goldin-Meadow S. 2014 Creating a communication system from scratch: gesture beats vocalization hands down. *Front. Psychol.* **5**, 354. (doi:10.3389/fpsyg.2014.00354)
 14. Goldin-Meadow S, Feldman H. 1977 The development of language-like communication without a language model. *Science* **197**, 401–403. (doi:10.1126/science.877567)
 15. Sandler W, Meir I, Padden C, Aronoff M. 2005 The emergence of grammar: systematic structure in a new language. *Proc. Natl Acad. Sci. USA* **102**, 2661–2665. (doi:10.1073/pnas.0405448102)
 16. Senghas A, Kita S, Özyürek A. 2004 Children creating core properties of language: evidence from an emerging sign language in Nicaragua. *Science* **305**, 1779–1782. (doi:10.1126/science.1100199)
 17. Lister CJ, Burtenshaw T, Walker B, Ohan JL, Fay N. 2021 A cross-sectional test of sign creation by children in the gesture and vocal modalities. *Child Dev.* **92**, 2395–2412. (doi:10.1111/cdev.13587)
 18. Kempe V, Gauvrit N, Panayotov N, Cunningham S, Tamariz M. 2021 Amount of learning and signal stability modulate emergence of structure and iconicity in novel signaling systems. *Cogn. Sci.* **45**, e13057. (doi:10.1111/cogs.13057)
 19. Bar-on D. 2013 Origins of meaning: must we 'go gricean'? *Mind Lang.* **28**, 342–375. (doi:10.1111/mila.12021)
 20. Fitch WT. 2010 *The evolution of language*. Cambridge, UK: Cambridge University Press.
 21. Pleyer M. 2017 Protolanguage and mechanisms of meaning construal in interaction. *Lang. Sci.* **63**, 69–90. (doi:10.1016/j.langsci.2017.01.003)
 22. Sterelny K. 2017 From code to speaker meaning. *Biol. Phil.* **32**, 819–838. (doi:10.1007/s10539-017-9597-8)
 23. Planer R, Sterelny K. 2021 *From signal to symbol: the evolution of language*. Cambridge, MA: MIT Press.
 24. Arbib MA, Liebal K, Pika S. 2008 Primate vocalization, gesture, and the evolution of human language. *Curr. Anthropol.* **49**, 1053–1076. (doi:10.1086/593015)
 25. Arnold K, Zuberbühler K. 2006 Semantic combinations in primate calls. *Nature* **441**, 303–303. (doi:10.1038/441303a)
 26. Boë L-J, Fagot J, Perrier P, Schwartz J-L. 2017 *Origins of human language: continuities and discontinuities with nonhuman primates*. Berlin, Germany: Peter Lang International Academic Publishers.
 27. Fedurek P, Slocombe KE. 2011 Primate vocal communication: a useful tool for understanding human speech and language evolution? *Hum. Biol.* **83**, 153–173. (doi:10.3378/027.083.0202)
 28. Zuberbühler K. 2005 The phylogenetic roots of language: evidence from primate communication and cognition. *Curr. Dir. Psychol. Sci.* **14**, 126–130. (doi:10.1111/j.0963-7214.2005.00357.x)
 29. Leavens DA, Russell JL, Hopkins WD. 2005 Intentionality as measured in the persistence and elaboration of communication by chimpanzees (*Pan troglodytes*). *Child Dev.* **76**, 291–306. (doi:10.1111/j.1467-8624.2005.00845.x)
 30. Bates E, Benigni L, Bretherton I, Camaioni L, Volterra V. 1979 *The emergence of symbols: cognition and communication in infancy*. New York, NY: Academic Press.
 31. Liebal K, Waller BM, Burrows AM, Slocombe KE. 2014 *Primate communication: a multimodal approach*. Cambridge, UK: Cambridge University Press.
 32. Townsend SW *et al.* 2017 Exorcising Grice's ghost: an empirical approach to studying intentional communication in animals. *Biol. Rev.* **92**, 1427–1433. (doi:10.1111/brv.12289)
 33. Poss SR, Kuhar C, Stoinski TS, Hopkins WD. 2006 Differential use of attentional and visual communicative signaling by orangutans (*Pongo pygmaeus*) and gorillas (*Gorilla gorilla*) in response to the attentional status of a human. *Am. J. Primatol.* **68**, 978–992. (doi:10.1002/ajp.20304)
 34. Cartmill EA, Byrne RW. 2007 Orangutans modify their gestural signaling according to their audience's comprehension. *Curr. Biol.* **17**, 1345–1348. (doi:10.1016/j.cub.2007.06.069)
 35. Roberts AI, Vick S-J, Buchanan-Smith HM. 2013 Communicative intentions in wild chimpanzees: persistence and elaboration in gestural signalling. *Anim. Cogn.* **16**, 187–196. (doi:10.1007/s10071-012-0563-1)
 36. Call J, Tomasello M. 2007 *The gestural communication of apes and monkeys*. New York, NY: Lawrence Erlbaum Associates.
 37. Liebal K, Pika S, Tomasello M. 2006 Gestural communication of orangutans (*Pongo pygmaeus*). *Gesture* **6**, 1–38. (doi:10.1075/gest.6.1.02lie)
 38. Crockford C, Wittig R, Mundry R, Zuberbühler K. 2012 Wild chimpanzees inform ignorant group members of danger. *Curr. Biol.* **22**, 142–146. (doi:10.1016/j.cub.2011.11.053)
 39. Crockford C, Wittig R, Zuberbühler K. 2017 Vocalizing in chimpanzees is influenced by social-cognitive processes. *Sci. Adv.* **3**, e1701742. (doi:10.1126/sciadv.1701742)
 40. Liebal K, Oña L. 2018 Mind the gap—moving beyond the dichotomy between intentional gestures and emotional facial and vocal signals of nonhuman primates. *Interact. Stud.* **19**, 121–135. (doi:10.1075/is.17040.lie)
 41. Fischer J, Noser R, Hammerschmidt K. 2013 Bioacoustic field research: a primer to acoustic analyses and playback experiments with primates. *Am. J. Primatol.* **75**, 643–663. (doi:10.1002/ajp.22153)
 42. Font E, Carazo P. 2010 Animals in translation: why there is meaning (but probably no message) in animal communication. *Anim. Behav.* **80**, e1–e6. (doi:10.1016/j.anbehav.2010.05.015)
 43. Smith WJ. 1965 Message, meaning, and context in ethology. *Am. Nat.* **99**, 405–409. (doi:10.1086/282382)
 44. Wheeler BC, Fischer J. 2012 Functionally referential signals: a promising paradigm whose time has passed. *Evol. Anthropol.: Issues News Rev.* **21**, 195–205. (doi:10.1002/evan.21319)
 45. Hobaiteer C, Byrne RW. 2014 The meanings of chimpanzee gestures. *Curr. Biol.* **24**, 1596–1600. (doi:10.1016/j.cub.2014.05.066)
 46. Graham KE, Hobaiteer C, Ounsley J, Furuichi T, Byrne RW. 2018 Bonobo and chimpanzee gestures overlap extensively in meaning. *PLoS Biol.* **16**, e2004825. (doi:10.1371/journal.pbio.2004825)
 47. Slocombe KE, Waller BM, Liebal K. 2011 The language void: the need for multimodality in primate communication research. *Anim. Behav.* **81**, 919–924. (doi:10.1016/j.anbehav.2011.02.002)
 48. Wilke C, Kavanagh E, Donnellan E, Waller BM, Machanda ZP, Slocombe KE. 2017 Production of and responses to unimodal and multimodal signals in wild chimpanzees, *Pan troglodytes schweinfurthii*. *Anim. Behav.* **123**, 305–316. (doi:10.1016/j.anbehav.2016.10.024)
 49. Hobaiteer C, Byrne RW, Zuberbühler K. 2017 Wild chimpanzees' use of single and combined vocal and

- gestural signals. *Behav. Ecol. Sociobiol.* **71**, 1–13. (doi:10.1007/s00265-017-2325-1)
50. Oña LS, Sandler W, Liebal K. 2019 A stepping stone to compositionality in chimpanzee communication. *PeerJ* **7**, e7623. (doi:10.7717/peerj.7623)
51. Fröhlich M, Bartolotta N, Frysns C, Wagner C, Momon L, Jaffrezic M, Mitra Setia T, van Noordwijk MA, van Schaik CP. 2021 Multicomponent and multisensory communicative acts in orang-utans may serve different functions. *Commun. Biol.* **4**, 1–13. (doi:10.1038/s42003-020-01566-0)
52. Liebal K, Slocombe KE, Waller BM. 2022 The language void 10 years on: multimodal primate communication research is still uncommon. *Ethol. Ecol. Evol.* **34**, 1–14. (doi:10.1080/03949370.2021.2015453)
53. Enfield NJ. 2009 *The anatomy of meaning: speech, gesture, and composite utterances*. Cambridge, UK: Cambridge University Press.
54. Holler J, Levinson SC. 2019 Multimodal language processing in human communication. *Trends Cogn. Sci.* **23**, 639–652. (doi:10.1016/j.tics.2019.05.006)
55. Vigliocco G, Perniss P, Vinson D. 2014 Language as a multimodal phenomenon: implications for language learning, processing and evolution. *Phil. Trans. R. Soc. B* **369**, 20130292. (doi:10.1098/rstb.2013.0292)
56. Cavicchio F, Dachkovsky S, Leemor L, Shamay-Tsoory S, Sandler W. 2018 Compositionality in the language of emotion. *PLoS ONE* **13**, e0201970. (doi:10.1371/journal.pone.0201970)
57. Frank MC, Goodman ND. 2012 Predicting pragmatic reasoning in language games. *Science* **336**, 998. (doi:10.1126/science.1218633)
58. Goodman ND, Frank MC. 2016 Pragmatic language interpretation as probabilistic inference. *Trends Cogn. Sci.* **20**, 818–829. (doi:10.1016/j.tics.2016.08.005)
59. Goodman ND, Stuhlmüller A. 2013 Knowledge and implicature: modeling language understanding as social cognition. *Top. Cogn. Sci.* **5**, 173–184. (doi:10.1111/tops.12007)
60. Kao JT, Wu JY, Bergen L, Goodman ND. 2014 Nonliteral understanding of number words. *Proc. Natl Acad. Sci. USA* **111**, 12 002–12 007. (doi:10.1073/pnas.1407479111)
61. Lassiter D, Goodman ND. 2017 Adjectival vagueness in a Bayesian model of interpretation. *Synthese* **194**, 3801–3836. (doi:10.1007/s11229-015-0786-1)
62. Tessler MH, Goodman ND. 2019 The language of generalization. *Psychol. Rev.* **126**, 395. (doi:10.1037/rev0000142)
63. Yoon EJ, Tessler MH, Goodman ND, Frank MC. 2020 Polite speech emerges from competing social goals. *Open Mind* **4**, 71–87. (doi:10.1162/opmi_a_00035)
64. Bohn M, Tessler MH, Merrick M, Frank MC. 2022 Predicting pragmatic cue integration in adults' and children's inferences about novel word meanings. *J. Exp. Psychol.: Gen.* (doi:10.1037/xge0001216)
65. Bohn M, Tessler MH, Merrick M, Frank MC. 2021 How young children integrate information sources to infer the meaning of words. *Nat. Hum. Behav.* **5**, 1046–1054. (doi:10.1038/s41562-021-01145-1)
66. De Boer B, Tecumseh Fitch W. 2010 Computer models of vocal tract evolution: an overview and critique. *Adapt. Behav.* **18**, 36–47. (doi:10.1177/1059712309350972)
67. Riede T, Bronson E, Hatzikirou H, Zuberbühler K. 2005 Vocal production mechanisms in a non-human primate: morphological data and a model. *J. Hum. Evol.* **48**, 85–96. (doi:10.1016/j.jhevol.2004.10.002)
68. Altmann SA. 1965 Sociobiology of rhesus monkeys. II: Stochastics of social communication. *J. Theor. Biol.* **8**, 490–522. (doi:10.1016/0022-5193(65)90024-X)
69. Arbib M, Ganesha V, Gasser B. 2014 Dyadic brain modelling, mirror systems and the ontogenetic ritualization of ape gesture. *Phil. Trans. R. Soc. B* **369**, 20130414. (doi:10.1098/rstb.2013.0414)
70. Gasser B, Arbib M. 2019 A dyadic brain model of ape gestural learning, production and representation. *Anim. Cogn.* **22**, 519–534. (doi:10.1007/s10071-018-1228-5)
71. Arbib MA. 2016 Towards a computational comparative neuroprimatology: framing the language-ready brain. *Phys. Life Rev.* **16**, 1–54. (doi:10.1016/j.plrev.2015.09.003)
72. Gasser B, Cartmill EA, Arbib MA. 2014 Ontogenetic ritualization of primate gesture as a case study in dyadic brain modeling. *Neuroinformatics* **12**, 93–109. (doi:10.1007/s12021-013-9182-5)
73. Yuille A, Kersten D. 2006 Vision as Bayesian inference: analysis by synthesis? *Trends Cogn. Sci.* **10**, 301–308. (doi:10.1016/j.tics.2006.05.002)
74. Griffiths TL, Tenenbaum JB. 2006 Optimal predictions in everyday cognition. *Psychol. Sci.* **17**, 767–773. (doi:10.1111/j.1467-9280.2006.01780.x)
75. Chater N, Oaksford M. 1999 Ten years of the rational analysis of cognition. *Trends Cogn. Sci.* **3**, 57–65. (doi:10.1016/S1364-6613(98)01273-X)
76. Lieder F, Griffiths TL. 2020 Resource-rational analysis: understanding human cognition as the optimal use of limited computational resources. *Behav. Brain Sci.* **43**, e1. (doi:10.1017/s0140525x1900061x)
77. Degen J, Hawkins RD, Graf C, Kreiss E, Goodman ND. 2020 When redundancy is useful: a Bayesian approach to 'overinformative' referring expressions. *Psychol. Rev.* **127**, 591. (doi:10.1037/rev0000186)
78. Ekman P, Friesen WV. 1978 *Facial action coding system: a technique for the measurement of facial movement*. Palo Alto, CA: Consulting Psychologists Press.
79. Parr LA, Waller BM, Vick SJ, Bard KA. 2007 Classifying chimpanzee facial expressions using muscle action. *Emotion* **7**, 172–181. (doi:10.1037/1528-3542.7.1.172)
80. Bohn M, Kordt C, Braun M, Call J, Tomasello M. 2020 Learning novel skills from iconic gestures: a developmental and evolutionary perspective. *Psychol. Sci.* **31**, 873–880. (doi:10.1177/0956797620921519)
81. Dezechache G, Bourgeois A, Bazin C, Schlenker P, Chemla E, Maille A. 2019 Orangutans' comprehension of zoo keepers' communicative signals. *Animals* **9**, 300. (doi:10.3390/ani9060300)
82. Hare B, Tomasello M. 2004 Chimpanzees are more skilful in competitive than in cooperative cognitive tasks. *Anim. Behav.* **68**, 571–581. (doi:10.1016/j.anbehav.2003.11.011)
83. Herrmann E, Tomasello M. 2006 Apes' and children's understanding of cooperative and competitive motives in a communicative situation. *Dev. Sci.* **9**, 518–529. (doi:10.1111/j.1467-7687.2006.00519.x)
84. Kirchofer KC, Zimmermann F, Kaminski J, Tomasello M. 2012 Dogs (*Canis familiaris*), but not chimpanzees (*Pan troglodytes*), understand imperative pointing. *PLoS ONE* **7**, e30913. (doi:10.1371/journal.pone.0030913)
85. Tomasello M, Call J, Gluckman A. 1997 Comprehension of novel communicative signs by apes and human children. *Child Dev.* **68**, 1067–1080. (doi:10.2307/1132292)
86. Tempelmann S, Kaminski J, Liebal K. 2013 When apes point the finger: three great ape species fail to use a conspecific's imperative pointing gesture. *Interact. Stud.* **14**, 7–23. (doi:10.1075/is.14.1.02tem)
87. Moore R, Call J, Tomasello M. 2015 Production and comprehension of gestures between orang-utans (*Pongo pygmaeus*) in a referential communication game. *PLoS ONE* **10**, e0129726. (doi:10.1371/journal.pone.0129726)
88. Margiotoudi K, Bohn M, Schwob N, Tagliatalata J, Pulvermüller F, Epping A, Schweller K, Allrutz M. 2022 Bo-NO-bouba-kiki: picture-word mapping but no spontaneous sound symbolic speech-shape mapping in a language trained bonobo. *Proc. R. Soc. B* **289**, 20211717. (doi:10.1098/rspb.2021.1717)
89. Mulcahy NJ, Call J. 2009 The performance of bonobos (*Pan paniscus*), chimpanzees (*Pan troglodytes*), and orangutans (*Pongo pygmaeus*) in two versions of an object-choice task. *J. Comp. Psychol.* **123**, 304–309. (doi:10.1037/a0016222)
90. Lyn H, Russell JL, Hopkins WD. 2010 The impact of environment on the comprehension of declarative communication in apes. *Psychol. Sci.* **21**, 360–365. (doi:10.1177/0956797610362218)
91. Behne T, Carpenter M, Tomasello M. 2005 One-year-olds comprehend the communicative intentions behind gestures in a hiding game. *Dev. Sci.* **8**, 492–499. (doi:10.1111/j.1467-7687.2005.00440.x)
92. Rütter J, Liszkowski U. 2020 Ontogenetic emergence of cognitive reference comprehension. *Cogn. Sci.* **44**, e12869. (doi:10.1111/cogs.12869)
93. Kachel G, Hardecker DJ, Bohn M. 2021 Young children's developing ability to integrate gestural and emotional cues. *J. Exp. Child Psychol.* **201**, 104984. (doi:10.1016/j.jecp.2020.104984)
94. Itakura S. 1996 An exploratory study of gaze-monitoring in nonhuman primates 1. *Jpn. Psychol. Res.* **38**, 174–180. (doi:10.1111/j.1468-5884.1996.tb00022.x)
95. Bohn M, Call J, Tomasello M. 2016 Comprehension of iconic gestures by chimpanzees and human

- children. *J. Exp. Child Psychol.* **142**, 1–17. (doi:10.1016/j.jecp.2015.09.001)
96. van Arkel J, Woensdregt M, Dingemans M, Blokpoel M. 2020 .SL: A simple repair mechanism can alleviate computational demands of pragmatic reasoning: simulations and complexity analysis. In *Proc. of the 24th Conf. on Computational Natural Language Learning, November 2020* (eds R Fernández, T Linzen), pp. 177–194. Stroudsburg, PA: Association for Computational Linguistics.
97. Scott-Phillips TC. 2015 *Speaking our minds*. London, UK: Palgrave Macmillan.
98. Piantadosi ST, Tily H, Gibson E. 2012 The communicative function of ambiguity in language. *Cognition* **122**, 280–291. (doi:10.1016/j.cognition.2011.10.004)
99. Hawkins RD, Franke M, Frank MC, Smith K, Griffiths TL, Goodman ND. 2021 From partners to populations: a hierarchical Bayesian account of coordination and convention. *arXiv*.
100. Woensdregt M, Cummins C, Smith K. 2020 A computational model of the cultural co-evolution of language and mindreading. *Synthese* **199**, 1–39. (doi:10.1007/s11229-020-02798-7)
101. Fröhlich M, Sievers C, Townsend SW, Gruber T, van Schaik CP. 2019 Multimodal communication and language origins: integrating gestures and vocalizations. *Biol. Rev.* **94**, 1809–1829. (doi:10.1111/brv.12535)