

The evolution of speech: vision, rhythm, cooperation

Asif A. Ghazanfar¹ and Daniel Y. Takahashi²

¹Princeton Neuroscience Institute, Departments of Psychology and Ecology & Evolutionary Biology, Princeton University, Princeton, NJ 08544, USA

²Princeton Neuroscience Institute, Department of Psychology, Princeton University, Princeton, NJ 08544, USA

A full account of human speech evolution must consider its multisensory, rhythmic, and cooperative characteristics. Humans, apes, and monkeys recognize the correspondence between vocalizations and their associated facial postures, and gain behavioral benefits from them. Some monkey vocalizations even have a speech-like acoustic rhythmicity but lack the concomitant rhythmic facial motion that speech exhibits. We review data showing that rhythmic facial expressions such as lip-smacking may have been linked to vocal output to produce an ancestral form of rhythmic audiovisual speech. Finally, we argue that human vocal cooperation (turn-taking) may have arisen through a combination of volubility and prosociality, and provide comparative evidence from one species to support this hypothesis.

Introduction

Believing, as I do . . . , that the possession of articulate speech is the grand distinctive character of man . . . , I find it very easy to comprehend that some . . . inconspicuous structural differences may have been the primary cause of the immeasurable and practically infinite divergence of the Human form from the simian strips. – Thomas Huxley [1] (p. 63).

The uniqueness of speech to humans is indisputable, but the question of how it came to be in humans and no other animal remains a source of contention. Did speech evolve gradually via communication precursors in the primate lineage or did it arise ‘spontaneously’ through a fortuitous confluence of genetic and/or neuroanatomical changes found only in humans? Some argue that, unlike traits such as opposable thumbs or color vision, where there is clear evidence for a gradual evolution, speech essentially arose suddenly, almost *de novo*. Even Thomas Huxley, Darwin’s irascible promoter of the theory of evolution by natural selection, found the idea that speech could evolve gradually – with many factors at play – through animal precursors too difficult to swallow. Huxley’s attitude is shared by modern scientists who continue to argue for ‘primary causes’ whereby key changes in one factor were pivotal for our ‘infinite divergence’ from other primates in the realm of

communication. There are advocates for human-specific expression of genes (e.g., *FOXP2* [2]), changes in anatomy (e.g., laryngeal descent [3]), increases in the size of the neocortex or particular neocortical areas [4,5], use of peculiar neural circuitry (e.g., mirror neurons [6], neocortical connections with brainstem nuclei [7]), or expression of unique behavioral precursors (e.g., gestures [8] and cooperation [9]).

Each of these factors may have played an important role in the evolution of human communication, but certainly none can be considered a lynch-pin. This is largely because the problem of speech evolution is one about how a whole suite of features integrate to produce uniquely human vocal output patterns and their perception. That is, similarly to language [10–12], speech is a complex adaptation that evolved in a piecemeal fashion. As such, determining the many substrates required for the evolution of human speech is a difficult task, particularly because most traits thought to give rise to it – the vocal production apparatus and the brain – do not fossilize. We are left with one robust method of inquiry: comparing our vocal behaviors and brain organization with those of other extant mammals, and of primates in particular. Humans have long had a fascination with the utterances of other animals and how their vocal signals may or may not relate to our speech [13]. Even the daring adventurer and master linguist, Sir Richard Burton (1821–1890), could not resist investigating whether monkeys communicated using speech-like vocalizations [14]. Our interest in monkey and other animal vocalizations, and their putative relation to human speech, continues unabated because it is our only path to understanding how human vocal communication evolved.

We will explore three complex phenotypes that are part and parcel of human speech and universal across all languages, but that are typically ignored when considering speech origins: its audiovisual nature, its rhythmicity, and its coordination during conversations. In brief, here are the motivations: (i) speech is produced by making different facial expressions which change the shape of the vocal tract. Not surprisingly, humans recognize the correspondence between vocalizations and the facial postures associated with them. Because speech is thus inherently ‘multisensory’, it is important to investigate the role of facial expressions in the vocalizations of other primates. (ii) One key characteristic of audiovisual speech is that the acoustic output and associated movements of the mouth are both rhythmic and tightly coordinated. Some monkey

Corresponding author: Ghazanfar, A.A. (asifg@princeton.edu).

1364-6613/

© 2014 Elsevier Ltd. All rights reserved. <http://dx.doi.org/10.1016/j.tics.2014.06.004>

vocalizations have similar acoustic rhythmicity but lack the concomitant rhythmic facial motion. This raises the question of how we evolved from a presumptive ancestral acoustic-only vocal rhythm to one that is audiovisual. (iii) Finally, speech is a behavior that occurs between individuals and is thus a cooperative endeavor. Humans take turns during a conversation to be heard clearly and to facilitate social interactions. Because of its importance and obvious communicative advantage, how vocal cooperation evolved is of great interest. We explore one possible evolutionary trajectory – a combination of prosociality and volubility – for the origin of vocal turn-taking, and use data from marmoset monkeys to explore this idea.

Before we begin we would like to address two caveats. First, speech and language are two separable phenomena that need not have evolved in parallel [12,15]. Speech is an audiovisual signaling system whereas language is a system for communicating complex concepts, irrespective of modality (e.g., writing and sign language as well as speech). In this review we focus on the evolution of speech. Nevertheless, because speech is the default signal system for language in all human cultures, its evolution may also have implications for linguistic evolution [12], but we do not explore these implications. The second caveat is that, as in any review on the evolutionary origins of a behavior, our arguments below are only as good as the amount of comparative evidence available (i.e., the number of species tested). Thus, we hope that even if what we suggest seems too speculative it will spur more experiments in other species (and potentially falsify our claims).

On the origins of multisensory speech

As with humans, many of the signals that nonhuman primates (hereafter, *primates*) exchange to mediate social interactions take the forms of facial expressions and vocalizations [16]. Indeed, in anthropoid primates, as social group size grows, the complexity of facial expressions [17] and vocal expressions grows in parallel [18,19]. Although facial and vocal expressions are typically treated separately in most studies, they are in fact often inextricably linked: a vocal expression typically cannot be produced without concomitant movements of the face. When we speak our face moves and deforms around the mouth and other regions [20,21]. These dynamics and deformations lead to a variety of visual motion cues related to the auditory components of speech. In noisy, real-world environments, these visual cues increase speech intelligibility [22,23], increase detection speeds [24], and are hard to ignore – visual cues integrate readily and automatically with auditory speech [25]. In light of this, audiovisual (or ‘multisensory’) speech is really the primary mode of speech perception, and is not a capacity that was simply piggy-backed onto auditory speech perception later in the course of our evolution [26].

If audiovisual speech is our default mode of communication, then this should be reflected in its evolution. Many species integrate audiovisual signals during communication, including frogs [27,28] and spiders [29]. Moreover, any vertebrate organism that produces vocalizations will have a simple, concomitant visual motion in the area of the mouth. However, in the primate lineage both the number

and diversity of muscles innervating the face [30] and the amount of neural control related to facial movement [31,32] increased over time relative to other mammals. This ultimately allowed the production of a greater diversity of facial and vocal expressions in primates [33], with different patterns of facial motion being uniquely linked to different vocal expressions [34,35]. Vocalizations are the result of coordinated movements of the lungs, larynx (vocal folds), and the vocal tract [36,37]. The vocal tract consists of the column of air that extends from the vocal folds to the mouth and nasal passages. Changing the shape of the vocal tract not only allows different sounds to be produced (by modifying the resonance frequencies of the vocal tract), but also results in the predictable deformation of the face around the mouth and other parts of the face [20,34]. Put another way, different facial expressions can result in different-sounding vocalizations.

Given that vocalizations are physically linked to different facial expressions, it is perhaps not surprising that many primates other than humans recognize the correspondence between the visual and auditory components of vocal signals. Both macaque monkeys (*Macaca mulatta*) and chimpanzees (*Pan troglodytes*) recognize auditory–visual correspondences between their vocalizations under various contextual and experiential constraints [38–44]. Although these ‘matching’ experiments show that monkeys and apes can recognize the correspondence between visual and auditory signals, they do not demonstrate directly whether such recognition leads to a behavioral advantage – one that would lead to the natural selection of multisensory processes. In a recent vocal detection study, macaque monkeys were trained to detect auditory, visual, or audiovisual vocalizations embedded in noise as quickly and accurately as possible [45] (Figure 1A). Monkeys exhibited greater accuracy and faster reaction times to audiovisual vocalizations than to unisensory events (Figure 1B), as also seen in humans (Figure 1C). Under these task conditions, monkeys truly integrated faces and voices; that is, they combined them in such a way that behavioral performance was significantly better than either of the unisensory conditions. This was the first evidence for a behavioral advantage for combining faces and voices in a primate [45].

There are also some very important differences in how humans versus primates produce their utterances [37], and these differences further enhance human multisensory communication above and beyond what monkeys can do. One universal feature of speech – typically lacking in at least macaque monkey vocalizations – is its bi-sensory rhythm. That is, when humans speak both the acoustic output and the movements of the mouth are highly rhythmic and tightly correlated with each other [21]. This enhances perception and the parsing of long-duration vocal signals [46]. How did this bisensory speech rhythm evolve?

On the origins of the speech rhythm

Across all languages studied to date, both mouth motion and the acoustic envelope of speech typically exhibit a 3–8 Hz rhythm that is, for the most part, related to the rate of syllable production [21,47]. This 3–8 Hz rhythm is crucial for speech perception. Disrupting the acoustic

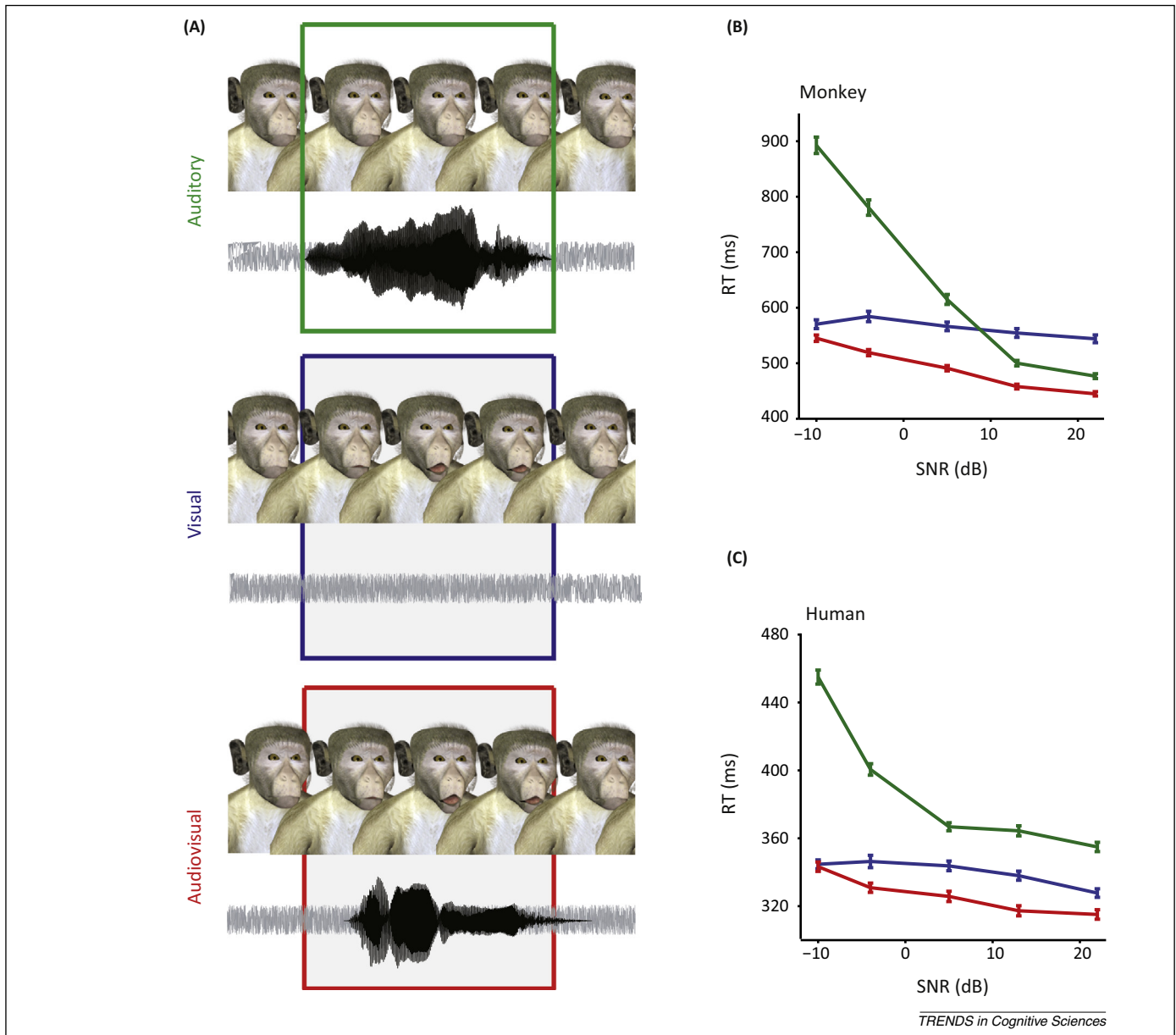


Figure 1. Auditory, visual, and audiovisual vocalization detection. **(A)** Monkeys were trained to detect auditory (green box), visual (blue box), or audiovisual (red box) vocalizations embedded in noise as quickly and as accurately as possible. An avatar and background noise were continuously presented. In the auditory condition, a coo call was presented. In the visual condition, the mouth of the avatar moved without any corresponding vocalization. In the audiovisual, a coo call with a corresponding mouth movement was presented. Each stimulus was presented with four different signal-to-noise ratios (SNR). **(B)** Mean reaction times as a function of SNR for the unisensory and multisensory conditions for one monkey. The color code is the same as in (A); x axes denote SNR in dB; y axes depict reaction times (RT) in ms. **(C)** An analogous experiment with human avatar and speech was carried out in humans. The graph represents the mean reaction times as a function of SNR for the unisensory and multisensory conditions for one individual. Conventions as in (B).

component [48–51] or the visual component arising from facial movements [52] decreases intelligibility. It is thought that the speech rhythm parses the signal into basic units from which information on a finer (faster) temporal scale can be extracted [46]. Given the importance of this rhythm in speech and its underlying neurophysiology [53,54], understanding how speech evolved requires investigating the origins of its bi-sensory rhythmic structure.

Unfortunately, not much is known about the rhythmicity of primate vocalizations. We do know that macaque monkey vocalizations have a similar acoustic rhythmicity to human speech, but without the concomitant and temporally correlated rhythmic facial motion [55]. Modulation spectra analyses of the acoustic rhythmicity of macaque

monkey vocalizations reveal that their rhythmicity is strikingly similar to that of the acoustic envelope for speech [55] (Figure 2A). Both signals fall within the 3–8 Hz range (see also [56] for evidence of shared low-frequency components between macaque monkey calls and speech). Figure 2B shows that, unlike human speech (top panel), macaque coo vocalizations (bottom panel) are typically produced with a single ballistic facial motion – a motion that does not correspond to the amplitude modulation of the produced sound beyond its onset and offset. Thus, one key evolutionary question is – how did we evolve from a presumptive ancestral unisensory, acoustic-only vocal rhythm (Figure 3A) to one that is audiovisual, with both mouth movements and acoustics sharing the same rhythmicity (Figure 3C)?

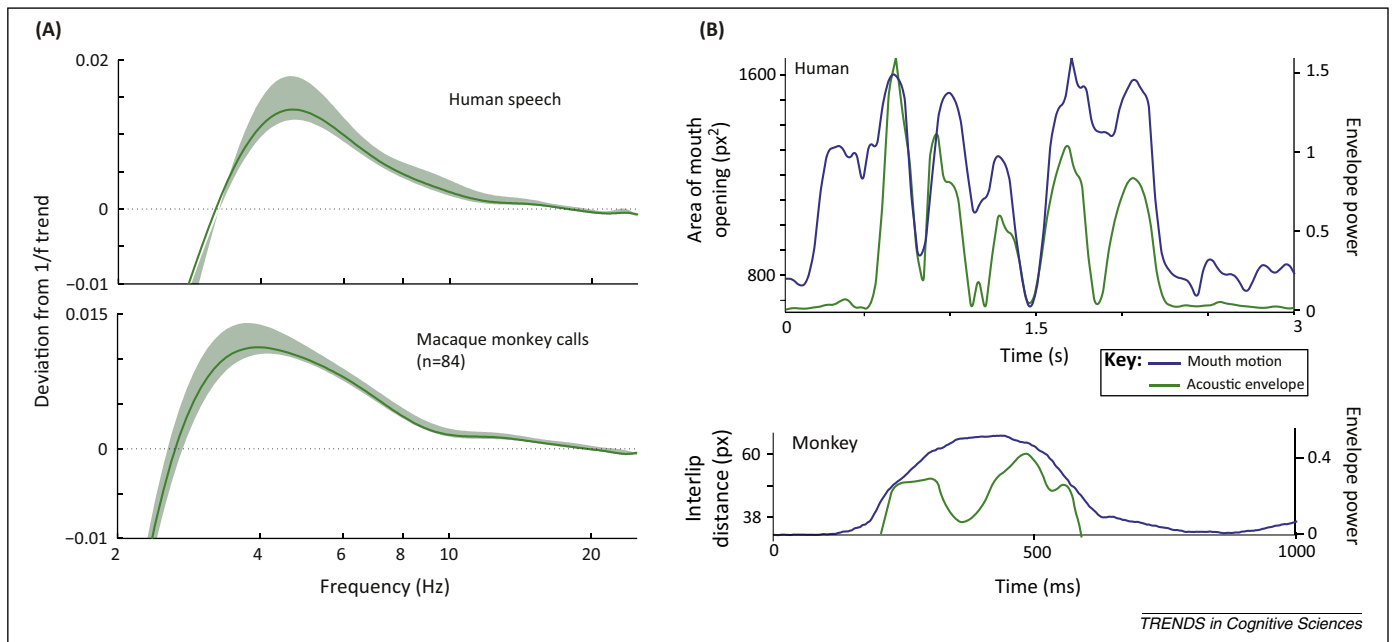


Figure 2. (A) Speech and macaque monkey calls have similar rhythmic structure in their acoustic envelopes. Modulation spectra for human speech and long-duration (>400 ms) macaque monkey calls; x axes represent the frequency in log Hz; y axes depict power deviations from a $1/f$ trend. (B) Mouth motion and auditory envelope for a single sentence produced by human (top panel; the x axis depicts time in s; the y axis on the left depicts the area of the mouth opening in pixel squared; the y axis on the right depicts the acoustic envelope in Hilbert units. The bottom panel shows mouth motion and the auditory envelope for a single coo vocalization produced by a macaque monkey; the x axis depicts time in ms; the y axis on the left depicts the distance between lips in pixels; the y axis on the right depicts the acoustic envelope power in Hilbert units.

One theory posits that the speech rhythm evolved through the modification of rhythmic facial movements in ancestral primates [57] (Figure 3B). In extant primates such facial movements are extremely common as visual communicative gestures. Lip-smacking, for example, is an affiliative signal commonly observed in many genera of primates including virtually every species of Old World monkey [58–61], chimpanzees [62], and in the few New World monkey species whose facial expressions have been studied (common marmosets, *Callithrix jacchus* [63], and capuchins, *Cebus apella* [64]). There are no reports of lip-smacking behavior in prosimian primates [65]. Lip-smacking is characterized by regular cycles of vertical jaw movement, often involving a parting of the lips, but sometimes occurring with closed, puckered lips. Although lip-smacking by both monkeys and chimpanzees is often produced during grooming interactions, macaque monkeys (at least) also exchange lip-smacking bouts during face-to-face interactions [61,66–68]. According to MacNeilage [57], during the course of speech evolution, such non-vocal rhythmic facial expressions were coupled with vocalizations to produce the audiovisual components of babbling-like (i.e., consonant-vowel-like) speech expressions in the human lineage (Figure 3C).

Although direct tests of such an evolutionary hypothesis are usually impossible, in this case one can use the 3–8 Hz rhythmic signature of speech as a foundation to explore its veracity. There are now many lines of evidence that demonstrate that the production of lip-smacking in macaque monkeys is similar to the production of orofacial rhythms during speech. First and foremost, lip-smacking exhibits a speech-like rhythm in the 3–8 Hz frequency range [69]. This rhythmic frequency range is distinct from that of chewing and teeth-grinding (an anxiety-driven expression), although all three rhythmic orofacial motions use

the same effectors. Nevertheless, it still may be that the 3–8 Hz range is large enough that the correspondence between the speech rhythm and the lip-smacking rhythm is not surprising. However, recent evidence from development, X-ray cineradiography, and perception dismiss the possibility that the similarities between lip-smacking and visual speech rhythm are coincidental.

Development

If the underlying mechanisms that produce the rhythm in monkey lip-smacking and human speech are homologous, then their developmental trajectories should be similar [70]. In humans, babbling – the earliest form of rhythmic and voluntary vocal behavior [71–73] – is characterized by the production of canonical syllables that have acoustic characteristics similar to those to adult speech and involves rhythmic sequences of mouth close–open alternation [74–76]. Babbling does not emerge with the same rhythmic structure as adult speech. It starts slower and is more variable. During development, the rhythmic frequency increases from ~ 3 Hz to ~ 5 Hz [21,47,77,78], and the variability of this rhythm starts out very high [77] and does not become fully adult-like until post-pubescence [72]. Importantly, this developmental trajectory from babbling to speech is distinct from that of another cyclical mouth movement, that of chewing. The frequency of chewing movements in humans is highly stereotyped and slow in frequency, remaining unchanged from early infancy into adulthood [79,80]. Chewing movements are often used as a reference movement in speech production studies because, again, both movements use the same effectors.

The developmental trajectory of macaque monkey lip-smacking parallels speech development [81,82]. Measurements of the rhythmic frequency and variability of lip-smacking across neonates, juveniles, and adults revealed

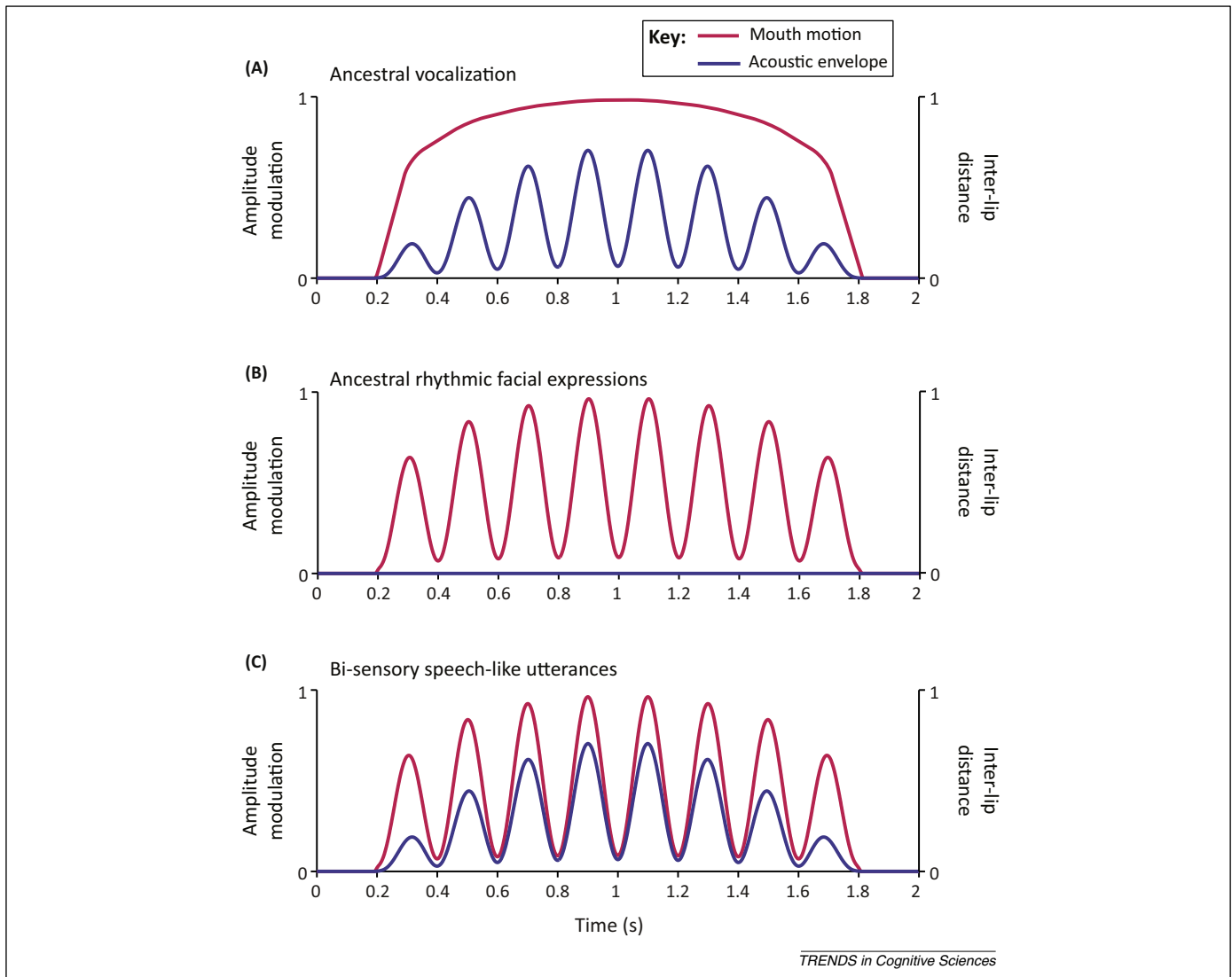


Figure 3. Hypothetical transition from an ancestral unisensory, acoustic-only vocal rhythm to one that is audiovisual, with both mouth movements and acoustics sharing the same rhythmicity. **(A)** Schematic of a presumptive ancestral vocalization with a rhythmic auditory component (blue line) and non-rhythmic visual component (red line). **(B)** Graphical representation of a presumptive ancestral rhythmic facial expression without any vocal component; convention as in (A). **(C)** Illustration of a speech-like utterance with rhythmic and coupled audiovisual components.

that young individuals produce slower and more variable mouth movements and, as they get older, these movements become faster and less variable [82]. Moreover, the developmental trajectory for lip-smacking is distinct from that of chewing. As in humans [79,80], macaque monkey chewing had the same slow frequency and consistent low variability across age groups [82]. Thus, in terms of rhythmicity, the trajectory of lip-smacking development is identical to that of babbling-to-consonant-vowel production in humans. The differences in the developmental trajectories between lip-smacking and chewing are also identical to those reported in humans for speech and chewing [77,83–85].

The coordination of effectors

If human speech and monkey lip-smacking have a shared neural basis, one would expect commonalities in the coordination of the effectors involved. During speech, different sounds are produced through the functional coordination between key vocal tract anatomical structures: the

jaw/lips, tongue, and hyoid. The hyoid is a bony structure to which the laryngeal muscles attach. These effectors are more loosely coupled during speech movements than during chewing movements [86–89]. X-ray cineradiography (X-ray movies) used to visualize the internal dynamics of the macaque monkey vocal tract during lip-smacking and chewing revealed that lips, tongue, and hyoid move during lip-smacking (as in speech) and do so with a speech-like 3–8 Hz rhythm. Relative to lip-smacking, movements during chewing were significantly slower for each of these structures. Importantly, the temporal coordination of these structures was distinct for each behavior. Partial directed coherence measures – an analysis that measures to what extent one time series can predict another [90] – revealed that although the hyoid moves continuously during lip-smacking there is no coupling of the hyoid with lips and tongue movements, whereas during chewing more coordination was observed between the three structures. These patterns are consistent with what is observed during human speech and chewing [86,87]: the effectors are more

Box 1. Vocal coordination: other forms in other species

Many species of animals exchange vocalizations often taking the form of a single 'call-and-response'. For example, naked mole-rats [117], squirrel monkeys [118], female Japanese macaques [119], large-billed crows [120], bottlenose dolphins [121], and some anurans [122,123] are all capable of simple call-and-response behaviors. It is not known how many animals engage in extended, structured sequences of vocal interactions. Instances of extended, coordinated vocal exchanges include the chorusing behaviors of male anurans and insects in the competitive context of mate attraction [124] and duetting between pair-bonded songbirds (e.g., [125,126]; [127] for review), titi monkeys [128], and gibbons (e.g., [129]; [130] for review). Duetting is usually associated with cooperative defense of territory and perhaps mate-guarding. Unlike vocal turn-taking in marmosets and humans, chorusing and duetting occur within the limited contexts of competitive interactions or pair-bonds, respectively. Marmosets and humans are able to flexibly coordinate extended vocal exchanges with any conspecific, regardless of pair-bonding status or relatedness [101].

One possibility is that 'call-and-response' behavior, duetting, and cooperative vocal turn-taking are evolutionarily related to one another [131]. For example, Yoshida and Okanoya [131] argue that the more general call-and-response behavior was derived from duetting behavior. Another possibility is that cooperative vocal turn-taking exhibited by marmoset monkeys and humans is like duetting, which has at its foundation a strong social bond between a mated pair. In the case of marmosets and humans, both of which exhibit stable social bonds with unrelated individuals, prosocial behaviors such as cooperative vocal turn-taking may have been driven by their

cooperative breeding strategy [132]. Thus, cooperative vocal turn-taking may be an extension of 'duetting-like' vocal coordination to any conspecific. More comparative data are necessary to distinguish between the most plausible evolutionary scenarios. Regardless of the initial conditions, cooperative vocal turn-taking in marmosets and humans is the result of convergent evolution because even call-and-response vocal exchanges are not consistently observed among Old World primates. Convergent evolution of vocal behaviors is not uncommon: both vocal learning [133] and duetting [134] evolved multiple times in birds. The evolution of duetting in birds is related to a decline in migration, which promotes the formation of more stable social bonds between mates [134]. The cooperative breeding strategy of marmosets and humans also produces more stable social bonds, but beyond the mated pair.

Importantly, convergent evolution of vocal behaviors does not mean that new mechanisms must be deployed at each instance. For example, coupled oscillatory mechanisms can explain the chorusing behaviors of frogs [124], duetting in birds [125], and vocal turn-taking in marmosets [101] and humans [114]. Of course, it is impossible that the specific neural instantiation (the central pattern generators, their connectivity and modulation) of the coupled oscillator mechanisms is the same across all species. However, it may be the case that convergent evolution of vocal turn-taking in marmosets and humans is the outcome of a homologous neural circuit [100]. This is for two reasons: developmental trajectories are highly constrained across related species [135] and radically different behaviors (e.g., turn-taking versus no turn-taking) can hinge on differential neuromodulation of the same circuit [136].

loosely coupled during lip-smacking than during chewing. Furthermore, the spatial displacement of the lips, tongue, and hyoid is greater during chewing than for lip-smacking [91], again similar to what is observed in human speech versus chewing [87].

Perceptual tuning

In speech, disrupting the auditory or visual component of the 3–8 Hz rhythm significantly reduces intelligibility [48–52]. To test whether macaque monkeys were differentially sensitive to lip-smacking produced with a rhythmic frequency in the species-typical range (mean 4–6 Hz [69,82,91]), a preferential-looking procedure was used [92]. Computer-generated monkey avatars were used to produce stimuli varying in lip-smacking frequency within (6 Hz) and outside (3 and 10 Hz) the species-typical range but with otherwise identical features [45,93]. Although there were at least four alternative outcomes in this experiment, monkeys showed a preference for the 6 Hz lip-smacking over the 3 and 10 Hz. This lends behavioral support for the hypothesis that perceptual processes are similarly tuned to the natural frequencies of communication signals as they are for the speech rhythm in humans.

Bridging the gap

How easy would it be to link vocalizations to a rhythmic facial expression during the course of evolution? Recent work on gelada baboons (*Theropithecus gelada*) proves to be illuminating. Geladas are a highly-specialized type of baboon. Their social structure and habitat are unique among baboons and other Old World primates, as are some of their vocalizations [18]. One of these unique vocalizations, known as a 'wobble', is produced only by males of this species and during close affiliative interactions with

females. Wobbles are essentially lip-smacking expressions produced concurrently with vocalization [94]. Moreover, their rhythmicity falls within the range of speech rhythm and lip-smacking by macaque monkeys. Given that gelada baboons are very closely related to yellow baboons (their taxa are separated by 4 million years), who do not produce anything resembling wobble vocalizations, it suggests that linking rhythmic facial expressions such as lip-smacking to vocal output may not be a complex evolutionary process. How geladas achieved this feat at the level of neural circuits is unknown, but finding out could reveal key information about the human transition to rhythmic audiovisual vocal output – and, more generally, to the production of consonants (another evolutionary puzzle [95]) – during the course of our evolution.

In humans, this rhythmic signal perception and production is often nested in another rhythm – the extended exchanges of speech across two individuals during a conversation. The evolution of such vocal cooperation between subjects is, of course, as important as the coupling between the visual and auditory modalities within a subject. Effective and efficient vocal communication is achieved by minimizing signal interference. Taking turns is one mechanism that reduces interference. To be conversation-like, such turn-taking would involve multiple exchanges, not simply a call-and-response (Box 1). Until recently humans were thought to be the only primates to exhibit vocal cooperation in the form of turn-taking.

On the origins of cooperative vocal communication

Cooperation is central to human communication [9,96]. Conversation, a form of vocal cooperation, proceeds smoothly because of turn-taking. Typically, speech exchanges between two individuals occur without any

explicit agreement on how the talk may flow [97]. A smooth speech interaction consists of vocal exchanges with gaps of silence and minimal overlaps. These features are universal, being present in the conversations of traditional indigenous peoples as well as those speaking any of the major world languages [98]. Given its central importance in everyday human social interactions, it is natural to ask how conversational, vocal turn-taking evolved. It has been argued that human cooperative vocal communication is unique and evolved in essentially three steps (put forth most cogently in [9]; see also [6,99] for similar scenarios). First, an ape-like ancestor used manual gestures to point and direct the attention of others. Second, later ancestors with prosocial tendencies used manual gestures in communications to mediate shared intentionality. Finally, and most mysteriously, a transition from primarily gestural to primarily vocal forms of cooperative communication formed, perhaps to more efficiently express shared intentionality. Implicit in this idea is that a large brain is required for these behaviors. Is this the only plausible scenario? Not necessarily. Perhaps vocal turn-taking evolved through a voluble and prosocial ancestor without the prior scaffolding of a manual gestures or big brains

[100]. The vocal exchanges of the common marmoset monkey provide evidence for this alternative route [101].

Marmoset monkeys are part of the Callitrichinae subfamily of the Cebidae family of New World primates. Marmosets display little evidence of shared intentionality nor do they produce manual gestures. Similarly to humans, they are cooperatively breeding and voluble. Marmosets and humans are among the very few primate species that form pair bonds and exhibit bi-parental and allo-parental care of infants [102]. These cooperative care behaviors are thought to scaffold prosocial motivational and cognitive processes – such as attentional biases toward monitoring others, the ability to coordinate actions, increased social tolerance, and increased responsiveness to the signals of others [103]. Apart from marmosets and humans, and perhaps to some extent bonobos [104], this suite of prosocial behaviors is not typically seen in other primate species. Importantly, when out of visual contact, marmoset monkeys and other callitrichid primates will participate in vocal exchanges with out-of-sight conspecifics [105–108].

In the laboratory and in the wild, marmosets typically use phee calls, a high-pitched call that can be monosyllabic

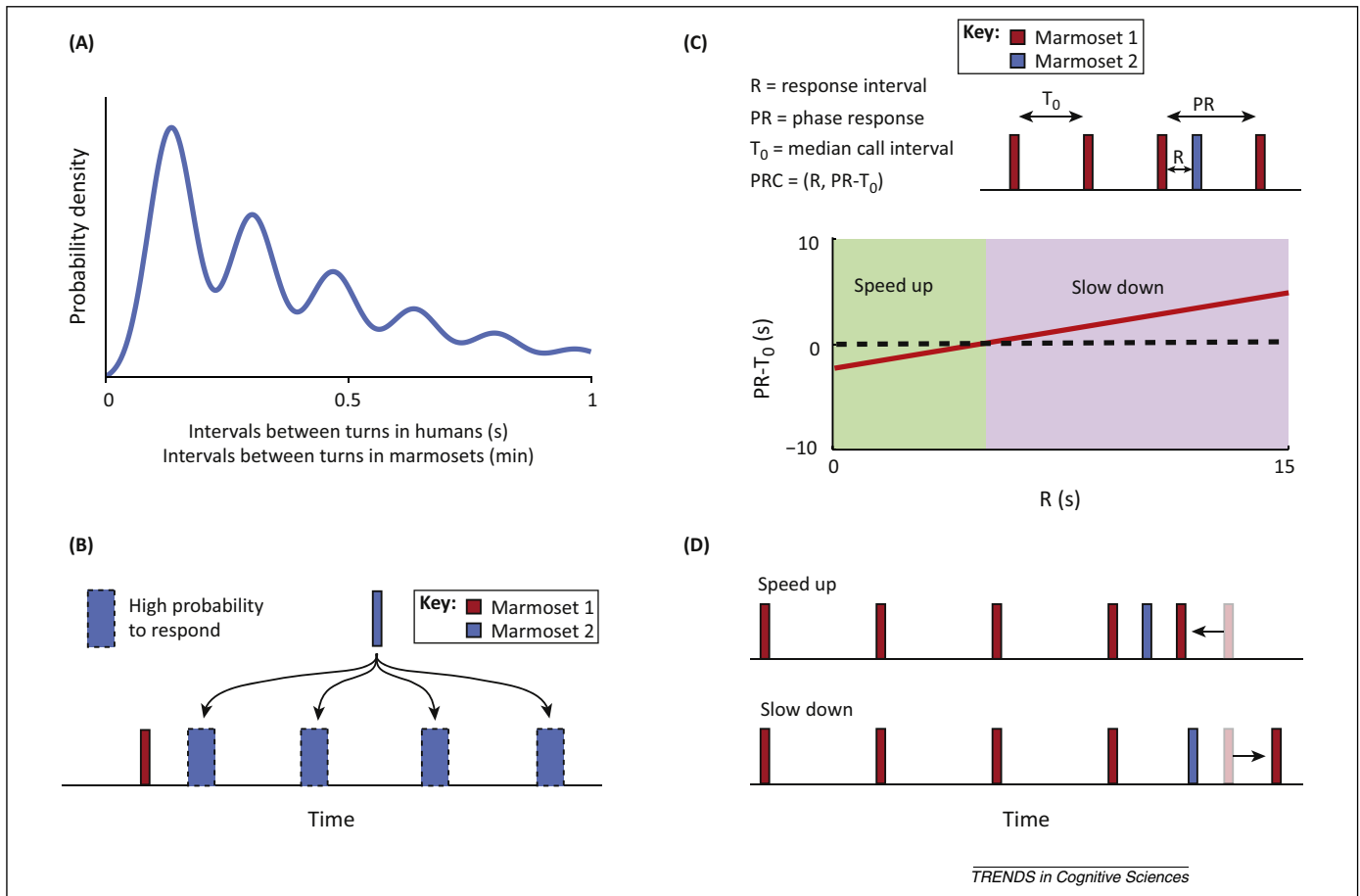


Figure 4. Coupled oscillator dynamics of vocal turn-taking. **(A)** Schematic of the probability distribution of the interval duration between turns during vocal exchanges in humans and marmosets. The same pattern of distribution of the intervals is observed in humans and marmosets, but with a difference in the timescale. **(B)** Coupled rhythmicity implies that once a marmoset calls (red rectangle), the responses from a second marmoset (blue rectangle) will arrive with high probability at one of the intervals regularly spaced from each other (blue rectangle with broken outline). **(C)** Illustration of the correlation between response interval (R) and phase response (PR) when there is an entrainment between call exchanges of marmoset 1 (red rectangle) and marmoset 2 (blue rectangle). When R is short (green area) PR is shorter than the median call interval (T_0), therefore there is a speed-up in the call interval of marmoset 1. When R is long (purple area), PR is longer than the median call interval (T_0), therefore there is a slow-down in the call interval of marmoset 1. **(D)** Schematic of the effect of short and long R on PR . Convention as in (C). The transparent red rectangle indicates where the call from marmoset 1 would be produced had marmoset 2 not responded.

or multisyllabic, as their contact call [109]. A phee call contains information about gender, identity, and social group [110,111]. Marmoset vocal exchanges can last as long as 30 minutes [101] and have a temporal structure that is strikingly similar to the turn-taking rules that humans use in informal, polite conversations [98]. First, there are rarely, if ever, overlapping calls (i.e., no interruptions, and thus no interference). Second, there is a consistent silent interval between utterances across two individuals. Importantly, as in human conversations, marmoset vocal turn-taking occurs spontaneously with another conspecific regardless of pair-bonding status or relatedness [101]. Thus, although some other animal species exhibit vocal coordination over an extended time-period (as opposed to a simple call-and-response), these behaviors are typically confined to competitive chorusing among males of the species or to duetting between pair-bonded mates (Box 1).

Dynamic system models incorporating coupled oscillator-like mechanisms are thought to account for the temporal structure of conversational turn-taking and other social interactions in humans [112,113] (Figure 4A). Such a mechanism would have two basic features: (i) periodic coupling in the timing of utterances across two interacting individuals (Figure 4A,B), and (ii) entrainment, where if the timing of one individual's vocal output quickens or slows, the other follows suit (Figure 4C,D). The vocal exchanges of marmoset monkeys share both of these features [101]. Thus, marmoset vocal communication, like human speech communication [114], can be modeled as loosely coupled oscillators. As a mechanistic description of vocal turn-taking, coupled oscillators are advantageous because they are consistent with the functions of brain oscillations underlying speech processing [54] and its evolution [55]. Further, such oscillations do not require any higher-order cognitive capacities to function [101]. In other words, a coupled oscillator can occur without the involvement of a big brain [100], something worth considering given the small encephalization quotient of the marmoset monkey compared to great apes and humans [115].

The split between the New World primate lineage and the Old World primate lineage occurred around 40 million years ago [116] and, because no other Old World monkey or ape has been observed to vocally cooperate with conspecifics outside of pair-bond, it is unlikely that the cooperative vocal behavior exhibited by both humans and marmosets is shared with a common ancestor. Thus, it is an example of convergent evolution. However, we argue that such convergent evolution of turn-taking behavior may occur through similar or identical modulation of a homologous neuronal circuit [100] (Box 1). Such modulation is driven by the two behavioral features shared by both humans and marmosets: prosociality and volubility. This hypothesis is consistent with the available data on cooperative vocal behaviors in other taxa, in which the strength of social bonds correlates with frequency and complexity of vocal interaction (Box 1). Given that marmosets engage in vocal cooperation in a manner similar to what we observe in humans, it suggests that cooperative vocal communication could have evolved in a manner very different than gestural-origins hypotheses predict

[6,9,99]. Instead of taking an evolutionary route that requires the elaboration of manual gestures and shared intentionality, cooperative vocal communication could have evolved in a more direct fashion. In this alternative scenario, existing vocal repertoires were used in a cooperative, turn-taking manner when prosocial behaviors in general emerged. They developed in both humans and callitrichid primates when they evolved a cooperative breeding strategy.

Concluding remarks and future directions

The default mode of communication in many primates is multisensory. Humans, apes and monkeys all recognize the correspondence between vocalizations and the facial postures associated with them. One striking dissimilarity between some monkey vocalizations and human speech is that the latter has a unique bi-sensory rhythmic structure in that both the acoustic output and the movements of the mouth are rhythmic and tightly correlated. According to one hypothesis, this bimodal speech rhythm evolved through the rhythmic facial expressions of ancestral primates. Developmental, cineradiographic, electromyographic, and perceptual data from macaque monkeys all support the notion that a rhythmic facial expression common among many primate species – lip-smacking – may have been one such ancestral expression. Further explorations of this hypothesis must include a broader comparative sample, especially investigations of the temporal dynamics of facial and vocal expressions in the great apes. Understanding the neural basis of both lip-smacking and speech production – and their similarities and differences – would also be illuminating (Box 2).

In parallel to the evolution of audiovisual coordination within a subject, the evolution of temporal coordination between subjects would need to take place to achieve speech-like behavior. One pragmatic underlying successful speech communication is the ability to take turns. Until recently no nonhuman primate had been observed to naturally take turns using vocalizations in an extended manner with any conspecific. However, such behavior was

Box 2. Outstanding questions

- Beyond the advantages that facial motion provides for vocal detection in noisy environments, do non-human primate species also use facial motion to discriminate between different call types?
- What neural mechanisms and/or biomechanical structures link rhythmic facial motion with rhythmic vocal acoustics?
- Is cooperative vocal turn-taking evident in species that are closely related to marmoset monkeys and humans but that lack prosocial tendencies and/or cooperative breeding strategies (e.g., squirrel monkeys and chimpanzees)?
- What are the neural bases for the coupled oscillator dynamics during vocal turn-taking, and are these mechanisms the same across, for example, marmoset monkeys and humans? Are the neural bases the same or similar to those exhibited by duetting birds?
- What changes in neural circuitry (or in its modulation) lead to changes in prosociality and/or cooperative vocal communication? Is this neural mechanism shared across all species that exhibit some form of vocal coordination (e.g., duetting) with conspecifics?

recently documented in the common marmoset. Because the common marmoset is distantly related to humans, we argue that turn-taking arose as an instance of convergent evolution and is part of a suite of prosocial behaviors. Such behaviors in both humans and marmosets may be, at least in part, the outcome of a cooperative breeding strategy. Here again more comparative evidence is needed either to bolster or falsify this claim (Box 2). Importantly, marmoset vocal turn-taking demonstrates that a large brain size and complex cognitive machinery is not needed for vocal cooperation to occur. Consistent with this idea, the structure of marmoset vocal exchanges can be described in terms of coupled oscillator dynamics that are similar to those used to describe human conversations.

Acknowledgments

We thank Diego Cordero, Lauren Kelly, and our two anonymous reviewers for their thoughtful comments on this manuscript. We thank David Logue for information on, and insights into, duetting in songbirds. This work was supported by National Institutes of Health grant R01NS054898 (A.A.G.), a James S. McDonnell Scholar Award (A.A.G.), a Pew Latin American Fellowship (D.Y.T.), and a Brazilian Science without Borders Fellowship (D.Y.T.).

References

- Huxley, T.H. (1863) *Evidences as to Man's Place in Nature*, Williams and Norgate
- Vargha-Khadem, F. et al. (2005) FOXP2 and the neuroanatomy of speech and language. *Nat. Rev. Neurosci.* 6, 131–138
- Lieberman, P. et al. (1969) Vocal tract limitations on the vowel repertoires of rhesus monkey and other non-human primates. *Science* 164, 1185–1187
- Deacon, T.W. (1997) *The Symbolic Species: The Coevolution of Language and the Brain*, W.W. Norton and Company
- Dunbar, R. (1998) *Grooming, Gossip and the Evolution of Language*, Harvard University Press
- Rizzolatti, G. and Arbib, M.A. (1998) Language within our grasp. *Trends Neurosci.* 21, 188–194
- Jarvis, E.D. (2004) Learned birdsong and the neurobiology of human language. *Ann. N. Y. Acad. Sci.* 1016, 749–777
- Arbib, M.A. et al. (2008) Primate vocalization, gesture, and the evolution of human language. *Curr. Anthropol.* 49, 1053–1076
- Tomasello, M. (2008) *Origins of Human Communication*, MIT Press
- Bates, E. (1979) The emergence of symbols: ontogeny and phylogeny. In *Children's Language and Communication* (Minnesota Symposia on Child Psychology, Vol. 12) (Collins, W.A., ed.), pp. 121–157, Lawrence Erlbaum Associates
- Pinker, S. and Jackendoff, R. (2005) The faculty of language: what's special about it? *Cognition* 95, 201–236
- Fitch, W.T. (2010) *The Evolution of Language*, Cambridge University Press
- Radick, G. (2007) *The Simian Tongue: The Long Debate About Animal Language*, University of Chicago Press
- Lowell, M.S. (1998) *A Rage to Live: A Biography of Richard and Isabel Burton*, Little, Brown and Company
- Fitch, W.T. (2000) The evolution of speech: a comparative review. *Trends Cogn. Sci.* 4, 258–267
- Ghazanfar, A.A. and Santos, L.R. (2004) Primate brains in the wild: the sensory bases for social interactions. *Nat. Rev. Neurosci.* 5, 603–616
- Dobson, S.D. (2009) Socioecological correlates of facial mobility in nonhuman anthropoids. *Am. J. Phys. Anthropol.* 138, 413–420
- Gustison, M.L. et al. (2012) Derived vocalizations of geladas (*Theropithecus gelada*) and the evolution of vocal complexity in primates. *Philos. Trans. R. Soc. Lond. B: Biol. Sci.* 367, 1847–1859
- McComb, K. and Semple, S. (2005) Coevolution of vocal communication and sociality in primates. *Biol. Lett.* 1, 381–385
- Yehia, H. et al. (2002) Linking facial animation, head motion and speech acoustics. *J. Phon.* 30, 555–568
- Chandrasekaran, C. et al. (2009) The natural statistics of audiovisual speech. *PLoS Comput. Biol.* 5, e1000436
- Sumby, W.H. and Pollack, I. (1954) Visual contribution to speech intelligibility in noise. *J. Acoust. Soc. Am.* 26, 212–215
- Ross, L.A. et al. (2007) Do you see what I am saying? Exploring visual enhancement of speech comprehension in noisy environments. *Cereb. Cortex* 17, 1147–1153
- van Wassenhove, V. et al. (2005) Visual speech speeds up the neural processing of auditory speech. *Proc. Natl. Acad. Sci. U.S.A.* 102, 1181–1186
- McGurk, H. and MacDonald, J. (1976) Hearing lips and seeing voices. *Nature* 264, 746–748
- Rosenblum, L.D. (2005) Primacy of multimodal speech perception. In *The Handbook of Speech Perception* (Pisoni, D.B. and Remez, R.E., eds), pp. 51–78, Blackwell Publishing
- Narins, P.M. et al. (2003) Bimodal signal requisite for agonistic behavior in a dart-poison frog, *Epipedobates femoralis*. *Proc. Natl. Acad. Sci. U.S.A.* 100, 577–580
- Taylor, R.C. et al. (2011) Multimodal signal variation in space and time: how important is matching a signal with its signaler? *J. Exp. Biol.* 214, 815–820
- Uetz, G.W. and Roberts, J.A. (2002) Multisensory cues and multimodal communication in spiders: insights from video/audio playback studies. *Brain Behav. Evol.* 59, 222–230
- Burrows, A.M. et al. (2009) Facial musculature in the rhesus macaque (*Macaca mulatta*): evolutionary and functional contexts with comparisons to chimpanzees and humans. *J. Anat.* 215, 320–334
- Sherwood, C.C. (2005) Comparative anatomy of the facial motor nucleus in mammals, with an analysis of neuron numbers in primates. *Anat. Rec. A: Discov. Mol. Cell. Evol. Biol.* 287A, 1067–1079
- Sherwood, C.C. et al. (2004) Cortical orofacial motor representation in old world monkeys, great apes, and humans. II. Stereologic analysis of chemoarchitecture. *Brain Behav. Evol.* 63, 82–106
- Andrew, R.J. (1962) The origin and evolution of the calls and facial expressions of the primates. *Behaviour* 20, 1–109
- Hauser, M.D. et al. (1993) The role of articulation in the production of rhesus monkey, *Macaca mulatta*, vocalizations. *Anim. Behav.* 45, 423–433
- Partan, S.R. (2002) Single and multichannel facial composition: facial expressions and vocalizations of rhesus macaques (*Macaca mulata*). *Behaviour* 139, 993–1027
- Fitch, W.T. and Hauser, M.D. (1995) Vocal production in nonhuman primates - acoustics, physiology, and functional constraints on honest advertisement. *Am. J. Primatol.* 37, 191–219
- Ghazanfar, A.A. and Rendall, D. (2008) Evolution of human vocal production. *Curr. Biol.* 18, R457–R460
- Ghazanfar, A.A. and Logothetis, N.K. (2003) Facial expressions linked to monkey calls. *Nature* 423, 937–938
- Parr, L.A. (2004) Perceptual biases for multimodal cues in chimpanzee (*Pan troglodytes*) affect recognition. *Anim. Cogn.* 7, 171–178
- Jordan, K.E. et al. (2005) Monkeys match the number of voices they hear with the number of faces they see. *Curr. Biol.* 15, 1034–1038
- Ghazanfar, A.A. et al. (2007) Vocal tract resonances as indexical cues in rhesus monkeys. *Curr. Biol.* 17, 425–430
- Sliwa, J. et al. (2011) Spontaneous voice-face identity matching by rhesus monkeys for familiar conspecifics and humans. *Proc. Natl. Acad. Sci. U.S.A.* 108, 1735–1740
- Adachi, I. and Hampton, R.R. (2011) Rhesus monkeys see who they hear: spontaneous crossmodal memory for familiar conspecifics. *PLoS ONE* 6, e23345
- Habbershon, H.M. et al. (2013) Rhesus macaques recognize unique multimodal face-voice relations of familiar individuals and not of unfamiliar ones. *Brain Behav. Evol.* 81, 219–225
- Chandrasekaran, C. et al. (2011) Monkeys and humans share a common computation for face/voice integration. *PLoS Comput. Biol.* 7, e1002165
- Ghitza, O. (2011) Linking speech perception and neurophysiology: speech decoding guided by cascaded oscillators locked to the input rhythm. *Front. Psychol.* 2, 130
- Greenberg, S. et al. (2003) Temporal properties of spontaneous speech – a syllable-centric perspective. *J. Phon.* 31, 465–485

- 48 Saberi, K. and Perrott, D.R. (1999) Cognitive restoration of reversed speech. *Nature* 398, 760
- 49 Smith, Z.M. *et al.* (2002) Chimaeric sounds reveal dichotomies in auditory perception. *Nature* 416, 87–90
- 50 Elliot, T.M. and Theunissen, F.E. (2009) The modulation transfer function for speech intelligibility. *PLoS Comput. Biol.* 5, e1000302
- 51 Ghitza, O. and Greenberg, S. (2009) On the possible role of brain rhythms in speech perception: Intelligibility of time-compressed speech with periodic and aperiodic insertions of silence. *Phonetica* 66, 113–126
- 52 Vitkovitch, M. and Barber, P. (1996) Visible speech as a function of image quality: effects of display parameters on lipreading ability. *Appl. Cogn. Psychol.* 10, 121–140
- 53 Ghazanfar, A.A. and Poeppel, D. (2014) The neurophysiology and evolution of the speech rhythm. In *The Cognitive Neurosciences V* (Gazzaniga, M.S., ed.), pp. 629–638, MIT Press
- 54 Giraud, A.-L. and Poeppel, D. (2012) Cortical oscillations and speech processing: emerging computational principles and operations. *Nat. Neurosci.* 15, 511–517
- 55 Ghazanfar, A.A. and Takahashi, D.Y. (2014) Facial expressions and the evolution of the speech rhythm. *J. Cogn. Neurosci.* 26, 1196–1207
- 56 Cohen, Y.E. *et al.* (2007) Acoustic features of rhesus vocalizations and their representation in the ventrolateral prefrontal cortex. *J. Neurophysiol.* 97, 1470–1484
- 57 MacNeilage, P.F. (1998) The frame/content theory of evolution of speech production. *Behav. Brain Sci.* 21, 499–511
- 58 Preuschoff, S. (2000) Primate faces and facial expressions. *Soc. Res.* 67, 245–271
- 59 Hinde, R.A. and Rowell, T.E. (1962) Communication by posture and facial expressions in the rhesus monkey (*Macaca mulatta*). *Proc. Zool. Soc. Lond.* 138, 1–21
- 60 Redican, W.K. (1975) Facial expressions in nonhuman primates. In *Primate Behavior: Developments in Field and Laboratory Research* (Rosenblum, L.A., ed.), pp. 103–194, Academic Press
- 61 Van Hooff, J.A.R.A.M. (1962) Facial expressions of higher primates. *Symp. Zool. Soc. Lond.* 8, 97–125
- 62 Parr, L.A. *et al.* (2005) Influence of social context on the use of blended and graded facial displays in chimpanzees. *Int. J. Primatol.* 26, 73–103
- 63 Kemp, C. and Kaplan, G. (2013) Facial expressions in common marmosets (*Callithrix jacchus*) and their use by conspecifics. *Anim. Cogn.* 16, 773–788
- 64 De Marco, A. and Visalberghi, E. (2007) Facial displays in young tufted capuchin monkeys (*Cebus apella*): appearance, meaning, context and target. *Folia Primatol. (Basel)* 78, 118–137
- 65 Newell, T.G. (1971) Social encounters in two prosimian species. *Psychon. Sci.* 24, 128–130
- 66 Ferrari, P.F. *et al.* (2009) Reciprocal face-to-face communication between rhesus macaque mothers and their newborn infants. *Curr. Biol.* 19, 1768–1772
- 67 Livneh, U. *et al.* (2012) Self-monitoring of social facial expressions in the primate amygdala and cingulate cortex. *Proc. Natl. Acad. Sci. U.S.A.* 109, 18956–18961
- 68 Shepherd, S.V. *et al.* (2012) Facial muscle coordination during rhythmic facial expression and ingestive movement. *J. Neurosci.* 32, 6105–6116
- 69 Ghazanfar, A.A. *et al.* (2010) Dynamic, rhythmic facial expressions and the superior temporal sulcus of macaque monkeys: implications for the evolution of audiovisual speech. *Eur. J. Neurosci.* 31, 1807–1817
- 70 Gottlieb, G. (1992) *Individual Development and Evolution: The Genesis of Novel Behavior*, Oxford University Press
- 71 Locke, J.L. (1993) *The Child's Path to Spoken Language*, Harvard University Press
- 72 Smith, A. and Zelaznik, H.N. (2004) Development of functional synergies for speech motor coordination in childhood and adolescence. *Dev. Psychobiol.* 45, 22–33
- 73 Preuschoff, K. *et al.* (2008) Human insula activation reflects risk prediction errors as well as risk. *J. Neurosci.* 28, 2745–2752
- 74 Davis, B.L. and MacNeilage, P.F. (1995) The articulatory basis of babbling. *J. Speech Hear. Res.* 38, 1199–1211
- 75 Lindblom, B. *et al.* (1996) Phonetic systems and phonological development. In *Developmental Neurocognition: Speech and Face Processing in the First Year of Life* (de Boysson-Bardies, B. *et al.*, eds), pp. 399–409, Kluwer Academic Publishers
- 76 Oller, D.K. (2000) *The Emergence of the Speech Capacity*, Lawrence Erlbaum
- 77 Dolata, J.K. *et al.* (2008) Characteristics of the rhythmic organization of vocal babbling: Implications for an amodal linguistic rhythm. *Infant Behav. Dev.* 31, 422–431
- 78 Nathani, S. *et al.* (2003) Final syllable lengthening (FSL) in infant vocalizations. *J. Child Lang.* 30, 3–25
- 79 Green, J.R. *et al.* (1997) Development of chewing in children from 12 to 48 months: longitudinal study of EMG patterns. *J. Neurophysiol.* 77, 2704–2727
- 80 Kiliaridis, S. *et al.* (1991) Characteristics of masticatory mandibular movements and velocity in growing individuals and young adults. *J. Dent. Res.* 70, 1367–1370
- 81 Locke, J.L. (2008) Lipsmacking and babbling: syllables, sociality and survival. In *The Syllable in Speech Production* (Davis, B.L. and Zajdo, K., eds), pp. 111–129, Lawrence Erlbaum Associates
- 82 Morrill, R.J. *et al.* (2012) Monkey lip-smacking develops like the human speech rhythm. *Dev. Sci.* 15, 557–568
- 83 Moore, C.A. and Ruark, J.L. (1996) Does speech emerge from earlier appearing motor behaviors? *J. Speech Hear. Res.* 39, 1034–1047
- 84 Steeve, R.W. (2010) Babbling and chewing: jaw kinematics from 8 to 22 months. *J. Phon.* 38, 445–458
- 85 Steeve, R.W. *et al.* (2008) Babbling, chewing, and sucking: oromandibular coordination at 9 months. *J. Speech Lang. Hear. Res.* 51, 1390–1404
- 86 Hiemae, K.M. and Palmer, J.B. (2003) Tongue movements in feeding and speech. *Crit. Rev. Oral Biol. Med.* 14, 413–429
- 87 Hiemae, K.M. *et al.* (2002) Hyoid and tongue surface movements in speaking and eating. *Arch. Oral Biol.* 47, 11–27
- 88 Ostry, D.J. and Munhall, K.G. (1994) Control of jaw orientation and position in mastication and speech. *J. Neurophysiol.* 71, 1528–1545
- 89 Matsuo, K. and Palmer, J.B. (2010) Kinematic linkage of the tongue, jaw, and hyoid during eating and speech. *Arch. Oral Biol.* 55, 325–331
- 90 Takahashi, D.Y. *et al.* (2010) Information theoretic interpretation of frequency domain connectivity measures. *Biol. Cybern.* 103, 463–469
- 91 Ghazanfar, A.A. *et al.* (2012) Cineradiography of monkey lip-smacking reveals putative precursors of speech dynamics. *Curr. Biol.* 22, 1176–1182
- 92 Ghazanfar, A.A. *et al.* (2013) Monkeys are perceptually tuned to facial expressions that exhibit a theta-like speech rhythm. *Proc. Natl. Acad. Sci. U.S.A.* 110, 1959–1963
- 93 Steckenfinger, S.A. and Ghazanfar, A.A. (2009) Monkey visual behavior falls into the uncanny valley. *Proc. Natl. Acad. Sci. U.S.A.* 106, 18362–18466
- 94 Bergman, T.J. (2013) Speech-like vocalized lip-smacking in geladas. *Curr. Biol.* 23, R268–R269
- 95 Lameira, A.R. *et al.* (2014) Primate feedstock for the evolution of consonants. *Trends Cogn. Sci.* 18, 60–62
- 96 Levinson, S.C. (2006) On the human interactional engine. In *Roots of Human Sociality* (Enfield, N.J. and Levinson, S.C., eds), pp. 39–69, Berg Publishers
- 97 Sacks, H. *et al.* (1974) Simplest systematics for organization of turn-taking for conversation. *Language* 50, 696–735
- 98 Stivers, T. *et al.* (2009) Universals and cultural variation in turn-taking in conversation. *Proc. Natl. Acad. Sci. U.S.A.* 106, 10587–10592
- 99 Hewes, G.W. (1973) Primate communication and the gestural origin of language. *Curr. Anthropol.* 14, 5–24
- 100 Borjon, J.I. and Ghazanfar, A.A. (2014) Convergent evolution of vocal cooperation without convergent evolution of brain size. *Brain Behav. Evol.* (in press)
- 101 Takahashi, D.Y. *et al.* (2013) Coupled oscillator dynamics of vocal turn-taking in monkeys. *Curr. Biol.* 23, 2162–2168
- 102 Zahed, S.R. *et al.* (2008) Male parenting and response to infant stimuli in the common marmoset (*Callithrix jacchus*). *Am. J. Primatol.* 70, 84–92
- 103 Burkart, J.M. and van Schaik, C.P. (2010) Cognitive consequences of cooperative breeding in primates? *Anim. Cogn.* 13, 1–19
- 104 Hare, B. *et al.* (2007) Tolerance allows bonobos to outperform chimpanzees on a cooperative task. *Curr. Biol.* 17, 619–623

- 105 Chen, H.-C. *et al.* (2009) Contact calls of common marmosets (*Callithrix jacchus*): influence of age of caller on antiphonal calling and other vocal responses. *Am. J. Primatol.* 71, 165–170
- 106 Ghazanfar, A.A. *et al.* (2001) The units of perception in the antiphonal calling behavior of cotton-top tamarins (*Saguinus oedipus*): playback experiments with long calls. *J. Comp. Physiol. A* 187, 27–35
- 107 Ghazanfar, A.A. *et al.* (2002) Temporal cues in the antiphonal long-calling behaviour of cottontop tamarins. *Anim. Behav.* 64, 427–438
- 108 Miller, C.T. and Wang, X. (2006) Sensory-motor interactions modulate a primate vocal behavior: antiphonal calling in common marmosets. *J. Comp. Physiol. A: Neuroethol. Sens. Neural Behav. Physiol.* 192, 27–38
- 109 Bezerra, B.M. and Souto, A. (2008) Structure and usage of the vocal repertoire of *Callithrix jacchus*. *Int. J. Primatol.* 29, 671–701
- 110 Miller, C.T. *et al.* (2010) The communicative content of the common marmoset phoe call during antiphonal calling. *Am. J. Primatol.* 72, 974–980
- 111 Norcross, J.L. and Newman, J.D. (1993) Context and gender-specific differences in the acoustic structure of common marmoset (*Callithrix jacchus*) phoe calls. *Am. J. Primatol.* 30, 37–54
- 112 Oullier, O. *et al.* (2008) Social coordination dynamics: measuring human bonding. *Soc. Neurosci.* 3, 178–192
- 113 Schmidt, R. and Morr, S. (2010) Coordination dynamics of natural social interactions. *Int. J. Sport Psychol.* 41, 105–106
- 114 O'Dell, M.L. *et al.* (2012) Modeling turn-taking rhythms with oscillators. *Linguist. Ural.* 48, 218–227
- 115 Jerison, H.J. (1973) *Evolution of the Brain and Intelligence*, Academic Press
- 116 Steiper, M.E. and Young, N.M. (2006) Primate molecular divergence dates. *Mol. Phylogenet. Evol.* 41, 384–394
- 117 Yosida, S. *et al.* (2007) Antiphonal vocalization of a subterranean rodent, the naked mole-rat (*Heterocephalus glaber*). *Ethology* 113, 703–710
- 118 Masataka, N. and Biben, M. (1987) Temporal rules regulating affiliative vocal exchanges of squirrel monkeys. *Behaviour* 101, 311–319
- 119 Sugiura, H. (1998) Matching of acoustic features during the vocal exchange of coo calls by Japanese macaques. *Anim. Behav.* 55, 673–687
- 120 Kondo, N. *et al.* (2010) A temporal rule in vocal exchange among large-billed crows *Corvus macrorhynchos* in Japan. *Ornithol. Sci.* 9, 83–91
- 121 Nakahara, F. and Miyazaki, N. (2011) Vocal exchanges of signature whistles in bottlenose dolphins (*Tursiops truncatus*). *J. Ethol.* 29, 309–320
- 122 Grafe, T.U. (1996) The function of call alternation in the African reed frog (*Hyperolius marmoratus*): precise call timing prevents auditory masking. *Behav. Ecol. Sociobiol.* 38, 149–158
- 123 Zelick, R. and Narins, P.M. (1985) Characterization of the advertisement call oscillator in the frog *Eleutherodactylus coqui*. *J. Comp. Physiol. A* 156, 223–229
- 124 Greenfield, M.D. (1994) Synchronous and alternating choruses in insects and anurans: common mechanisms and diverse functions. *Am. Zool.* 34, 605–615
- 125 Laje, R. and Mindlin, G.B. (2003) Highly structured duets in the song of the South American hornero. *Phys. Rev. Lett.* 91, 258104
- 126 Logue, D.M. *et al.* (2008) The behavioural mechanisms underlying temporal coordination in black-bellied wren duets. *Anim. Behav.* 75, 1803–1808
- 127 Hall, M.L. (2009) A review of vocal duetting in birds. *Adv. Stud. Behav.* 40, 67–121
- 128 Caselli, C.B. *et al.* (2014) Vocal behavior of black-fronted titi monkeys (*Callicebus nigrifrons*): acoustic properties and behavioral contexts of loud calls. *Am. J. Primatol.* 76, 788–800
- 129 Mitani, J.C. (1985) Gibbon song duets and intergroup spacing. *Behaviour* 92, 59–96
- 130 Geissmann, T. (2002) Duet-splitting and the evolution of gibbon songs. *Biol. Rev.* 77, 57–76
- 131 Yoshida, S. and Okanoya, K. (2005) Evolution of turn-taking: a bio-cognitive perspective. *Cogn. Stud.* 12, 153–165
- 132 Burkart, J.M. *et al.* (2009) Cooperative breeding and human cognitive evolution. *Evol. Anthropol.* 18, 175–186
- 133 Jarvis, E.D. (2006) Selection for and against vocal learning in birds and mammals. *Ornithol. Sci.* 5, 5–14
- 134 Logue, D.M. and Hall, M.L. (2014) Migration and the evolution of duetting in songbirds. *Proc. Biol. Sci.* 281, 20140103
- 135 Krubitzer, L.A. and Seelke, A.M. (2012) Cortical evolution in mammals: the bane and beauty of phenotypic variability. *Proc. Natl. Acad. Sci. U.S.A.* 109 (Suppl. 1), 10647–10654
- 136 Marder, E. (2012) Neuromodulation of neuronal circuits: back to the future. *Neuron* 76, 1–11