

Universals and cultural variation in turn-taking in conversation

Tanya Stivers^{a,1}, N. J. Enfield^a, Penelope Brown^a, Christina Englert^b, Makoto Hayashi^c, Trine Heinemann^d, Gertie Hoymann^a, Federico Rossano^a, Jan Peter de Ruiter^{a,e}, Kyung-Eun Yoon^f, and Stephen C. Levinson^a

^aLanguage and Cognition Group, Max Planck Institute for Psycholinguistics, 6525XD Nijmegen, The Netherlands; ^bCenter for Language and Cognition, University of Groningen, 9172TS Groningen, The Netherlands; ^cDepartment of East Asian Languages and Cultures, University of Illinois at Urbana-Champaign, Urbana, IL 61801; ^dSønderborg Participatory Innovation Research Center and The Institute of Business Communication and Information Science, University of Southern Denmark, 6400 Sønderborg, Denmark; ^eFaculty for Linguistics and Literary Sciences, Bielefeld University, D-33501 Bielefeld, Germany; and ^fDepartment of African and Asian Languages and Literatures, University of Florida, Gainesville, FL 32611

Edited by Paul Kay, International Computer Science Institute, Berkeley, CA, and approved April 28, 2009 (received for review April 2, 2009)

Informal verbal interaction is the core matrix for human social life. A mechanism for coordinating this basic mode of interaction is a system of turn-taking that regulates who is to speak and when. Yet relatively little is known about how this system varies across cultures. The anthropological literature reports significant cultural differences in the timing of turn-taking in ordinary conversation. We test these claims and show that in fact there are striking universals in the underlying pattern of response latency in conversation. Using a worldwide sample of 10 languages drawn from traditional indigenous communities to major world languages, we show that all of the languages tested provide clear evidence for a general avoidance of overlapping talk and a minimization of silence between conversational turns. In addition, all of the languages show the same factors explaining within-language variation in speed of response. We do, however, find differences across the languages in the average gap between turns, within a range of 250 ms from the cross-language mean. We believe that a natural sensitivity to these tempo differences leads to a subjective perception of dramatic or even fundamental differences as offered in ethnographic reports of conversational style. Our empirical evidence suggests robust human universals in this domain, where local variations are quantitative only, pointing to a single shared infrastructure for language use with likely ethological foundations.

cooperation | response speed | social interaction

Crucial to understanding the nature and origins of human language, perhaps our most distinctive trait, is understanding the social-interactive matrix in which it is used. Informal conversation is where language is learned and where most of the business of social life is conducted. A fundamental part of the infrastructure for conversation is turn-taking, or the apportioning of who is to speak next and when (1). Previous research on turn-taking has examined cues used in recognizing opportunities for turn transition (1–4), the time course of a turn in an exchange (5), and the timing of turn transitions (1, 6–10). In English conversation speakers do not wait for pauses to begin their turn but avoid gaps and overlaps. To achieve this they use grammar, prosody, and pragmatics to project when they can start a next turn, suggesting that turn-taking is specifically organized to achieve this close timing. Here, we consider whether this organization varies across human cultures or is reflective of a universal system of rules for turn-taking in conversation. To our knowledge, no previous study has set out to test the robustness of a turn-taking system for informal interaction across the diversity of human cultures.

In the anthropological literature there are frequent claims that cultures differ radically in the timing of conversational turn-taking, and thus that the findings for English are culture-specific. Nordic cultures, for example, are said to relish long delays between one turn and the next. As the report goes, “Two brothers of Häme (Finland) were on their way to work in the morning. One says, ‘It is here that I lost my knife’. Coming back

home in the evening, the other asks, ‘Your knife, did you say?’” (11). Or receiving visitors in the North of Sweden: “We would offer coffee. After several minutes of silence the offer would be accepted. We would tentatively ask a question. More silence, then a ‘yes’ or a ‘no’” (12). Compare this preference for silence between turns with the reported “fast rate of turn-taking” and “preference for simultaneous speech” in New York Jewish conversation (13) or the “anarchic” conversation of an Antiguan village, in which there is said to be “no regular requirement for 2 or more voices not to be going on at the same time” (12). Although there are many such claims in the anthropological literature of cultures where substantial overlap is the norm (14–16) or where long silences are said to be the rule (11, 12, 17), no broad-ranging, quantitative comparison has been made. These claims suggest that there are culturally variable turn-taking systems.

In contrast to these claims of diversity, there are arguments in favor of a universal system for turn-taking, that, as in English, follows a norm of “minimal-gap minimal-overlap” (18). First, there is a functional basis for turns to be immediately adjacent (rather than overlapping or overly separated): a timely response makes clear its link to another speaker’s prior utterance (19), displaying that it is directly contingent on that utterance (20), and showing how the prior utterance was understood, allowing rapid correction if necessary (1, 21, 22). Second, there is evidence for a human ethological basis for adjacent sequences of communicative action and response, for example in very early “proto-conversation” between newborns and caregivers (23–26). Systems in which turn transitions occur with minimal delay or overlap have been described for several languages (1, 8, 27, 28), but no systematic cross-linguistic comparison has been undertaken.

Here, we test these opposing hypotheses: (i) a universal system hypothesis, by which turn-taking is a universal system with minimal cultural variability, and (ii) a cultural variability hypothesis, by which turn-taking is language and culture dependent. The universal system hypothesis predicts a unimodal distribution of turn transitions with most transitions occurring ≈ 0 in all languages, whereas the cultural variability hypothesis predicts that overlap is more common in some languages and gaps more common in others.

If a community of speakers shows a highly regular target for the timing of turn transition, deviations will come to have a natural communicative significance (e.g., delays implying prob-

Author contributions: T.S. and N.J.E. designed research; T.S., N.J.E., P.B., C.E., M.H., T.H., G.H., F.R., J.P.d.R., K.-E.Y., and S.C.L. performed research; T.S. and N.J.E. analyzed data; and T.S., N.J.E., and S.C.L. wrote the paper.

The authors declare no conflict of interest.

This article is a PNAS Direct Submission.

¹To whom correspondence should be addressed. E-mail: tanya.stivers@mpi.nl.

This article contains supporting information online at www.pnas.org/cgi/content/full/0903616106/DCSupplemental.

lems with the prior utterance), so giving rise to implicit norms of timely response that will be maintained to avoid such added implications (29). Research on questions in English conversation has shown that speakers display inhibition in producing responses that in some way fail to conform with the terms of the question or with the questioner's agenda: thus, responses are often delayed by up to 1 s if, for example, they do not answer the question (e.g., *I don't know* or *I can't remember*) (30, 31) or if they give a response that runs against the bias of the question (e.g., A: *Is that your car?* B: *No*) (32, 33).

Two further explanations for variation in turn transition speed are associated with nonverbal behavior such as head movements (e.g., nodding) and gaze. Although the rules for turn-taking may discourage overlap in the vocal channel, they may nevertheless leave other channels exempt. If nonverbal signals are viewed as less intrusive upon speech, they may come earlier than purely verbal responses. Additionally, if questioners fix their eye gaze on their addressees, this may be expected to elicit faster responses. Research on conversation in European languages suggests that a speaker's gaze toward a listener may increase the pressure to respond and to respond quickly: eye gaze does this by indicating who is addressed (1), by providing early possible cues that the speaker's turn is now coming to an end (4, 6, 34) and signaling the speaker's heightened expectancy for a response (35). However, gaze behavior may show substantial cultural variation (36).

With respect to these 4 accounts for delayed turn transition (nonanswering responses, disconfirmations, vocal-only responses, and nongazing questions), the 2 hypotheses make different predictions. The universal system hypothesis predicts that the languages will all show the same pattern of slower turn transitions when these factors are present. By contrast, the cultural variability hypothesis predicts that delayed turn transition will be explained by different factors in different languages and that the 4 factors just mentioned are unlikely to account for variation in the same way cross-linguistically.

To test these competing hypotheses, we compared data from video recordings of informal natural conversation in 10 languages from 5 continents, e.g., from Southeast Asia, Mexico, Namibia, and Papua New Guinea (see Table S1). The languages vary fundamentally in type (e.g., in word order, sound structure, grammatical options) and are drawn from cultures of quite different kinds (from hunter-gatherer groups to peasant societies to large-scale postindustrial nations). To achieve a natural control over the discourse environment to be compared, we took advantage of a universal context for turn transition, namely that between questions and their responses. For optimal comparability we restricted the comparison to polar questions (questions that expect a yes or no answer). These are the most common type of questions in 9 of the 10 languages (67% of total questions in our 10-language sample were of this type), and they are also logically the simplest type: unlike responses to WH- questions (see Table S2), the desired response to a polar question comes from a small, closed set, usually yes or no. Although not all languages have precise equivalents of English yes and no, they all do have ways of asking polar questions and ways of conveying the basic functions of yes and no. For example, yes can be conveyed by repeating the key information in the question [e.g., Q: *Is John going?*, A: *He's going* (= yes)] or the use of nonstandard expressions like *uh huh* or *yep*. To determine whether question-response sequences are representative of turn-taking in general, we examined a corpus of Dutch conversation (8) for timing across all types of turns and responses and found no difference between response times after questions and nonquestions (see Fig. S1). This suggests that the use of question-answer sequences is a reasonable proxy for turn-taking more generally.

Results

Distribution of Turn Transitions. The temporal relation between a turn and its response we will call the response offset, measured in milliseconds, when there is a gap we have a positive offset, when there is an overlap we have a negative offset. As Fig. 1 shows, we find that the response timings for each language, although slightly skewed to the right, have a unimodal distribution with a mode offset for each language between 0 and +200 ms, and an overall mode of 0 ms (see Fig. 1 and Table S3). The medians are also quite uniform, ranging from 0 ms (English, Japanese, Tzeltal, and Yéli-Dnye) to +300 ms (Danish, Ṃkhoe Hai|om, Lao) (overall cross-linguistic median +100 ms).

The means display somewhat more variation, as shown in Fig. 2. Danish has the slowest response time on average (+469 ms) and Japanese has the fastest (+7 ms). The mean response offset for the full dataset is +208 ms, and the language-specific means fall within ≈ 250 ms either side of this cross-language mean, approximately the length of time it takes to produce a single English syllable (37).

The Implications of Turn Delay. Answering vs. not. Speakers of all of the languages provide answers significantly faster than nonanswer responses to questions (Fig. 3). In all of the languages we also found a greater proportion of answers than nonanswer responses (ranging from 64% of all responses in Korean to 87% in Dutch and Yéli-Dnye) (see Table S4).

Confirming vs. disconfirming. Within the set of answers, those that are confirmations are delivered faster than disconfirmations in all languages, between 100 and 500 ms faster on average (see Fig. 4). This difference reaches significance in 7/10 languages. In all of the languages, we also found a greater proportion of confirmations than disconfirmations (ranging from 70% of all answers in Danish to 89% of all answers in Yéli-Dnye; see Table S4). This advantage for affirmation also holds, incidentally, even if the affirming response is negative in form (as in "You're not coming?" and "No, I'm not"), showing that it is not simply a side-effect of the greater processing costs of negative responses (38) (confirmations using no are not significantly slower than confirmations using yes; 90 vs. 36 ms; $t[693] = -1.1$).

The Implications of Nonverbal Channels. Visible responses vs. vocal-only responses. Visible responses were most commonly head nods, but we also found shrugs and head shakes, and in some languages like Yéli-Dnye conventionalized extended blinks and eyebrow flashes in response to questions. When visible responses occurred in response to a question, they were faster than speech in every language (see Fig. 5). This reached significance in 7/10 of the languages even though there was substantial variation in how frequently visible responses were included in a response (from 21% of responses including a visible component in Ṃkhoe Hai|om to 60% in Italian) (see Table S4).

Questioner gaze vs. no gaze. We found in 9 of the 10 languages that responses were delivered earlier if the speaker was looking at the recipient while the question was asked (Fig. 6). The differences reach statistical significance in only 5 languages. That Danish shows the opposite timing trend, although nonsignificant, combined with known differences in reliance on interactional gaze in different languages, suggests that gaze may be more culturally variable than other behaviors (36). This is also supported by the range of frequencies of gaze to addressee (from 21% in Ṃkhoe Hai|om to 88% in Japanese) (see Table S4). This is incidentally not the expectation in the literature, where addressee gaze rather than speaker gaze has been argued to be the norm (4, 6, 39).

Multivariate Results. The results so far show broadly similar patterns of response timing across the languages, with the same 4 factors each independently accounting for faster or slower than

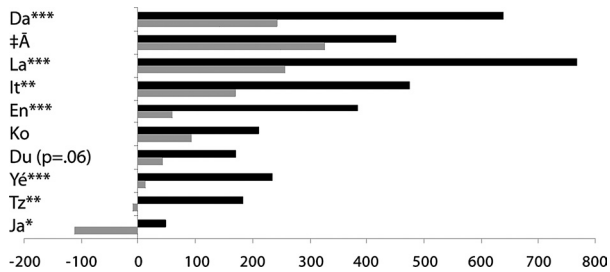


Fig. 4. The mean time of turn transition for responses coded as confirmations versus responses coded as disconfirmations in each of the languages. Speakers of all languages produced confirmations (gray) faster, on average, than they produced disconfirmations (black). *, $P \leq 0.05$; **, $P \leq 0.01$; ***, $P \leq 0.001$. Milliseconds are shown on the x axis. Languages are arrayed along the y axis. Da, Danish; ꞤĂ, ꞤĂkhoe Hai||om; La, Lao; It, Italian; En, English; Ko, Korean; Du, Dutch; Yé, Yéli-Dnye; Tz, Tzeltal; Ja, Japanese.

offset of next turn in each language departing no more than a quarter-second from the overall mean, is not of the kind that would imply fundamentally different types of turn-taking systems in the different languages, as the cultural variability hypothesis would suggest.

Language structure does not explain the variance we observe. Languages that mark questions using a sentence-final marker might plausibly have been associated with slower responses because the fact that the utterance is a question may not be evident until the very end of the turn (28). However, Japanese, Korean, and Lao all use sentence-final marking for questions, yet they do not cluster together within the cross-language range of mean turn offsets (Fig. 2). A converse prediction, that languages like Danish, Dutch, and English, which tend to mark questions at the beginning of a turn, would allow faster responses, also turns out not to hold up. These 3 languages similarly do not cluster together (Fig. 2). Finally, note that this failure of Dutch, English, and Danish to cluster within the cross-language range of mean turn offsets is also evidence that linguistic and cultural kinship (in this case, West Germanic) does not predict interactional tempo.

We suggest that the differences involve a different cultural “calibration” of delay, thus constituting minor variation in the local implementation of a universal underlying turn-taking system, in which speakers aim to minimize the perceived gap before producing a following turn at talk. This target for ideal turn transition remains in a narrow window within each language, with each of 4 factors predisposing a response to be slower (or faster in the case of gaze) than the mean and having similar effects for all of the languages. These differences could either be

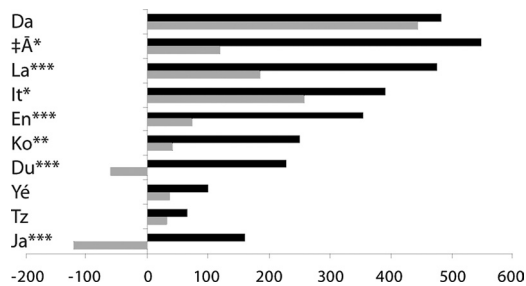


Fig. 5. The mean time of turn transition for responses coded as including a visible response versus responses coded as vocal only in each of the languages. Speakers of all languages produced responses with a visible component (gray) faster, on average, than they produced vocal only responses (black). *, $P \leq 0.05$; **, $P \leq 0.01$; ***, $P \leq 0.001$. Milliseconds are shown on the x axis. Languages are arrayed along the y axis. Da, Danish; ꞤĂ, ꞤĂkhoe Hai||om; La, Lao; It, Italian; En, English; Ko, Korean; Du, Dutch; Yé, Yéli-Dnye; Tz, Tzeltal; Ja, Japanese.

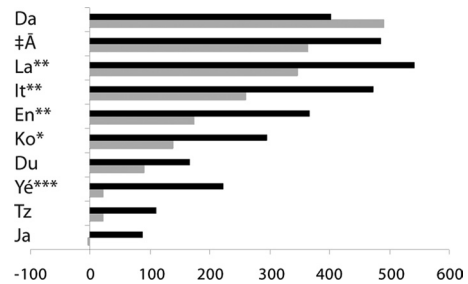


Fig. 6. The mean time of turn transition for questions coded as with speaker gaze versus questions coded as without speaker gaze in each of the languages. Speakers of 9/10 languages produced responses to questions with speaker gaze (gray) faster, on average, than they produced responses to questions without speaker gaze (black). *, $P \leq 0.05$; **, $P \leq 0.01$; ***, $P \leq 0.001$. Milliseconds are shown on the x axis. Languages are arrayed along the y axis. Da, Danish; ꞤĂ, ꞤĂkhoe Hai||om; La, Lao; It, Italian; En, English; Ko, Korean; Du, Dutch; Yé, Yéli-Dnye; Tz, Tzeltal; Ja, Japanese.

a because of a specific cultural interactional pace or follow from more general differences in the overall tempo of social life (40). This would mean that speakers of all languages aim at minimizing significant delays relative to the specific rhythm of that language in conversation (e.g., ref. 41), a perspective that is supported by existing studies of some non Indo-European languages (27, 28, 42). To address this hypothesis we coded the offset of our responses for whether or not, when a relative subjective measure of the conversation’s rhythm was taken into account, responses were coded as late versus on time. Mean response times for subjectively on-time responses are much longer in Danish and Lao (203 and 202 ms, respectively) than in Japanese and Tzeltal (36 and 83 ms, respectively) and comparing the 3 languages with longest response offsets to all others, the difference is significant [$t(847) = -10.97, P < 0.001$]. Thus, a silence of 200 ms, judged as a delay in most languages, was still considered on time. Such a silence is thus not phenomenologically salient within a speech community (but may be to an outside observer). In short, what constitutes a subjectively notable delay involves greater absolute duration in some languages than in others. This is consistent with the presence of a universal, stable system of turn taking avoiding overlap and

Table 1. Mixed-level multiple linear regression model predicting response time

Level 1 variables	Estimate	95% CI
Response variables		
Nonanswer response	131.78***	59.34, 204.23
Confirmation	-206.87***	-268.61, -145.12
Visible response component	-86.93***	-136.76, -37.10
Question variables		
Information request only	129.38***	79.30, 179.46
Questioner gaze	-69.28**	-123.48, -15.08
Context variables		
Level 2: Variance at language level	19555.05*	7342.23, 57304.20
Level 3: Variance at interaction level	14091.24***	7715.57, 25735.37

The multivariate model shows that nonanswer responses are slower than answer responses and responses to information questions are slower than responses to other sorts of questions. Confirmations, responses with a visible component, and responses to questions that are delivered with speaker gaze are shown to be faster than disconfirmations, vocal responses and other sorts of questions, respectively. Language (i.e., the language being spoken) and conversation (i.e., the conversation from which a data point was taken) were treated as levels and thus the results are language independent. *, $P \leq 0.05$; **, $P \leq 0.01$; ***, $P \leq 0.001$.

