

DOES A ROBOT'S GAZE BEHAVIOR AFFECT ENTRAINMENT IN HRI?

Jay KEJRIWAL

*Institute of Informatics
Slovak Academy of Sciences
Bratislava, Slovakia*

✉

*Faculty of Informatics and Information Technology
Slovak Technical University
Bratislava, Slovakia
e-mail: kejriwal.jay@gmail.com*

Chinmaya MISHRA

*Furhat Robotics AB
Stockholm, Sweden*

✉

*Max Planck Institute for Psycholinguistics
Nijmegen, Netherlands
e-mail: chinmaya.mishra@mpi.nl*

Gabriel SKANTZE

*Furhat Robotics AB
Stockholm, Sweden*

✉

*KTH Royal Institute of Technology
Stockholm, Sweden
e-mail: skantze@kth.se*

Tom OFFREDE

*Institut für deutsche Sprache und Linguistik
Humboldt-Universität zu Berlin*

Berlin, Germany

e-mail: offredet@hu-berlin.de

Štefan BEŇUŠ

*Institute of Informatics
Slovak Academy of Sciences
Bratislava, Slovakia*

✉

*Constantine the Philosopher University
Nitra, Slovakia
e-mail: sbenus@ukf.sk*

Abstract. Speakers tend to engage in adaptive behavior, known as entrainment, when they reuse their partner's linguistic representations, including lexical, acoustic prosodic, semantic, or syntactic structures during a conversation. Studies have explored the relationship between entrainment and social factors such as likeability, task success, and rapport. Still, limited research has investigated the relationship between entrainment and gaze. To address this gap, we conducted a within-subjects user study ($N = 33$) to test if gaze behavior of a robotic head affects entrainment of subjects toward the robot on four linguistic dimensions: lexical, syntactic, semantic, and acoustic-prosodic. Our results show that participants entrain more on lexical and acoustic-prosodic features when the robot exhibits well-timed gaze aversions similar to the ones observed in human gaze behavior, as compared to when the robot keeps staring at participants constantly. Our results support the predictions of the computers as social actors (CASA) model and suggest that implementing well-timed gaze aversion behavior in a robot can lead to speech entrainment in human-robot interactions.

Keywords: Entrainment, alignment, HRI, linguistic

1 INTRODUCTION

Entrainment in spoken interaction is a ubiquitous and multi-faceted phenomenon observed in Human-Human Interaction (HHI), whereby people adjust their speaking behavior in response to the speech patterns of their interlocutors. Several studies have examined this phenomenon using diverse approaches and referred to it with various names, such as alignment [1], accommodation [2], 'the Chameleon Effect' [3], convergence [2, 4], coordination [5], coupling [6, 7], mimicry [8], mirroring [9], priming [10], and synchrony [11], among many more. According to the psycholinguistic literature, entrainment happens on various linguistic dimensions, such as acoustic-prosodic features [12], lexical choice [13], syntactic structure [14],

or semantic [15, 16]. A comprehensive discussion on the types of entrainment, classification criteria, and terminology can be found in [17]. All these approaches assess the level of similarity of different linguistic features and attribute the similarities to internal (social) cognitive mechanisms or external social factors.

Entrainment in HHI has been studied extensively, and several theories have been proposed to explain it. The Interactive Alignment Model (IAM) [1] and Communication Accommodation Theory (CAT) [18] are two major theoretical frameworks that address entrainment. CAT suggests that speakers dynamically adapt their communication behaviors based on their interaction with their partner. This process involves either converging toward their interlocutors' communication behaviors to reduce social distance or diverging from them to increase it. On the other hand, IAM suggests that entrainment is an automatic process triggered by a priming mechanism that operates on linguistic representations and is based on a direct link between perception and production in conversation. Although differing in their perspectives, both models agree that entrainment plays a crucial role in HHI. The theories suggest that social processes and automatic cognitive mechanisms can coexist and vary in significance across individuals, which may help explain why speakers exhibit different degrees of entrainment.

The development of spoken dialogue systems (SDS) that can accurately recognize and understand social cues and behaviors is a complex and ongoing process. Despite significant progress, researchers have not yet been able to satisfactorily model the intricate dynamics involved in human conversations. One line of research in this domain is to explore the application of entrainment findings from HHI to HMI. Entrainment functionality at various linguistic levels has shown the potential to improve the naturalness and effectiveness of SDS, which could increase the number of potential applications. Several studies have reported encouraging results, such as [19], who proposed a model for lexical entrainment that uses a data-driven approach to identify the most appropriate terms for system prompts, leading to improved SDS performance. Similarly, [20] reported accuracy improvements in speech recognition through speech rate induction, and [21] reported students' increased knowledge gains when a tutoring SDS entrained to their pitch and intensity. Similarly, in [22], authors found that adjusting the conversational agent's mean pitch to match that of its human interlocutor resulted in more rapport and natural communication. This suggests that advanced methods of implementing entrainment may improve the efficiency and effectiveness of HMI.

Researchers have also investigated the relationship between entrainment and various social factors. They found that entrainment is associated with different social aspects of a conversation, such as naturalness [23, 24], rapport [25, 22], task success [26], liking [27], and cooperation [28]. Further, researchers have explored non-verbal aspects of communication. Eye gaze behavior, such a non-verbal cue and the focus of this paper, has proven vital in facilitating smooth communication. Studies have explored the relationship between gaze and various social factors in HHI, such as conversational feedback [29], trust, rapport and shared attention [30], and turn-taking [31]. It has been observed that lack of eye contact during video-

conferencing can negatively impact turn-taking efficiency [32]. In addition, the presence of eye contact during spoken interaction can significantly enhance performance in word acquisition tasks [33]. However, the relationship between gaze and entertainment in HMI has received less attention. In a recent study conducted by [34], speech entrainment was analyzed by measuring the mean pitch of speech collected from 33 participants subjected to two modes of robot's gaze behavior (fixed vs. variable) described in [35]. However, the results indicated no significant differences between the two conditions.

In this work, we extend the study in [34] by focusing on other linguistic dimensions, i.e., lexical, syntactic, and semantic levels. Further, at the acoustic-prosodic level, we extend the entertainment analysis to eight acoustic-prosodic features: mean and max pitch, mean and max intensity, jitter, shimmer, noise-to-harmonics ratio (NHR), and speech rate, which in previous studies showed entrainment in HHI. The current study thus aims to investigate entrainment in Human-Robot Interaction (HRI) on four linguistic levels (lexical, syntactic, semantic, and acoustic) under two different gaze conditions. Entrainment was measured using the entrainment metrics proposed by [36] for acoustic-prosodic features and [16] for text-based features extracted from transcripts.

The contributions of this work can be summarized as follows:

- It explores the relationship between gaze behavior and entrainment in HRI.
- It investigates entrainment in four linguistic dimensions, i.e., lexical, syntactic, semantic, and acoustic-prosodic.

2 RELATED WORK

The relationship between gaze and social factors in HHI has been extensively studied. For instance, [31] explored the relationship between interlocutors' eye gaze and spoken utterances and how it affects entrainment. They used their own corpus [37], which consisted of three-party conversations, to train data-driven models to classify turn-taking. The data was annotated with dialogue acts, eye-gaze, and turn-taking features for analysis. The results showed that combining dialogue act features with eye-gaze features resulted in higher classification accuracy. Moreover, the study found that eye-gaze features were more important than speech signals for turn management. Similarly, [38] explored the relationship between visual cues and speech entrainment by investigating whether speakers entrain more when they see each other as opposed to when they only hear each other. In an interactive search task, pairs of participants were given a set of keywords to say repeatedly. While one half could only hear each other, the other half could see and hear each other. The study results indicated that the speakers entrained more towards each other when they could see each other, suggesting that visual information enhances speech alignment.

[39] conducted a study to explore the relationship between gaze and gestural alignment during face-to-face interactions. The latter was operationalized as a de-

gree of similarity between adjacent representational hand gestures from two interlocutors in terms of finger and palm orientation, handedness, gesture type, and hand shape. They used the InSight Interaction Corpus [40], which consists of 15 recordings of face-to-face conversations that last about 20 minutes each. The study revealed that the listener's gaze significantly affects gestural alignment, whereas the speaker's gaze does not significantly impact gestural alignment. It was also found that individuals tend to mimic similar gestures in their next turn when they concentrate their visual attention on the speaker's movements. The study highlights the importance of gaze behavior in gestural alignment. In a recent study [41], the authors investigated the characteristics of gaze and its relation to speech behavior during video-mediated face-to-face interactions between parents and their children. The study involved 81 parent-child dyads who interacted with each other in two scenarios, namely, cooperative and conflictive family topics.

The study's findings showed that children spoke more in the cooperation scenario, whereas parents spoke more in the conflict scenario. Additionally, parents gazed slightly more at their children's eyes in the conflict scenario compared to the cooperation scenario. Both parents and children looked more at the other's mouth region while listening than speaking. Overall, this study contributes to the literature on the importance of non-verbal communication cues in HHI.

When it comes to HRI, however, studies exploring the relationship between gaze and social factors are limited. Few researchers have started addressing this gap. For example, in a recent study [35] examined the relationship between robot and human gaze behavior. The study involved a within-subjects design where 33 participants interacted with a Furhat robot in two experimental conditions: Fixed Gaze and Gaze Aversion. In the Fixed Gaze condition, the robot maintained constant eye contact with the participant; in the Gaze Aversion condition, it produced gaze aversions throughout the conversations, more similar to how humans behave. The study found that participants tended to avert their gaze more often and for longer when the robot maintained constant eye contact than when it produced gaze aversion. This shows the significance of incorporating well-timed gaze aversions in robotic conversational agents. If robots do not exhibit gaze aversions, then users may have to put in extra effort to avoid frequent mutual gaze with the robot, which can make the interaction more difficult. On the contrary, the same study also reported subjective evaluations of the perceived interaction, where participants preferred the robot under the fixed gaze condition to be more human-like. In subsequent work, [34] utilized data collected in [35] and explored the relationship between gaze and speech entrainment. PRAAT toolkit was employed to extract mean pitch values of the participants' and robots' speech at each turn exchange. It was found that speakers tend to entrain to the mean pitch of the robot. However, no significant differences in mean pitch entrainment between the Fixed Gaze and Gaze Aversion conditions were reported.

This work aims to add to the existing studies on speech entrainment and gaze behavior in HRI by adopting a more comprehensive approach. While previous work [34] measured entrainment on the mean pitch, entrainment may have happened on other

features. Empirical evidence on speech entrainment has shown that speakers entrain and dis-entrain on different prosodic features [36, 42, 43, 44, 45]. Thus, more acoustic-prosodic features should be examined to assess speech entrainment. Further, the linear regression models used in the previous study did not consider the order effect – the sequence in which the conditions were presented to the participants (Fixed Gaze followed by Gaze Aversion or vice versa). This could have influenced the results and should be taken into account. Additionally, the study only investigated prosodic entrainment. We believe that a more comprehensive understanding of speech entrainment can be achieved by complementing the acoustic-prosodic evaluation of entrainment with also analyzing text-based features extracted from the transcripts at different linguistic levels, such as lexical, syntactic, and semantic. We expand the scope of the original study by examining eight acoustic-prosodic features and four linguistic dimensions, including lexical, syntactic, semantic, and acoustic-prosodic. Furthermore, we use linear mixed-effect models to compare entrainment in two different gaze conditions while considering the order effect and its interaction with the gaze condition. Our study thus provides a comprehensive understanding of entrainment in HRI and how gaze behavior affects entrainment on various linguistic levels.

3 HYPOTHESIS

The computers are social actors (CASA) theory [46] suggests that humans interact with media and computers as if they were real individuals. This theory proposes that individuals subconsciously apply scripts for interacting with humans to social interactions when they detect social cues of humanity. While this may no longer apply to old technology such as desktop computers, as per a recent study [47], the authors conclude that the CASA theory would apply to emergent technologies. We argue that HRI is one such emergent technology to which the CASA theory would apply. When a robot exhibits human-like behavior, it is perceived by humans as having agency, which in turn encourages them to treat the robot as a social actor/agent. Few studies further support this observation. For instance, when a robot exhibits appropriate emotions, it tends to be perceived as more intelligent [48] and trustworthy [49].

Entrainment is a phenomenon that reflects the degree of social closeness among speakers during an interaction. It suggests that the closer speakers get, the lesser the social distance between them [2]. In the context of HRI, research has shown that people tend to avert their gaze less when a robot exhibits well-timed gaze aversion behavior [35]. It indicates that human-like gaze aversion behavior in robots can have a positive influence on overall experience and ease of communicating with the robot. We assume that speech entrainment might be one of the computationally accessible indicators that can, in part, inform us about the cognitive states of the human interlocutors and their perceived agency of the robot. We thus expect that human-like gaze behavior by a robot during HRI would also have a positive influ-

ence on the entrainment exhibited by human interlocutors. Although the original study [34] examining only a single feature of the mean pitch did not find support for this expectation, entrainment has been established as a complex and multi-faceted phenomenon [17, 50]. Therefore, we believe that employing comprehensive features and extensive analysis could shed additional light on the relationship between speech entrainment and gaze behavior. Hence, in a similar experimental setup, as the one used in [35], we hypothesize that participants will entrain more with the robot when it exhibits well-timed gaze aversions during an interaction (**H1**).

4 METHOD

In this section, we provide a brief overview of the study design, procedure, and data collection. A more detailed description of the method can be found in [35].

Participants interacted with two robot characters in a within-subjects design under two experimental conditions: Fixed Gaze (FG) and Gaze Aversion (GA). The robot was a Furhat robot [51], which has a back-projected face capable of exhibiting human-like gaze behaviors. In the FG condition, the robot maintained a fixed gaze at the participant, while in the GA condition, the robot averted its gaze away at appropriate timings using the GCS proposed in [52], mimicking human-like gaze behavior. This gaze aversion behavior was designed to emulate conversational gaze cues related to turn-taking, intimacy regulation, and joint attention. The order in which the conditions were presented to the participants (Order 1: FG \rightarrow GA; and Order 2: GA \rightarrow FG) were alternated. For instance, if Participant 1 was presented with Order 1; then Participant 2 was presented with Order 2. The participant’s speech, eye gaze behavior, and subjective perception of the conversation with the robot were recorded during the study.

4.1 Participants

The study involved 33 participants assigned male at birth, with ages ranging between 21 and 56 years ($M = 30.55$; $SD = 8.07$). Most participants were L2 speakers of English, with only five being L1 speakers of English. Based on their LexTALE language proficiency scores [53], 16 participants were classified at the C1 to C2 level, 15 at B2, and two at B1. Each participant in the study was compensated with a voucher valued at 100 SEK.

4.2 Procedure

The robot began by introducing itself and explaining the purpose of the conversation, see Figure 1. It then asked the participant six questions, giving the participant as much time as needed in between questions. The robot also answered its own question after the participant had finished answering it, before asking the next question. This made the conversation feel more interactive rather than a one-sided interview. The

procedure was the same for both conditions. After each interaction, the participant completed a questionnaire regarding their perception of the interaction with the robot. They also completed the [54]'s version of the Big Five personality inventory and the LexTALE English proficiency test [53] between the experimental conditions, which served as a distractor task.



Figure 1. Snapshot of the data collection setup where a participant interacts with the Furhat robot

4.3 Data Collection

We recorded three types of data during the study: gaze data, speech data, and subjective responses. The gaze behavior of the participants during the interactions was recorded using a Tobii Pro Glasses 2 eye-tracker. Audio recordings of the conversation between the participants and the robot were made using a Zoom H5 multi-track microphone. Finally, the subjective responses to a 9-point Likert scale questionnaire about the participants' perception of their interaction with the robot were collected at the end of each interaction. The participants were informed about the data being collected and gave their informed consent at the beginning of the experiment. The study was approved by the Ethics Committee of the Faculty of Language, Literature and Humanities of the Humboldt-Universität zu Berlin.

5 MEASURES AND ANALYSIS

Hypothesis **H1** proposed that the participants would entrain more towards the robot when it exhibits human-like gaze aversions (GA condition) as compared to the FG

condition. To test this, we investigated entrainment on four different linguistic dimensions: lexical, syntactic, semantic, and acoustic-prosodic. Textual data extracted from the audio data was used to assess the entrainment at the lexical, syntactic, and semantic levels. Various acoustic-prosodic features were extracted from the audio data to analyze entrainment at each of the extracted feature levels. The following subsections discuss the feature extraction process, the measures of entrainment used in this study, and the annotation and analysis of the data.

5.1 Feature Extraction

We extracted lexical, syntactic, semantic, and acoustic features from each turn. As a first step for text-based features, we pre-processed each utterance by removing numbers, punctuation, and other special symbols.

Lexical and syntactic features: We utilized the methodology proposed in the ALIGN toolkit [55], an entrainment analysis tool, to extract both lexical and syntactic features from each utterance in the dialog. The tool employs n-gram sequences to extract these features. For lexical feature extraction, we first tokenized each word in every utterance and then converted them into their lemma form using the Stanza toolkit [56]. By doing so, we were able to reduce all inflectional and derived forms of words to a common base form. Subsequently, we measured the term frequency (TF) for each lemmatized word in the text and generated vectors for each turn based on the TF of every lemmatized word. Lastly, each utterance was encoded into 512 one-dimensional lexical features.

To extract syntactic features from utterances, we tokenized each word and transformed them into their respective parts of speech (POS) tags using the Stanza toolkit [56]. The POS tagging process helps to classify words in an utterance based on their associated parts of speech (e.g., noun, verb, adjective), which is essential for understanding their grammatical structure and meaning. We then converted POS sequences into bi-gram units for each utterance, as these units contain crucial information about the grammatical structure and the relationship between adjacent words in an utterance. The frequency of the bi-gram sequence, representing syntactic units within each turn, was calculated and represented as a vector with 512 one-dimensional syntactic features.

We also utilized CASSIM (ConversAtion level Syntax SIMilarity Metric) [57] to extract syntactic features. This tool allowed us to compare structural differences utilizing parse trees. We generated parse trees for each utterance using CASSIM and measured conversational syntax similarity using edit distance¹.

Semantic features: We utilized a neural network-based DistilBERT model (ms-marco-distilbert-base-v4) [58] pre-trained on the MS MARCO (Microsoft Machine Reading Comprehension) dataset. The dataset comprises a large-scale

¹ A lower edit distance indicates closeness, while a greater distance indicates the opposite.

information retrieval corpus based on real user search queries using the Bing search engine. Each turn in the dialog was encoded into a set of fixed-length vectors known as embeddings. Each turn is represented by 768 one-dimensional semantic features, which enables us to capture the meaning and context of the conversation efficiently.

Acoustic features: Using the PRAAT toolkit [59], we extracted eight acoustic-prosodic features for each turn, namely the mean and max pitch, the mean and max intensity, jitter, shimmer, the noise-to-harmonics ratio (NHR), and the speaking rate. The speaking rate was computed by counting the number of syllables per second from the orthographic transcriptions of the data. Additionally, we normalized all the extracted features by the speaker using the z -score.

5.2 Quantifying Entrainment

Various entrainment metrics have been proposed by researchers that capture different aspects of entrainment and employ different methodologies (for a review, see [17, 50]). For measuring entrainment in acoustic and textual features, we employed two different metrics.

We utilized the approach suggested in [36] to measure acoustic-prosodic proximity. To determine the entrainment distance between dyads, we measured the absolute distance between each adjacent turn of the speakers on each feature, as shown in Equation (1):

$$Ent_{acoustic} = |SpeakerA_{feat} - SpeakerB_{feat}|. \quad (1)$$

Here, $feat$ denotes the corresponding speaker’s feature and $Ent_{acoustic}$ computation was performed 8 times, separately for each prosodic feature. Entrainment distance represents the similarity of a prosodic feature over these adjacent turn transitions uttered by speakers A and B in a conversation. A lower distance indicates closeness, while a greater distance indicates the opposite.

To measure entrainment on the text-based features (i.e., at the lexical, syntactic, and semantic levels), we used cosine similarity as a distance measure. Specifically, we calculated the cosine similarity between a speaker’s embedding² and the adjacent embedding of their interlocutor, as shown in Equation (2):

$$Ent_{text-based} = \cos(\vec{A}, \vec{B}) = \frac{\vec{A} \cdot \vec{B}}{|\vec{A}| |\vec{B}|}. \quad (2)$$

Here, $Ent_{text-based}$ computation was performed 3 times, separately for each textual feature. In contrast to acoustic-prosodic entrainment distance, a greater textual entrainment distance indicates closeness, while a lower distance indicates the opposite.

² Embeddings are dense numerical representations of textual/acoustic features expressed as vectors in a low-dimensional space.

5.3 Annotation

The beginning and end of the turns of each speaker (human and robot) were annotated manually using Praat [59]. Text transcription of the audio signal for each turn was automatically obtained using the fairseq model *facebook/mms-1b-all* by [60].

5.4 Analysis

The analysis compared entrainment on four linguistic dimensions, i.e., lexical, syntactic, semantic, and acoustic. This analysis aimed to determine on which linguistic levels speakers are closest to each other under different gaze conditions. Entrainment distance was measured from each session at each linguistic level using Equations (1) and (2). We analyzed entrainment distance across four distinct linguistic levels, with separate linear mixed models (LMMs) developed for each linguistic feature, using the *lmerTest* R package [61]. Specifically, we trained eleven models, eight for each acoustic-prosodic feature and three for the lexical, syntactic, and semantic features, respectively. Each model considers entrainment distance as a dependent variable. The fixed effects for each model included

- the experimental condition (Fixed Gaze, FG, and Gaze Aversion, GA),
- the order in which the conditions were presented to the participants, and
- the interaction between condition and order.

We included participant as a random effect variable. Formula (3) was used to fit each model. We fit each LMM by REML t-tests and used Satterthwaite approximations to determine the degrees of freedom. Finally, the p-values were derived from the output of each model. The post-hoc testing of each model was carried out by adjusting multiple comparisons using Tukey’s Multiple Contrasts (part of R package “emmeans”) [62].

$$\text{Entrainment distance} \sim \text{condition} + \text{Order} + \text{condition} * \text{Order} + (1 | \text{Participant}). \quad (3)$$

6 RESULTS

6.1 Text-Based Entrainment Models

Figure 2 shows the mean entrainment distance (see Section 5.2) of participants in the two experimental conditions GA and FG for a) lexical, b) syntactic, and c) semantic linguistic levels. Table 1 summarizes the results of the LMM fits for all three levels and post-hoc comparison for the significant models.

In lexical entrainment, the results indicated no significant main effect of the experimental condition on the entrainment of the participant towards the robot. However, a significant main effect of the order was observed, indicating that speakers

a) Lexical level

Fixed effects:					
Variable	Estimate	SE	df	<i>t</i>	<i>p</i>
(Intercept)	0.276	0.017	50.243	16.244	< .001
ConditionFG	0.034	0.023	52.656	1.443	0.155
Order 2	0.055	0.023	53.713	2.367	0.022
ConditionFG:Order 2	-0.085	0.041	31.495	-2.066	0.047

Post-hoc comparison:					
	Contrast	β	<i>t ratio</i>	<i>p</i>	
ConditionGA	Order 1-Order 2	-0.055	-2.367	0.022	
ConditionFG	Order 1-Order 2	0.029	1.25	0.217	
Order 1	GA-FG	-0.034	-1.44	0.155	
Order 2	GA-FG	0.051	2.173	0.034	

b) Syntax level

Fixed effects:					
Variable	Estimate	SE	df	<i>t</i>	<i>p</i>
(Intercept)	0.549	0.010	50.516	56.964	< 0.001
ConditionFG	-0.016	0.013	54.400	-1.161	0.251
Order 2	0.008	0.013	53.281	0.572	0.57
ConditionFG:Order 2	0.004	0.023	28.336	0.168	0.868

c) Semantic level

Fixed effects:					
Variable	Estimate	SE	df	<i>t</i>	<i>p</i>
(Intercept)	0.205	0.011	68.976	19.472	< 0.001
ConditionFG	0.037	0.015	73.682	2.529	0.014
Order 2	-0.027	0.015	76.141	-1.818	0.073
ConditionFG:Order 2	-0.052	0.024	33.079	-2.206	0.034

Post-hoc comparison:					
	Contrast	β	<i>t ratio</i>	<i>p</i>	
ConditionGA	Order 1-Order 2	0.027	1.818	0.073	
ConditionFG	Order 1-Order 2	0.079	5.415	< 0.001	
Order 1	GA-FG	-0.037	-2.528	0.013	
Order 2	GA-FG	0.016	1.06	0.293	

Table 1. LMM model output comparing entrainment at the lexical, syntactic, and semantic levels in two different conditions with Gaze Aversion (GA) condition as the reference value. Significant *p*-values are shown in bold with *p* < 0.05 with post-hoc comparisons for significant models.

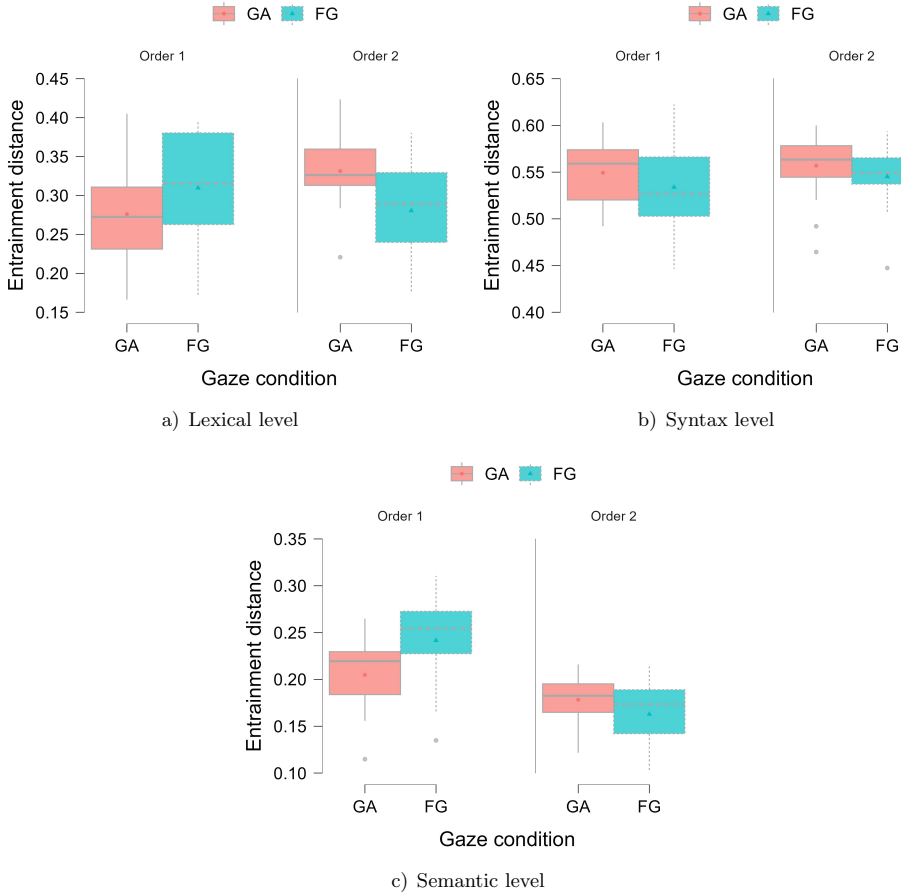


Figure 2. Entrainment in two different gaze conditions, Fixed Gaze (FG) and Gaze Aversion (GA), and two orders, i.e., Order 1 (FG → GA) and Order 2 (GA → FG) at lexical, syntactic, and semantic features

entrained more in Order 2, i.e., (GA → FG). The significant interaction and its subsequent post-hoc analysis revealed that the greater lexical entrainment in Order 2 is driven by the significantly higher lexical entrainment for the GA. This implied that speakers entrained lexically more under the GA condition only in Order 2. This partially supports hypothesis **H1**, which predicted that participants would entrain more under the GA condition. At the semantic level, we found a main effect of the experimental condition whereby speakers entrained more in the FG condition as compared to the GA condition. Further, the significant interaction and its subsequent analysis showed that greater semantic entrainment in Order 1, i.e., (FG → GA), is only significant for the FG condition and that speakers entrained

semantically more under the FG condition only in Order 1. This finding does not support **H1**.

We did not find any significant main effects of experimental conditions or the orders on syntactic entrainment using the ALIGN toolkit. However, previous empirical studies have demonstrated that the choice of methodology can significantly influence entrainment results [50]. In measuring syntactic entrainment, two commonly used methods include n-gram sequence [55] and parse-tree comparison [57]. To further test the degree of syntactic entrainment, we used CASSIM (ConversAtion level Syntax SImilarity Metric) [57] (section 5.1) and compared syntactic entrainment distance in both experimental conditions using the LMM model. Table 2 shows the results where we found no significant difference across both conditions, order, and their interactions, which indicates that participants used similar syntactic structures in both conditions.

a) Syntax (CASSIM)

Fixed effects:					
Variable	Estimate	SE	df	<i>t</i>	<i>p</i>
(Intercept)	0.564	0.009	49.226	56.969	< 0.001
ConditionFG	-0.005	0.013	52.903	-0.43	0.669
Order 2	0.017	0.013	51.108	1.315	0.194
ConditionFG:Order 2	-0.017	0.023	31.470	-0.752	0.457

Table 2. LMM model output comparing entrainment at syntactic level using CASSIM in two different gaze conditions and order with Gaze Aversion (GA) condition as the reference value

6.2 Acoustic-Prosodic Based Entrainment Models

Figure 3 shows the mean entrainment distance of participants under the two experimental conditions: GA and FG for a) Mean pitch and b) NHR. It needs to be kept in mind that for acoustic-prosodic features, the lower the mean entrainment distance the more the entrainment (see Section 5.2). We observed that only the LMM models for mean pitch and NHR, out of the eight models fit for acoustic-prosodic features, showed significant effects. The outcomes of the LMM fits and post-hoc comparisons for these two features are summarized in Table 3.

For the mean pitch model, we found a main effect of the experimental condition whereby speakers aligned significantly more on mean pitch with the robot in the GA condition as compared to the FG condition. This supported hypothesis **H1**. For the NHR model, we observed significant main effects of both experimental conditions and order. Participants entrained significantly more in the GA condition, which was in line with **H1**. Additionally, it was observed that participants entrained more in Order 1, where they interacted with the robot under the FG condition first followed by the GA condition. Since the interaction yielded the *p*-value of 0.051, we also

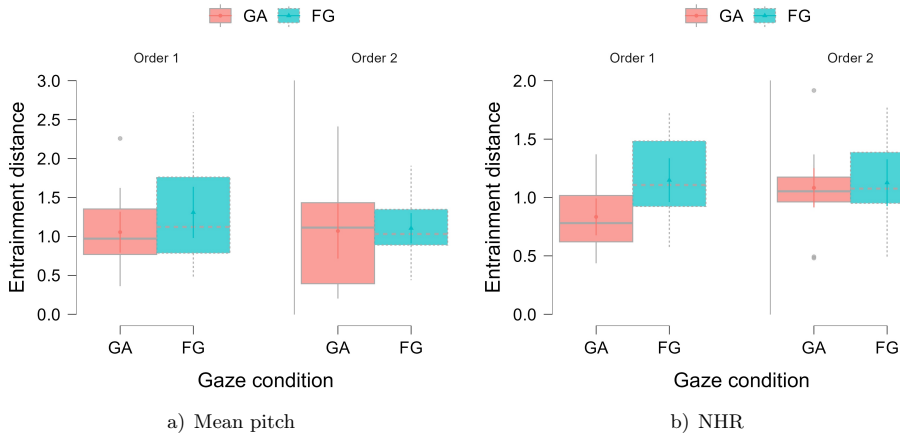


Figure 3. Entrainment in two different gaze conditions, Fixed Gaze (FG) and Gaze Aversion (GA), and two orders, i.e., Order 1 (FG → GA) and Order 2 (GA → FG) on mean pitch and NHR

a) Mean pitch

Fixed effects:					
Variable	Estimate	SE	df	<i>t</i>	<i>p</i>
(Intercept)	1.056	0.075	50.317	14.081	< 0.001
ConditionFG	0.224	0.106	52.866	2.116	0.039
Order 2	0.041	0.106	54.079	0.387	0.701
ConditionFG:Order 2	-0.227	0.185	30.903	-1.229	0.228

b) NHR

Fixed effects:					
Variable	Estimate	SE	df	<i>t</i>	<i>p</i>
(Intercept)	0.828	0.061	73.411	13.558	< 0.001
ConditionFG	0.324	0.087	78.741	3.726	< 0.001
Order 2	0.255	0.088	82.338	2.91	0.039
ConditionFG:Order 2	-0.28	0.138	31.528	-2.036	0.051
Post-hoc comparison:					
	Contrast	β	<i>t ratio</i>	<i>p</i>	
ConditionGA	Order 1–Order 2	-0.251	-2.113	0.039	
ConditionFG	Order 1–Order 2	0.022	0.188	0.851	
Order 1	GA–FG	-0.032	-2.68	0.009	
Order 2	GA–FG	-0.044	-0.371	0.712	

Table 3. LMM model output and post-hoc comparisons for significant acoustic-prosodic models mean pitch and NHR in two different gaze conditions and order with Gaze Aversion (GA) condition as the reference value. Significant *p*-values are shown in bold with $p < 0.05$

performed a post-hoc analysis. It was observed that people entrained more in GA condition only in Order 1, further supporting **H1**.

7 DISCUSSION

Based on the gaze behavior of the robotic interlocutor, entrainment was measured in two different experimental conditions (FG & GA) to examine how participants aligned on lexical, syntactic, semantic, and acoustic-prosodic levels. Potential differences (or the lack thereof) in entrainment under the two different gaze conditions of the robot can inform us about the underlying relationship between gaze and entrainment during HRI. We predicted that the participants would exhibit more entrainment towards the robot in the GA condition as compared to the FG condition, across different linguistic levels (**H1**). Significant differences between conditions were observed across lexical, semantic, and acoustic-prosodic levels. We found that participants entrained more in GA condition at the lexical and acoustic-prosodic levels (specifically at mean pitch and NHR), which was in line with **H1**. Additionally, we found that the order of the experimental conditions to which the participants were exposed had a significant effect on entrainment at the lexical level. This meant the participants lexically entrained more with a robot depending on whether they first interacted under the GA or FG conditions.

We observed no significant differences between the experimental conditions, order, or interactions at the syntactic level using both the bi-gram and parse-tree methodologies. This might be because of the specific role assigned to the participants, where they always had to answer the open-ended questions asked by the robot across both conditions. This restricted the syntactic structure of the participants' responses to be similar across the conditions, as answering questions entails a similar syntactic structure. Since participants' responses lacked variability, this might have resulted in similar syntactic entrainment distance in both gaze conditions. On the other hand, if the conversation were free-flowing, interlocutors would freely alternate between asking questions and answering. It might result in more variability in the syntactic structure of the interlocutors' responses. The lack of a free-flowing conversation with a robot, thereby, the restricted syntactic structure of the responses by the participants across the conditions, might have led to finding no significant differences between the conditions at the syntactic level.

Contrary to our predictions, it was observed that participants entrained more in the FG condition as compared to the GA condition at the semantic level. This may have arisen due to the erratic gaze behavior by the robot under the GA condition as reported in [35]. It was observed that during the GA condition, the robot directed its gaze away from the participants even when they were seated in front of them until the confederate initiated the interaction. This unnatural robot gaze behavior could have negatively influenced the perception of the robot's abilities by the participants (the robot could have been perceived as having less agency). As a result, in subjective evaluations of the perceived interaction outlined in [35], participants rated the robot

in the FG condition as more human-like than the GA condition, which aligns with the semantic entrainment results obtained in our study. This could suggest that the perception of a robot’s capabilities could have a direct influence on the entrainment at the semantic level during an HRI.

Secondly, we examined acoustic-prosodic entrainment on eight prosodic features across the two experimental conditions. We found that only two features, mean pitch, and NHR, displayed significant differences in entrainment across conditions. We observed that speakers entrained more with the robot when in the GA condition as compared to the FG condition. Mean pitch is often related to naturalness and rapport [25, 22] between interlocutors. Since participants entrained significantly more on the mean pitch in the GA condition, we can infer that the robot is perceived as more natural and has a better rapport with the participants. Further, we also found that the order of gaze conditions significantly affected NHR, where participants interacted more acoustically with the robot under the GA condition when they first interacted with the robot under the FG condition. This could highlight that the participants were able to perceive the difference in the gaze behavior of the robot across the conditions. The more human-like gaze aversion behavior in the GA condition after being exposed to the unnatural fixed gaze had a positive influence on the entrainment in the NHR level. Empirical evidence on acoustic-prosodic entrainment suggests people entrain and dis-entrain on different acoustic-prosodic features depending on a variety of social factors such as gender and personality of the interlocutors [63], the emotional state of the speaker [64], the relationship between the speakers [65], the context of the conversation, and the interaction between all these factors [42]. Therefore, we only found a significant difference in entrainment in two acoustic-prosodic features.

The current study does not corroborate results in [34], where we found no significant difference across conditions on the mean pitch. There are two potential reasons for distinct results. First, the mean pitch extracted in [34] was not z-score normalized. Second, we used different entrainment metrics. For instance, [34] used metrics proposed by [66], whereas, in the current study, we utilized the methodology proposed by [36]. Empirical evidence has shown entrainment results are affected by utilizing different methodologies [42, 50]. Further, our analysis included the Order effect. We observed this effect on lexical and acoustic-prosodic levels, with different orders showing varying degrees of entrainment. Specifically, participants entrained more in Order 2 at the lexical level and Order 1 at the NHR level. We are unable to explain this finding at present, and further research is needed. Lastly, our results show that human-like gaze aversion facilitates entrainment on the acoustic and lexical levels, whereas the semantic level shows the facilitatory effect of the FG condition. We speculate that lexical and acoustic-prosodic levels might be considered more “automatic” or low-level when it comes to priming-based entrainment ([1]). On the other hand, the semantic level can be construed as more high-level and potentially affected more by various social and attitudinal factors. Thus, various aspects of the assumed robot’s agency might affect entrainment at linguistic dimensions differently.

To sum up, the current study’s findings suggest that people entrain more at lexical and acoustic-prosodic levels in the GA condition compared to the FG condition. This finding of the current study is in line with the *computers are social actors (CASA)* theory proposed by [46] as described in Section 3. In the GA condition, the robot’s gaze behavior emulated human-like gaze aversion behavior, which made the participants feel more comfortable during the interaction. This suggests that endowing human-like behavior in robots can be beneficial in HRI.

8 LIMITATIONS AND FUTURE WORK

The results reported in the paper should be interpreted with caution. Among several limitations, we mention four. Firstly, we utilized neural network-based BERT models to extract semantic features from each utterance. These models are trained explicitly on a conversational corpus that allows us to assess semantic entrainment. Our previous study [16] demonstrated that using different neural-based models can influence the results. We compared entrainment behavior in the Columbia games corpus [67] using BERT [68], trained explicitly on conversational data, and the Universal Sentence Encoder (USE) model [69], trained on multiple languages. Our findings indicate that the utilization of features extracted from BERT and USE has a significant impact on the results of entrainment. It is worth noting that USE does not offer any insights into the dataset it is trained on, whereas BERT is trained specifically on the English language dataset. Secondly, we employed Facebook’s fair sequence model for extracting text transcriptions. However, as with all speech-to-text (STT) models, errors can occur during the process of extracting textual features from speech. High word error rates can negatively impact entrainment results. To address this, manual annotation with inter-annotator agreement can be used as a solution. Thirdly, scarcity of data may affect entrainment results. In the current study, participants were asked 6 questions, each in two different conditions. If there were several turn-exchanges in HRI, then the accuracy of the entrainment analysis could be strengthened. Lastly, the robot’s speech was fixed across conditions. As the robot’s questions and responses were pre-determined and fixed, there was no variation in the interaction between each participant, which might affect entrainment outcomes.

We utilized two different entrainment measures, i.e., absolute distance and cosine similarity, for measuring entrainment at acoustic and textual levels, respectively. Hence, only an indirect comparison between acoustic and textual entrainment is possible. In future work, auditory features can be extracted from each turn using TRIPlet Loss network (TRILL) [70], and entrainment distance can be measured using cosine similarity. Since all four entrainment distances will be identical, i.e., measured using cosine similarity, we can further compare the entrainment distance across different linguistic levels more directly by constructing a single (LMM) model. It can allow us to determine which linguistic levels speakers are closest to each other under different gaze conditions.

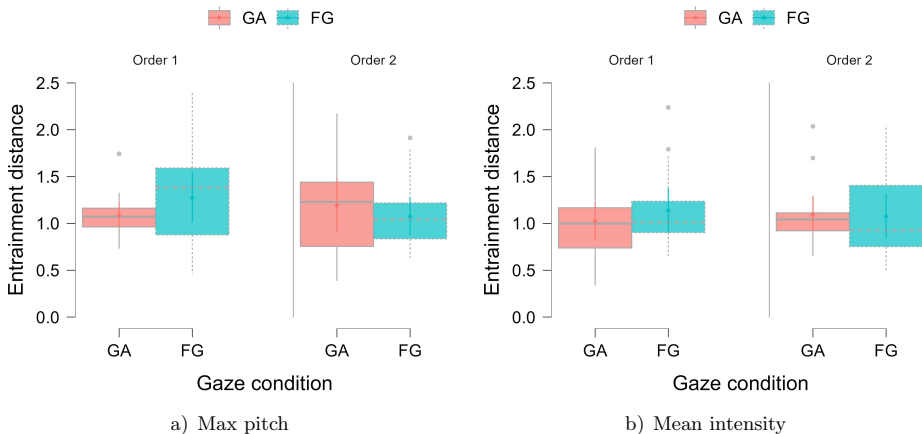
9 CONCLUSION

Our study analyzed entrainment across four linguistic dimensions in HRI with a Furhat robot and revealed interesting findings. Our study found that speakers entrained more in the Gaze Aversion and Fixed Gaze conditions at the lexical and semantic levels, respectively. Furthermore, we observed that the order of interaction had a significant effect on lexical entrainment. At the acoustic level, speakers entrained more in the GA condition on mean pitch and NHR. The results suggest that entrainment can be influenced by various factors, such as the robot’s gaze behavior, the order of the robots one interacts with, and linguistic dimensions. Overall, this study provides valuable insights into the nature of entrainment in HRI and highlights the importance of considering multiple factors in understanding the phenomenon.

Acknowledgements

We would like to thank Susanne Fuchs (Leibniz – Centre General Linguistics (ZAS), Berlin, Germany) and Christine Mooshammer (Humboldt-Universität zu Berlin, Berlin, Germany) for their insights and support.

This project has received funding from the European Union’s Framework Programme for Research and Innovation Horizon 2020 (2014–2020) under the Marie Skłodowska-Curie Grant Agreement No. 859588 and in part from the Slovak Granting Agency grant VEGA2/0165/21 and Slovak Research and Development Agency grant APVV-21-0373.



10 APPENDIX

Figure 4 shows the mean entrainment distance of participants under the two experimental conditions: GA and FG for a) Max pitch, b) Mean intensity, c) Max

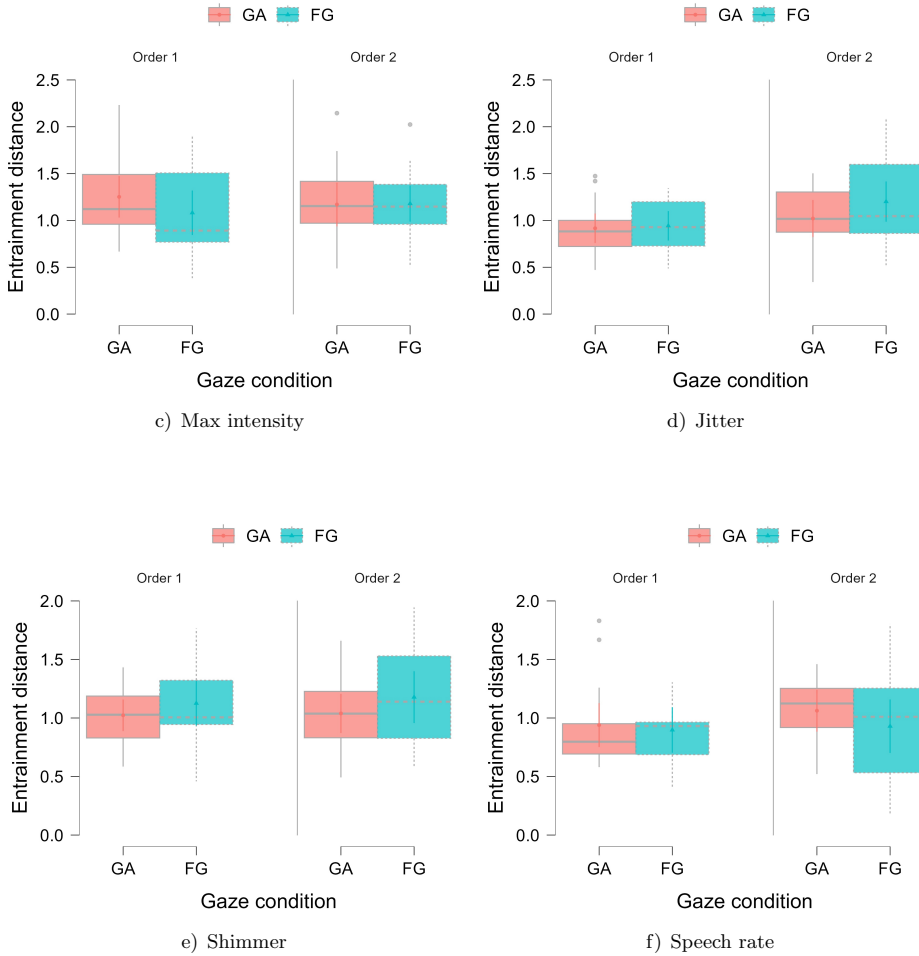


Figure 4. Entrainment in two different gaze conditions, Fixed Gaze (FG) and Gaze Aversion (GA), and two orders, i.e., Order 1 (FG → GA) and Order 2 (GA → FG) on Max pitch, Mean and Max intensity, Jitter, Shimmer, and Speech rate

a) Max pitch

Fixed effects:

Variable	Estimate	SE	df	<i>t</i>	<i>p</i>
(Intercept)	1.085	0.075	52.993	14.435	< .001
ConditionFG	0.174	0.106	55.828	1.638	0.107
Order2	0.124	0.107	57.226	1.164	0.249
ConditionFG:Order2	-0.308	0.184	31.365	-1.676	0.104

b) Mean intensity

Fixed effects:					
Variable	Estimate	SE	df	<i>t</i>	<i>p</i>
(Intercept)	1.027	0.081	48.916	12.667	< 0.001
ConditionFG	0.109	0.114	51.357	0.953	0.345
Order2	0.078	0.115	52.509	0.677	0.501
ConditionFG:Order2	-0.142	0.200	30.334	-0.707	0.485

c) Max intensity

Fixed effects:					
Variable	Estimate	SE	df	<i>t</i>	<i>p</i>
(Intercept)	1.253	0.084	46.676	14.950	< 0.001
ConditionFG	-0.161	0.118	48.616	-1.367	0.178
Order2	-0.081	0.118	49.457	-0.681	0.499
ConditionFG:Order2	0.165	0.212	31.887	0.777	0.443

d) Jitter

Fixed effects:					
Variable	Estimate	SE	df	<i>t</i>	<i>p</i>
(Intercept)	0.915	0.065	66.502	14.021	< 0.001
ConditionFG	0.029	0.093	70.851	0.311	0.757
Order2	0.103	0.093	73.399	1.107	0.272
ConditionFG:Order2	0.148	0.152	32.875	0.973	0.338

e) Shimmer

Fixed effects:					
Variable	Estimate	SE	df	<i>t</i>	<i>p</i>
(Intercept)	1.016	0.064	76.009	15.787	< 0.001
ConditionFG	0.124	0.091	81.553	1.354	0.18
Order2	0.019	0.092	85.340	0.204	0.839
ConditionFG:Order2	0.023	0.144	32.298	0.159	0.875

(f) Speech rate

Fixed effects:					
Variable	Estimate	SE	df	<i>t</i>	<i>p</i>
(Intercept)	0.942	0.086	70.439	10.915	< 0.001
ConditionFG	-0.042	0.122	75.137	-0.340	0.734
Order2	0.121	0.123	77.963	0.980	0.33
ConditionFG:Order2	-0.088	0.200	33.952	-0.443	0.661

Table 4. LMM model output for insignificant acoustic-prosodic models Max pitch, Mean and Max intensity, Jitter, Shimmer, and Speech rate in two different gaze conditions and orders with Gaze Aversion (GA) condition as the reference value

intensity, d) Jitter, e) Shimmer, and f) Speech rate. Table 4 summarizes the results of the LMM fits for the models that were not significant.

REFERENCES

- [1] PICKERING, M. J.—GARROD, S.: Toward a Mechanistic Psychology of Dialogue. *Behavioral and Brain Sciences*, Vol. 27, 2004, No. 2, pp. 169–190, doi: 10.1017/S0140525X04000056.
- [2] GILES, H.—COUPLAND, N.—COUPLAND, J.: Accommodation Theory: Communication, Context, and Consequence. In: Giles, H., Coupland, J., Coupland, N. (Eds.): *Contexts of Accommodation: Developments in Applied Sociolinguistics*. Cambridge University Press, Studies in Emotion and Social Interaction, 1991, pp. 1–68, doi: 10.1017/CBO9780511663673.001.
- [3] CHARTRAND, T. L.—BARGH, J. A.: The Chameleon Effect: The Perception–Behavior Link and Social Interaction. *Journal of Personality and Social Psychology*, Vol. 76, 1999, No. 6, pp. 893–910, doi: 10.1037/0022-3514.76.6.893.
- [4] PARDO, J. S.: On Phonetic Convergence During Conversational Interaction. *The Journal of the Acoustical Society of America*, Vol. 119, 2006, No. 4, pp. 2382–2393, doi: 10.1121/1.2178720.
- [5] BERNIERI, F. J.—ROSENTHAL, R.: Interpersonal Coordination: Behavior Matching and Interactional Synchrony. In: Feldman, R. S., Rimé, B. (Eds.): *Fundamentals of Nonverbal Behavior*. Cambridge University Press, Studies in Emotion and Social Interaction, 1991, pp. 401–432.
- [6] GARROD, S.—PICKERING, M. J.: Joint Action, Interactive Alignment, and Dialog. *Topics in Cognitive Science*, Vol. 1, 2009, No. 2, pp. 292–304, doi: 10.1111/j.1756-8765.2009.01020.x.
- [7] FUSAROLI, R.—RĄCZASZEK-LEONARDI, J.—TYLÉN, K.: Dialog as Interpersonal Synergy. *New Ideas in Psychology*, Vol. 32, 2014, pp. 147–157, doi: 10.1016/j.newideapsych.2013.03.005.
- [8] SCISSORS, L. E.—GILL, A. J.—GERGLE, D.: Linguistic Mimicry and Trust in Text-Based CMC. *Proceedings of the 2008 ACM Conference on Computer Supported Cooperative Work (CSCW '08)*, 2008, pp. 277–280, doi: 10.1145/1460563.1460608.
- [9] DIJKSTERHUIS, A.—BARGH, J. A.: The Perception-Behavior Expressway: Automatic Effects of Social Perception on Social Behavior. In: Zanna, M. (Ed.): *Advances in Experimental Social Psychology*. Elsevier, Vol. 33, 2001, pp. 1–40, doi: 10.1016/S0065-2601(01)80003-4.
- [10] BOCK, J. K.: Syntactic Persistence in Language Production. *Cognitive Psychology*, Vol. 18, 1986, No. 3, pp. 355–387, doi: 10.1016/0010-0285(86)90004-6.
- [11] RICHARDSON, D.—DALE, R.—SHOCKLEY, K.: Synchrony and Swing in Conversation: Coordination, Temporal Dynamics, and Communication. In: Wachsmuth, I., Lenzen, M., Knoblich, G. (Eds.): *Embodied Communication in Humans and Machines*. Oxford University Press, 2008, pp. 75–94.

- [12] LEVITAN, R.—HIRSCHBERG, J.: Measuring Acoustic-Prosodic Entrainment with Respect to Multiple Levels and Dimensions. *Proceedings of Interspeech 2011*, 2011, pp. 3081–3084, doi: 10.21437/Interspeech.2011-771.
- [13] BRENNAN, S. E.—CLARK, H. H.: Conceptual Pacts and Lexical Choice in Conversation. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, Vol. 22, 1996, No. 6, pp. 1482–1493, doi: 10.1037/0278-7393.22.6.1482.
- [14] REITTER, D.—MOORE, J.—KELLER, F.: Priming of Syntactic Rules in Task-Oriented Dialogue and Spontaneous Conversation. In: Sun, R., Miyake, N. (Eds.): *Proceedings of the 28th Annual Conference of the Cognitive Science Society*. 2006, pp. 685–690.
- [15] LIU, Y.—LI, A.—DANG, J.—ZHOU, D.: Semantic and Acoustic-Prosodic Entrainment of Dialogues in Service Scenarios. *Companion Publication of the 2021 International Conference on Multimodal Interaction (ICMI '21)*, ACM, 2021, pp. 71–74, doi: 10.1145/3461615.3491105.
- [16] KEJRIWAL, J.—BEŇUŠ, Š.: Relationship Between Auditory and Semantic Entrainment Using Deep Neural Networks (DNN). *Proceedings of INTERSPEECH 2023*, 2023, pp. 2623–2627, doi: 10.21437/Interspeech.2023-1947.
- [17] WYNN, C. J.—BORRIE, S. A.: Classifying Conversational Entrainment of Speech Behavior: An Expanded Framework and Review. *Journal of Phonetics*, Vol. 94, 2022, Art. No. 101173, doi: 10.1016/j.wocn.2022.101173.
- [18] GILES, H.: Accent Mobility: A Model and Some Data. *Anthropological Linguistics*, Vol. 15, 1973, No. 2, pp. 87–105, <https://www.jstor.org/stable/30029508>.
- [19] LOPES, J.—ESKENAZI, M.—TRANCOSO, I.: From Rule-Based to Data-Driven Lexical Entrainment Models in Spoken Dialog Systems. *Computer Speech & Language*, Vol. 31, 2015, No. 1, pp. 87–112, doi: 10.1016/j.csl.2014.11.007.
- [20] BELL, L.—GUSTAFSON, J.—HELDNER, M.: Prosodic Adaptation in Human-Computer Interaction. *15th International Congress of Phonetic Sciences (ICPhS-15)*, 2003, pp. 2453–2456, https://www.internationalphoneticassociation.org/icphs-proceedings/ICPhS2003/papers/p15_2453.pdf.
- [21] THOMASON, J.—NGUYEN, H. V.—LITMAN, D.: Prosodic Entrainment and Tutoring Dialogue Success. In: Lane, H. C., Yacef, K., Mostow, J., Pavlik, P. (Eds.): *Artificial Intelligence in Education (AIED 2013)*. Springer, Berlin, Heidelberg, *Lecture Notes in Computer Science*, Vol. 7926, 2013, pp. 750–753, doi: 10.1007/978-3-642-39112-5_104.
- [22] LUBOLD, N.—PON-BARRY, H.—WALKER, E.: Naturalness and Rapport in a Pitch Adaptive Learning Companion. *2015 IEEE Workshop on Automatic Speech Recognition and Understanding (ASRU)*, 2015, pp. 103–110, doi: 10.1109/ASRU.2015.7404781.
- [23] NENKOVA, A.—GRAVANO, A.—HIRSCHBERG, J.: High Frequency Word Entrainment in Spoken Dialogue. *Proceedings of the 46th Annual Meeting of the Association for Computational Linguistics on Human Language Technologies: Short Papers (HLT-Short '08)*, 2008, pp. 169–172, <https://dl.acm.org/doi/pdf/10.5555/1557690.1557737>.
- [24] COHN, M.—RAVEH, E.—PREDECK, K.—GESSINGER, I.—MÖBIUS, B.—

- ZELLOU, G.: Differences in Gradient Emotion Perception: Human vs. Alexa Voices. *Proceedings of Interspeech 2020*, 2020, pp. 1818–1822, doi: 10.21437/Interspeech.2020-1938.
- [25] LUBOLD, N.—PON-BARRY, H.: Acoustic-Prosodic Entrainment and Rapport in Collaborative Learning Dialogues. *Proceedings of the 2014 ACM Workshop on Multimodal Learning Analytics Workshop and Grand Challenge (MLA '14)*, 2014, pp. 5–12, doi: 10.1145/2666633.2666635.
- [26] REITTER, D.—MOORE, J. D.: Predicting Success in Dialogue. In: Zaenen, A., van den Bosch, A. (Eds.): *Proceedings of the 45th Annual Meeting of the Association of Computational Linguistics (ACL 2007)*. 2007, pp. 808–815, <https://aclanthology.org/P07-1102>.
- [27] IRELAND, M. E.—SLATCHER, R. B.—EASTWICK, P. W.—SCISSORS, L. E.—FINKEL, E. J.—PENNEBAKER, J. W.: Language Style Matching Predicts Relationship Initiation and Stability. *Psychological Science*, Vol. 22, 2011, No. 1, pp. 39–44, doi: 10.1177/0956797610392928.
- [28] MANSON, J. H.—BRYANT, G. A.—GERVAIS, M. M.—KLINE, M. A.: Convergence of Speech Rate in Conversation Predicts Cooperation. *Evolution and Human Behavior*, Vol. 34, 2013, No. 6, pp. 419–426, doi: 10.1016/j.evolhumbehav.2013.08.001.
- [29] GOODWIN, C.: Action and Embodiment Within Situated Human Interaction. *Journal of Pragmatics*, Vol. 32, 2000, No. 10, pp. 1489–1522, doi: 10.1016/S0378-2166(99)00096-X.
- [30] TREVARTHEN, C.: Emotions in Infancy: Regulators of Contact and Relationships with Persons. In: Scherer, K. R., Ekman, P. (Eds.): *Approaches to Emotion*. Psychology Press, 1984, pp. 129–157.
- [31] JOKINEN, K.—HARADA, K.—NISHIDA, M.—YAMAMOTO, S.: Turn-Alignment Using Eye-Gaze and Speech in Conversational Interaction. *Proceedings of Interspeech 2010*, 2010, pp. 2018–2021, doi: 10.21437/Interspeech.2010-571.
- [32] VERTEGAAL, R.—WEEVERS, I.—SOHN, C.—CHEUNG, C.: GAZE-2: Conveying Eye Contact in Group Video Conferencing Using Eye-Controlled Camera Direction. *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems (CHI '03)*, ACM, 2003, pp. 521–528, doi: 10.1145/642611.642702.
- [33] QU, S.—CHAI, J. Y.: The Role of Interactivity in Human-Machine Conversation for Automatic Word Acquisition. In: Healey, P., Pieraccini, R., Byron, D., Young, S., Purver, M. (Eds.): *Proceedings of the SIGDIAL 2009 Conference*. 2009, pp. 188–195, <https://aclanthology.org/W09-3928>.
- [34] OFFREDE, T.—MISHRA, C.—SKANTZE, G.—FUCHS, S.—MOOSHAMMER, C.: Do Humans Converge Phonetically When Talking to a Robot? In: Skarnitzl, R., Volín, J. (Eds.): *Proceedings of the 20th International Congress of Phonetic Sciences (ICPhS 2023)*. International Phonetic Association, 2023, pp. 3507–3511, https://www.internationalphoneticassociation.org/icphs-proceedings/ICPhS2023/full_papers/712.pdf.
- [35] MISHRA, C.—OFFREDE, T.—FUCHS, S.—MOOSHAMMER, C.—SKANTZE, G.: Does a Robot's Gaze Aversion Affect Human Gaze Aversion? *Frontiers in Robotics and AI*, Vol. 10, 2023, Art.No. 1127626, doi: 10.3389/frobt.2023.1127626.

- [36] LEVITAN, R.—GRAVANO, A.—HIRSCHBERG, J.: Entrainment in Speech Preceding Backchannels. In: Lin, D., Matsumoto, Y., Mihalcea, R. (Eds.): Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies. 2011, pp. 113–117, <https://aclanthology.org/P11-2020>.
- [37] JOKINEN, K.—YAMAMOTO, S.—NISHIDA, M.: Collecting and Annotating Conversational Eye-Gaze Data. In: Kipp, M., Martin, J., Paggio, P., Heylen, D. (Eds.): The Proceedings of the Workshop on Multimodal Corpora: Advances in Capturing, Coding and Analyzing Multimodality. European Language Resources Association, 2010, pp. 125–130.
- [38] DIAS, J. W.—ROSENBLUM, L. D.: Visual Influences on Interactive Speech Alignment. *Perception*, Vol. 40, 2011, No. 12, pp. 1457–1466, doi: 10.1068/p7071.
- [39] OBEN, B.—BRÔNE, G.: What You See Is What You Do: On the Relationship Between Gaze and Gesture in Multimodal Alignment. *Language and Cognition*, Vol. 7, 2015, No. 4, pp. 546–562, doi: 10.1017/langcog.2015.22.
- [40] BRÔNE, G.—OBEN, B.: InSight Interaction: A Multimodal and Multifocal Dialogue Corpus. *Language Resources and Evaluation*, Vol. 49, 2015, No. 1, pp. 195–214, doi: 10.1007/s10579-014-9283-2.
- [41] HOLLEMAN, G. A.—HOOGE, I. T. C.—HUIJDING, J.—DEKOVIĆ, M.—KEMNER, C.—HESSELS, R. S.: Gaze and Speech Behavior in Parent–Child Interactions: The Role of Conflict and Cooperation. *Current Psychology*, Vol. 42, 2023, No. 14, pp. 12129–12150, doi: 10.1007/s12144-021-02532-7.
- [42] WEISE, A.—LEVITAN, R.: Looking for Structure in Lexical and Acoustic-Prosodic Entrainment Behaviors. Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers), 2018, pp. 297–302, doi: 10.18653/v1/N18-2048.
- [43] REICHEL, U. D.: CoPaSul Manual – Contour-Based Parametric and Superpositional Intonation Stylization. *CoRR*, 2016, doi: 10.48550/arXiv.1612.04765.
- [44] OSTRAND, R.—CHODROFF, E.: It’s Alignment All the Way Down, But Not All the Way Up: Speakers Align on Some Features But Not Others Within a Dialogue. *Journal of Phonetics*, Vol. 88, 2021, Art.No. 101074, doi: 10.1016/j.wocn.2021.101074.
- [45] PATEL, S. P.—COLE, J.—LAU, J. C. Y.—FRAGNITO, G.—LOSH, M.: Verbal Entrainment in Autism Spectrum Disorder and First-Degree Relatives. *Scientific Reports*, Vol. 12, 2022, No. 1, Art.No. 11496, doi: 10.1038/s41598-022-12945-4.
- [46] NASS, C.—STEUER, J.—TAUBER, E. R.: Computers Are Social Actors. Proceedings of the SIGCHI Conference on Human Factors in Computing Systems (CHI ’94), ACM, 1994, pp. 72–78, doi: 10.1145/191666.191703.
- [47] HEYSELAAR, E.: The CASA Theory No Longer Applies to Desktop Computers. *Scientific Reports*, Vol. 13, 2023, No. 1, Art.No. 19693, doi: 10.1038/s41598-023-46527-9.
- [48] GONSIOR, B.—SOSNOWSKI, S.—MAYER, C.—BLUME, J.—RADIG, B.—WOLLHERR, D.—KÜHNLENZ, K.: Improving Aspects of Empathy and Subjective Performance for HRI Through Mirroring Facial Expressions. 2011 RO-MAN, IEEE,

- 2011, pp. 350–356, doi: 10.1109/ROMAN.2011.6005294.
- [49] COMINELLI, L.—FERI, F.—GAROFALO, R.—GIANNETTI, C.—MELÉNDEZ-JIMÉNEZ, M. A.—GRECO, A.—NARDELLI, M.—SCILINGO, E. P.—KIRCHKAMP, O.: Promises and Trust in Human–Robot Interaction. *Scientific Reports*, Vol. 11, 2021, No. 1, Art.No. 9687, doi: 10.1038/s41598-021-88622-9.
- [50] KRUYT, J.—DE JONG, D.—D'AUSILIO, A.—BEŇUŠ, Š.: Measuring Prosodic Entrainment in Conversation: A Review and Comparison of Different Methods. *Journal of Speech, Language, and Hearing Research*, Vol. 66, 2023, No. 11, pp. 4280–4314, doi: 10.1044/2023.JSLHR-23-00094.
- [51] AL MOUBAYED, S.—BESKOW, J.—SKANTZE, G.—GRANSTRÖM, B.: Furhat: A Back-Projected Human-Like Robot Head for Multiparty Human-Machine Interaction. In: Esposito, A., Esposito, A. M., Vinciarelli, A., Hoffmann, R., Müller, V. C. (Eds.): *Cognitive Behavioural Systems*. Springer, Berlin, Heidelberg, *Lecture Notes in Computer Science*, Vol. 7403, 2012, pp. 114–130, doi: 10.1007/978-3-642-34584-5_9.
- [52] MISHRA, C.—SKANTZE, G.: Knowing Where to Look: A Planning-Based Architecture to Automate the Gaze Behavior of Social Robots. 2022 31st IEEE International Conference on Robot and Human Interactive Communication (RO-MAN), 2022, pp. 1201–1208, doi: 10.1109/RO-MAN53752.2022.9900740.
- [53] LEMHÖFER, K.—BROERSMA, M.: Introducing LexTALE: A Quick and Valid Lexical Test for Advanced Learners of English. *Behavior Research Methods*, Vol. 44, 2012, No. 2, pp. 325–343, doi: 10.3758/s13428-011-0146-0.
- [54] THOMPSON, E. R.: Development and Validation of an International English Big-Five Mini-Markers. *Personality and Individual Differences*, Vol. 45, 2008, No. 6, pp. 542–548, doi: 10.1016/j.jpaid.2008.06.013.
- [55] DURAN, N. D.—PAXTON, A.—FUSAROLI, R.: ALIGN: Analyzing Linguistic Interactions with Generalizable techNiques – A Python Library. *Psychological Methods*, Vol. 24, 2019, No. 4, pp. 419–438, doi: 10.1037/met0000206.
- [56] QI, P.—ZHANG, Y.—ZHANG, Y.—BOLTON, J.—MANNING, C. D.: Stanza: A Python Natural Language Processing Toolkit for Many Human Languages. In: Celikyilmaz, A., Wen, T. H. (Eds.): *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations (ACL 2020)*. 2020, pp. 419–438, doi: 10.18653/v1/2020.acl-demos.14.
- [57] BOGHRATI, R.—HOOVER, J.—JOHNSON, K. M.—GARTEN, J.—DEGHANI, M.: Conversation Level Syntax Similarity Metric. *Behavior Research Methods*, Vol. 50, 2018, No. 3, pp. 1055–1073, doi: 10.3758/s13428-017-0926-2.
- [58] REIMERS, N.—GUREVYCH, I.: Sentence-BERT: Sentence Embeddings Using Siamese BERT-Networks. *CoRR*, 2019, doi: 10.48550/arXiv.1908.10084.
- [59] BOERSMA, P.—WEENINK, D.: Praat: Doing Phonetics by Computer (Version 5.1.13). 2009, <http://www.praat.org>.
- [60] PRATAP, V.—TJANDRA, A.—SHI, B.—TOMASELLO, P.—BABU, A. et al.: Scaling Speech Technology to 1,000+ Languages. *Journal of Machine Learning Research*, Vol. 25, 2024, Art.No. 97, <http://jmlr.org/papers/v25/23-1318.html>.
- [61] KUZNETSOVA, A.—BROCKHOFF, P. B.—CHRISTENSEN, R. H. B.: lmerTest Pack-

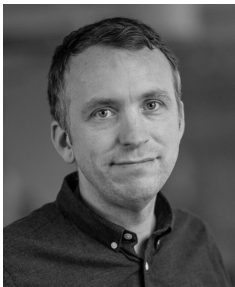
- age: Tests in Linear Mixed Effects Models. *Journal of Statistical Software*, Vol. 82, 2017, No. 13, pp. 1–26, doi: 10.18637/jss.v082.i13.
- [62] LENTH, R. V.—BUERKNER, P. et al.: Emmeans: Estimated Marginal Means, Aka Least-Squares Means. 2024, <https://cran.r-project.org/web/packages/emmeans/emmeans.pdf>.
- [63] REICHEL, U. D.—BEŇUŠ, Š.—MÁDY, K.: Entrainment Profiles: Comparison by Gender, Role, and Feature Set. *Speech Communication*, Vol. 100, 2018, pp. 46–57, doi: 10.1016/j.specom.2018.04.009.
- [64] HE, S.—ZHENG, X.—ZENG, D.—XU, B.—TIAN, G.—HAO, H.: Emotion Evolution Under Entrainment in Social Media. In: Huang, H., Liu, T., Zhang, H. P., Tang, J. (Eds.): *Social Media Processing*. Springer, Berlin, Heidelberg, Communications in Computer and Information Science, Vol. 489, 2014, pp. 155–163, doi: 10.1007/978-3-662-45558-6_14.
- [65] TEMPLETON, E. M.—CHANG, L. J.—REYNOLDS, E. A.—CONE LEBEAUMONT, M. D.—WHEATLEY, T.: Long Gaps Between Turns Are Awkward for Strangers But Not for Friends. *Philosophical Transactions of the Royal Society B: Biological Sciences*, Vol. 378, 2023, No. 1875, Art.No. 20210471, doi: 10.1098/rstb.2021.0471.
- [66] SCHWEITZER, A.—LEWANDOWSKI, N.: Convergence of Articulation Rate in Spontaneous Speech. *Proceedings of Interspeech 2013*, 2013, pp. 525–529, doi: 10.21437/Interspeech.2013-148.
- [67] HIRSCHBERG, J.—GRAVANO, A.—BEŇUŠ, Š.—WARD, G.—GERMAN, E. S.: Columbia Games Corpus LDC2021S02. Web Download. Linguistic Data Consortium, 2021, doi: 10.35111/ayn3-sp31.
- [68] DEVLIN, J.—CHANG, M. W.—LEE, K.—TOUTANOVA, K.: BERT: Pre-Training of Deep Bidirectional Transformers for Language Understanding. In: Burstein, J., Doran, C., Solorio, T. (Eds.): *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers) (NAACL 2019)*. 2019, pp. 4171–4186, doi: 10.18653/v1/N19-1423.
- [69] CER, D.—YANG, Y.—KONG, S.—HUA, N.—LIMTIACO, N.—ST. JOHN, R.—CONSTANT, NOAHAND GUAJARDO-CESPEDES, M.—YUAN, S.—TAR, C.—STROPE, B.—KURZWEIL, R.: Universal Sentence Encoder for English. In: Blanco, E., Lu, W. (Eds.): *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations (EMNLP 2018)*. Association for Computational Linguistics, 2018, pp. 169–174, doi: 10.18653/v1/D18-2029.
- [70] SHOR, J.—JANSEN, A.—MAOR, R.—LANG, O.—TUVAL, O.—CHAUMONT QUITRY, F.—TAGLIASACCHI, M.—SHAVITT, I.—EMANUEL, D.—HAVIV, Y.: Towards Learning a Universal Non-Semantic Representation of Speech. *Proceedings of Interspeech 2020*, 2020, pp. 140–144, doi: 10.21437/Interspeech.2020-1242.



Jay KEJRIWAL received his Bachelor's degree in computer science and engineering from the Maharishi Dayanand University, India, in 2012, his Master's degree in computational linguistics from the University of Tübingen in 2019, and his Ph.D. in computer science from the Slovak University of Technology in Bratislava, Slovakia in 2024. His primary interest is applying artificial intelligence to human-computer interaction to improve naturalness in conversation.



Chinmaya MISHRA joined the Multimodal Language Department at the Max Planck Institute for Psycholinguistics, Nijmegen as a Postdoctoral Researcher in 2024. He obtained his Ph.D. in human-robot interaction (HRI) from the Radboud University, Nijmegen in 2024 and has his M.Sc. in intelligent systems from the Technical University of Kaiserslautern, Germany (2020). In his research he uses an interdisciplinary approach to find novel solutions to model/automate multimodal robot behaviors (verbal and non-verbal) that are needed to facilitate a seamless HRI and study how these robot behaviors are perceived by humans.



Gabriel SKANTZE is Professor in speech communication and technology at the Department of Speech, Music and Hearing at KTH Royal Institute of Technology, Stockholm. He is also the co-founder and chief scientist at the Furhat Robotics, a Stockholm based company which develops the Social Robotics platforms: Furhat and Misty. He received his M.Sc. in cognitive science from the Linköping University, Linköping in 2000 and obtained his Ph.D. in speech communication from KTH, Stockholm in 2007. In his research, he investigates computational models of conversation that allow computers and robots to have

face-to-face conversations with humans. This includes both verbal and non-verbal aspects of communication, including turn-taking, alignment, joint-attention, and language grounding, and how they can be applied to human-robot interaction.



Tom OFFREDE obtained his Bachelor degree in English studies (linguistics, literature and teaching) at the University of Brasilia, Brazil, in 2017. He obtained his Research Master's degree in language and cognition at the University of Groningen, Netherlands, in 2020. He concluded his Ph.D. at the Humboldt-Universität zu Berlin, Germany, in 2024, where he worked on the intersection between phonetic adaptation and psychophysiology. His research interests are language production, bilingualism, and embodiment.



Štefan BEŇUŠ is Professor in the Department of English and American Studies, Constantine the Philosopher University in Nitra, Slovakia, and senior researcher at the Department of Speech Analysis and Synthesis, Institute of Informatics, Slovak Academy of Sciences, Bratislava. He obtained his Ph.D. in linguistics at the New York University in 2005 and did post-docs with the Spoken Language Processing Group at Columbia University and with the Institute for Phonetics and Speech Processing at Ludwig-Maximilians-Universität München. He has done experimental research in articulatory and acoustic phonetics and

recently focuses on better understanding of the alignment between interlocutors in terms of various aspects of speech prosody. He is also interested in how ideas from these areas might be applied in speech processing and foreign language acquisition.