



Universiteit  
Leiden  
The Netherlands

## Effect of prosody awareness training on the quality of consecutive interpreting between English and Farsi

Yenkimaleki, M.

### Citation

Yenkimaleki, M. (2017, June 7). *Effect of prosody awareness training on the quality of consecutive interpreting between English and Farsi*. LOT dissertation series. LOT, Utrecht. Retrieved from <https://hdl.handle.net/1887/49507>

Version: Not Applicable (or Unknown)

License: [Licence agreement concerning inclusion of doctoral thesis in the Institutional Repository of the University of Leiden](#)

Downloaded from: <https://hdl.handle.net/1887/49507>

**Note:** To cite this publication please use the final published version (if applicable).

Cover Page



Universiteit Leiden



The handle <http://hdl.handle.net/1887/49507> holds various files of this Leiden University dissertation

**Author:** Yenkimaleki, Mahmood

**Title:** Effect of prosody awareness training on the quality of consecutive interpreting between English and Farsi

**Issue Date:** 2017-06-07

## Chapter seven

# Objective correlates of the quality of interpreting

### Abstract

This study attempts to relate the intersubjective expert judgments to objective measures that can be expected to correlate with the judgments. If such correlates can be found, the expert judgment can be predicted by some combination of objective correlates. If the prediction is sufficiently accurate, expert judgments could be dispensed with in the future and be replaced by objective measurements. I have investigated the relationships between the expert judgments of the quality of the participants' interpreting performance on the one hand and objective correlates of their performance on the other. Somewhat surprisingly, the results show that the intersubjective ratings of the students' interpreting performance can be quite adequately predicted from objective measures for members of the control group through multiple linear regression analysis but that such predictions are less successful in the case of the experimental group. Crucially, the members of the control group were given more favorable ratings by the expert judges as their speed of delivery was faster. Such a linear relationship was absent in the case of the experimental group. For the latter group it seemed as though the relationship between speed of delivery of the interpretation was most favorably rated if it was in the middle of the range. Less favorable ratings were obtained not only for slow delivery (as in the case of the control group) but also for excessively fast delivery by some members of the experimental group.

**Keywords:** expert judgments, objective correlates, interpretation, pace, prosodic features

## 7.1 Introduction

In Chapter 3 the quality of consecutive interpreting by control and experimental groups was determined by an intersubjective evaluation procedure. Three experts rated the participants' performance as interpreters on a number of subjective scales, which together aimed to capture all relevant quality aspects of interpreting. The three experts, moreover, agreed strongly in their ratings. In fact, they were so much in agreement that it was decided to simply average the three scores given by the three raters. The analysis showed that the experimental group did better on each rating scale than the control group, and that the gain of the experimental group over the control group was specifically greater on those quality aspects that relate to prosody.

In the present chapter I will make an attempt to relate the intersubjective expert judgments to objective measures that can be expected to correlate with the judgments. If such correlations can be found, the expert judgment can be predicted by some combination of objective correlates. If the prediction is sufficiently accurate, expert judgments could be dispensed with in the future and be replaced by objective measurements.

Ten rating scales were used. For the sake of convenience, they are repeated here in Table 7.1, which is a copy of Table 3.1.

*Table 7.1. Ten evaluation criteria subdivided into three domains used in the quality judgment of interpreting performance. Weights add up to 100. After Sanyer (2004).*

Meaning		Language use		Presentation	
Accuracy	20	Grammar	7	Pace	10
Omissions	15	Expression	7	Voice	10
Overall coherence	10	Word choice	7		
		Terminology	7		
		Accent	7		

Unfortunately, in the set of evaluation criteria not all scales can be grounded in experimental measures. I will not try, for instance, to come up with objective measures that might predict overall coherence of the interpretation into Farsi relative to the original English text, nor will I attempt to define an objective measure for 'Expression'. However, omissions – i.e., failure to translate an important word or concept – can be counted, and the number of grammatical anomalies can be determined by analyzing transcripts of the interpretations. I will also try to establish correlates of at least some of the prosodic evaluation criteria such as accent and pace. Especially 'Pace' (or fluency) would seem to be amenable to objective testing. At least two correlates of pace will be considered, viz. speaking rate and articulation rate. Speaking rate is traditionally defined as the number of linguistic units, i.e., words or syllables, produced per unit time (per minute or per second). Here the total speaking time includes all pauses, whether silent or filled (*eh*, *ehm*). Articulation rate is computed the same way as is speaking rate but the

total time does not include pauses and hesitations. Defined this way, obviously, speaking rate and articulation rate are strongly correlated. When trying to predict judgments on a rating scale from objective measures it is better to work with independent predictors, i.e., predictors that do not or only weakly correlate with each other. It seemed to us that a feasible way to disentangle speaking rate and articulation rate would be to use articulation rate only and supplement this parameter with a more direct measure of the incidence of pauses and hesitations. This latter aspect can be adequately captured by computing the percentage of the total speaking time that is taken up by pauses. I will call this latter parameter ‘%-pause’.

I note, in passing, that it will not necessarily be the case that pace (fluency) is monotonically related to either %-pause or to articulation rate. It would seem more likely that the relationship between the judgment and the acoustic measure will be U-shaped, i.e., judgments may well be most favorable for values in the middle of the range, when the speaker does not insert a great many pauses (indicative of difficulties in producing the interpretation) nor speaks with very few pauses (which would create a burden on the listener). Similarly, articulation rates in the middle of the range are expected to receive the most favorable judgments.

## 7.2 Objective measures used

In the next few sections I will outline the procedures followed to quantify the objective correlates I used. Here I will distinguish between counts of phenomena that can be established by analyzing written transcripts of the interpreter’s performance (and comparing it to the original text), and measurement of acoustic properties, which, of course, cannot be done from a written transcript.

### 7.2.1 Count measures

Generally, the norm is that interpreters should have a complete transfer of the source text to their audience without any omission of ideas or changes of meaning. This issue has received a lot of attention in typology and error analysis of translation and interpreting performance. But we know that in some cases omission of some aspects in interpretation enhances the quality of interpreting and as a result communication of message is done properly (‘less is more’). Jones (2014: 139) pointed out that interpreters in some situations are not in position to render the exact and complete message. So, in these situations the interpreters omit part of the source text in order to have a coherent message for the audience. Therefore in some cases, the interpreters intentionally omit part of the source language because they want to transfer the gist of the message so that the audience may perceive the message more easily. So, when it happens, the communication of the message between the audience and interpreters can be achieved comprehensively. We should know that in interpreting the important aspects and essentials are preferred over the completeness of the message.

It is an open question, in the present study, whether the judgment of accuracy and omissions is monotonically related to the number of words (or concepts) incorrectly

translated or left out altogether. One hypothesis would be that the more accurate and complete the interpreting is, the better the accuracy and omission judgments. I leave room, however, for a more sophisticated possibility; viz., the relationship between the objective counts and the global judgments is U-shaped. In the latter case, keeping in all details would detract from the judged adequacy or optimality of the interpreting job.

The number of omissions was established by comparing an optimal translation of the original English texts into Farsi with transcripts of the student's interpretation. The unit of measurement was the content word. I checked for every content word in the model translation whether it occurred in some adequate or at least acceptable form (identical, synonym or paraphrase) in the student's transcript. When the word or concept was not an acceptable stand-in for the original, it was counted as an inaccuracy or meaning error. When the word or concept was absent from the student's interpretation altogether it was scored as an omission. The total number of errors was then equal to the number of inaccuracies and the number of omissions added together.

### 7.2.2 Acoustic measures of pace

The sound recordings of each of the 30 speakers were segmented into interpausal units. An interpausal unit, or IPU, is defined as a stretch of speech not interrupted by a silent or filled pause (Koiso et al. 1998, Buhmann et al. 2002).<sup>1</sup> In order to qualify as a pause, a silence in the spoken utterance must be longer than 100 ms. If shorter silences would also be considered, the occlusion phases of voiceless plosives would be counted as pauses, which would be undesirable.

The recordings were recoded from mp3 format to .wav-format. Normally, lossy coding such as mp3 would be ill-advised for the analysis of speech but in the present case, where only duration, fundamental frequency and intensity will be measured, measurements will be quite faithful. The segmentation of the recordings was done semi-automatically with PRAAT speech processing software (Boersma & Weenink 1996, 2017). As a first approximation, the recordings for a given speaker were automatically split up into stretches of uninterrupted speech and pauses using the annotation module with automatic speech/silence detection. For male speakers the bottom pitch was set at 70 Hz, for females at 120 Hz. For all other parameters the default setting was used (both speech and silences should exceed 100 ms, silence threshold at -25 dB). The resulting annotation grids plus waveforms were inspected by ear and eye. The procedures laid down by Buhmann et al. (2002) were followed. Filled pauses, which are not detected as such by the algorithm, were set by hand, and misplaced segmentation boundaries were corrected when necessary. Each speaker produced three fragments. Time intervals preceding and following fragments were discarded. Only pauses within each of the three fragments were included in the computations. Filled pauses were separately labelled. A filled pause, by definition, is not coarticulated with whatever precedes it. As a result, a filled pause is always preceded by a short stretch of silence. It occurred regularly that a speaker fell silent for several hundreds of milliseconds, then produced an *eh* or *ehm*

---

<sup>1</sup> IPU's are sometimes also referred to as 'fluent runs' (e.g., De Jong & Perfetti 2011).

filled pause, which could or could not be followed fluently by the onset of the next fragment. In such cases two or even three pauses were distinguished, one of which was filled and the others were silent. As a result of this procedure the number of pauses found could be greater than the number of IPUs. In a number of cases the speaker lengthened a word-final vowel, which was clearly indicative of a hesitation. In such cases we did not mark a pause; lengthened vowels lead to slower speaking rates. The occurrences of such lengthened vowels were also marked and counted.

The transcripts of the students' interpretations were automatically converted from the Arabic script to a Western transliteration. This transliteration is close enough to a broad phonemic transcription of what was said to enable correct syllabification. Word boundaries were checked and corrected by hand. A list of word types was extracted from the transcripts. In each word in the list, syllable boundaries were inserted by hand. Syllable boundaries were then inserted automatically in the materials by applying a series of find-and-replace commands using the words and their hyphenation in the list of types. The number of syllables as well as the number of words was then counted automatically for each IPU and stored in the database.

For each speaker the following speech rate-related measures were computed from the duration data and the syllable and word counts:

- Total articulation time: i.e., the duration of all the IPUs added together
- Total pause time: the duration of all the intervals, whether silent or filled, added together
- Total filled pause time: the duration of all filled pauses (*eh*, *ehm*, *mm*, *mmm*) added together
- Number of IPUs
- Number of silent pauses
- Number of filled pauses
- Standard deviation of IPU duration
- Standard deviation of pause duration
- Speaking rate in words/s: (total articulation time + total pause time) / number of words
- Speaking rate in syll/s: (total articulation time + total pause time) / number of syllables
- Articulation rate in words/s: total articulation time / number of words
- Articulation rate in syll/s: total articulation time / number of syllables
- %pause: total pause duration / (total articulation time + total pause duration).

### 7.3 Results

#### 7.3.1 Count measures

The number of inaccuracies and omissions were counted by comparing each individual student's written transcript with the ideal, model interpretation. Note that the model interpretation (see appendix 5 in chapter 3) contained a rendition of all words and concepts that occurred in the English source text.

Table 7.2 presents the number of inaccurately translated words as well as the number of omissions, for the members of the control group and the experimental group separately. Moreover, the individuals in the two groups were matched pairwise on the basis of their performance on the overall TOEFL score obtained in the pre-test (the same matching criterion I used in chapter 3).

*Table 7.2. Number of incorrectly translated words, omitted words and total number of errors for individual subjects in control and experimental groups. Subjects are matched on TOEFL score in pre-test, and listed from best to poorest.*

Subject	Control group			Subject	Experimental group		
	Wrong words	Omissions	Total		Wrong words	Omissions	Total
1	15	15	30	1	5	13	18
2	12	18	30	2	6	30	36
3	20	10	30	3	22	8	30
4	24	10	34	4	15	8	23
5	20	12	32	5	19	8	27
6	34	6	40	6	22	8	30
7	13	19	32	7	10	17	27
8	21	29	50	8	26	7	33
9	32	15	47	9	22	10	32
10	22	30	52	10	25	15	40
11	23	56	79	11	19	12	31
12	26	41	67	12	18	22	40
13	37	43	80	13	32	32	64
14	31	50	81	14	23	27	50
15	24	110	134	15	39	40	79
Mean	23.6	30.9	54.5		20.2	17.1	37.3

Interestingly, although the two groups did not differ from one another on the pre-test, there is a substantial difference in the number of word errors counted in the transcripts of the subjects' interpreting performance in the post-test. The number of word errors is significantly smaller for the experimental group for both wrong words 24 versus 20) and for omissions (31 versus 17),  $t(14) = 1.8$  ( $p = .045$ , one-tailed) and  $t(14) = 2.6$  ( $p =$



.016, one-tailed), respectively. The effect is most clearly seen in the total number of errors (55 versus 37 errors),  $t(14) = 4.0$  ( $p < .001$ ).

The crucial question is if the experts' global judgment of the accuracy of the interpretations can be understood from the objective post-hoc error counts. To answer this question, I computed the correlation coefficients between the objective counts and the expert judgments. The correlations are shown in Table 6.3.

Table 7.3. Correlation matrix for objective error counts and global expert judgments. The lower triangle contains Pearson's  $r$ , the upper triangle shows the non-parametric Spearman's  $\rho$ .

	Objective error counts			Judgments	
	wrong words	word omissions	total errors	accuracy	omissions
wrong words		.218	.732	-.696	-.718
omissions	.255		.755	-.712	-.698
total errors	.555	.946		-.921	-.904
Accuracy	-.691	-.704	-.838		.973
omissions	-.694	-.718	-.851	.976	

Note:  $r > .555$ :  $p < .001$

It can be observed, first of all, that the global accuracy and omission judgments are extremely strongly correlated ( $r = .976$ ). This means that the judges did not differentiate between these two aspects of the interpreting performance. This seems understandable, given that leaving out words or concepts that occurred in the English source text from the interpretation is to all intents and purposes a form of inaccuracy. In the objective post-hoc error counts the numbers of incorrectly translated words and omitted words are not significantly correlated, so that these two types of error might have contributed to the global judgments separately and independently. Note that the number of omission errors was much larger than the number of inaccurately translated words, which explains the much higher correlation between the former ( $r = .946$ ) and the total error score than the latter ( $r = .555$ ). Observe, finally, that the non-parametric  $\rho$  coefficients tend to be somewhat better than their parametric counterparts  $r$ . This suggests that the relationships between the objective error counts and the global judgments are non-linear. I will come back to this issue presently.

The most important conclusion to be drawn from Table 7.3 is that the global accuracy and omission judgments (which are virtually identical anyway) can be predicted with great precision by the objective error counts, especially when the total number of errors is used as the predictor, with  $\rho$ -values in excess of .9. Clearly, then, the experts' global judgments have a high concurrent validity in that they lead to the same ranking of students as can be achieved on the basis of laborious error counts.

To conclude this part of the analysis, Figure 7.1 plots the mean of the global accuracy and omission judgments as a function of the total error number for each of the 30 students. The x-axis of the plot, however, is not linear but logarithmic. A preliminary

check revealed that the percentage of the judgment scores accounted for by a logarithmic model was appreciably better than by a linear model, with  $R^2 = .828$  and  $.720$ , respectively.

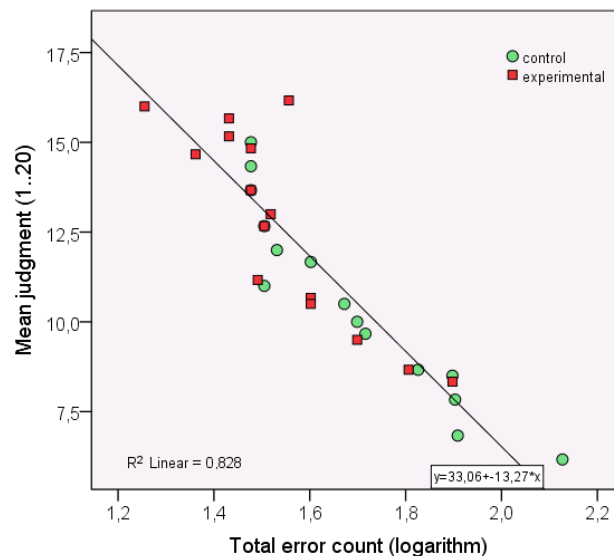


Figure 7.1. Mean of global word accuracy and word omission judgments as a function of the logarithm of the total count of errors.

It can be seen in Figure 7.1 that accuracy judgments for the experimental group are better than those for the control group, which reiterates what was observed in chapter 5. We now understand that the difference between the two groups is related in a perfectly straightforward manner to the difference in numbers of incorrectly translated words and words omitted during interpretation from English into Farsi. Moreover, the relationship works the same way for both groups of student interpreters. What we do not know is how this difference in performance can be explained. The experimental group received ample explanation of prosodic differences between English and Farsi, and did practical exercises emphasizing these prosodic differences, but this in and by itself does not explain why the accuracy of the translation of the contents should improve.

### 7.3.2 Acoustic measures

A total of 15 speech rate related parameters were computed (see above). Some of these were measured from the acoustic signal, other were counted in written transcripts of the interpreting performance by the participants. Compound measures were derived by computing ratios or percentages based on raw measurements. For instance, Articulation rate was defined as the Total articulation time divided by the Total number of syllables

counted. Table 7.4 presents the summary statistics for these 15 parameters, for the experimental and control groups separately. Independent t-tests indicate that the small differences between the two groups never reach statistical significance for any of the 15 parameters, with p-values ranging between .187 and .950.

*Table 7.4. Mean and standard deviation of 15 fluency-related acoustical correlates for control group and experimental group. The difference between the two means ( $\Delta$ ) and the t and p-values are given (df = 28 for each parameter).*

Parameters	Control group		Exper. Group		All		$\Delta$	t	P
	Mean	SD	Mean	SD	Mean	SD			
Total articulation time (s)	71.2	14.2	73.6	10.7	72.4	12.4	-2.4	-.5	.606
Total pause time (silent + filled)	21.9	10.5	21.7	9.0	21.8	9.6	.2	.1	.950
Total N words	220.3	33.8	231.5	17.3	225.9	27.0	-11.1	-1.1	.266
Total N syllables	444.5	72.9	470.9	36.6	457.7	58.3	-26.4	-1.3	.221
Percent pause (silent + filled)	22.9	6.0	22.2	5.9	22.6	5.9	.7	.3	.742
Speech rate (words/s)	2.4	.5	2.5	.4	2.5	.4	-.1	-.3	.756
Speech rate (syllables/s)	4.9	1.0	5.1	.8	5.0	.9	-.2	-.4	.657
Articulation rate (words/s)	3.1	.5	3.2	.4	3.2	.4	.0	-.3	.781
Articulation rate (syllables/s)	6.3	1.0	6.5	.8	6.4	.9	-.2	-.5	.640
SD IPU duration (s)	1.2	.3	1.6	1.3	1.4	1.0	-.5	-1.4	.187
SD pause duration (s)	.9	.4	1.0	.7	.9	.5	-.1	-.5	.644
SD N words in IPU	4.0	1.1	5.0	3.6	4.5	2.7	-1.0	-1.0	.311
SD N syllables in IPU	8.0	2.5	10.3	7.7	9.1	5.8	-2.4	-1.1	.270
N IPUs	33.9	11.4	33.1	11.4	33.5	11.2	.7	.2	.862
N pauses (silent + filled)	34.1	15.9	31.7	15.8	32.9	15.6	2.3	.4	.690

Figure 7.2 presents the relationship between percentage of pause and the judged pace of delivery, shown separately for the experimental and control groups.

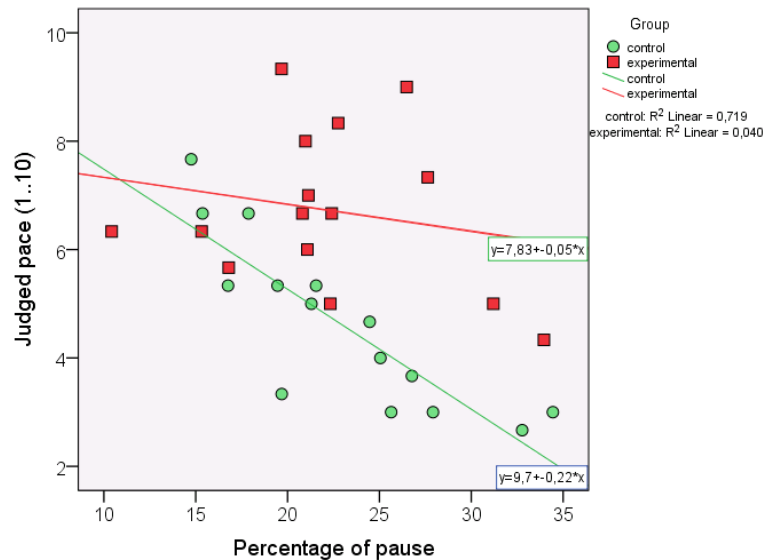


Figure 7.2. Judged pace as a function of percentage of pause (silent and filled) in spoken text, shown separately for members of experimental and control groups.

The figure reiterates that the judged pace is better for the experimental group than for the control group, as was seen earlier in Chapter 3. However, the figure also shows, quite clearly, that the relationship between percentage of pause and the judged pace of delivery is strong and linear, as far as the control group is concerned. The greater the percentage of pausing, the poorer the judged pace, where the objective measure accounts for 72% of the variance the judged pace score,  $R^2 = .719$ . The relationship is much weaker, in fact almost absent, for the experimental group. It is not the case that the experimental group has no variability in percentage of pause: the distribution of this objective parameter is roughly the same for experimental and control group alike, with a spread between 10 and 35%. In order to shed light on this curious asymmetry, let us now examine the relationship between articulation rate (words/s) and judged pace. The expectation, of course, is that a faster articulation rate should correlate with better pace judgments. The results are in Figure 7.3.

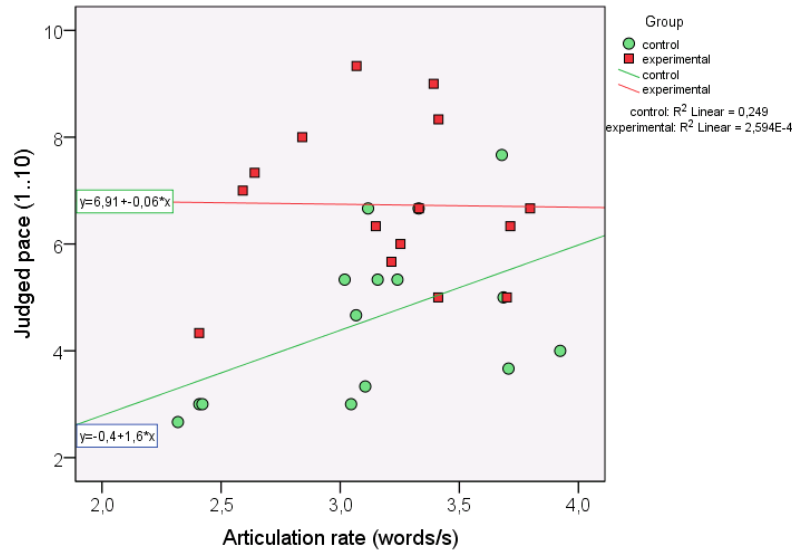


Figure 7.3. Judged pace as a function of percentage of articulation rate (words/s) in spoken text, shown separately for members of experimental and control groups.

Again, it can be observed that the relationship between the objective measure and judged pace is fairly strong for the control group,  $R^2 = .249$ , and explains a quarter of the variance in the judgments. It can also be noticed that there was no correlation at all for the experimental group.

In order to understand the asymmetry in the results of the experimental and control groups, at first it is needed to examine the relationship between the predictor variables used here, viz. percentage pause and articulation rate (in words/s and in syllables/s). It turns out that it is rather immaterial whether articulation rate is expressed in words/s or in syllables/s. The intercorrelation between these two measures is almost perfect at  $r = .991$  for the control group and  $r = .994$  for the experimental group ( $N = 15$ ,  $p < .001$  in both cases). The intercorrelation between articulation rate and %-pause shows the same remarkable discrepancy between the two groups we met before, such that the correlation is relatively strong and significant for the control group,  $r = -.564$  ( $p = .014$ ) and  $-.618$  ( $p = .007$ ) for words/s and sylls/s, respectively, but weaker and insignificant for the experimental group,  $r = -.340$  ( $p = .107$ , one-tailed) and  $-.367$  ( $p = .089$ , one-tailed) for words and syllables per second, respectively (see also Figure 7.3).

These results suggest that articulation rate and %-pause in the control group are both indices of cognitive difficulty in task performance: when these participants find it difficult to interpret the incoming message, they tend to speak more slowly, leading at the same time to fewer syllables (or words) per second and to more and/or longer pauses. These would be pauses for the sake of the speaker rather than for the sake of the listener. The speaker needs more time to find appropriate words and formulations to

get the message across. The speaker does not insert pauses to help the listener by clearly marking off processing units (be they clauses or constituents). A reasonable prediction here would be that these speakers also produce relatively many filled pauses, which are the hallmark of problems with finding words or formulations. In contradistinction to this we would expect pauses in the experimental group to be used as structure markers for the sake of the listener. These would be relatively short pauses, which are planned by the speaker to guide the listener. Additionally, fewer filled pauses and other over markers of planning difficulty should be observed for the experimental group.

These hypotheses can be tested by examining the number of disfluency markers, which is what it is done in the following section.

### 7.3.3 Disfluencies

In order to understand the discrepancy between the experimental and control group, let us now consider the number of disfluencies marked for each. I distinguish the following four categories:

1. Long silent pause, indicative of extra planning time needed. Assuming that pauses between IPUs in fluent speech normally do not exceed a duration of 1000 ms, any silent pause longer than 1000 ms was considered a (potential) disfluency.
2. Filled pause, indicative of hesitation, i.e., any instance of *eh*, *ehm* or *mm* that is not fluently coarticulated with whatever precedes it.
3. Lengthened vowel, i.e., a word-final vowel that is lengthened and is indicative of hesitation
4. Repetition, i.e., the repetition of something that was said in the immediately preceding IPU, then broken off, and repeated in a second attempt. In a number of cases there was no break (no silent or filled pause after the false start); the repetition followed seamlessly after the false start. I decided to count the repetitions only (and only if the repetition was not an instance of stuttering – which happened on two occasions).

Table 7.5 lists the disfluencies found, for the experimental and control groups, together with the number of regular IPUs and short silent pauses. The latter two categories are indices of fluent speech, whereas the other four categories point to planning difficulties on the part of the speaker.

Table 7.5. Mean duration (in seconds) and number of IPUs, regular silent pauses, long pauses and filled pauses produced by experimental and control groups.

Disfluency	Control group		Experimental group		$\Delta$ (exp - cont)	
	Duration	N	duration	N	duration	N
regular IPU	2.14886	456	2.22998	485	0.08112	29
silent pause short	.41173	394	.43771	392	0.02598	-2
silent pause long	3.22033	43	3.21933	43	-0.00100	0
filled pause	.36330	73	.38007	42	0.01677	-31
lengthened vowel	1.44213	30		0		-30
Repetition	2.01930	23	1.97880	11	-0.04050	-12

Table 7.5 shows no systematic differences in the behavior of the experimental and control participants in terms of the duration and number of disfluencies, with three notable exceptions. The number and mean duration of regular IPUs, as well as those of both short and long silent pauses (the latter would be indicative of planning problems on the part of the speaker) are virtually identical between the two groups. This also goes for the duration of the remaining categories of disfluencies but, remarkably, the number of disfluencies in the latter three categories differs between the groups such that the control group shows many more disfluencies in the categories filled pause, excessive prepausal vowel length and repetitions after a false start. These three categories, obviously, are indicative of planning problems. Before drawing any conclusions from these observations let us first see how the numbers are distributed over the 15 participants in each group.

Table 7.6 presents the numbers of disfluencies in the categories filled pause, excessive prepausal vowel lengthening and IPUs that repeat materials after a false start, broken down by the two groups of participants. In order to make the comparison maximally sensitive, the participants in the two groups have been matched for their TOEFL test scores.

Inspection of Table 7.6 reveals, first of all, that the TOEFL pre-test predicts the number of disfluencies observed in the interpretation tasks rather well. The correlations are negative, of course, since high TOEFL scores (indicating good proficiency in English) should lead to better performance, with fewer hesitations in the interpreting task. The best fit was obtained when the TOEFL scores were used to predict the logarithm of the number of disfluencies. Quite a few participants fulfilled their interpreting task without any disfluency. Since the logarithm of 0 is undefined, I remedied this by incrementing the overall disfluency count for each of the 30 participants by 1. I then find the same asymmetry in the predictability that we met before. The interpreting performance of the control can be predicted from objective measures much better than the scores of the experimental group. The correlation coefficients are  $r = -0.742$  ( $N = 15$ ,  $p = .001$ , one-tailed) for the control group and  $r = -0.440$  ( $N = 15$ ,  $p = 0.050$ , one-tailed). Across all participants  $r = -0.612$  ( $N = 30$ ,  $p < .001$ , one-tailed).

Table 7.6. Number of over disfluencies in three categories (excessive pre-pausal vowel lengthening, filled pause, repetition of words after a false start) for participants in control and experimental groups. Participants are rank ordered within their group on the basis of their pre-test overall TOEFL score.

Nr.	Student	Gender	Lengthen	Filled pause	Repeat	Total	TOEFL
Control Group							
1.	SaM	Female	0	1	0	1	610.0
2.	PoP	Female	0	0	0	0	586.7
3.	DaD	Male	0	0	1	1	563.3
4.	EiM	Male	0	0	2	2	553.3
5.	KhR	Female	0	0	2	2	540.0
6.	ZaN	Female	0	0	0	0	530.0
7.	EiM	Female	0	0	0	0	513.3
8.	AtH	Female	0	0	0	0	510.0
9.	ReR	Male	0	1	9	10	506.7
10.	AlA	Male	0	0	1	1	503.3
11.	MaN	Male	0	4	0	4	500.0
12.	LeK	Female	0	10	2	12	490.0
13.	JaR	Male	0	0	5	5	473.3
14.	AsH	Male	20	35	1	56	446.7
15.	NeF	Female	10	22	0	32	446.7
Total			30	73	23	126	
Experimental group							
1.	AlR	Male	0	0	1	1	613.3
2.	MaH	Female	0	1	1	2	603.3
3.	RaM	Male	0	0	0	0	566.7
4.	MoH	Male	0	0	0	0	563.3
5.	NaN	Female	0	0	1	1	553.3
6.	SaK	Female	0	4	1	5	553.3
7.	ArA	Male	0	3	0	3	550.0
8.	ZoM	Female	0	0	0	0	550.0
9.	PaN	Female	0	2	0	2	546.7
10.	BaN	Male	0	1	0	1	523.3
11.	KiK	Female	0	0	0	0	516.7
12.	MaR	Female	0	0	0	0	493.3
13.	NaH	Male	0	0	4	4	480.0
14.	HaM	Male	0	0	2	2	476.7
15.	TaB	Female	0	31	1	32	446.7
Total			0	42	11	53	

Although the total number of overt disfluencies in the performance of the control group (126) is more than twice as large as for the experimental group (53), the difference falls short of significance. A sign test on the counts (11 pairs matched on within-group TOEFL rank, excluding 4 tied scores) yields  $p = .114$  (one-tailed), which is a (weak) trend at best.



I may also normalize the number of overt disfluencies by speech time. After all, when a speaker produces more speech materials (words, syllables) during a longer stretch of time, there is more opportunity to produce errors and disfluencies. I therefore divided the total number of disfluencies per speaker by the duration of all his/her IPUs added together.

To conclude this part of the analysis, I will now try to establish a possible relationship between the incidence of overt disfluency markers and the pace of the interpreting performance as judged by the expert raters. The correlation between the raw number of disfluencies and judged pace is slightly poorer than when the logarithm of the number of disfluencies used, but even then  $r$  is rather weak at  $-.526$  ( $N = 30$ ,  $p < .001$ ). Moreover, similar correlation coefficients are obtained between the disfluency counts and all other judged aspects of the interpreting performance (which tend to be strongly correlated, see Table 3.5). When I compute the correlations separately for experimental and control groups, we observe the same asymmetry as before: correlations are appreciably better for the control group than for the experimental group, not just for pace but for all judged aspects.

#### 7.3.4 Predicting pace from multiple correlates of fluency

In the preceding sections we have seen that the prosodic parameter with the most tangible measureable correlates, i.e., pace of delivery, correlates with a large number of variables. These variables can be located in the acoustical domain, e.g., articulation rate (syllables per second) and percent pause. However, pace also correlated with the number of disfluencies per unit time as counted in the transcripts of the interpreting performances obtained from the participants. Interestingly, the intercorrelations between the disfluency counts and the acoustic correlates of pace were relatively modest, so that there is reason to try to predict judged pace from acoustic and count parameters together. Table 7.7 presents the correlation matrix for judged pace (dependent) and the acoustic and count parameters of (dis)fluency. Only the non-redundant lower triangle of the matrix is shown.

Table 7.7. Correlation matrix of judged pace (dependent) and five predictors: Percent pause, Articulation rate (syllables/second), Standard deviation of interpausal units (ms), Standard deviation of (filled and silent) pauses (ms) and the Relative number of disfluencies per unit time.  $N = 30$  for each cell.

	Pace	Perc pause	Art rate	SD speech	SD pause
Percent pause	-.469**				
Artic. rate (syll/s)	.314*	-.504**			
SD speech	.503**	-.181	.121		
SD pause	.075	.503**	.112	.561**	
Rel. disfluencies	-.543**	.646**	-.623**	-.267	-.035

\*  $r > .300$ :  $p < .05$ , \*\*  $r > .460$ :  $p < .01$  (one-tailed).

This table summarizes the information presented earlier in § 7.3.2 with one exception: we now see that the variability in the duration of the inter-pausal units (or fluent runs) is, in fact, fairly good predictor of judged pace, better, for instance, than articulation rate or percent pause, though still weaker than the relative number of disfluencies. This is somewhat unexpected, especially since the correlation is positive. One would expect competent speakers to divide their delivery into chunks of roughly equal size, which should yield a negative correlation with judged pace: the smaller the variability in the chunk size, the better the fluency. Variability in the pause duration, however, does not correlate with judged pace.

Table 7.8 A-B contains the same correlation matrix as in Table 7.7 but now the data are presented separately for the experimental and control groups.

Table 7.8 A-B. Correlation matrix of judged pace (dependent) and five predictors. Further see Table 6.7,  $N = 15$  per cell.

		Pace	Perc pause	Art rate	SD speech	SD pause
A. Control	Percent pause	-.848**				
	Artic. rate (syll/s)	.592**	-.618**			
	SD speech	.741**	-.773**	.337		
	SD pause	-.454*	.593**	.073	-.397	
	Rel. disfluencies	-.583*	.727**	-.727**	-.615**	-.007
B. Experimental	Percent pause	-.200				
	Artic. rate (syll/s)	-.009	-.367			
	SD speech	.487*	-.056	.064		
	SD pause	.338	.487*	.137	.732**	
	Rel. disfluencies	-.470*	.601**	-.434	-.210	-.034

\*  $r > .450$ ;  $p < .05$ , \*\*  $r > .590$ ;  $p < .01$  (one-tailed).

Breaking the correlations down separately for the experimental and control groups shows the by now familiar result that the correlations are clearly stronger for the control group than for the experimental group. There is, however, one parameter that behaves differently between the two groups. The variability in duration of the (filled and silent) pauses correlates negatively with judged pace in the control group ( $r = -.454$ ,  $p = .045$ , one-tailed) but positively in the experimental group ( $r = -.338$ ,  $p = .109$ , ins.). Variability in pause duration in the control group is typically caused by long silences and hesitations when the student interpreter is stuck for words. The better participants in this group have fewer of these long pauses, so that the variability in their pause durations is reduced. The experimental group, however, has fewer long pauses and disfluencies as a general characteristic; their pause variability is probably conditioned by the grammatical structure of their utterances such that light prosodic boundaries (at the phrase and clause level) have relatively short pauses and deeper boundaries (at the sentence level) are marked by longer pauses – as is typically found in other languages such as English (e.g., Grosjean, Grosjean & Lane 1979, Selkirk 1984)

and Dutch (e.g., Swerts 1997). In that case, the more variable the pause duration, the more competently does the speaker use prosodic markers. Note also that for the experimental group longer pauses tend to go together with longer IPUs, whereas the correlation is reversed for the control group.

Multiple regression analyses were then conducted for the two groups combined ( $N = 30$ ) and for the experimental and control groups separately. All five predictors mentioned in the correlation matrix were entered simultaneously in one analysis and in step-wise mode in another. The results are shown in Table 7.9 A-B-C.

Table 7.9. Summary of multiple regression analysis with judged pace as the dependent from five predictors: Percent pause, Articulation rate (syllables/second), Standard deviation of interpausal units (ms), Standard deviation of (filled and silent) pauses (ms) and the Relative number of disfluencies per unit time. Analysis were run with simultaneous entry (left part of table) and in stepwise mode (right part of table) for all participants combined (panel A,  $N = 30$ ) as well as for the control (panel B,  $N = 15$ ) and experimental groups (panel C,  $N = 15$ ) separately.

Simultaneous entry						Stepwise mode						
Predictors	Beta	R <sup>2</sup>	F	df	P	Predictors	Beta	R <sup>2</sup>	$\Delta R^2$	F	df	P
A. Combined groups												
Rel.disfl.	-.443					Rel. disfl.	-.440	.295	.295	11.7	1, 28	
SD speech	.617					SD speech	.385	.433	.138	6.6	1, 27	.016
SD pause	-.380											
% Pause	.166											
Artic. rate	.089	.473	4.3	5, 24	.006							
B. Experimental group												
SD pause	.642											
% Pause	-.444											
Artic. rate	-.419											
Rel.disfl.	-.376											
SD speech	-.059	.457	1.5	5, 9	.227							
C. Control group												
% Pause	-.751					% Pause	-.221	.719	.719	33.3	1, 13	<.001
Rel.disfl.	.376											
SD speech	.342											
Artic. rate	.277											
SD pause	.110	.785	6.6	5, 9	.008							

For the total group of participants combined we find an  $R^2$  of .473 when all five predictors are entered simultaneously. In the stepwise mode it turns out that only two predictors make a sufficient contribution to be included in an optimal model, which then accounts for 43.3 percent of the variance.

When the analysis is performed for the experimental group separately, no model is produced that is better than chance. As can be seen in Table 7.8 B two predictors correlate significantly (but only just) with the criterion when studied as single predictors, viz. SD

speech and Relative number of disfluencies but they lose significance in the simultaneous entry of five predictors because of the increased degrees of freedom.

Judged pace can be best predicted for the control group. Entering all five predictors simultaneously yields an  $R^2$  of .785, i.e., the model accounts for 79 percent of the variance. However, as was shown earlier, one single predictor, i.e., percent pause duration, accounts for 72 percent of the variance; none of the remaining four predictors makes a further contribution that reached significance.

#### 7.4 Conclusions and discussion

In this chapter I have examined the relationships between the expert judgments of the quality of the participants' interpreting performance on the one hand and objective correlates of their performance on the other. In the quality judgments a rating instrument was used that was comprised of ten scales. Seven of these pertain to aspects of quality that can be (and actually were) established by examining written transcripts of the interpreting tasks. These aspects relate to abstract linguistic properties of the interpretations, such as the accuracy with which the ideas in the source text were expressed, number of words omitted, appropriateness of choice of words and terminology, number of grammatical errors, and overall coherence of the interpretation. The remaining three scales were meant to capture the phonetic aspects of the delivery of the interpretation, i.e., the degree of accentedness, the pace (or fluency) of the delivery and the pleasantness of the voice. These three phonetic aspects all relate to relatively long-term aspects of speech, i.e., are not properties of specific vowels or consonants (see the definition of prosody in § 3.1), and are therefore prosodic features.

It was shown in Chapter 3 that the seven textual/linguistic scales intercorrelate very strongly, as do the three prosodic scales, but the correlations between scales in different categories are relatively low. The possibility to divide the ten scales into one group of seven non-prosodic and three prosodic scales was borne out by a factor analysis, which showed opposite factor loadings by the two groups of scales on the two principal components extracted in the analysis.

The results presented in this chapter bear out, quite clearly, that the expert judgments of the non-phonetic aspects can be related in a rather straightforward manner to a number of structural properties that could be quantified or counted in written transcripts of the interpreters' deliveries. Since the seven rating scales are very strongly intercorrelated there is little point in trying to predict each of these scales separately from objective counts. It would be sufficient, therefore, to summarize the most striking correlations found.

It turned out, then, the total number of errors in the interpreted passages (i.e., wrong words and number of omitted words added together) afford excellent prediction of the accuracy (and omissions) rating, with correlation coefficients in excess of .900. The actual numbers of wrong words and omissions were quite disparate, however. The conclusion follows, therefore, that the expert judges were not able to differentiate between these two aspects of accuracy even though they were clearly different in the interpreters' productions. This conclusion does not reflect negatively on the quality of

the raters – it just shows that these two closely related aspects are extremely difficult to distinguish when asked to give an on-the-spot evaluation of an interpreter's performance. Proper differentiation between the two types of inaccuracy in interpreting can only be achieved when a written transcript is available for a detailed and more time-consuming analysis.

These lexical accuracy parameters (words incorrectly translated or omitted altogether) are the two most important aspects of the overall rating of the students' interpreting performance. Incorrect words were weighted by 20, omissions by 15, so that together they represent 35 percent of the overall score. The other eight aspects together, with weights of either 7 or 10, represent the remaining 65 percent.

It should be noted in this context that the objective measures that predict the judged accuracy of the interpretation performance so well, are also the quantitative measures that optimally differentiate between the participants in the experimental and the control groups. The experimental group produced a very significantly smaller number of (lexical) inaccuracies than the control group (with a mean of 55 versus 37 lexical inaccuracies per speaker). It remains unclear at this stage why the experimental group would produce fewer inaccuracies than the control group. Why would a 36-hour training module emphasizing prosody and prosodic differences between English and Farsi, which is what differentiates the experimental groups from the control group, lead to a reduction in number of lexical errors?

The total number of disfluencies counted in a participant's delivery proved to be a reasonable predictor for the rated adequacy of the speaker's expression and coherence, explaining between a quarter and a third of the variance in the ratings. Interestingly, the ratings could be better predicted by a relative than an absolute count of the number of disfluencies. In the relative measure the number of disfluencies were related to the duration of the interpretation. So the expert judges did not just keep track of the number of disfluencies they heard in the interpreter's delivery but normalized for the length of the delivery.

There is no point in trying to predict the ratings of grammatical correctness of the interpretations. Since the interpretation was from English into Farsi, all participants spoke the target language as their native language. Although numerous disfluencies were found in the Farsi utterances produced, no ungrammatical structures were observed.

Turning now to the prosodic rating scales, it appeared that the pace of the delivery is clearly related to a number of objective parameters. The three phonetic-prosodic evaluation scales are very highly intercorrelated, even if the correlation coefficients are computed for the experimental and control groups combined ( $.888 < r < .976$ ). I decided to concentrate on the prediction of pace (fluency) as this parameter has rather straightforward acoustical correlates. The results show that the pace rating for the control group can be predicted most successfully by a single parameter, i.e., percent pause duration, which by itself explains 72 percent of the variance in the pace rating. Curiously enough, no predictive model is possible for the experimental group and only

two single predictors yield marginally significant correlations with pace, i.e., the variability in the duration of the interpausal units and the relative number of disfluencies.

In the overall prediction of pace for the group of 30 participants combined a regression model was found that explains 43 percent of the variance. The best predictor here was the number of disfluencies (normalized for the total duration of the interpretation), followed by the variability in the duration of (filled and silent pauses).

It remains unclear at this time why the pace (or fluency) judgments can be predicted in a rather straightforward fashion from a number of objective properties of the speech produced by the student interpreters in the control group, whereas no convincing relationships could be found between the acoustic measurements and counts of errors and disfluencies for the experimental group. Part of the solution of this problem may be that the assumption underlying the analysis I applied is that the relationships between the predictors and the criterion should be linear. In § 7.1, however, I briefly speculated that it might be more reasonable to assume a U-shaped (i.e., quadratic or parabolic) relationship between such parameters as speech rate and percent pause on the one hand and judged pace on the other. Obviously, when there is excessive pausing or an exceedingly slow speaking rate, which would cause poor judgments of pace (or fluency). However, a speaker may also speak so fast and with so few pauses that the listener suffers from information overload – which would yield unfavorable ratings of pace. I argued that speech rates and speech pause ratio (i.e., percent pause duration) should ideally be somewhere in the middle of the range, neither too slow nor too fast.

No signs of non-linearity can be observed in the results obtained for the control group. For this group the overall tendency is: the faster the better. However, when we examine the results of the experimental group more closely, we may observe a tendency in the scatterplots (Figures 7.2-3) to reveal non-linear, in fact, parabolic relationships between the acoustic predictors and judged pace. Table 7.10 lists side-by-side the correlation coefficients between the acoustic predictors and judged pace obtained for linear and quadratic (U-shaped) regression functions for the experimental and the control group separately.

*Table 7.10. Correlation coefficients (Pearson's  $r$ ) between acoustic predictors and judged pace for experimental and control groups, assuming linear versus quadratic relationships.*

Acoustic predictor	Experimental group			Control group		
	Linear	Quadratic	$\Delta$	Linear	Quadratic	$\Delta$
% Pause	.200	.531	.331	.848	.887	.039
Articulation rate	.001	.430	.429	.592	.642	.050

Table 7.10 shows that the U-shaped function fits the data much better (by 33 to 43 points) than a linear function. For the control group, however, the difference between linear and quadratic functions is almost negligible (5 points or less). I am inclined to interpret this difference as an indication that some speakers in the experimental group

speak so fast and pause so little that the raters judge this speed of delivery (or pace) uncomfortable.

### References

- Boersma, P. & Weenink, D. J. M. (1996). *Praat, a system for doing phonetics by computer, version 3.4*, report 132, Institute of Phonetic Sciences University of Amsterdam.
- Boersma, P. & Weenink, D. J. M. (2017). *Praat: doing phonetics by computer* [Computer program]. Version 6.0.28, retrieved 23 March 2017 from <http://www.praat.org/>
- Buhmann, J., Caspers, J., Heuven, V. J. van, Hoekstra, H., Martens, J.-P. & Swerts, M. (2002). Annotation of prominent words, prosodic boundaries and segmental lengthening by non-expert transcribers in the spoken Dutch corpus. In *Proceedings of LREC 2002 Paris*: ELRA, 779–785.
- Grosjean, F., Grosjean, L. & Lane, H. (1979). The patterns of silence: performance structures in sentence production. *Cognitive Psychology*, 11, 58–81.
- Jones, R. (2014). *Conference interpreting explained*. Routledge: New York.
- Jong, N. de & Perfetti, C. A. (2011). Fluency training in the ESL classroom: An experimental study of fluency development and proceduralization. *Language Learning*, 61, 533–568.
- Koiso, H., Horiuchi, Y., Tutiya, S., Ichikawa, A. & Den, Y. (1998). An analysis of turn-taking and backchannels based on prosodic and syntactic features in Japanese Map Task dialogs. *Language and Speech*, 41, 295–321.
- Selkirk, E. O. (1984). *Phonology and syntax. The relation between sound and structure*. Cambridge: MIT Press.
- Swerts, M. (1997). Prosodic features at discourse boundaries of different strength. *Journal of the Acoustical Society of America*, 101, 514–521.

