

## AFFECTIVE SPEECH GATING

*Ioulia Grichkovtsova<sup>1</sup>, Anne Lacheret<sup>2</sup>, Michel Morel<sup>1</sup>,  
Virginie Beaucousin<sup>3</sup>, Nathalie Tzourio-Mazoyer<sup>3</sup>*

<sup>1</sup>CRISCO, Université de Caen, France, <sup>2</sup>MoDyCO, Université Paris X, Nanterre, France,

<sup>3</sup>GIN, CNRS UMR 6194, GIP Cyceron, France

*ioulia.grichkovtsova@unicaen.fr*

### ABSTRACT

This study tested the hypothesis that emotions may be identified earlier than attitudes in the flow of speech. The gating paradigm was chosen to investigate if such differentiation between emotions and attitudes was possible. Perception test results included the following variables: the identification point, the isolation point and the confusion matrix. Acoustic analysis was conducted and linked to the perception results. Anger and sadness were separated from the other studied affective states on the basis of the results analysis. Interestingly, happiness followed the identification pattern found for attitudes. The future directions of work are presented.

**Keywords:** Affective speech, speech perception, gating paradigm.

### 1. INTRODUCTION

The scientific context of the present study is placed in the multi-disciplinary perspective of psycholinguistics, phonetics and neurolinguistics. The tested hypothesis originates from the collaborative study on neural bases of affective speech comprehension between two laboratories: linguistic and neuroimaging. In the frames of the study, two categories of affective states were studied and compared: emotions (anger, sadness and happiness) and attitudes (obviousness, irony and doubt). According to the definitions proposed by Scherer [11], emotions represent a synchronized organismic response to the event of major significance, while attitudes like obviousness, doubt, may be described as affective colouring of interpersonal interaction.

The results of the identification psycholinguistic test, reported in [2], showed no significant differences between the two categories on the bases of response analysis. It was hypothesized that the difficulty to distinguish emotions from attitudes in the results of the identification test may be explained by the fact that subjects were permitted to answer only after listening to the whole utterance. The presence of specific prosodic characteristics, occurring from the very beginning of emotional utterances,

may allow early identification. Listeners may be able to recognize emotions before the end of the utterance. The identification of attitudes, requiring an integrated analysis of lexical, syntactic and prosodic structures of the utterance [13], may happen later than for emotions.

The gating paradigm was chosen to test the hypothesis that emotions may be identified earlier than attitudes. The gating paradigm was originally developed by Grosjean [6], and it was used in spoken word recognition research [8]. More recently, it was applied in intonation research [1, 12]. The gating task is based on the principal that an audio stimulus is presented in segments of increasing duration, and the subjects are asked to identify what was said at the end of each segment or “gate”. The gating paradigm allows to understand how much acoustic-phonetic information is needed to identify a stimulus.

### 2. METHOD

#### 2.1. Stimuli description

Two utterances were pronounced by each of the two recorded speakers (an actor and an actress) for six studied affective states (anger, sadness, happiness, obviousness, irony and doubt), thus the total of 24 utterances was used in the experiment. The lexical meaning of the utterances was not neutral. They were designed to carry natural lexical meaning. The beginning of the utterances was potentially possible for several affective states, and the meaning was disambiguated by prosodic and lexical means by the end of the utterance. Ironic utterances were marked by a deliberate contrast between their apparent and intended meaning. It was achieved in the used utterances by opposing the literal meaning of the words between themselves and reinforced by the tone of voice. Examples of the corpus are given in the Appendix.

#### 2.2. Participants

13 subjects were recruited in University of Caen (5 females and 8 males). They were students and pro-

professionals working in the university. Average subject age was 31 years old, with standard deviation of 13.2. All subjects were native speakers of French and none reported having any hearing difficulty.

### 2.3. Procedure

Special software for psycholinguistic perception tests *Perceval* was used for the experiment design [4]. It allows the programming of stimuli presentation according to the gating paradigm without cutting utterances at the gating points. Increment size of gates was fixed at 200 msec. The duration-blocked presentation format was chosen: first all the stimuli of the particular segment size were presented to the listeners in a randomized order, then all the stimuli of the following segment size, and so on. Each block contained 24 stimuli. The experiment was run on a computer in a quiet laboratory room. The whole experiment was run by the *Perceval* software, responses and response times were automatically recorded in the data file. In average, the experiment took about 40 minutes.

### 2.4. Analyzed variables

Responses proposed after each gate were analyzed, both correct responses and observed patterns of confusion. Two variables have been selected for the analysis of correct responses: identification point and isolation point. *Identification point* refers to the gate of the stimulus where correct identification achieves 50%. *Isolation point* is the gate of the stimulus where the highest identification is achieved and maintained without any change in response thereafter.

### 2.5. Perception test results

**Table 1:** Results for the identification point and the isolation point. Two values show the observed range of gating points for each affective state across the four utterances. Anger - A, Doubt - D, Obviousness - O, Happiness - H, Irony - I, Sadness - S.

Affect	Identification point	Isolation point
A	200/400 msec	800/1600 msec
S	400/1000 msec	800/1600 msec
O	400/1200 msec	800/complete
D	600/1400 msec	800/complete
H	800/1800 msec	1000/complete
I	1600/1800 msec	1600/complete

Results for the identification point are shown in Table 1. Anger stands out of the other affective states, as it is the earliest to get recognized at 50%. Angry utterances are successfully recognised at the first or maximum at the second gate. Sadness

**Table 2:** Confusion matrix calculated on the whole response data (in percentage). Anger - A, Doubt - D, Obviousness - O, Happiness - H, Irony - I, Sadness - S.

	A	S	H	D	O	I
A	<b>74,9</b>	1,1	10,5	2,5	7,7	3,1
S	1,3	<b>70,1</b>	1,1	5,8	18,3	3,1
H	11,9	3,8	<b>54,1</b>	4,9	10,0	15,1
D	8,4	3,3	6,8	<b>57,6</b>	8,4	15,3
O	10,8	7,7	1,3	12,4	<b>56,9</b>	10,6
I	2,8	8,9	29,2	4,0	15,4	<b>39,6</b>

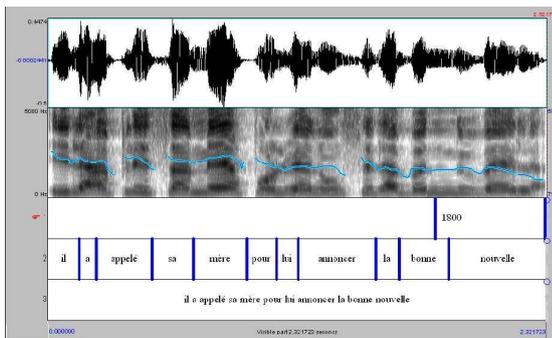
also can be recognised as early as the second gate, but more variability in the identification of sad utterances was observed. Interestingly enough, obviousness may be also identified early. Other affective states, and especially irony, were identified later.

Results for the isolation point are presented in Table 2. It shows the range of observed gates for each affective state where the highest identification is achieved and maintained without any change in response thereafter. All the studied affective states have some utterances which may be recognized before the end of the utterance. At the same time, only anger and sadness have *all* their utterances recognized before hearing the complete form. This observation allowed to separate the affective states into two groups. Interestingly, based on the analysis of the identification and isolation points, happiness was classified not with emotions, but with attitudes.

A confusion matrix was used to study the percentage of correct responses and the pattern of confusions, see Table 3. The confusion matrix is based on the principle to plot encoded affective states against the decoded affective states. Table 3 shows that confusion patterns are observed for all the affective states, nevertheless anger and sadness have the lowest level of confusion. Obviousness, doubt, happiness, and especially irony have high level of confusion.

### 2.6. Acoustic analysis

The main objective of the acoustic analysis was to understand the relationship between the level of identification and the used pattern of acoustic correlates. A large number of acoustic correlates (voice quality, intensity, initial F0, pitch span, pitch level, peak steepness, speech rate) may be used for the encoding and decoding of affective states. The acoustic analysis was conducted with PRAAT software [3] and linked to the identification and isolation points. The auditory analysis of voice quality was based on definitions, proposed by Laver [7], and it involved only evaluation of the main phonation setting and the presence of smiling (lip spreading). A linguistic



**Figure 1:** The graph represents an example of the acoustic analysis realized in PRAAT. Utterance 15 was produced with happiness by the male actor. The identification point is at 1800 msec, and the isolation point is at the end.

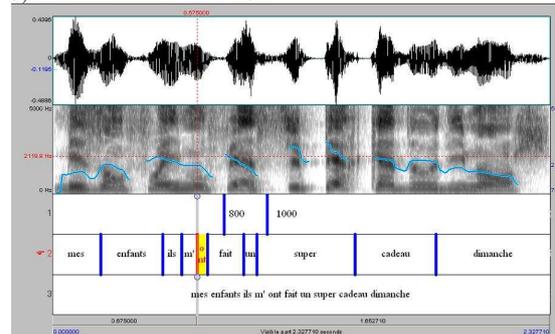
**Table 3:** Results of acoustic analysis for utterances 15 and 16 (happiness by male actor)

Measure	Utterance 15	Utterance 16
Initial $F_0$	6.8 ERB	5.0 ERB
Pitch span	40.7 st	42.6 st
Final $F_0$	5.1 ERB	4.8 ERB
Speech rate	6.1 syl/sec	5.6 syl/sec
Voice quality	-	smiling

approach to the pitch range measurement [10] was used: level, “height” of a speaker’s voice, was measured as the final  $F_0$  value; span, the width of pitch frequencies covered by a speaker, was measured as the difference between  $F_0$  peaks and valleys of the utterance. Examples of acoustic analysis are shown in Figures 1 and 2. Due to the lack of space, individual results for all 24 analyzed utterances are not presented in this paper.

Some difficulty to generalize the acoustic results was encountered due to the fact that the two speakers used different strategies to express the same affective states. Moreover, variability within the same speaker was also observed. This variability on the production level had its impact on the identification level. For example, utterances 15 and 16 (see Figures 1 and 2) said by the male speaker with happiness differ in that utterance 16 has its isolation point at the 5th gating point (1000 msec), while utterance 15 is identified only at 50 % even in its complete form. Table 5 gives results of the analysed acoustic measures. It is possible to see that utterance 15 is pronounced fast with high pitch level, but narrow pitch span. On the contrary, utterance 16 has a lower pitch level, but a wider pitch span, moreover the actor is smiling. The isolation point for the utterance 16 occurs before the lexical meaning is disambiguated. Apparently, the realisation pattern allowed

to recognize happiness earlier and better in utterance 16, than in utterance 15.



**Figure 2:** Utterance 16 was produced with happiness by the male actor. The identification point is at 800 msec, and the isolation point is at 1000 msec.

### 3. DISCUSSION

The hypothesis that emotions may be identified earlier than attitudes was tested in the present study by the gating paradigm. Based on the results for the recognition point, it is possible to see that all angry utterances were identified very early (1-2 gating point), while other affective states show much more variability. Due to the isolation point, the studied affective states were separated into two groups. The first group comprises anger and sadness, as all angry and sad utterances were recognized before the end of the utterance. Anger and sadness showed much less confused responses than the other affective states. Based on these results, the differentiation of emotions from attitudes was realized. Happiness makes an exception from the other studied emotions, as its recognition follows the pattern identified for attitudes. It is a very interesting observation, it may be explained by a particular communicative role, played by happiness. Happiness may be better controlled by speakers, and it can be used intentionally to colour interpersonal communication, in a way attitudes are used.

Acoustic results showed variability both in the production of the two speakers, and in the production of the same speaker for the studied affective states. It has been acknowledged in the previous studies [9] that affective states are communicated through voice by a variable combination of acoustic parameters. The search for stable associations between a group of particular acoustic parameters and specific affective states has not yet been successful. It has been recently suggested [5] that several different strategies may be successfully used for the expression of the same affective state. The acoustic analysis of the studied corpus does not allow to understand in detail which correlations between the us-

age of acoustic correlates and their perceptual value are possible, as the corpus has only two speakers, moreover they had different utterances for the same affective state. Nevertheless, it was shown that the chosen pattern of acoustic correlates had influence on the level of recognition, some strategies were more successful than others.

#### 4. CONCLUSION

The gating paradigm allowed to discover some differences in the identification of emotions and attitudes. The hypothesis that emotions may be identified earlier than attitudes was confirmed for anger and sadness, at the same time it was not possible to separate happiness from attitudes in the perception test results. The study has also shown inter- and intra-speaker variability in the realization of affective states, and its influence on the perception. Some strategies to encode the same affective state were more successful and easier for the identification than others. The used corpus of affective speech included only two speakers, it cannot give deep understanding of possible strategies in the production of affective states.

At the present, a new corpus is being developed for the perception tests with the gating paradigm. A large number of speakers will be recorded encoding several affective states on the same neutral utterance. It will help to investigate the width of possible variability on the production level and to search for correlations between the usage of acoustic correlates and their perception value.

#### 5. APPENDIX

**Happiness:** Il a appelé sa mère pour lui annoncer la bonne nouvelle. (He called his mother to tell her the good news.) **Irony:** J'ai réussi ma chute en pleine rue brillament. (I managed to fall in the middle of the street very nicely.)

#### 6. REFERENCES

- [1] Aubergé, V., Grépillat, T., Rilliard, A. 1997. Can we perceive attitudes before the end of sentences? the gating paradigm for prosodic contours. *Proceedings of the European Conference on Speech Communication and Technology* Rhodes, Greece. 871–874.
- [2] Beaucousin, V. 2006. *Bases neurales de la compréhension de la phrase affective: Des fonctions ortholinguistiques à la prosodie affective*. PhD thesis University of Caen.
- [3] Boersma, P. 2001. Praat, a system for doing phonetics by computer. *Glott International* 5, 341–345.

- [4] Ghio, A., André, B., C. and Teston, C., C. 2003. Perceval: une station automatisée de tests de perception et d'évaluation auditive et visuelle. *TIPA Aix-en-Provence, France*. 115–133.
- [5] Grichkovtsova, I. 2007. *A cross-linguistic study of affective prosody production by monolingual and bilingual children: Scottish English and French*. PhD thesis Queen Margaret University.
- [6] Grosjean, F. 1996. Gating. *Language and Cognitive Processes* 11, 597–604.
- [7] Laver, J. 1994. *Principles of phonetics*. Cambridge: Cambridge University Press.
- [8] Lickley, R., McKelvie, D., Bard, E. 1999. Comparing human and automatic speech recognition using word gating. *Proceedings of the ICPhS Satellite meeting on Disfluency in Spontaneous Speech* UC Berkeley. 23–26.
- [9] Murray, I., Arnott, J. 1993. Towards the simulation of emotion in synthetic speech: A review of the literature on human vocal emotion. *Journal of Acoustic Society of America* 93, 1097–1108.
- [10] Patterson, D., Ladd, R. D. 1999. Pitch range modelling: New demensions of variation. *Proceedings of the ICPhS* San Francisco, USA.
- [11] Scherer, K. R. 2003. Vocal communication of emotions: A review of research paradigms. *Speech Communication* 40, 227–256.
- [12] Vion, M., Colas, A. 2006. Pitch cues for the recognition of yes-no questions in french. *Journal of Psycholinguistic Research* 35, 427–445.
- [13] Wichmann, A. 2000. The attitudinal effects of prosody, and how they relate to emotion. *Proceedings of SpeechEmotion-2000* Belfast. 143–148.