

Identification des émotions en voix naturelle et synthétique : paradigme d'ancrage

Ioulia Grichkovtsova¹, Michel Morel¹ et Anne Lacheret²

¹Laboratoire CRISCO, Université de Caen
Esplanade de la Paix, 14032 Caen Cedex, France

²Laboratoire MODYCO, Institut universitaire de France, Paris
Université Paris X, Nanterre
ioulia.grichkovtsova@unicaen.fr

ABSTRACT

The two main objectives of the study were to identify the perceptive role of intonation and voice quality in the identification of emotions and to determine gating points for each studied emotion. A multi-speaker corpus was used for the development of a new perception test on the basis of gating paradigm and transplantation paradigm. The results of the study are discussed in the light of their relevance for the voice synthesis of affective speech.

Keywords: affective prosody, emotions, speech synthesis, gating paradigm.

1. INTRODUCTION

La valeur distinctive de l'intonation et de la qualité vocale n'est toujours pas établie de manière fiable pour l'identification des émotions. L'intérêt de cette problématique n'est pas seulement théorique : une meilleure modélisation des émotions permettra des avancées importantes en TAP, en particulier en synthèse de la parole.

La difficulté d'obtenir des résultats conclusifs est en partie au moins due aux types de corpus utilisés généralement pour l'analyse de l'encodage des émotions. Il en existe deux types principaux : les phrases avec ou sans contenu lexical.

Une combinaison inventée des syllabes présentes dans la langue [1], ou un énoncé dans une langue étrangère, inconnue des locuteurs et des auditeurs [2], ont pu être utilisés dans les expériences perceptives sur la base de phrases non lexicales. Ces études n'ont pas pu établir une association claire de l'intonation avec les émotions étudiées, contrairement à la qualité vocale dont l'importance pour l'identification des émotions a été parfaitement démontrée. [3]

L'utilisation des phrases dépourvues de sens lexical peut être problématique dans l'étude de l'intonation, car elles ne permettent pas de découpage syntaxique et pragmatique. Notre hypothèse est que sans la distinction mots pleins/mots vides, sans les frontières syntaxiques et les frontières informationnelles, les points d'ancrage pour les patrons prosodiques qui encodent l'expression des émotions sont absents.

Les études [4,5] utilisant des phrases dotées d'un sens lexical mettent en lumière l'importance de l'intonation dans l'identification des émotions, même si cette importance est variable. La présence de l'intonation naturelle pour l'identification peut même être cruciale pour certains états affectifs. Néanmoins, cette approche provoque un autre biais : la neutralisation complète du lexique émotionnel n'est pas possible. L'utilisation d'un patron prosodique, même "neutre", préserve toujours l'amplification des mots importants du lexique et de ce fait facilite l'identification des émotions.

Une autre faiblesse méthodologique est présente dans les expériences sur l'identification des émotions : l'insuffisance du nombre de locuteurs (1 ou 2). Or, la parole émotionnelle est caractérisée par une variabilité inter- et intra- individuelle significative [6]. La question se pose de savoir si cet échantillon des locuteurs est représentatif de l'ensemble des stratégies phonostylistiques possibles pour l'encodage des émotions.

2. METHODOLOGIE

2.1. Établissement du corpus

L'analyse des problèmes méthodologiques soulevés en *supra* nous a conduits à créer un corpus multi-locuteur sur la base d'une phrase émotionnellement neutre : *Je vais rentrer à la maison maintenant*. 22 locuteurs natifs du français (11 hommes et 11 femmes) ont été enregistrés. La taille du corpus représente environ 7 heures d'enregistrement. Six émotions (colère, tristesse, joie, dégoût, peur, chagrin) et une affirmation neutre ont été sélectionnés pour ce travail.

Pour chaque émotion, un texte spécifique devait être oralisé mais à chaque fois avec l'insertion de la même phrase neutre au milieu du texte. Il a été ainsi possible de récolter la même séquence lexicalement neutre mais dans des contextes émotionnels variés. L'hypothèse sous-jacente était la suivante : l'expression des émotions induites dès les premiers mots du texte, au contenu lexical émotionnellement marqué, se propage dans la production de cette phrase dite neutre, autrement dit, la phrase en question porte les traces prosodiques de l'émotion globale véhiculée dans le texte, sans que le locuteur en soit conscient.

Chaque texte a été lu trois fois par chaque locuteur pour que celui-ci imagine mieux le contexte de la situation et exprime l'émotion désirée le plus naturellement possible¹.

2.2. Validation du corpus

Le corpus enregistré (154 stimuli présentés dans un ordre aléatoire, 20 minutes par passation) a été validé par un test d'évaluation avec le logiciel PERCEVAL [7]. Dix auditeurs français (moyenne d'âge 30 ans) ont participé au test.

Nous avons fixé le seuil d'acceptation des énoncés identifiés à 50% (au moins 50% des auditeurs ont pu identifier l'émotion exprimée). Les résultats du test perceptif sont présentés dans le Tableau 1.

Table 1 : Taux du nombre d'énoncés identifiés par les auditeurs avec le seuil d'acceptation de 50%.

Emotions	Identification réussite
Chagrin	32%
Colère	73%
Dégoût	9%
Joie	36%
Peur	14%
Tristesse	68%

A l'issue de ce test préliminaire, nous avons mis à l'écart le dégoût et la peur, et retenu pour l'analyse les 4 émotions suivantes : chagrin, tristesse, colère et joie.

2.3. Test perceptif

Le corpus multi-locuteur ainsi validé a été utilisé pour le développement d'un nouveau test perceptif qui inclut deux paradigmes méthodologiques : le *paradigme d'ancrage* et le *paradigme de transplantation*.

Paradigme d'ancrage

Le paradigme d'ancrage [8] permet de comprendre quel type et quelle quantité d'information phonétique est nécessaire pour une reconnaissance optimale des émotions. Ce paradigme est fondé sur l'hypothèse selon laquelle un point d'ancrage correspond à une zone de l'énoncé qui déclenche la reconnaissance de l'émotion. Pour mesurer ces points, la méthodologie est la suivante : le stimulus vocal est coupé en fragments qui augmentent progressivement par pas de une ou deux syllabes jusqu'à l'énoncé entier² :

je vais | ren | trer | à la | mai | son | main | tenant |

Il est ensuite demandé à chaque auditeur d'identifier l'état affectif à la fin de chaque segment.

¹ Un exemple de texte est donné en appendice.

² Nous n'avons pas retenu la méthode temporelle classique (augmentation par rapport à un seuil de durée déterminé) de manière à contourner les variations inter-locuteur de débit.

Paradigme de transplantation

Le paradigme de transplantation [9] a été utilisé pour évaluer les rôles perceptifs respectifs de la qualité vocale et de l'intonation pour l'identification des 4 émotions (chagrin, tristesse, colère et joie) ainsi que de la production neutre.

Cette méthode repose sur une extraction et un échange de la prosodie entre deux énoncés avec le même contenu segmental, l'une naturelle, l'autre synthétique avec intonation neutre. La transformation prosodique agit sur le contour intonatif, l'intensité, le rythme et les paramètres temporels, mais elle ne modifie pas les caractéristiques de la qualité vocale.

Quatre versions de chaque phrase ont été produites selon cette méthode de la transplantation avec le logiciel de synthèse vocale Kali [10]: version 1 « naturelle » (prosodie naturelle et qualité vocale naturelle de l'acteur), version 2 « qualité vocale » (qualité vocale naturelle de l'acteur et prosodie de Kali), version 3 « prosodie » (prosodie naturelle de l'acteur et qualité vocale de Kali). Version 4 : énoncé entièrement synthétique dépourvu de tout contenu phonétique affectif.

Au total, 1560 stimuli ont été construits. Compte tenu du nombre de stimuli, le test est basé sur le principe du carré latin qui permet de réduire le nombre de stimuli présentés à chaque sujet. Un **carré latin** est un tableau carré de n lignes et n colonnes rempli de n éléments distincts dont chaque ligne et chaque colonne ne contient qu'un seul exemplaire. Sur ce principe, huit versions du test ont été préparées. Pour sa validation statistique, le test a nécessité 16 sujets (auditeurs français, âge moyen 25 ans, 20 minutes de passation par sujet).

Les stimuli ont été présentés par groupes de segments en commençant par le plus petit, l'ordre des stimuli étant aléatoire dans chaque groupe.

3. RESULTATS

Le *chagrin* est caractérisé par une qualité vocale spécifique, Figure 1. La montée régulière dans l'identification de la version « qualité vocale » est observée, malgré la dernière valeur étonnante. Une dégradation à la fin peut être expliquée par le dernier mot « maintenant » moins chargé émotionnellement qui, en une fraction de seconde, réduit le taux de l'identification.

Ce premier résultat n'est pas étonnant quand on sait que le chagrin est une émotion forte et peu contrôlée, qui, en conséquence, modifie d'abord les paramètres physiologiques non codés en langue et influence essentiellement la qualité vocale.

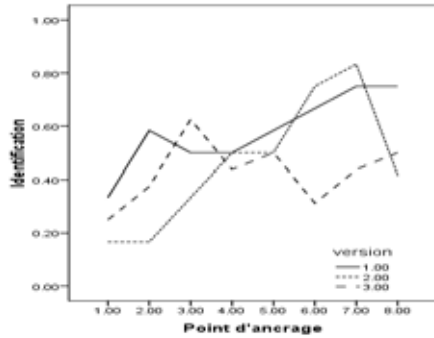


Figure 1 : Chagrin. Version 1 « naturelle », version 2 « qualité vocale », version 3 « intonation ».

La *colère* met en lumière une bonne identification pour la version naturelle, Figure 2. La qualité vocale joue un rôle faible en comparaison avec l'intonation. L'importance de l'intonation fait de la colère un mécanisme socialement et donc linguistiquement contrôlé : relevant plus de l'attitude que de l'émotion.

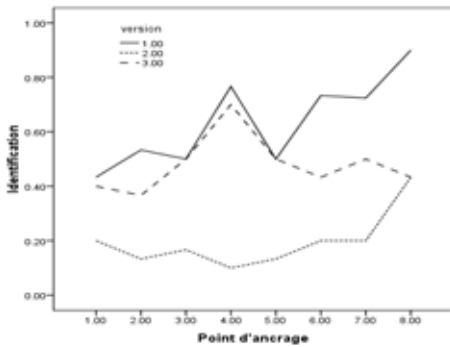


Figure 2 : Colère. Version 1 « naturelle », version 2 « qualité vocale », version 3 « intonation ».

L'identification de la *joie* naturelle est très supérieure aux versions modifiées, Figure 3. L'intonation est tout juste significative, la qualité vocale négligeable. Le résultat négatif pour la qualité vocale résulte de la transplantation prosodique qui altère la voix souriante dans la version 2 « qualité vocale » et masque la joie. Si ce résultat montre les limites de la méthode par transplantation, elle met néanmoins en valeur le rôle de l'intonation dans la synthèse d'une émotion comme la joie.

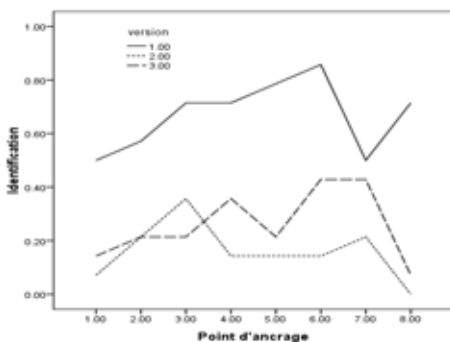


Figure 3 : Joie. Version 1 « naturelle », version 2 « qualité vocale », version 3 « intonation ».

Les résultats pour la *tristesse* mettent en lumière un fait intéressant : l'identification de la version 2 « qualité vocale » est aussi bonne que celle de la version naturelle. L'intonation, en revanche, joue un rôle significativement moins important.

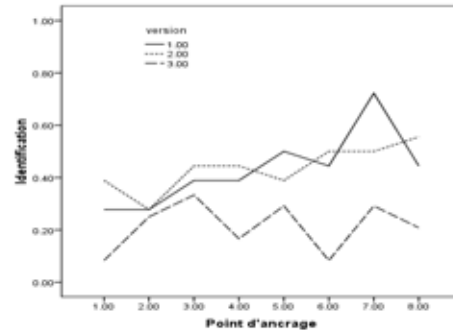


Figure 4 : Tristesse. Version 1 « naturelle », version 2 « qualité vocale », version 3 « intonation ».

L'identification du neutre est équivalente pour la version « naturelle » et la version « intonation ». La participation de la qualité vocale est peu significative. Le neutre se distingue donc bien des modalités émotionnelles, surtout par l'intonation.

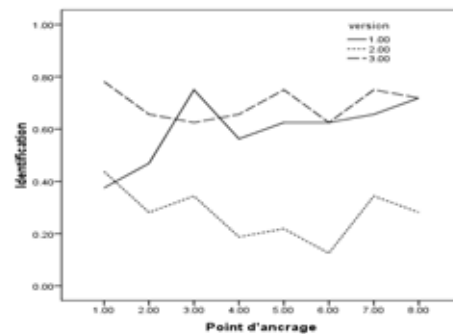


Figure 5 : Neutre. Version 1 « naturelle », version 2 « qualité vocale », version 3 « intonation ».

4. DISCUSSION

Le choix méthodologique d'utiliser le même énoncé pourvu d'un sens global purement référentiel pour une étude contextualisée de l'encodage des émotions a permis de déterminer de façon plus fiable la valeur perceptive de l'intonation et de la qualité vocale. Une telle méthode, en effet, nous a permis d'éviter deux biais. Le premier consiste à analyser des énoncés non écologiques dépourvus de signification, dans lesquels l'absence de structure syntaxique réduit la valeur perceptive de l'intonation. Le second, à l'inverse, repose sur l'analyse d'énoncés déjà très connotés par leur matériel verbal. Si pour ces derniers, on a pu montrer que la seule présence du lexique ne suffisait pas pour une identification optimale des émotions [4,5], il restait très difficile de quantifier objectivement l'apport des paramètres suprasegmentaux dans les processus de reconnaissance.

La méthode présentée ici conforte l'hypothèse qui est la

nôtre concernant le double statut de l'intonation dans l'activité de langage : 1) elle amplifie les mots informationnellement importants dans un contexte émotionnel spécifique; 2) elle facilite l'identification d'une émotion par un contour prosodique qui lui est propre. L'utilisation d'un énoncé lexicalement neutre permet justement d'évaluer le rôle d'un contour prosodique dans l'identification de l'émotion en l'isolant de l'amplification. Il est ensuite possible de caractériser ses propriétés acoustiques sur le plan qualitatif et quantitatif.

5. CONCLUSION

L'utilisation des patrons intonatifs naturels par la synthèse vocale donne accès, au moins partiellement, à l'expression des états affectifs. Certes, la qualité vocale joue parfois un rôle important, mais c'est un paramètre difficile à faire varier en temps réel. Le contour prosodique (hauteur, intensité, durée) est beaucoup plus facile à implémenter. Nos résultats montrent que le contour prosodique joue un rôle non négligeable dans l'expression globale des émotions. L'analyse de la variabilité inter-individuelle du corpus et la sélection des stratégies les plus adaptées pour la reproduction dans la synthèse vocale constituent l'étape suivante de notre projet.

BIBLIOGRAPHIE

- [1] T. Bänziger and K.R. Scherer. The role of intonation in emotional expressions. *Speech Communication*, volume 46, pages 252-267, 2005.
- [2] C. Gobl and A. Ní Chasaide The role of voice quality in communicating emotion, mood and attitude. *Speech Communication*, volume 40, issues 1-2, pages 189-212, 2003.
- [3] C. d'Alessandro. Voice source parameters and prosodic analysis. In: S. Sudhoff, D. Lenertová, R. Meyer, S. Pappert, P. Augurzy, I. Mleinek, N. Richter & J. Schließer (eds): *Methods in Empirical Prosody Research*. Berlin, New York: De Gruyter, pages 63-87, 2006.
- [4] I. Grichkovtsova, A. Lacheret and M. Morel. The role of intonation and voice quality in the affective speech perception. In *Proc. of Interspeech*, 2007.
- [5] I. Grichkovtsova, A. Lacheret, M. Morel, V. Beaucousin and N. Tzourio-Mazoyer. Affective speech gating. In *Proc. of ICPHS*, pages 805-808, 2007.

6. APPENDICE

Chagrin : « Tu sais comme j'aimais mon chien ? Hier, quand je suis revenu(e) de voyage, j'ai appris qu'il était mort. Je suis bouleversé(e). *Je vais rentrer à la maison maintenant.* J'vais jamais pouvoir le dire aux enfants. »

Colère : « Vous appelez ça une chambre d'hôtel ? Regardez un peu ces draps ! Ils sont ignobles. Vous ne croyez quand même pas que je vais dormir ici ! C'est révoltant ! *Je vais rentrer à la maison maintenant !* Ce n'est pas un hôtel ici, c'est un élevage de cafards ! »

Joie : « Mon frère reviendra demain ! Quelle joie ! Je suis si content(e) ! *Je vais rentrer à la maison maintenant !* Je vais annoncer cette super nouvelle à ma famille ! »

Neutre : « J'ai fini de ranger les boîtes. Elles sont classées et numérotées. *Je vais rentrer à la maison maintenant.* Je reviendrai demain à dix heures. »

Tristesse : « Ce que tu m'as appris m'a fichu le moral à zéro. C'est vraiment déprimant... *Je vais rentrer à la maison maintenant.* J'ai l'impression que c'est une situation sans issue. »

7. REMERCIEMENTS

Le projet était financé par le Conseil Régional de Basse-Normandie. Les auteurs remercient tous les participants aux expériences.

- [6] I. Grichkovtsova. A cross-linguistic study of affective prosody production by monolingual and bilingual children: Scottish English and French. PhD thesis. Queen Margaret University. 2007.
- [7] A. Ghio, C. André, B. Teston and C. Cavé. PERCEVAL: une station automatisée de tests de PERCEption et d'EVALuation auditive et visuelle. In *TIPA*, vol. 22, pages 115-133, 2007.
- [8] F. Grosjean. Gating. *Language and Cognitive Processes*, volume 11, pages 597-604, 1996.
- [9] M.-N. Garcia, C. d'Alessandro, G. Bailly, P. Boula de Mareuil and M. Morel. A joint prosody evaluation of French text-to-speech systems. In *Proc. LREC*, pages 307-310, 2006.
- [10] M. Morel and A. Lacheret-Dujour. Kali, synthèse vocale à partir du texte : de la conception à la mise en oeuvre. *Traitement Automatique des Langues*, volume 42, pages 1-29, 2001.