



**HAL**  
open science

## Synthèse concaténative de la voix chantée

Luc Ardaillon, Axel Roebel

► **To cite this version:**

Luc Ardaillon, Axel Roebel. Synthèse concaténative de la voix chantée. Journées des Jeunes Chercheurs en Audition, Acoustique musicale et Signal , Jul 2014, Lyon, France. 2014. hal-01251892

**HAL Id: hal-01251892**

**<https://hal.science/hal-01251892>**

Submitted on 7 Jan 2016

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Les recherches présentées ici concernent la synthèse de voix chantée. Le but est d'obtenir, à partir d'un texte et d'une partition, une voix synthétique à la fois naturelle et expressive chantant la partition et le texte donnés en entrée.  
La méthode utilisée pour cela est la concaténation et transformation d'unités [Bonada 2006] : une base de données est enregistrée par un chanteur, de laquelle des segments sont extraits puis assemblés afin de générer la synthèse. Les discontinuités entre ces segments doivent alors être supprimées. Afin de rendre la synthèse naturelle et expressive, un module de contrôle doit permettre d'en générer les paramètres (hauteur, durées, intensité et timbre) en fonction des contextes donnés en entrée par la partition (mélodie, paroles) et certains choix esthétiques (style, accents expressifs, ...). Les paramètres ainsi générés permettent alors de transformer les segments concaténés afin d'obtenir le résultat final de la synthèse.

1

## Base de données

### Enregistrements :

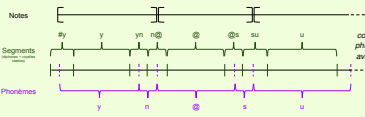
- Série de courtes phrases chantées lentement à **hauteur et intensité fixes** couvrant l'ensemble des **diphones de la langue française**

- Bases différentes enregistrées pour **2 styles "lyrique" et "variété"**, afin de capturer les **qualités timbrales propres à chaque style**

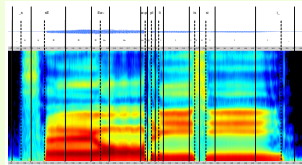
### Stratégie de segmentation :

- Double segmentation :

- (1) en **phonèmes** => gestion précise des durées et placement des phonèmes
- (2) en **diphones** "courts" + parties stables => gestion indépendante des parties co-articulées (transitions entre phonèmes) et stables (tenues de notes)



- **Segmentation automatique** en phonèmes et parties stables => permet de déduire les segments des diphones à utiliser pour la synthèse



3

## Contrôle

### Objectifs :

- Synthèse naturelle et expressive
- Modélisation du style
- Contrôle paramétrique intuitif
- Indépendant du système de synthèse

### Etude :

- Constitution d'un corpus (lyrique et variété française) avec 2 chanteurs professionnels pour étude et apprentissage des paramètres **en contexte**

5

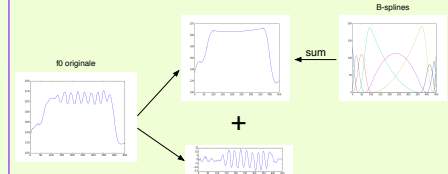
## Fréquence fondamentale

- 2 composantes principales :

- courbe "lisse" avec transitions entre notes
- vibrato

(+ inflexions liées au contexte phonétique)

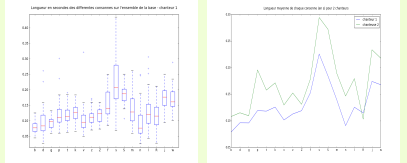
- Modélisation par B-splines avec approximation par méthode des moindres carrés => modèle paramétrique



## Durées et alignement des phonèmes

**Consonnes** : liées à l'articulation

- durée moyenne propre à chaque consonne pour un débit lent (dépendant du locuteur)



- durées varient en fonction du contexte phonétique (début de phrase, enchaînement de consonnes, ...) et du tempo => à apprendre

**Voyelles** : liées aux notes

- => début des notes sur les onsets des voyelles
- => durée\_voyelle = durée\_note - durée\_consonnes

## autres paramètres :

- intensité / effort vocal
- enveloppe spectrale (formants et source)
- => corrélation avec hauteur et intensité à établir

## Concaténation

- Dans le **domaine spectral**

- 2 types de **discontinuités** à supprimer

### 1). Phase

**Problème** : alignement "vertical" des phases entre harmoniques différent d'un segment à l'autre



**Solution** : étaler ces différences dans le temps, pour chaque harmonique, avec un écart maximum imperceptible entre 2 trames



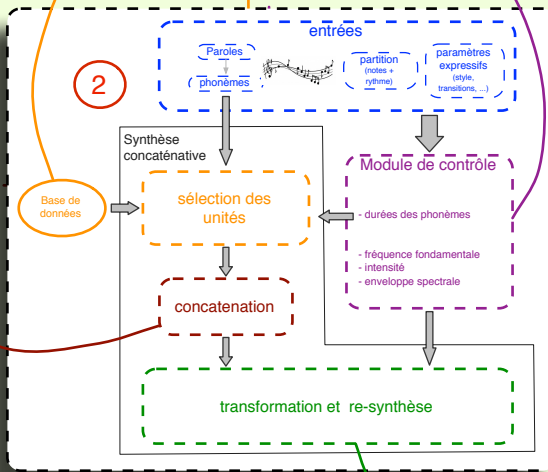
### 2). Enveloppe spectrale

**Problème** : différences d'enveloppe d'un segment à l'autre pour un même phonème

**Solution** : interpoler les enveloppes sur les parties stables (tenues) et autour des jonctions entre segments



4



Temps entre centre de la trame et phase 0 de la fondamentale :

$$dt = -\varphi_{\text{centre}}(\omega) / \omega \text{ avec } \omega = f_0 \times 2\pi$$

Différence de phase entre l'harmonique i et la fondamentale :

$$\varphi_i^0 = \text{princarg}(\varphi_{\text{centre}} + 2\pi \times f_0 \times dt)$$

Ecart entre les trames n et n+1 :

$$\Delta\varphi_i^0(n) = \text{princarg}(\varphi_i^0(n) - \varphi_i^0(n-1))$$

Si  $\text{abs}(\Delta\varphi_i^0(n)) > \Delta\varphi_{\text{max}}$  :

$$\Rightarrow \varphi_i^0(n) = \text{princarg}(\varphi_i^0(n-1) + \Delta\varphi_{\text{max}} \times \text{sign}(\Delta\varphi_i^0(n)))$$

## Expériences de synthèse et résultats

### Objectifs :

- Evaluation du moteur de synthèse
- Identification des composantes à améliorer

### Moyen :

- Re-synthétiser un extrait de chant enregistré dont on peut extraire certains paramètres (fréquence fondamentale et durées des phonèmes) et les imposer lors de la synthèse, en s'affranchissant du module de contrôle. Cela permet d'isoler l'influence de chaque paramètre sur le résultat perçu.

### Conclusion :

- Synthèse de bonne qualité (écoute informelle), ce qui semble valider l'ensemble du processus (contrôle mis à part)
- Durées des phonèmes et F0 essentielles pour un résultat naturel => importance d'un module de contrôle adapté
- Timbre "statique", sans variations expressives => adaptation de l'enveloppe spectrale nécessaire

7

6

## Transformation et re-synthèse

### Pitch-shifting :

- dans le domaine spectral
- Algorithme **SHIP** [Roebel 2010] : Conservation des relations de phases entre harmoniques d'une trame à l'autre => haute qualité pour la voix (pas d'effet de "phasiness")

### Time-stretching :

- Sur les parties stables pour notes tenues
- Sur les parties co-articulées pour contrôler la vitesse d'articulation

### Re-synthèse :

IFFT + overlap-add

## Résultats préliminaires :

Synthèse de bonne qualité, avec jonctions transparentes sans discontinuités, et peu d'artefacts  
Identification des points à améliorer

8

## Deux principaux axes de recherche :

- (1) contrôle des paramètres en fonction des contextes et du style, avec apprentissage
- (2) adaptation de l'enveloppe spectrale (source et conduit vocal) en fonction des autres paramètres (hauteur, intensité) et de la qualité vocale souhaitée