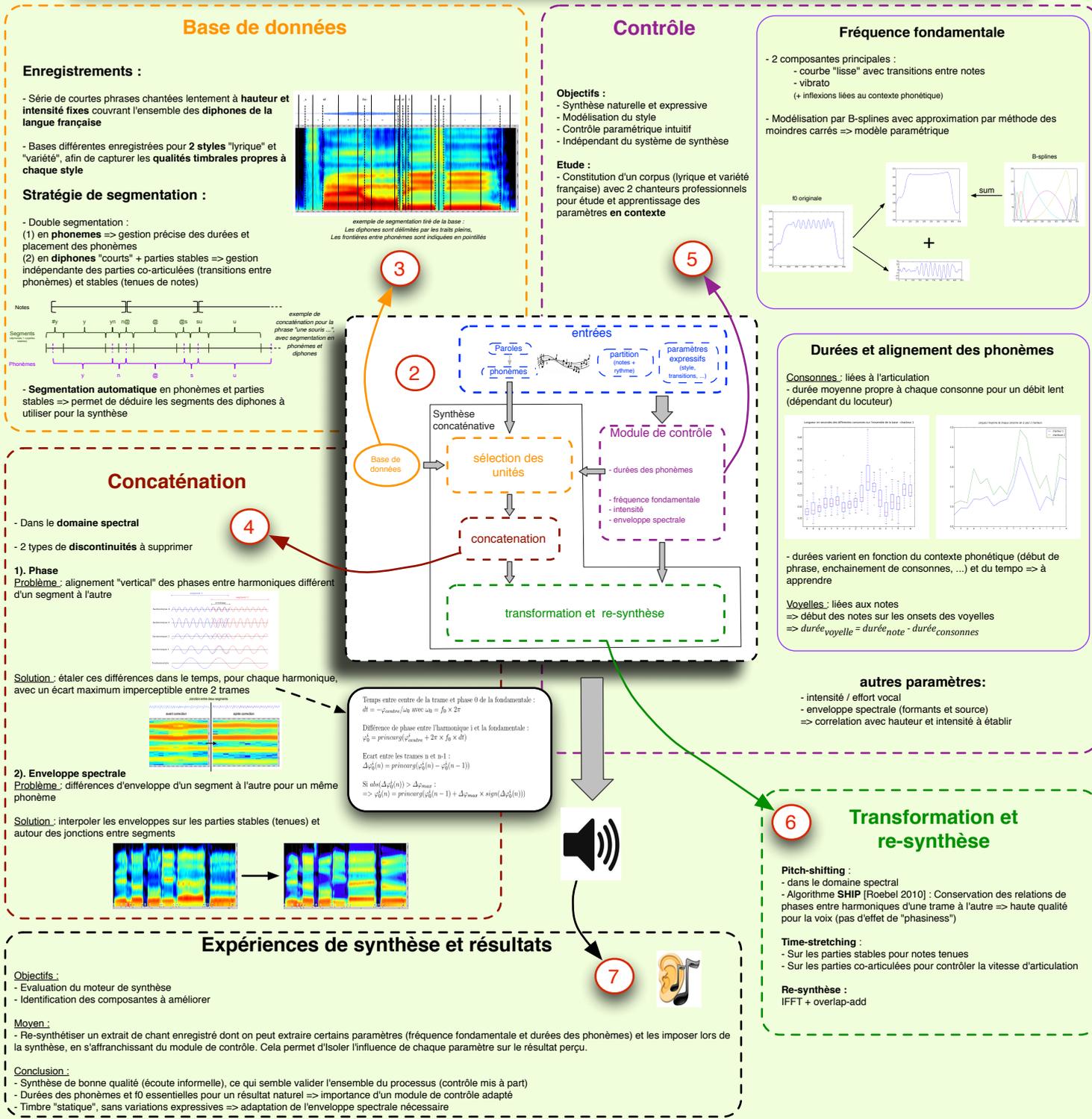


Les recherches présentées ici concernent la synthèse de voix chantée. Le but est d'obtenir, à partir d'un texte et d'une partition, une voix synthétique à la fois naturelle et expressive chantant la partition et le texte donnés en entrée.
La méthode utilisée pour cela est la concaténation et transformation 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.



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

IO originale + B-splines = sum

Durées et alignement des phonèmes

Consonnes : liées à l'articulation

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

Voyelles : liées aux notes

- => début des notes sur les onsets des voyelles
- => durée_voyelle = durée_note - durée_consonnes

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

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

$$dt = -\varphi_{\text{centre}}(\omega) / \omega = -\varphi_0 + 2\pi$$

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

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

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

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

autres paramètres :

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

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

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 souhaité