

1 [hc21](#)

# 2 Synthèse de la parole à partir 3 du texte

## 4 Synthèse de la parole à partir 5 du texte

6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18 par **Christophe D'ALESSANDRO**

19 *Directeur de recherches*  
20 *LIMSI-CNRS, Orsay, France*

21  
22 et **Gaëli RICHARD**

23 *Professeur*  
24 *Institut Mines-Télécom, Télécom ParisTech, CNRS-LTCl, Paris France*  
25  
26  
27  
28  
29

30	<b>1. Histoire de la synthèse de la parole .....</b>	H 7 288v2 - 2
31	<b>2. Texte – Analyses et traitements linguistiques .....</b>	– 3
32	2.1 Normalisation et prétraitement .....	– 3
33	2.2 Analyse lexicale et morpho-syntaxique .....	– 4
34	2.3 Analyse morpho-syntaxique .....	– 5
35	2.4 Analyse syntaxique.....	– 5
36	2.5 Transcription graphème-phonème.....	– 6
37	<b>3. Signal de parole – Modèle source-filtre .....</b>	– 9
38	3.1 Modèle paramétrique de synthèse de parole.....	– 9
39	3.2 Caractéristiques du filtre .....	– 9
40	3.3 Caractéristiques de la source .....	– 10
41	<b>4. Prosodie .....</b>	– 11
42	4.1 Prosodie et syntaxe.....	– 11
43	4.2 Calcul du rythme.....	– 12
44	4.3 Calcul de l'intonation .....	– 13
45	<b>5. Synthèse acoustique .....</b>	– 13
46	5.1 Synthèse à formants par règles.....	– 14
47	5.2 Synthèse non paramétrique par concaténation d'unités acoustiques ...	– 15
48	5.3 Synthèse par diphtonges.....	– 15
49	5.4 Synthèse par sélection et concaténation .....	– 18
50	5.5 Synthèse paramétrique statistique.....	– 20
51	5.6 Construction du corpus textuel et sonore.....	– 22
52	<b>6. Applications de la synthèse de parole .....</b>	– 23
53	6.1 Exemples d'applications.....	– 23
54	6.2 Interfaces de programmation .....	– 24
55	6.3 Produits.....	– 24
56	<b>7. Évaluation de la synthèse .....</b>	– 25
57	7.1 Boîte noire ou boîte de verre .....	– 25
58	7.2 Évaluation de qualité globale.....	– 25
59	<b>8. Conclusion.....</b>	– 26
60	8.1 Bilan .....	– 26
61	8.2 Perspectives.....	– 27
62	<b>Pour en savoir plus.....</b>	Doc. H 7 288v2

L'objet de la synthèse de la parole à partir du texte (ou TTS, Text-To-Speech) est de calculer automatiquement le signal de parole correspondant à un texte donné. Le texte lui-même peut provenir de diverses sources : journaux, livres, systèmes de réponse vocale, de dialogue ou traduction automatique (borne interactive, assistant personnel), base de données d'un système d'information, jeu vidéo, courriers électroniques, SMS, documents butinés sur la toile, ou tout simplement texte saisi au clavier d'un ordinateur.

La réponse vocale sous sa forme la plus simple peut être un ensemble de messages préenregistrés (ou « prompts »). L'ambition de la synthèse de la parole à partir du texte est plus grande : il s'agit de calculer automatiquement les échantillons sonores correspondant à un énoncé écrit quelconque, qui n'est pas connu d'avance et qui peut être de grande taille.

Les deux versants de la synthèse de parole sont d'une part, l'analyse et l'interprétation du texte, d'autre part, la prédiction des paramètres acoustico-phonétiques du son et la synthèse du signal proprement dite :

- analyse du texte. La première étape de la transformation d'un texte en parole implique la capacité d'analyser, de comprendre le texte écrit, ses nuances et ses connotations, la situation du discours et l'acte de parole à effectuer. En plus du texte, le contexte peut être spécifié (style de parole, émotion, attitude, type de personnage, voix spécifique...);

- synthèse du signal. Une fois le texte analysé, il s'agit de calculer le signal acoustique qui interprète au mieux le contenu linguistique, avec une voix aussi naturelle que possible, ressemblant à un locuteur particulier, et avec les nuances d'attitude, voire d'émotion que le texte réclame. En plus du signal audio, le synthétiseur peut fournir des indications pour synchroniser le mouvement des lèvres d'un avatar ou personnage vidéo, ou les mouvements d'un robot.

## 1. Histoire de la synthèse de la parole

L'histoire de la synthèse à partir du texte est déjà longue : par exemple le premier système autonome de synthèse automatique de la parole en français, l'Icophone V du LIMSI, date de 1974. Cependant, cela reste encore un domaine de recherche très actif. Les travaux actuels portent à la fois sur la compréhension des textes et sur la restitution d'une parole naturelle, personnalisée et expressive. Les applications se multiplient. Le lecteur désireux d'approfondir les notions présentées dans cet article pourra consulter les références portées dans la partie documentation.

L'architecture générale d'un système de synthèse se compose ainsi de ces deux parties principales (figure 1). Les principaux modules correspondant à ces traitements sont décrits dans cet article. Bien que tous les exemples cités dans la suite soient tirés du français, il est important de souligner que les problèmes posés sont similaires pour toutes les langues. Cependant, des différences notables existent en fonction des spécificités linguistiques de chaque langue particulière et notamment en ce qui concerne leur :

- notation graphique [alphabétique (romaine, cyrillique, hébraïque, sanskrite, arabe...), syllabique (coréenne, japonaise...), idéogrammatique (chinoise, japonaise...)] ;
- grammaire (agglutinante, flexionnelle, niveaux de langue...);
- phonologie et phonétique (système de phonèmes, langues à ton, langues à clics...);
- prosodie (durées, intonation, qualité vocale).

Les systèmes actuellement développés pour toutes les langues s'inspirent de principes identiques, même s'ils diffèrent pour les corpus, lexiques, analyses et heuristiques linguistiques.

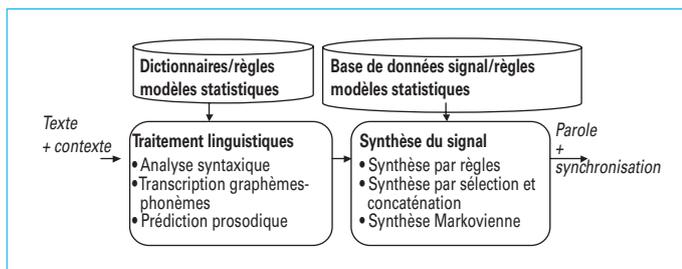


Figure 1 – Architecture générale d'un système de synthèse à partir du texte

## 2. Texte – Analyses et traitements linguistiques

La première étape d'un système TTS comprend les modules de traitements linguistiques qui permettent de transformer le texte à synthétiser en une chaîne de symboles représentant les sons distinctifs de la langue, ou « phonèmes », et un ensemble d'indications « prosodiques » caractérisant l'élocution (durée des différents sons et des pauses, évolution de la mélodie). Cette représentation phonético-prosodique est ensuite utilisée par l'étape de synthèse sonore qui assure la génération du signal de parole numérisé.

Les traitements linguistiques se décomposent en plusieurs étapes, qui obéissent en général à la séquence suivante :

- **normalisation et prétraitement du texte brut** : cherche les phrases, les anomalies du texte et produit une séquence lexicale (séquence d'unités lexicales ou « mots ») ;
- **analyse lexicale et morpho-syntaxique** : enrichit la séquence lexicale avec des étiquettes lexicales, des marques morpho-syntaxiques et produit un découpage en composantes syntaxiques ;
- **transcription graphème-phonème** : transforme la suite lexicale de sa forme orthographique à la forme phonétique.

### 2.1 Normalisation et prétraitement

Les sources du texte à synthétiser peuvent être variées : interfaces de programmation (API), textes issus de la toile, SMS, journaux ou livres électroniques, courriels, etc.

Confrontés à des sources de textes réelles et variées, le système de synthèse doit normaliser les nombreuses « anomalies » des textes au niveau graphémique (c'est-à-dire des « graphèmes » ou caractères écrits). Cette étape de prétraitement a pour objet de retranscrire en toutes lettres les chaînes de caractères non lexicales (hors dictionnaire), ou inconnues.

Il peut s'agir de nombres et de chiffres (3 999 356 à transcrire « trois millions neuf cent quatre vingt dix neuf mille trois cent cinquante six »), de dates (« 24/01/63 », « 24 janvier 1963 »), ou plus généralement de mots composés de caractères orthographiques et numériques (« vol AF102 », « référence SD44 »), de symboles spéciaux (« € », « © »).

Les diverses formes d'abréviations (par exemple, « c.-à-d. », « Pr. ») doivent aussi être repérées et traitées, comme les sigles (suite de lettres initiales non prononçable, qui est épelée, « CNRS », « SNCF »), les acronymes (suite de lettres initiales prononçable, « UNESCO », « ENA »).

L'analyse de la ponctuation est aussi un élément de prétraitement. L'étape de normalisation utilise à la fois un système de

1 règles de transcription (pour le traitement des quantités numé-  
2 riques ou des dates, des abréviations standards, ou des  
3 acronymes), et un lexique paramétré par l'utilisateur, spécifique de  
4 chaque domaine d'application de la synthèse.

5 Il est important de noter que beaucoup d'éléments sont  
6 ambigus. À titre d'exemple, il faut distinguer le rôle du point  
7 comme élément syntaxique (point final) ou comme symbole  
8 numérique (« 1.3 »). Seule une procédure de désambiguïsation,  
9 utilisant des heuristiques contextuelles permet de décider de la  
10 transcription appropriée.  
11

12 L'étape de prétraitement permet de former une suite de  
13 « lexèmes », ou unités lexicales à partir du texte d'entrée. Les  
14 signes de ponctuation et les caractères typographiques non litté-  
15 raux sont également analysés, afin de former la suite lexicale (par  
16 exemple en séparant les lexèmes comme dans des formes  
17 composées comme « l'y » ou « 3,5 »).

18 L'apparition de textes électroniques peu structurés, voire entachés  
19 de fautes plus ou moins systématiques (courriers électroniques, tex-  
20 tes sans accent, orthographe approximative, etc.) pose des problè-  
21 mes nouveaux que l'on se doit de traiter automatiquement.  
22

## 23 2.2 Analyse lexicale et morpo-syntaxique

24 Par le prétraitement, le texte d'entrée a été transformé en une  
25 suite d'unités lexicales, séparées, et encadrées par la ponctuation.  
26 L'étape d'analyse lexicale consiste à rechercher dans un lexique  
27 les informations associées aux différents lexèmes.  
28

29 ■ Le lexique contient les **classes lexicales** (ou parties du discours)  
30 associées au lexème. Suivant les systèmes, la classe lexicale peut  
31 être très simple, comme par exemple « mot outil/mot plein », qui  
32 sépare les mots qui n'ont qu'une fonction syntaxique de  
33 connexion et ceux qui ont une fonction sémantique (c'est-à-dire  
34 qui portent du sens). D'autres systèmes peuvent utiliser des caté-  
35 gories grammaticales « fines » qui combinent les grandes classes  
36 grammaticales (« nom », « adjectif », « verbe », « pronom »...) et  
37 des propriétés grammaticales de type (« genre », « nombre »,  
38 « infinitif », « verbe d'état », « verbe transitif »...).

39 ■ Il est fréquent que la catégorie grammaticale d'un lexème hors  
40 contexte soit ambiguë, c'est-à-dire que le même mot ait **plusieurs**  
41 **catégories grammaticales**.  
42

43 Par exemple, un mot tel que « voile » peut être soit un nom mas-  
44 culin (« porter la voile »), un nom féminin (« larguer les voiles »),  
45 un verbe transitif à la première ou troisième personne de l'indicatif  
46 présent (« je voile », « il voile ») ou du subjonctif présent (« qu'il  
47 voile », « que je voile »), verbe pronominal (« le ciel se voile »).  
48

49 Ainsi, plus de 70 % des lexèmes peuvent avoir plusieurs catégories  
50 grammaticales différentes.  
51

52 ■ La recherche lexicale est dans certains cas précédée d'une étape  
53 d'analyse morphologique, qui a pour objet de décomposer le  
54 lexème en composantes élémentaires, les **morphèmes**, correspon-  
55 dant aux **préfixes**, **suffixes**, **désinences** (marque du féminin ou du  
56 pluriel pour les noms et les adjectifs, temps, personne et mode  
57 pour les verbes), **racines**. On compte en français environ  
58 500 préfixes, suffixes et désinences. Les désinences permettent  
59 généralement de déterminer des catégories grammaticales  
60 précises. Cela est particulièrement vrai pour les formes verbales  
61 conjuguées qui sont généralement composées d'une racine et  
62 d'indices grammaticaux, propres à la conjugaison du verbe,  
63 porteurs d'information sur le temps, le mode, et la personne. Dans  
64 les systèmes mettant en œuvre un traitement morphologique, le  
65 lexique se compose principalement de formes sources, ou radi-  
66 caux (verbes, noms et adjectifs non fléchis, éventuellement privés  
67 de certains éléments de formation et/ou d'indices grammaticaux).  
68 Le lexique comporte aussi les exceptions de décomposition  
69  
70  
71

1 morphologique, c'est-à-dire les mots qui ne se décomposent pas  
2 suivant les règles morphologiques du français standard.

3  
4 ■ La taille des lexiques varie significativement suivant les  
5 systèmes de synthèse. Certains systèmes utilisent des lexiques  
6 restreints, par exemple, aux mots outils et verbes (de quelques  
7 centaines à quelques milliers de mots). D'autres systèmes plus fins  
8 utilisent des lexiques de l'ordre de 50 000 à 100 000 mots (à titre  
9 de comparaison, un dictionnaire de langue comporte environ  
10 50 000 articles). Les lexiques contenant toutes les formes fléchies,  
11 qui comprennent les mots et l'ensemble de leurs dérivations  
12 morphologiques se composent de 400 000 à 1 000 000 de mots et  
13 de locutions.

## 14 15 16 2.3 Analyse morpho-syntaxique

17  
18 À l'issue du prétraitement des éléments non lexicaux et de la  
19 recherche lexicale, chaque lexème se trouve affecté à une ou  
20 plusieurs catégories grammaticales. Le choix de la catégorie de  
21 chaque mot s'effectue au moyen de règles contextuelles, ou de  
22 modèles statistiques, prenant en compte les catégories grammat-  
23icales des mots adjacents.

24 Les contextes sont réduits à quelques mots précédents ou  
25 suivants, on parle alors d'analyse morpho-syntaxique en  
26 micro-contextes. L'analyse des dépendances syntaxiques globales  
27 est beaucoup plus difficile à mettre en œuvre.

28  
29 L'analyse syntaxique peut faire appel à des **règles heuristiques**  
30 **issues des règles de la grammaire de la langue** (par exemple « on  
31 ne peut observer la succession de deux verbes conjugués »).

32 Une autre stratégie consiste à suivre une approche statistique,  
33 exploitant des modèles probabilistes du langage. Ces modèles  
34 sont fondés sur l'observation que toutes les successions de caté-  
35 gories grammaticales dans une langue donnée ne sont pas équi-  
36 probables.

37 On peut donc chercher à résoudre les ambiguïtés en recherchant  
38 dans l'ensemble des successions possibles de catégories grammat-  
39icales (chaque mot est *a priori* porteur de plusieurs catégories  
40 possibles et l'on considère l'ensemble des transitions entre ces  
41 différentes catégories) la succession de catégories la plus  
42 probable.

43  
44 Ces modèles probabilistes présentent l'avantage de ne requérir  
45 qu'une connaissance sommaire de la langue à traiter, à l'inverse  
46 des approches heuristiques qui compilent des connaissances et  
47 des observations très fines sur la structure des dépendances  
48 fonctionnelles des mots dans chaque langue. Cette caractéristique  
49 présente un avantage décisif, dès que l'on s'intéresse aux  
50 systèmes multilingues, les mêmes paradigmes d'apprentissage  
51 pouvant être déclinés pour plusieurs langues, sans remettre en  
52 question la structure du système.

## 53 54 55 2.4 Analyse syntaxique

56  
57 L'analyse morpho-syntaxique permet la désambiguïsation des  
58 parties du discours, et associe à chaque unité lexicale une catégo-  
59rie grammaticale.

60 Cette suite de catégories grammaticales permet de réaliser une  
61 analyse syntaxique de la phrase, c'est-à-dire de la découper en  
62 constituants syntaxiques et de grouper les mots.

63  
64 Les groupes de mots permettent de structurer la phrase, en  
65 proposant un « parenthésage ». La constitution des groupes  
66 syntaxiques est effectuée par des règles heuristiques ou des  
67 modèles statistiques sur les successions de catégories pour former  
68 des groupes syntaxiques. Le groupement des mots est une étape  
69 essentielle pour le calcul de la prosodie (§ 4), des contours mélo-  
70 diques et rythmiques de l'énoncé. Voici un **exemple** d'analyse  
71 syntaxique simple « en tronçon ».

1 Considérons la phrase :

2  
3 Maintenant, un peu de voix féminine.

4 La suite de catégories grammaticales, ou parties du discours, asso-  
5 ciée est :

6  
7 ADV PMK ART ADV ART NOM ADJ PMK

8 Avec : ART article, ADV adverbe NOM nom, ADJ adjectif, PMK  
9 marque de ponctuation. La liste des mot-outils et des mots pleins  
10 (nom, adjectif, verbe et adverbe) et les règles heuristiques de  
11 groupement permettent de regrouper les mots de la façon suivante.

12 (ART ADV ART ADV) (ART NOM ADJ) PMK

13  
14 Grâce aux analyses lexicale, morpho-syntaxique, et syntaxique « en  
15 tronçons », le texte brut initial est maintenant structuré en sept mots  
16 avec une marque de ponctuation et deux groupes syntaxiques.

17  
18 Le lecteur désireux d'approfondir ces notions consultera avec  
19 profit les références [1] [2].

## 22 2.5 Transcription graphème-phonème

23  
24 Après la normalisation, étape qui transforme un texte brut en  
25 suite de mots sous la forme orthographique, et l'analyse lexicale et  
26 morpho-syntaxique, il s'agit de calculer la prononciation du texte.  
27 Cette « **orthographe inversée** » qui permet de passer des lettres aux  
28 sons a été abordée par le biais de systèmes de règles, de lexiques  
29 spécialisés, ou de techniques d'apprentissage automatique.

30  
31 La transcription graphème-phonème, ou phonétisation, associe  
32 à la forme orthographique une chaîne de signes phonétiques qui  
33 spécifie la prononciation du mot.

34 Pour cela, est utilisé un « alphabet phonétique », sous-ensemble  
35 issu de l'alphabet phonétique international (API), et qui spécifie les  
36 sons élémentaires de la langue. Pour le français, un tel alphabet  
37 comporte 16 sons vocaliques, 20 sons consonantiques. Le codage  
38 informatique SAMPA est souvent utilisé en synthèse pour l'alpha-  
39 bet phonétique (tableau 1).

40  
41  
42 ■ Tout comme pour les catégories grammaticales, la même chaîne  
43 orthographique peut être associée à différentes transcriptions  
44 phonétiques : on parle alors d'« **homographes hétérophones** »,  
45 mots qui s'orthographient de la même façon mais se prononcent  
46 différemment.

47 Le français standard comprend environ 150 homographes  
48 hétérophones ; il s'agit dans la plupart des cas d'ambiguïtés entre  
49 un verbe conjugué et un adjectif ou un adverbe formé sur la même  
50 racine, comme par exemple un président (nom)/ils président  
51 (verbe) ; somnolent (adjectif formé sur le verbe somnoler)/ils  
52 somnolent (verbe).

53 Le français comporte des cas, beaucoup plus rares, d'homopho-  
54 nie mettant en jeu des mots de racines différentes : les portions  
55 (nom)/nous portions (verbe) ; est (nom)/est (verbe).

56 Dans les cas précédents, la connaissance de la catégorie  
57 grammaticale du mot permet de choisir la prononciation correcte.  
58 Dans certains cas, plus rares, des homographes hétérophones ont  
59 la même catégorie grammaticale, comme par exemple : fils (du  
60 père)/fils (de coton). Seule une analyse sémantique (une étude du  
61 sens de la phrase) ou l'analyse du contexte élargi, permet alors  
62 d'effectuer la phonétisation correcte.

63  
64 ■ La **transcription phonétique** proprement dite est effectuée à  
65 l'aide d'un lexique d'exceptions et de règles. Après analyse des  
66 exceptions aux règles de phonétisation standards et des homo-  
67 graphes hétérophones, les **règles générales de phonétisation**  
68 s'appliquent.

69 Un système de transcription orthographique-phonétique est un  
70 automate paramétré appliquant un ensemble de règles de réécriture,

**Tableau 1 – Alphabet phonétique du français, avec les symboles issus du projet SAMPA**

Description	Code phonétique SAMPA
archiphoneme / A /	A
[VANtardise, tEMPS]	A <sup>-</sup>
schwa	@
closed / oe / [crEUser, dEUx]	2
open / E / [pERdu, modEle]	E
closed / e / [Emu, otE]	e
[pEINture, matIN]	E <sup>-</sup>
[malhEUreux, pEUr]	9
[Idée, amI]	i
open / O / [Obstacle, cOrps]	O (o Majuscule)
closed / o / [AUditeur, bEAU]	O
[rONdeur, bON]	O <sup>-</sup>
[cOUpable, lOUp]	u
[pUnir]	y
[OUi, OIseau]	w
[hUIe]	H
[plétiner, paiLLe]	j
	p
	t
	k
	b
	d
	g
	f
	s
[CHanter, maCHine]	S
	v
	z
[Jardin, manGer]	Z
	l
uvular / R /	R
	m
	n
[αGNeau, rèGNe]	J
[campING]	N

qui permettent d'associer un phonème (ou un groupe de phonèmes) à un caractère (ou un groupe de caractères) orthographique(s), en prenant en compte le contexte gauche (caractères ou groupes de caractères précédant le segment à transcrire) et le contexte droit

(caractères ou groupes de caractères suivant le segment à transcrire). Ces règles sont organisées de façon hiérarchique, des règles les plus particulières aux règles les plus générales.

Le nombre de règles nécessaires pour effectuer la transcription orthographique phonétique dépend de la langue considérée ; par exemple, moins de 100 règles sont requises pour une langue comme l'espagnol, la forme orthographique étant très proche de la forme phonétique.

Même si les exceptions ne sont pas prises en compte, le cas du français est beaucoup plus complexe, car la forme orthographique est, pour des raisons historiques, très éloignée de la forme phonétique. Un système minimal de description des règles de phonétisation du français standard se compose environ de 500 règles. Les meilleurs systèmes associent un lexique important avec de plus de 2 000 règles.

Pour donner un **exemple**, le mot « oiseau » se transcrit phonétiquement en /wazo/, par application des règles suivantes :

1. la chaîne de caractères orthographiques « oi » se transcrit par la succession des phonèmes /wa/, parce qu'elle est précédée d'un séparateur de mot et qu'elle n'est pas suivie de la chaîne « gn » comme dans « oignon », ou d'un « n » comme dans « oindre » ;

2. la lettre « s » se transcrit par le phonème /z/ car cette lettre est entourée par deux voyelles et que « oiseau » ne fait pas partie d'une liste d'exceptions à cette règle, stockée dans le lexique (comme « paraSol » ou « vraiSemblance ») ;

3. la chaîne de caractère « eau » se transcrit par le phonème /o/, indépendamment du contexte.

La chaîne phonétique obtenue après transcription peut être enrichie de marques syllabiques, à l'aide d'un petit ensemble de règles. La syllabe est souvent considérée comme l'unité rythmique de base, utile pour le calcul de la prosodie.

La question de la **phonétisation**, qui semble en apparence assez régulière, se révèle en fait difficile lorsqu'il faut traiter de problèmes comme les noms propres (notamment ceux d'origine étrangère), les mots nouveaux et inconnus, les variantes de prononciation, les différents dialectes, idiolectes ou sociolectes. Cela pose d'importantes questions de phonologie, comme celles du « e » muet, de la coupe syllabique, des liaisons, de l'harmonie vocalique, de l'emprunt de phonèmes d'autres langues, etc. Même avec des taux de phonétisation correcte très élevés (plus de 95 ou 98 %), les erreurs effectives de phonétisation restent fréquentes dans un texte de quelques secondes, puisque nous prononçons environ 10 phonèmes par seconde.

Voici un **exemple** de chaînes obtenues à la **sortie du module de phonétisation** :

La légende veut que, en rêvant devant Grenade au cours de vacances espagnoles, lord Sydney Bernstein ait choisi le nom de son entreprise.

Texte :

Cinquante ans plus tard, le groupe pèse un peu plus de 8 milliards de francs et affiche un bénéfice de 900 millions pour 1986.

Phonèmes et mots :

slkAt A ply tar, lx grup pEz Ip@ ply dx hi miljar dx frA e afiS I benefis dx nXf sA miljO pur mil nXf sA katrx vI sis.

Catégories grammaticales des mots :

ADJ| NOM| ADV| ADV| VRG| PPL| VCO| VCO| ADV| ADV| PRP| ADJ| NOM| PRP| NOM| CCO| VCO| DET| NOM| PRP| ADJ| ADJ| NOM| PRP| ADJ| ADJ| ADJ| ADJ| ADJ| PMK|

### 3. Signal de parole – Modèle source-filtre

Cette approche est fondée sur un modèle source/filtre de production du signal vocal, **modèle commandé par un nombre restreint de paramètres**. La synthèse se décompose en deux étapes :

1. à l'aide de règles contextuelles, les informations phonético-prosodiques sont transformées en commandes permettant de spécifier l'évolution temporelle des paramètres du modèle de synthèse ;

2. les valeurs des paramètres ainsi déterminées sont utilisées par le vocodeur pour synthétiser le signal acoustique.

Historiquement, dès la fin des années 1950, ce type de technique a été la première à émerger avec les vocodeurs de l'époque, et des règles élaborées explicitement par l'analyse de corpus phonético-acoustique. Elle reste encore utilisée, après une éclipse en faveur des techniques par concaténation d'unités (§ 4), grâce aux techniques statistiques qui remplacent avantageusement les règles explicites par des techniques d'apprentissage automatique.

#### 3.1 Modèle paramétrique de synthèse de parole

Le signal de parole peut être modélisé par un système source-filtre. La parole résulte de l'excitation des **cavités acoustiques supra-glottiques** (conduit oral, conduit nasal) par des impulsions acoustiques créées par le flux d'air en provenance des poumons et modulé par les cordes vocales. Le modèle source-filtre de la parole représente la production vocale en distinguant deux éléments : une source de phonation, qui représente les impulsions ou le bruit généré à la glotte, et un filtre acoustique qui modélise la contribution de la partie articulaire.

#### 3.2 Caractéristiques du filtre

Les cavités supra-glottiques jouent le rôle d'un **résonateur acoustique**. Les caractéristiques acoustiques de ce résonateur (les fréquences amplifiées et atténuées par ce dispositif) dépendent de la **géométrie du conduit oral et de son degré de couplage avec le conduit nasal**. Ces caractéristiques géométriques sont elles-mêmes contrôlées par les **articulateurs**, à savoir les lèvres, la langue, et la mâchoire. Le couplage entre la cavité orale et la cavité nasale est contrôlé par l'intermédiaire du voile du palais (ou **velum**), qui peut s'abaisser ou se relever, réglant ainsi le débit du flux d'air dévié dans les fosses nasales.

Les caractéristiques acoustiques des cavités supra-glottiques peuvent être, en première approximation, modélisées à l'aide d'un filtre linéaire dont la fonction de transfert varie au cours du temps. Les variations de la fonction de transfert reflètent les modifications de la géométrie des cavités supra-glottiques au cours de l'articulation des différents phonèmes constituant l'énoncé.

Du fait de l'inertie mécanique de ces articulateurs et de la nature de leur contrôle musculaire, les variations sont relativement lentes : le temps typique de stabilité des caractéristiques spectrales est la dizaine de millisecondes (10 millisecondes pour les événements les plus brefs, 40-50 millisecondes pour les segments les plus stables). En pratique, il suffit de spécifier ces fonctions de transfert toutes les 10 à 20 millisecondes.

Les paramètres les plus fréquemment utilisés pour contrôler les caractéristiques de ce filtre sont ceux des « formants » spectraux, ou maxima de la fonction de transfert du conduit vocal, à savoir la fréquence centrale, la bande passante et l'amplitude des maxima (figure 2). Pour obtenir une parole intelligible, il suffit de spécifier les valeurs des trois à quatre premiers formants.

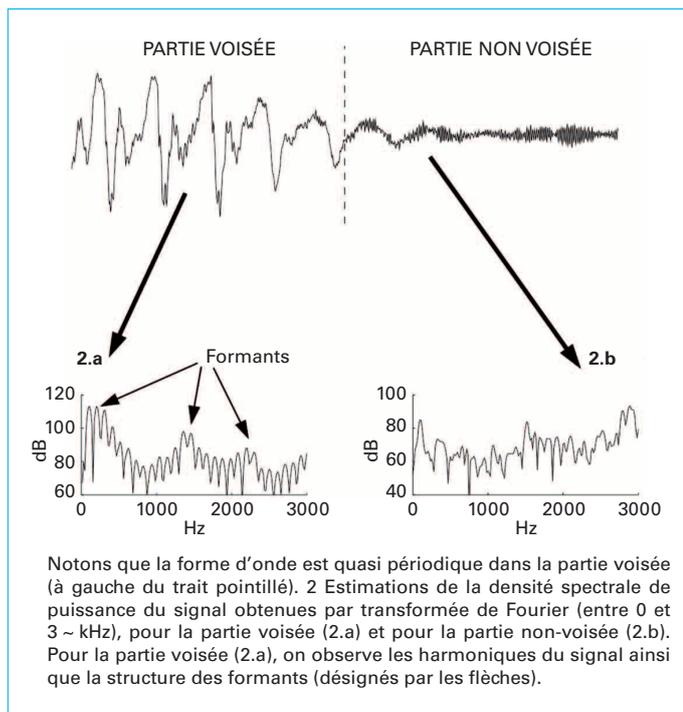


Figure 2 – Forme d'onde d'un signal de parole (environ 100 ms)

### 3.3 Caractéristiques de la source

L'excitation du conduit vocal, considéré comme un filtre acoustique avec plusieurs résonances, dépend du type de phonème considéré (voisé ou non voisé), du mode de phonation (chuchoté, crié...) et de l'effort vocal.

■ Pour les **phonèmes voisés, les voyelles et les consonnes voisées**, les plis vocaux, un peu comme les lèvres du trompettiste, modulent de façon périodique le débit d'air s'écoulant à travers la glotte. La fréquence de cette modulation est la fréquence de voisement, ou fréquence fondamentale. Pour des raisons aérodynamiques, cette onde de débit glottique est généralement très dissymétrique, comprenant une phase (dite d'**ouverture**) où le débit augmente lentement, suivie d'une phase (dite de **fermeture**) plus abrupte. La fréquence de vibration, ainsi que la forme de l'onde de débit est contrôlée par les muscles et cartilages dont est constitué le larynx. L'onde de débit glottique est représentée à l'aide d'un modèle paramétrique contrôlé par la période de voisement, l'amplitude, et des paramètres de forme.

■ Pour les **phonèmes non voisés, les consonnes non voisées**, les cordes vocales restent ouvertes pour permettre le passage d'un flux d'air en provenance des poumons ; l'excitation est due, soit au relâchement rapide d'une occlusion complète du conduit vocal (**plosive**), soit aux turbulences du flux d'air créées au passage d'une constriction du conduit vocal (**fricatives**). Ces signaux sont modélisés par des sources de bruit, transitoires ou continues, réparties dans le conduit vocal, sources de bruit dont on contrôle à la fois la position et la puissance.

■ Notons finalement que, dans la plupart des segments mais plus particulièrement dans les fricatives et plosives voisées (en français /b/ /d/ /g/, /v/ /z/ /j/), ces deux modes d'excitation (voisée/non

voisée) coexistent. On synthétise alors le signal d'excitation en combinant l'onde de débit glottique et les sources de bruit.

## 4. Prosodie

La prosodie est la « **musique** » de la parole, c'est-à-dire sa composante mélodique, rythmique et dynamique. Une même chaîne de phonèmes peut être prononcée sur des « tons » très différents, par la variation prosodique. Du point de vue de la synthèse, le calcul prosodique consiste à modéliser et prédire :

- les contours mélodiques, par l'évolution temporelle de la fréquence fondamentale de vibration des cordes vocales ;
- le rythme syllabique, par les durées des syllabes et des phonèmes ;
- le rythme des groupes de mots, par les positions et les durées des pauses ;
- la dynamique, l'effort vocal, l'intensité relative des syllabes ;
- la qualité vocale incluant le murmure, le chuchotement, les différents mécanismes de voisement ;
- les éléments extralinguistiques, comme les soupirs, rires, bruits de bouche, respiration, raclement de gorge, pauses vocalisées (« hummm »).

La prosodie joue un rôle important pour le naturel de la voix, mais aussi pour son intelligibilité. Certains systèmes de synthèse actuels offrent des styles vocaux différents, des voix dynamiques, enjouées, ou au contraire accueillantes, de proximité. Ces effets sont essentiellement des effets prosodiques. La modélisation prosodique est donc une **composante tout à fait essentielle d'un système de synthèse de parole** (figure 3).

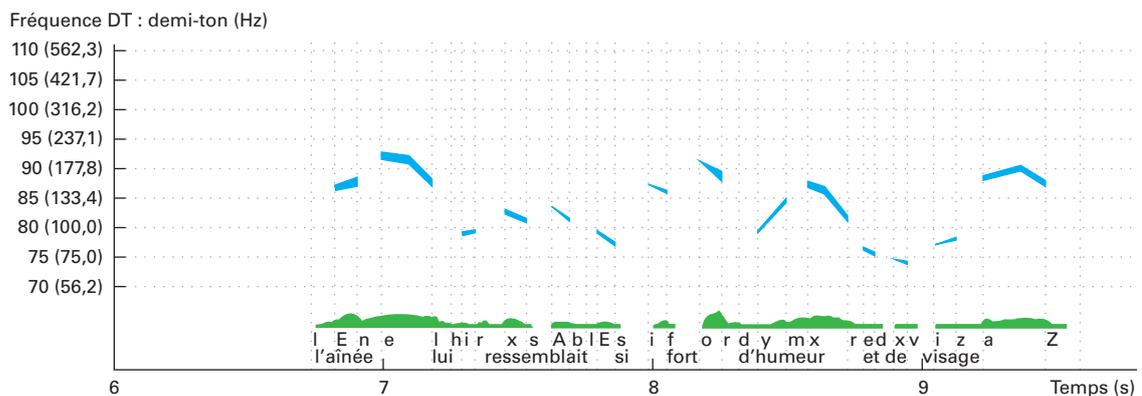


Figure 3 - Prosodie d'une phrase

### 4.1 Prosodie et syntaxe

Tout comme la phrase syntaxique s'ordonne de façon hiérarchique en groupes syntaxiques (groupe sujet, groupe verbal, groupe complément), la phrase prosodique s'organise en groupes prosodiques, groupes de mots le plus souvent séparés par des pauses.

Pour chaque groupe, les paramètres prosodiques (à savoir intonation, durées, dynamique, qualité vocale) suivent une évolution particulière, dépendant du rôle du groupe prosodique dans la

phrase, de sa position et de ses dépendances fonctionnelles avec les groupes adjacents, du nombre de syllabes, mais aussi du sens de la phrase, du style et de l'intention du locuteur.

Les frontières des groupes prosodiques ne reflètent pas systématiquement les frontières des groupes syntaxiques. Une certaine congruence peut toutefois être notée, surtout pour ce qui concerne les frontières syntaxiques majeures (frontière de phrase ou de clause, mais aussi frontière entre le groupe sujet et le groupe verbal associé).

■ La **structuration prosodique de la phrase** obéit à des règles moins bien formalisées que la structuration syntaxique, puisqu'elle ne dépend pas uniquement de la syntaxe, mais aussi du sens (sémantique linguistique), de l'interaction (pragmatique linguistique), du style de parole. Alors qu'il n'existe, une fois les règles grammaticales définies, qu'une seule façon d'analyser syntaxiquement une phrase, il existe plusieurs façons de décomposer la phrase en groupes prosodiques et plusieurs façons de structurer l'évolution des paramètres prosodiques sur chacun de ces groupes.

Par exemple, le sens du message dont est porteur la phrase et l'intention du locuteur (le ou les points clefs de la phrase sur lesquels le locuteur cherche à donner de l'emphase) jouent un rôle tout aussi important que la structure syntaxique pour calculer la prosodie.

■ Il est actuellement difficile, sauf dans des domaines restreints, d'analyser automatiquement les aspects sémantiques ou pragmatiques des énoncés. Pour cette raison, le calcul de la prosodie s'appuie sur le **découpage de la phrase en groupes prosodiques en fonction du découpage syntaxique**. Les contours intonatifs et rythmiques s'appuient sur la syntaxe et la ponctuation. Ce type de prosodie semble convenable pour la lecture de textes, de messages (consultation vocale de messagerie, par exemple) ou d'informations générales (journal téléphonique, bulletin météo...). La prosodie obtenue est néanmoins vite lassante, et peu sensible au contexte d'élocution.

## 4.2 Calcul du rythme

Les pauses sont les temps de silence, de durée variable (de 100 ms à plusieurs secondes), qui s'insèrent à la fin de chacun des groupes prosodiques. L'importance de la coupure syntaxique liée à un marqueur prosodique détermine la durée de la pause à insérer. Ce facteur est particulièrement important pour le **naturel de l'élocution**.

■ Pour le rythme, en plus des pauses, est associée à chaque segment (phonème, syllabe) une durée. Cette durée est déterminée en prenant en compte différents facteurs, en particulier, la **durée intrinsèque des sons** constituant le segment et le contexte. La durée intrinsèque correspond à la durée moyenne du segment tout contexte confondu (ou dans un contexte neutre). Cette durée intrinsèque est généralement déterminée en analysant un **corpus prosodique**, c'est-à-dire un ensemble de phrases qui ont été segmentées et annotées. En plus de la durée intrinsèque, le contexte influence la durée d'un phonème : il peut s'agir de la nature des phonèmes adjacents (certains phonèmes ont tendance à allonger les phonèmes adjacents, d'autres à les raccourcir), de la position de la syllabe porteuse dans le groupe prosodique (en français par exemple, la syllabe finale des mots est généralement allongée, d'un facteur d'autant plus important que le groupe précède une frontière syntaxique majeure), de la nature du groupe prosodique (sa fonction dans la phrase), de la longueur du groupe prosodique, etc.

Les éléments de contexte sont souvent combinés par des « sommes de produits », qui permettent de prendre en compte les allongements ou raccourcissements contextuels. Par exemple, la durée d'une voyelle V, suivie d'une consonne C dans une syllabe finale de phrase sera calculée par :

1 Durée (V, contexte droit C, syllabe finale) = (durée intrinsèque V)  
 2 + (consonne C) × (contexte droit) + [allongement final × (consonne  
 3 C et voyelle V)]

4  
 5 ■ Une bonne détermination des durées segmentales est cruciale  
 6 pour assurer le naturel de l'élocution (des durées erronées  
 7 produisent une parole heurtée, chaotique et, parfois, difficilement  
 8 intelligible).

9 Les durées des phonèmes et des syllabes peuvent être calculées  
 10 par des systèmes de règles, par des tables de durée, ou par des  
 11 méthodes d'apprentissage, souvent des arbres de classification et  
 12 de régression (CART : *Classification And Regression Trees*).

### 15 4.3 Calcul de l'intonation

17 Le calcul de l'intonation, ou contour mélodique, est particuliè-  
 18 rement important pour la **qualité de la synthèse**. L'intonation est la  
 19 variation de fréquence de voisement, la fréquence de vibration des  
 20 plis vocaux. Plusieurs modèles ont été proposés pour décrire cette  
 21 fonction continue du temps en termes discrets, tout comme le flux  
 22 continu de la parole est décrit par une chaîne de phonèmes :

23  
 24 – pour l'anglais et d'autres langues, le système ToBI (*Tones and*  
 25 *Boundary Indexes*) discrétise la courbe intonative, en la représen-  
 26 tant par un petit ensemble de tons hauts, bas, accentués (To), et  
 27 des niveaux de frontière (BI) ;

28 – le modèle tonal perceptif ou Prosogram est un système de  
 29 notation basé sur un modèle de perception de la hauteur tonale,  
 30 prenant comme unité de base la syllabe. Les courbes mélodiques  
 31 sont réduites par stylisation, et des tons statiques ou dynamiques  
 32 obtenus sont affectés à chaque syllabe ;

33 – INTSINT (*INTERNational Transcription System for INTonation*)  
 34 est un système de notation intonatif à visée multilingue, basé sur  
 35 un système de stylisation mélodique par points cibles (MoMel).

36 Tout comme les durées, **l'évolution de la fréquence fondamen-**  
 37 **tales** pour chaque phonème dépend de facteurs, reflétant le  
 38 contexte local et global. On peut donc la calculer par règles ou par  
 39 des méthodes d'apprentissage. Au niveau local, l'évolution de la  
 40 fréquence fondamentale est essentiellement influencée par la  
 41 nature du phonème, par sa position dans la syllabe et par son  
 42 environnement phonétique immédiat (certains phonèmes, comme  
 43 les occlusives voisées, contribuent à abaisser la fréquence  
 44 fondamentale ; d'autres ont tendance à l'augmenter). Au niveau  
 45 global (groupe prosodique, phrase), les facteurs importants sont la  
 46 position de la syllabe dans le groupe prosodique (en français, la  
 47 première et la dernière syllabe d'un groupe prosodique obéissent  
 48 à des règles spécifiques ; dans les langues à accent fixe, les  
 49 syllabes accentuées sont définies pour chaque mot dans le diction-  
 50 naire), la fonction du groupe prosodique dans la phrase, et le  
 51 mode de la phrase (interrogatif, déclaratif...).

52 Il est important de noter que la modélisation prosodique est  
 53 étroitement liée à la technique de synthèse acoustique utilisée.  
 54 Pour la synthèse par règles ou par concaténation de diphones, la  
 55 prosodie est en général calculée par règle. Pour la sélection et  
 56 concaténation d'unités non uniformes, la prosodie est obtenue  
 57 sans calcul explicite, puisque la sélection puis concaténation des  
 58 unités choisies dans un corpus de grande taille conserve la prosodie  
 59 originale des segments sélectionnés.

60 Dans la synthèse par modèles statistiques, la prosodie est  
 61 calculée par apprentissage sur un gros corpus, et générée par  
 62 exemple, à l'aide des chaînes de Markov.

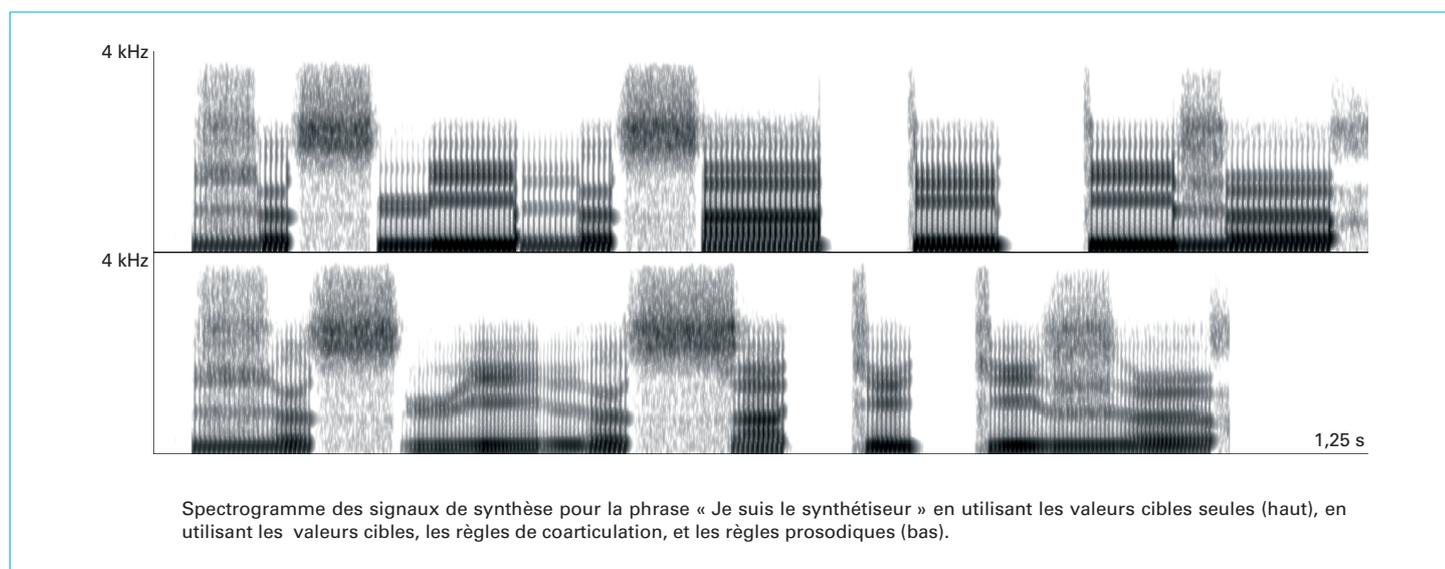
## 67 5. Synthèse acoustique

70 Après la phase d'analyse du texte à produire, intervient la phase  
 71 de synthèse acoustique, qui consiste à **transformer la suite de**

1 **symboles obtenus lors de l'analyse linguistique en suite d'échan-**  
 2 **tilions de signal.** Deux grandes classes de techniques existent pour  
 3 la synthèse acoustique : la synthèse paramétrique et la synthèse  
 4 non paramétrique. En synthèse paramétrique, le signal est calculé  
 5 en utilisant un modèle du signal de parole, le modèle source-filtre :  
 6 le signal passe par un « vocodeur » (pour « *voice coder* », ou  
 7 système d'analyse-synthèse). En synthèse non paramétrique, les  
 8 échantillons de signal sont concaténés avec des modifications  
 9 minimales, sans passer par un vocodeur.

## 5.1 Synthèse à formants par règles

15 Le synthétiseur par règles calcule l'évolution (les trajectoires)  
 16 des paramètres de contrôle du modèle de production à partir de la  
 17 représentation phonético-prosodique, qui spécifie la chaîne des  
 18 sons à prononcer, leur durée et le contour mélodique. La stratégie  
 19 généralement utilisée consiste à spécifier tout d'abord des points  
 20 cibles sur les segments stables du signal de parole (par exemple,  
 21 valeurs de la fréquence centrale, de la bande passante et de  
 22 l'amplitude de chaque formant au centre des voyelles) (figure 4).



49 **Figure 4 – Synthèse à formants par règles**

52 Des règles d'interpolation des paramètres entre les différents  
 53 points cibles sont ensuite mises en œuvre, modélisant les  
 54 phénomènes de coarticulation, c'est-à-dire les interactions  
 55 acoustiques entre phonèmes adjacents. Ces phénomènes de coar-  
 56 tulation sont la traduction acoustique des contraintes articula-  
 57 toires, c'est-à-dire de la dynamique des articulateurs, le conduit  
 58 vocal évoluant relativement lentement.

60 ■ La synthèse proprement dite est réalisée à l'aide d'un **synthé-**  
 61 **tiseur à formants** qui comprend :

63 1. un **module source**, qui comprend un générateur d'impulsions,  
 64 pour la parole voisée, et un générateur de bruit blanc gaussien  
 65 pour la parole non voisée. Les impulsions sont mises en forme par  
 66 un modèle de l'onde de débit glottique (ou alternativement d'un  
 67 modèle de la dérivée de l'onde de débit glottique), fonctions  
 68 construites à l'aide de sinusoides, d'exponentielles ou de poly-  
 69 nômes par morceaux. Bruits et impulsions peuvent se mélanger ;

70 2. un **module de filtre**, qui comprend en général cinq ou six  
 71 résonateurs du second ordre, combinés en série ou en parallèle,

1 pilotés par les valeurs des formants, calculées par les règles de  
2 synthèse.

3  
4 ■ **La qualité obtenue en synthèse par règles est limitée.** D'une  
5 part, la mise au point de règles de synthèse performantes est une  
6 tâche longue et experte : l'obtention de résultats convaincants  
7 nécessite plusieurs années d'efforts pour un expert phonéticien,  
8 un handicap certain lorsque l'on cherche à mettre en œuvre des  
9 systèmes multilingues ou multi locuteurs. D'autre part, malgré de  
10 nombreux efforts en ce sens, aucun système de règle n'est suffi-  
11 samment complexe et précis pour obtenir une qualité de voix  
12 complètement naturelle. Ainsi, dans le meilleur des cas, cette  
13 parole peut être très intelligible, mais assez peu naturelle.

14  
15 Les lecteurs désirant approfondir les aspects liés aux modèles  
16 ou aux stratégies de commande les plus couramment utilisés dans  
17 le cadre de la synthèse par règles peuvent se reporter à la  
18 référence [3].

## 21 5.2 Synthèse non paramétrique 22 par concaténation d'unités 23 acoustiques 24

25  
26 L'idée de restituer de la parole enregistrée en recombinaison des  
27 échantillons est apparue dès les premiers systèmes de réponse  
28 vocale. L'avantage est la qualité naturelle de la parole obtenue,  
29 puisque ce sont des échantillons. Les inconvénients sont l'espace  
30 mémoire requis (ce qui n'est plus vraiment un problème  
31 aujourd'hui), et la nécessité d'enregistrer un même locuteur si on  
32 veut compléter la base de données.

33  
34 Cette seconde approche n'utilise pas de modèle de production  
35 de la parole et ne dépend pas de paramètres d'un synthétiseur.  
36 Elle consiste à synthétiser le signal par concaténation d'unités  
37 acoustiques, c'est-à-dire de **segments de parole préenregistrés**.  
38 Cette technique, reposant sur l'utilisation de segments de signaux  
39 extraits de la parole naturelle, est la seule qui permette à ce jour  
40 de synthétiser des voix dont le timbre est proche, voire identique à  
41 celui d'un locuteur humain.

42  
43 Les premiers systèmes de réponse vocale de haute qualité utili-  
44 saient de la « **synthèse à parties manquantes** », c'est-à-dire des  
45 phrases porteuses, dont certaines parties seulement étaient  
46 variables. Certaines applications, comme lire des bulletins météo  
47 ou bien des itinéraires, fonctionnent très bien de cette façon.  
48 Cependant, il ne s'agit pas vraiment de synthèse à partir du texte,  
49 puisque le vocabulaire est forcément limité.

## 52 5.3 Synthèse par diphtongues 53

54  
55 Pour la synthèse par concaténation, il faut donc utiliser des uni-  
56 tés plus petites que le mot. L'analyse linguistique a montré que le  
57 phonème est l'unité minimale de base de la parole. Mais ces uni-  
58 tés, en petit nombre, sont en fait trop courtes et inappropriées car  
59 elles ne permettent pas de capturer la dynamique du processus de  
60 production de parole : la parole est **essentiellement un processus**  
61 **temporel continu** et la coarticulation entre sons voisins joue un  
62 rôle fondamental.

63  
64 ■ Ainsi, l'**unité minimale permettant d'obtenir une synthèse de**  
65 **qualité acceptable est le « diphtongue »**, qui est défini comme la  
66 portion du signal de parole comprise entre les noyaux stables de  
67 deux phonèmes consécutifs. Le diphtongue, à l'inverse du phonème,  
68 capture la transition entre les différentes cibles articulatoires asso-  
69 ciées aux phonèmes, transitions qui sont cruciales pour la percep-  
70 tion des différents sons. En théorie, le nombre de diphtongues est  
71 égal au carré du nombre de phonèmes, c'est-à-dire environ à 1 300

1 (36 × 36) pour le français (en sachant que certaines transitions  
2 entre phonèmes sont en fait impossibles en français).

3 En pratique, pour la synthèse par diphtongues, le nombre d'unités  
4 utilisées est légèrement plus important (de l'ordre de 1 500-2 000)  
5 pour tenir compte des différentes variantes contextuelles des  
6 phonèmes composant le diphtongue (dans certains systèmes, comme  
7 plusieurs représentants de chaque diphtongue sont disponibles,  
8 l'algorithme de concaténation choisi à chaque instant le  
9 « meilleur » représentant de façon à minimiser une fonction  
10 d'objectif). Le volume de stockage nécessaire est de l'ordre de 5 à  
11 10 Mo (2 à 6 min de parole numérisée avec une fréquence  
12 d'échantillonnage de 16 kHz). Cette quantité de données, qui a  
13 longtemps été considérable par rapport aux tailles mémoire dispo-  
14 nibles, est petite en regard des possibilités de stockage offertes  
15 par les systèmes informatiques actuels.

16  
17 ■ Pour accroître la qualité, il est possible de considérer des **unités**  
18 **plus longues que le diphtongue**, aptes à prendre en compte des phé-  
19 nomènes de coarticulation à plus long terme (disons, pour simpli-  
20 fier, à l'échelle de la syllabe). Parmi celles-ci, les unités de la forme  
21 voyelle-consonne-voyelle (V-C-V), ou de façon plus générale du  
22 type V-C-... -C-V (deux voyelles séparées par un nombre quel-  
23 conque de consonnes) permettent de n'avoir à effectuer des  
24 concaténations que dans les zones les plus stables du signal de  
25 parole, à savoir le centre des noyaux vocaliques. Elles capturent  
26 d'autre part la coarticulation de voyelle à voyelle à travers la (ou  
27 les) consonne(s), coarticulation qui joue un rôle important à la fois  
28 pour l'intelligibilité et l'agrément de la voix de synthèse. Le  
29 problème est que le nombre d'unités ainsi obtenues est beaucoup  
30 plus important (de l'ordre de 10 000-15 000, en ne retenant que les  
31 unités apparaissant effectivement). Un grand nombre de ces unités  
32 sont peu fréquentes et peuvent être éliminées pour satisfaire aux  
33 contraintes de taille.

34  
35 ■ La constitution du dictionnaire d'unités acoustiques se fait en  
36 enregistrant un **corpus de logatomes** (successions élémentaires de  
37 sons de parole n'ayant pas nécessairement de signification), qui  
38 servent de contexte aux unités choisies. L'extraction des unités  
39 acoustiques nécessite la segmentation des logatomes enregistrés.  
40 Celle-ci se fait de manière automatique en alignant, à l'aide de  
41 méthodes statistiques dérivées de la reconnaissance de parole, la  
42 transcription phonétique du mot et la forme acoustique. L'édition  
43 manuelle des résultats de segmentation, longue et fastidieuse, est  
44 nécessaire pour obtenir une synthèse de bonne qualité.

45 La synthèse proprement dite comprend trois étapes distinctes :

#### 46 1. Sélection des unités acoustiques

47 Cette première étape consiste à choisir dans le répertoire d'uni-  
48 tés acoustiques les unités qui seront effectivement utilisées pour  
49 synthétiser la succession de sons désirée. À partir de la représen-  
50 tation phonétique de l'énoncé, il s'agit de rechercher la suite de  
51 segments correspondants. Si les diphtongues sont utilisés, seule la  
52 présence de plusieurs versions pour le même segment est à  
53 prendre en considération. Cette étape est en revanche plus délicate  
54 pour les systèmes à base d'unités de taille variable. Pour une suite  
55 de sons donnée, plusieurs choix d'unités sont en général  
56 possibles. Il faut alors arbitrer entre les différentes décompositions  
57 avec des critères composites.

#### 58 2. Ajustement des paramètres prosodiques

59 Les unités acoustiques préenregistrées possèdent une prosodie  
60 intrinsèque. Cette prosodie intrinsèque doit être neutralisée, afin  
61 d'appliquer la prosodie de synthèse spécifiée par le module proso-  
62 dique. Dans ce cas, il est nécessaire d'utiliser une technique de  
63 traitement de signal pour ajuster aux valeurs cibles définies les  
64 paramètres prosodiques des unités de synthèse. Un tel système  
65 est décrit dans le paragraphe suivant.

#### 66 3. Concaténation des unités

67 Les unités acoustiques, quelles que soient les précautions prises  
68 lors de la sélection et de l'enregistrement des unités, ne possèdent

pas exactement à leur frontière les mêmes caractéristiques acoustiques (en particulier énergétiques). En l'absence de traitement, ces discontinuités vont engendrer des artefacts perceptibles et gênants. Il est donc important de lisser ces discontinuités au moment de la concaténation.

■ Dans un système de synthèse par diphones, les variations prosodiques déterminées lors de la phase d'analyse linguistique doivent être appliquées aux unités acoustiques. Un système de **synthèse par concaténation** (par opposition à la synthèse par règles) n'implique pas l'utilisation d'un modèle de production du signal de parole. Les modifications de durée et de fréquence fondamentale du signal utilisent des techniques non paramétriques.

Une des premières techniques non paramétriques de modification prosodique est l'algorithme PSOLA (pour *Pitch-Synchronous Overlap-Add*). La caractéristique la plus remarquable de la méthode PSOLA est qu'elle opère directement sur la forme d'onde du signal de parole. L'idée de base est d'extraire du signal des grains de sons élémentaires, représentant les caractéristiques locales du signal, et de déplacer ces grains de sons pour réaliser les modifications désirées. De façon plus précise, l'algorithme procède en trois étapes :

### 1. Analyse

Cette étape consiste à extraire du signal une suite de grains de sons élémentaires (signaux à court terme). Ces grains de sons sont obtenus en multipliant le signal par une fenêtre d'analyse centrée autour d'instant d'analyse. Les instants d'analyse sont disposés de façon synchrone à la fréquence fondamentale dans les segments de parole voisés. Ils sont répartis de façon non uniforme et arbitraire sur les segments non voisés.

### 2. Transformation

Cette étape consiste à calculer une suite d'instant de synthèse et une fonction d'association des instants de synthèse et des instants d'analyse pour réaliser la conversion désirée de durée et de fréquence fondamentale. On synchronise ensuite les signaux à court terme d'analyse sur les instants de synthèse, en utilisant la fonction d'association. On définit ainsi une suite de signaux à court terme de synthèse synchronisés sur les instants de synthèse. En l'absence de modification, les instants de synthèse correspondent aux instants d'analyse, et les signaux à court terme de synthèse sont égaux aux signaux à court terme d'analyse. La figure 5 illustre les caractéristiques de la fonction d'association dans le cas d'une modification simple de la fréquence fondamentale (abaissement par un facteur constant).

### 3. Synthèse

Cette dernière étape consiste à recombinaison des signaux à court terme de synthèse. On procède en additionnant les échantillons des signaux à court terme de synthèse associés au même instant de synthèse. Le facteur de normalisation variable tient compte des variations d'énergie liées à la cadence irrégulière de l'analyse et de la synthèse. On remarque qu'en l'absence de modification, le signal de synthèse correspond exactement au signal d'analyse.

■ Le **coût de calcul associé à la méthode PSOLA** est très raisonnable (typiquement moins de 10 multiplications-additions par échantillon de signal). Il faut cependant noter que préalablement à toute modification prosodique, il est nécessaire de déterminer la période du signal de parole (ainsi que son caractère voisé ou non voisé), opération qui est en général plus coûteuse que l'algorithme PSOLA lui-même ; en synthèse par concaténation d'unités, cela n'est pas très gênant, cette opération étant effectuée une fois pour toute lors de l'enregistrement des unités.

Des variantes de l'algorithme PSOLA, comme MBROLA ont connu également un grand succès pour le développement de synthèse multilingue. Une autre approche de modification non paramétrique est le vocodeur HNM (*Harmonic + Noise Model*) qui, par l'utilisation de la représentation sinusoïdale exploite la pério-

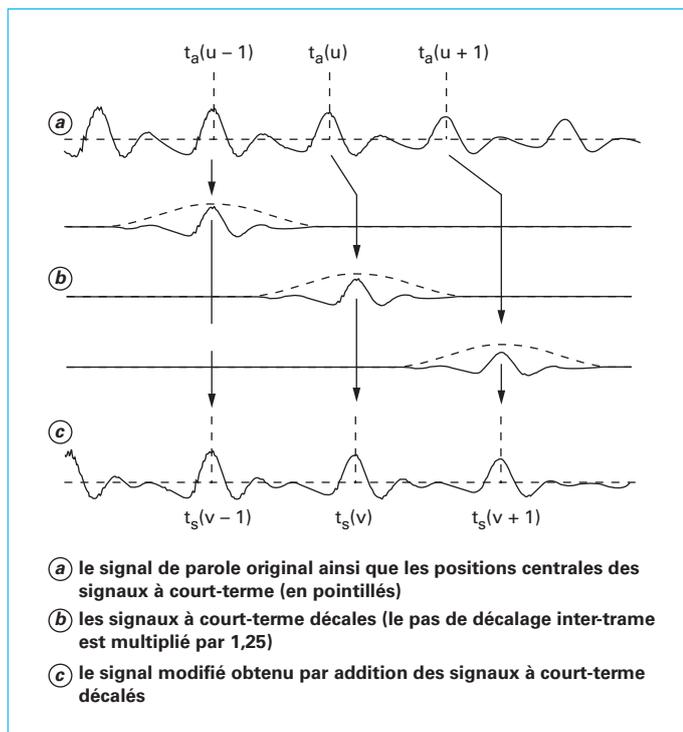


Figure 5 - Modification de la fréquence fondamentale par un facteur 0,8 avec la méthode PSOLA

dicité de la parole dans le domaine fréquentiel plutôt que dans le domaine temporel comme PSOLA.

## 5.4 Synthèse par sélection et concaténation

Pour améliorer le naturel de la voix de synthèse, les systèmes de dernière génération utilisent des **unités de taille variable, plus grandes que les diphones** : l'idée est d'enregistrer un corpus de parole de grande taille et d'aller y puiser l'ensemble optimal de segments ou unités de synthèse. Ces unités peuvent être au choix, des segments de phrase, des mots ou des fragments de mots, des syllabes, des diphones, ou même des sons isolés. Dans un tel système, il existe ainsi de nombreuses possibilités pour le choix des segments pour une même chaîne phonétique. Ces approches sont ainsi couramment dénommées : synthèse par sélection dynamique d'unités acoustiques, ou synthèse par sélection/concaténation.

■ La **synthèse par sélection/concaténation est une extension de la synthèse par unités** concaténées, comme les diphones, ou la synthèse à parties manquantes. Le principe est simple, sans grande théorie, mais avec un souci d'efficacité : à partir d'une grande base de données de parole, contenant une ou plusieurs heures de signal, il s'agit de rechercher les plus longs segments contigus de diphones, de demi-phones, ou de segments plus petits, qui correspondent à la phrase à synthétiser. En ce sens, c'est une extension de la synthèse à parties manquantes, puisque l'on utilisera si possible des mots, voire des membres de phrase entiers. Mais comme le vocabulaire est illimité, il faut s'appuyer sur des unités comme les diphones, afin de compléter les énoncés si l'on ne parvient pas à trouver des tronçons de signal longs.

Étant donné une phrase, la première étape consiste comme pour les autres types de synthèse à transcrire son contenu phonétique, ainsi que des informations sur la constitution des énoncés : ponc-

1 tuation, position des mots dans la phrase, des syllabes dans les  
 2 mots, des syllabes accentuées par exemple. Ces informations vont  
 3 permettre de rechercher à la fois des suites de phonèmes, mais  
 4 aussi des segments qui partagent certaines propriétés prosodiques  
 5 avec la phrase source.

6 ■ La sélection de la suite optimale de segments dans la base est  
 7 effectuée à l'aide de fonctions de coût. En général, deux fonctions  
 8 de coût, le « coût de cible » et le « coût de concaténation » sont  
 9 utilisées :

10  
 11 – **coût de cible** : il mesure et pondère l'avantage associé aux  
 12 différents segments possibles, en termes de longueur maximale  
 13 des segments, et de critères tels que la place du segment sélectionné  
 14 dans la syllabe, le mot ou la phrase afin, d'obtenir une voix  
 15 de synthèse avec une prosodie plus naturelle ;

16 – **coût de concaténation** : il mesure (en les pondérant) les quantités  
 17 acoustiques (distorsion de concaténation, fréquence  
 18 fondamentale moyenne) associées à la concaténation des unités.

19 Les unités optimales sont sélectionnées par une procédure  
 20 d'optimisation des coûts, souvent avec un algorithme de programmation  
 21 dynamique (ou un algorithme de Viterbi). La figure 6  
 22 illustre cette procédure sur un exemple simple pour la synthèse du  
 23 mot livre.

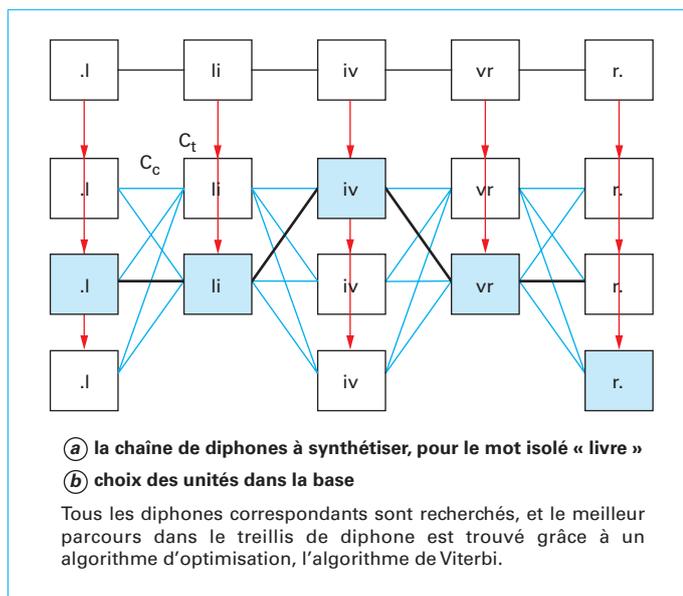


Figure 6 – Exemple de sélection des unités non uniformes  
 (algorithme type Viterbi)

Un exemple plus complet est donné dans le tableau 2, le texte à synthétiser est « La chasse aux papillons », et l'algorithme de synthèse des segments de longueur variable puisés dans différentes phrases de la base de données de très grande taille.

■ Dans ce type de synthèse, la **prosodie n'est pas calculée explicitement**. Les segments sélectionnés sont longs, en général au moins des syllabes, parfois des mots ou des groupes de mots (les locutions courantes se retrouvent souvent entières : « bonjour », « c'est-à-dire » etc.). Ils sont porteurs de leur propre prosodie, qui est réutilisée directement, sans passer par un modèle.

L'algorithme de sélection va ainsi non seulement reconstituer la chaîne de phonèmes, mais aussi la structure prosodique (groupes de syllabes, début, milieu, fin de phrases, ponctuation, début, milieu, fin des mots, etc.). Cette prosodie implicite est en général

Tableau 2 – Exemple de sélection

« La chasse aux papillons. »	
Phrases de la base de données	Segments de synthèse concaténés
/ .la Sas o papijO. /	
... d'autre part, le ...	/ .l /
... <b>la charte</b> ...	/ la Sα /
... <b>face aux</b> voisins ...	/ as o /
... participe <b>au premier</b> ...	/ o p /
... le <b>papy</b> boom	/ papi /
... publié ...	/ ij /
... <b>rayon</b> chaud ...	/ jO /
... 1 000 <b>chansons.</b>	/ O. /

Les segments de synthèse concaténés sont ceux extraits (texte en gras) des différentes phrases de la base de données.

d'un naturel étonnant, d'autant plus qu'elle offre beaucoup plus de variétés et moins de stéréotypes que la prosodie calculée par règles.

■ La plupart des systèmes de synthèse commerciaux actuels sont basés sur la synthèse par concaténation. La voix est souvent très naturelle et intelligible. Ce type de synthèse éprouve cependant des limites. C'est une synthèse assez coûteuse en termes de place mémoire, et la construction de bases de données de qualité demande beaucoup de soins et de temps. Le contrôle sur la voix synthétique est presque nul : c'est une méthode strictement non paramétrique, pour laquelle la qualité obtenue dépend uniquement du contenu de la base. Il n'est pas possible par exemple de changer la prosodie pour faire passer un contenu expressif, ou de changer la qualité de voix. Pour faire cela, il faut enregistrer et étiqueter une nouvelle base sonore. Cette synthèse est peu flexible : elle restitue bien ce qui est dans la base mais ne permet guère d'aller au-delà. Enfin, la qualité de synthèse peut être très variable d'une phrase à l'autre : dans certains cas parfaits, si les bons segments dans les bons contextes sont présents dans la base, dans d'autres cas très défavorables, si des combinaisons de segments manquent (ce qui est rare), ou si les contextes sont absents (ce qui est plus fréquent).

## 5.5 Synthèse paramétrique statistique

Une technique plus récente et en plein essor actuellement vise à combiner la flexibilité de la synthèse paramétrique et la qualité de la synthèse utilisant des gros corpus de parole : la synthèse paramétrique statistique.

■ Ce type de synthèse tire également avantage des dizaines d'années de recherche en reconnaissance de parole par des méthodes statistiques, et des outils disponibles dans ce cadre. La synthèse paramétrique statistique utilise le cadre général des modèles de Markov cachés (HMM pour *Hidden Markov Models*). Les travaux de Markov, au siècle dernier, ont montré que la succession des lettres dans les romans obéit à des lois de probabilité particulières. C'est l'origine des chaînes de Markov, ou automates probabilistes, particulièrement bien adaptées à la prédiction des enchaînements d'unités linguistiques, à tous les niveaux (phonèmes, syllabes, mots, suites de mots).

Contrairement aux approches par concaténation d'éléments sonores, la synthèse paramétrique par modèles statistiques repose

1 sur un modèle de signal et générera un nouveau signal à travers  
 2 un modèle de synthèse, pouvant être vu, en première approxima-  
 3 tion, comme le signal le plus probable, une moyenne d'un  
 4 ensemble de signaux de parole similaires.

5 Tout comme la synthèse par règles, la synthèse paramétrique  
 6 statistique repose sur un **modèle paramétrique du signal**, un  
 7 modèle source-filtre. Mais, contrairement à la synthèse par règles,  
 8 les paramètres de ce modèle ne sont pas analysés et définis expli-  
 9 citemment par un expert, et de façon déterministe. Ils sont appris  
 10 grâce à des modèles statistiques, sur un gros corpus de parole,  
 11 sous forme de chaînes de Markov, et générés de façon probabiliste  
 12 (la succession la plus probable de paramètres est émise) étant  
 13 donné un texte d'entrée à synthétiser. Ce principe est illustré sur la  
 14 figure 7.

15  
 16 ■ La **construction d'un système de synthèse** consiste, tout comme  
 17 en synthèse par sélection-concaténation, à enregistrer un gros  
 18 corpus de parole, et à l'étiqueter. Les étiquettes utilisées sont du  
 19 même type composite que pour la synthèse par concaténation :  
 20 phonème, phonèmes adjacents à droite et à gauche, syllabe,  
 21 position du phonème dans la syllabe, de la syllabe dans le mot,  
 22 mots adjacents, ponctuation, catégories syntaxiques, etc. Tout  
 23 type d'étiquette que l'on peut apposer sur des segments de façon  
 24 régulière est susceptible d'un apprentissage statistique, par  
 25 comptage des occurrences.

26  
 27 Le signal de parole enregistré est analysé par un modèle  
 28 source-filtre, avec d'un côté les paramètres du filtre, sous forme en  
 29 général non pas de formants mais de paramètres spectraux plus éla-  
 30 borés, et de l'autre les paramètres de la source, c'est-à-dire les para-  
 31 mètres prosodiques (intonation, durée, qualité vocale, intensité).

32 Les paramètres spectraux utilisés pour représenter le filtre sont  
 33 ceux utilisés en reconnaissance de parole. En effet, les synthé-  
 34 tiseurs par HMM sont directement issus des outils robustes et  
 35 sophistiqués développés pour la reconnaissance depuis plusieurs  
 36 décennies (par exemple, les MFCC, *Mel Frequency Cepstral Coeffi-*  
 37 *cients*, leurs dérivées et dérivées secondes). Ces paramètres se  
 38 prêtent bien à l'apprentissage automatique et leur extraction est  
 39 très robuste, contrairement aux formants par exemple.

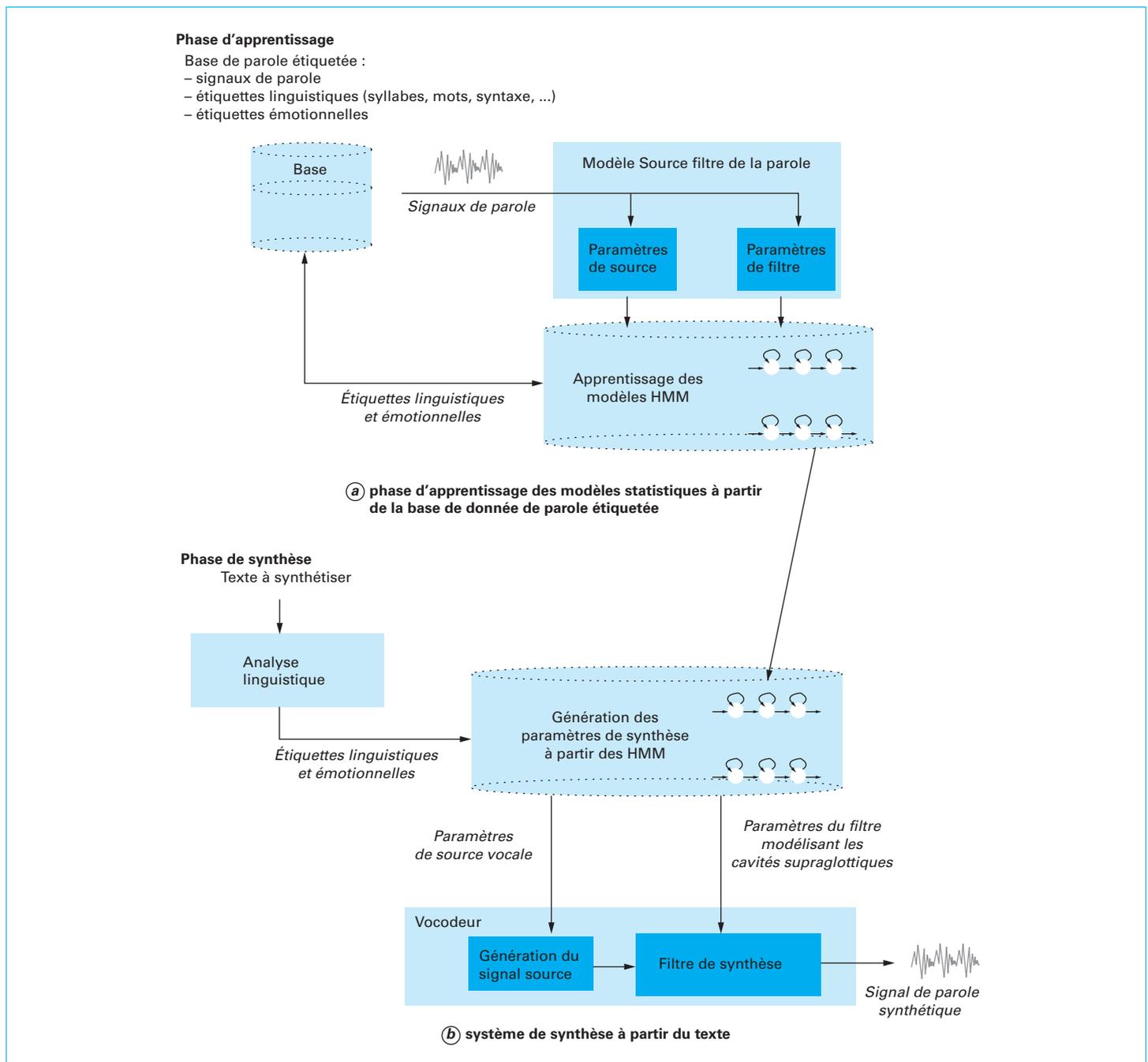
40  
 41 Ainsi, lors de la phase d'apprentissage, des modèles probabi-  
 42 listes représentant la base de données sonore sont construits. La  
 43 base de données elle-même n'est pas conservée dans le système  
 44 de synthèse, mais seulement sa représentation, ainsi l'empreinte  
 45 mémoire du système est réduite.

46  
 47 ■ Au **moment de la synthèse**, comme pour les autres modes de  
 48 synthèse, la phrase d'entrée est analysée par le module linguisti-  
 49 que et enrichie d'étiquettes linguistiques nombreuses. Ces éti-  
 50 quettes sont transmises aux modèles de Markov cachés  
 51 (sous-jacents) construits lors de l'apprentissage, qui génèrent les  
 52 paramètres les plus probables.

53 Par exemple, les modèles de phonème génèrent les paramètres  
 54 spectraux les plus probables du phonème en fonction de son iden-  
 55 tité et de son contexte (phonèmes voisins, syllabe, mot, etc.). Le  
 56 modèle d'intonation génère des points cibles pour les phonèmes  
 57 voisins en fonction des étiquettes d'entrée. Les durées sont  
 58 générées par d'autres modèles statistiques appris sur le corpus.

59 Lorsque les paramètres ont été générés, ils pilotent un vocodeur  
 60 de synthèse qui génère le signal acoustique proprement dit.  
 61 Plusieurs types de vocodeurs différents existent dans les différents  
 62 systèmes.

63  
 64 Cette approche est très prometteuse et commence à rivaliser en  
 65 termes de qualité avec les systèmes par sélection/concaténation  
 66 même si elle ne donne pas encore lieu, à ce jour, à des systèmes  
 67 commerciaux d'une qualité comparable à ceux par concaténation.  
 68 En effet, l'utilisation d'un vocodeur entraîne toujours une dégra-  
 69 dation de qualité des échantillons de parole. Par contre, la  
 70 flexibilité totale et les progrès des modèles statistiques, grâce à la  
 71 reconnaissance de parole rendent ces systèmes très attractifs.



57  
58 **Figure 7 – Schéma de principe de la synthèse paramétrique statistique**  
59

## 60 61 **5.6 Construction du corpus textuel et sonore** 62 63

64  
65 Pour la synthèse par sélection/concaténation, comme pour la synthèse paramétrique statistique, **la taille de la base de données est nécessairement importante**. Au moins une heure de parole (à comparer avec les 2 ou 3 min d'une base de diphones), parfois beaucoup plus, semble incontournable pour couvrir l'ensemble des combinaisons possibles de phonèmes et de contextes prosodiques.  
66  
67  
68  
69  
70  
71

1 La construction de la base de données est un sujet de recherche  
2 en soi. En effet, il faut qu'elle soit d'une taille optimale, c'est-à-dire  
3 que la plupart des segments qu'elle contient soient potentiel-  
4 lement utilisés, mais aussi qu'elle soit aussi complète que  
5 possible.

6 La base doit bien entendu contenir tous les phonèmes, tous les  
7 diphtongues, mais aussi des contextes plus longs, si possible toutes  
8 les syllabes. En plus de ces éléments phonétiques, l'algorithme de  
9 sélection utilise dans le calcul des coûts des éléments syntaxiques  
10 ou prosodiques, comme la place de la syllabe dans le mot, la place  
11 du mot dans la phrase, la ponctuation, la catégorie du mot, la  
12 longueur de la phrase, la durée de la syllabe, la valeur de  
13 fréquence fondamentale, etc.

14 Ainsi, la segmentation et l'étiquetage de cette base est réalisée à  
15 l'aide d'outils de reconnaissance vocale qui permettent d'obtenir  
16 des marques de frontières entre les différentes unités de base  
17 (certains systèmes utilisent des unités de base très petites comme  
18 le demi-phonème). Des outils d'analyse linguistique et d'analyse du  
19 signal complètent la construction de la base. Cette automatisation  
20 permet de construire des nouvelles voix de synthèse assez  
21 simplement puisqu'il suffit pour cela d'enregistrer un nouveau  
22 locuteur (ou nouvelle locutrice).

23 ■ La **qualité de synthèse** obtenue dépend principalement de celle  
24 du corpus sonore enregistré. La construction de ce corpus est donc  
25 fondamentale. Ces méthodes permettent d'ailleurs de reconstituer  
26 en synthèse de parole la voix d'un locuteur pour lequel on  
27 possède beaucoup d'enregistrements, comme celle d'un homme  
28 politique, d'un acteur, voire sa propre voix en s'enregistrant  
29 soi-même. La qualité de la synthèse dépend également du  
30 locuteur : même si exactement les mêmes outils et les mêmes  
31 phrases pour la base de données ont été utilisés, deux locuteurs  
32 différents ne donneront pas des systèmes de même qualité. Cela  
33 dépend de la netteté d'articulation, de la constance de prononciation,  
34 de la résistance à la fatigue vocale pour de longues séances  
35 d'enregistrement, et de la sonorité même de la voix.

## 40 6. Applications 41 de la synthèse de parole

### 46 6.1 Exemples d'applications

47 De nombreuses applications commerciales intègrent des  
48 systèmes de synthèse de parole. À l'heure actuelle, le marché prin-  
49 cipal de ce type de technique est celui des services de télécommu-  
50 nications. Ces services constituent l'exemple typique de situations  
51 dans lesquelles la synthèse de parole est le seul moyen par lequel  
52 un système informatique peut transmettre des informations à ses  
53 utilisateurs. Parmi les applications de la synthèse à partir du texte  
54 dans le domaine des services de télécommunications, citons :

55 - les services de réservation ou de prise de commandes  
56 téléphoniques ;

57 - les services d'information téléphonique pour lesquels le  
58 recours à la synthèse de parole se justifie, surtout lorsque l'infor-  
59 mation est amenée à évoluer vite, ce qui est notamment le cas  
60 pour les services bancaires (avec la fourniture, entre autres de  
61 l'état des comptes), les annonces météorologiques et routières, la  
62 lecture de méls ou de pages Internet. La synthèse de parole est  
63 aussi utilisée dans des contextes où le nombre des réponses  
64 potentielles du système est très important comme dans les appli-  
65 cations de renseignements téléphoniques ;

66 - les majordomes, assistants personnels, pour les téléphones  
67 mobiles ou autres terminaux, qui peuvent lire des messages reçus  
68 ou des courriers électroniques ;

69 - une application ambitieuse est envisagée à l'heure actuelle  
70 avec la téléphonie interprétée qui devrait permettre à deux corres-  
71

1 pondants ne parlant pas la même langue de dialoguer par télé-  
2 phone. Cette application fait intervenir plusieurs des grandes  
3 problématiques du traitement de la parole – reconnaissance, syn-  
4 thèse –, et bien sûr traduction automatique.

5 La synthèse de parole est aussi couramment employée dans des  
6 situations où l'utilisateur d'un système informatique n'a pas le  
7 loisir de consulter un écran, ou bien en complément de l'écran  
8 (cabine de pilotage d'un avion, systèmes industriels de fabrication,  
9 appareillage médical, etc.). Dans ce type d'applications, le rôle de  
10 la synthèse de parole consiste principalement à faire passer des  
11 informations brèves comme les messages d'erreurs du système.  
12 Les applications dans les systèmes d'information, fixes ou mobiles  
13 sont également nombreuses :

- 14 – portail vocaux d'application libre service ou de sites Internet ;
- 15 – systèmes de navigations ;
- 16 – systèmes de renseignement ;
- 17 – accessibilité des services ;
- 18 – vocalisation de journaux et de livres électroniques ;
- 19 – lecteurs d'écran ;
- 20 – jouets, robots et autres systèmes embarqués ;
- 21 – jeux vidéo ;
- 22 – jeux sérieux, éducation, *edutainment*.

24 La qualité accrue des systèmes de synthèse permet maintenant  
25 de développer des applications d'apprentissage des langues  
26 étrangères qui deviendront les évolutions naturelles des appli-  
27 cations actuelles de dictionnaire électronique de poche avec leur  
28 capacité à synthétiser des mots ou des phrases dans plusieurs  
29 langues.

30 Un autre aspect important des applications de la synthèse de  
31 parole à partir du texte concerne les services pour personnes han-  
32 dicapées. Dans ce domaine, le couplage de la synthèse de parole  
33 avec les techniques de reconnaissance automatique de caractères  
34 a permis la mise au point de véritables « machines à lire » pour les  
35 mal ou non-voyants.  
36

## 37 6.2 Interfaces de programmation

38 La synthèse de parole est intégrée dans les systèmes d'exploit-  
39 ations, les bibliothèques de logiciels ou les services sous la forme  
40 d'interface de programmation (API : *Application Programming*  
41 *Interface*), comme par exemple le *Speech API* (SAPI) de Microsoft.

42 Le *World Wide Web Consortium* (W3C) recommande le  
43 protocole SSML (*Speech Synthesis Markup Language*), protocole  
44 basé sur XML, afin d'annoter ou d'enrichir la synthèse avec des  
45 marqueurs prosodiques comme la fréquence fondamentale, le  
46 débit, les pauses, le volume, etc.  
47

## 48 6.3 Produits

49 L'offre en produits commerciaux de synthèse est aujourd'hui  
50 répandue. La plupart de ces systèmes sont multilingues,  
51 c'est-à-dire capables de produire des voix de synthèse dans  
52 plusieurs langues différentes. Ces systèmes incluent tous la syn-  
53 thèse de l'anglais (généralement, américain), et de la plupart des  
54 langues européennes et des grandes langues orientales.

55 Les configurations logicielles diffèrent suivant le type de produits  
56 et les applications. La plupart du temps cependant, l'obtention  
57 d'une voix de synthèse ne nécessite plus un matériel spécifique (si  
58 ce n'est une carte de restitution du son, disponible en standard sur à  
59 peu près toutes les plates-formes), la synthèse proprement dite ne  
60 requérant en fait qu'une fraction de la puissance de calcul d'un pro-  
61 cesseur moderne. Pour certaines applications spécifiques (serveurs  
62 vocaux ou applications embarquées), des implantations matérielles  
63 sont encore souvent nécessaires.

64 Comme pour beaucoup de produits dans le domaine des tech-  
65 nologies de l'information et de la communication, les offres  
66

1 évoluent très vite, les compagnies fusionnent ou sont rachetées, et  
2 le paysage est plutôt instable, à part pour les structures de grande  
3 taille.

4 Le paysage actuel est caractérisé par la concentration de l'offre  
5 commerciale sur peu d'acteurs. On note aussi la disparition des  
6 compagnies de téléphonie dans l'offre de produit, au profit de  
7 petites *startups* qui en sont issues, et des compagnies de logiciel.  
8 Nous donnons ci-dessous une liste non exhaustive des acteurs du  
9 domaine et fournissons pour certains d'entre eux une description  
10 plus détaillée de leur offre.  
11

12 **Nota** : se reporter à la partie documentation de cet article pour une présentation plus  
13 détaillée des acteurs commerciaux des produits de synthèse.  
14

## 15 7. Évaluation de la synthèse

### 16 7.1 Boîte noire ou boîte de verre

17 L'évaluation de la synthèse est nécessaire pour comparer les  
18 systèmes entre eux et pour mesurer les progrès réalisés. L'éva-  
19 luation de la parole de synthèse peut être globale (le signal  
20 généré) ou porter sur un des aspects du système.  
21

22 L'évaluation globale est celle en général du client, qui écoute et  
23 évalue le synthétiseur sur la parole produite. L'évaluation des  
24 composantes intéresse plus le chercheur ou le développeur, qui  
25 cherche à améliorer un aspect du système.  
26

27 La synthèse à partir du texte est en effet une chaîne de  
28 traitements, depuis le texte jusqu'au signal acoustique. C'est donc  
29 le maillon le plus faible de la chaîne qui en limitera la qualité, d'où  
30 l'importance d'évaluer chaque module (évaluation analytique ou  
31 interne, ou « boîte de verre »), en plus de l'évaluation globale (éva-  
32 luation externe, ou « boîte noire »). En général, la référence de  
33 l'évaluation des systèmes est la parole naturelle, ou de la parole  
34 naturelle avec un certain niveau de bruit.  
35

36 L'évaluation de la parole peut porter sur les aspects suivants :

#### 37 1. Évaluation globale :

- 38 a. intelligibilité (netteté à l'aide de syllabes, phrases dépourvues  
39 de sens, parole téléphonique, etc.),
- 40 b. qualité globale, suivant plusieurs critères (test d'opinion  
41 moyenne : *Mean Opinion Score* ou MOS) ;

#### 42 2. Évaluation analytique :

- 43 a. transcription graphème-phonème (taux de transcription  
44 correct, en particulier pour les noms propres),
- 45 b. prosodie (expressivité, agrément, etc.),
- 46 c. vocodeur - synthétiseur acoustique (qualité du signal,  
47 artefacts).  
48

### 49 7.2 Évaluation de qualité globale

50 Pour la qualité globale, une procédure multidimensionnelle  
51 d'évaluation a été normalisée en 1994 par l'Union internationale  
52 des télécommunications (UIT-T P.88). Les échantillons de parole  
53 doivent durer de 10 à 30 secondes, et il est recommandé d'utiliser  
54 une référence de parole naturelle (éventuellement dégradée).  
55

56 Les sujets écoutent chaque échantillon deux fois, pour une  
57 durée d'environ une heure (par exemple quatre stimuli pour quatre  
58 systèmes et trois références), en incluant les instructions et  
59 l'apprentissage.  
60

61 Les huit dimensions d'analyse sont :

- 62 1. Acceptabilité (pensez-vous que cette voix convienne pour tel  
63 service ?) ; 1 : oui, 2 : non.  
64

2. Impression d'ensemble (comment évaluez-vous la qualité de ce que vous venez d'entendre ?) ; 1 : excellente, 2 : bonne, 3 : correcte, 4 : faible, 5 : mauvaise.

3. Effort d'écoute (comment décririez-vous l'effort nécessaire pour comprendre le message ?) ; 1 : relaxation complète, aucun effort nécessaire, 2 : attention nécessaire, pas d'effort notable, 3 : effort modéré nécessaire, 4 : effort nécessaire, 5 : rien de compréhensible, quel que soit l'effort.

4. Compréhension (avez-vous trouvé des mots difficiles à comprendre ?) ; 1 : jamais, 2 : rarement, 3 : parfois, 4 : souvent, 5 : constamment.

5. Articulation (peut-on distinguer les sons avec netteté ?) ; 1 : oui, très net, 2 : oui, assez net, 3 : relativement net, 4 : non, pas très net, 5 : non, pas du tout net.

6. Prononciation (avez-vous remarqué des anomalies de prononciation ?) ; 1 : non, 2 : oui, mais pas gênantes, 3 : oui, gênantes, 4 : oui, très gênantes.

7. Débit de parole (que pensez-vous de la vitesse moyenne de prononciation ?) ; 1 : beaucoup trop rapide, 2 : trop rapide, 3 : correcte, 4 : trop lente, 5 : beaucoup trop lente.

8. Agrément de la voix (comment décririez-vous cette voix ?) ; 1 : très agréable, 2 : agréable, 3 : correcte, 4 : déplaisante, 5 : très déplaisante.

Le test d'opinion moyenne MOS est le plus souvent utilisé pour évaluer la synthèse. Ce test utilise généralement une échelle de cinq points, entre « très mauvaise qualité » et « excellente qualité ». La parole naturelle obtient généralement des scores de l'ordre de 4,5, et la parole de synthèse actuellement des scores de l'ordre de 3,5 pour les meilleurs systèmes, dès lors que les phrases sont longues.

## 8. Conclusion

### 8.1 Bilan

Bien que la liste des applications actuelles des systèmes de synthèse de parole soit assez conséquente, il serait faux de croire que celle-ci constitue une technique entièrement maîtrisée. Cependant, les travaux de recherche menés depuis des années ont permis d'atteindre une qualité qui se rapproche de celle de la voix naturelle.

La **qualité de synthèse est un problème crucial**. Il se manifeste essentiellement par le fait que la compréhension de la parole synthétique exige, de la part de l'auditeur, un effort plus important que pour la parole naturelle. Cet effort supplémentaire est rendu nécessaire par les artefacts éventuels du traitement, les erreurs de prononciation, la prosodie insuffisamment expressive, ou de manière plus générale, par le manque de naturel de l'élocution. Actuellement, si la qualité des systèmes de synthèse par sélection d'unités non uniformes est suffisante pour de nombreuses applications (lecture de méls, de messages d'information météo ou de navigation...), elle reste encore trop faible pour permettre une utilisation « prolongée » de la parole de synthèse ou pour des tâches véritablement expressive (lecture de livres par exemple).

Un exemple du type de problèmes qui restent associés à la synthèse de parole est le pourcentage important de prononciation incorrecte pour les noms propres. Ce problème n'est pas totalement inconnu dans le cas de la parole naturelle (il est même familier pour les enseignants confrontés à l'épreuve de l'appel en début d'année) ; toutefois, un locuteur humain est capable d'éliminer une proportion importante des erreurs potentielles en faisant appel à ses connaissances culturelles (notamment à celles qui concernent l'origine géographique du nom).

## 8.2 Perspectives

Les travaux en cours s'attachent ainsi à améliorer la flexibilité des systèmes suivant plusieurs axes :

- variabilité de la voix de synthèse au cours du temps, en fonction de l'énoncé ;
- possibilités d'expressivité accrue (attitudes, émotions, personnages particuliers) ;
- création rapide de nouvelles voix de synthèse.

La décennie passée a vu le développement de nouvelles techniques de synthèse (par **méthodes statistiques paramétriques**) proches de celles utilisées en reconnaissance de la parole. Ces techniques ne sont pas encore passées dans les produits mais sont très près de franchir ce pas.

Les offres de produit se sont diversifiées, avec en particulier le déploiement d'applications sur toutes les plates-formes, du serveur distribué au mobile, la possibilité de « studios » de synthèse pour préparer des messages d'une qualité équivalente à celle d'un locuteur enregistré, le développement d'assez nombreuses langues.

Le principal enjeu est aujourd'hui de rendre la synthèse interactive, capable de répondre et d'évoluer avec son environnement, capable d'exprimer des affects sociaux ou des émotions, et d'employer toute la gamme expressive de la voix humaine. Pour cela, des progrès sur la modélisation acoustique aussi bien que sur la compréhension des textes seront nécessaires.