



**HAL**  
open science

## Utilisation de la linguistique en reconnaissance de la parole : un état de l'art

Stéphane Huet, Pascale Sébillot, Guillaume Gravier

► **To cite this version:**

Stéphane Huet, Pascale Sébillot, Guillaume Gravier. Utilisation de la linguistique en reconnaissance de la parole : un état de l'art. [Rapport de recherche] RR-5917, INRIA. 2006, pp.72. inria-00077386v2

**HAL Id: inria-00077386**

**<https://inria.hal.science/inria-00077386v2>**

Submitted on 1 Jun 2006

**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.



INSTITUT NATIONAL DE RECHERCHE EN INFORMATIQUE ET EN AUTOMATIQUE

*Utilisation de la linguistique en reconnaissance de la  
parole : un état de l'art*

Stéphane Huet, Pascale Sébillot et Guillaume Gravier

N° 5917

Mai 2006

Thèmes SYM et COG

*R*apport  
de recherche





## Utilisation de la linguistique en reconnaissance de la parole : un état de l'art

Stéphane Huet, Pascale Sébillot et Guillaume Gravier \*

Thèmes SYM et COG — Systèmes symboliques et Systèmes cognitifs  
Projets TexMex et Metiss

Rapport de recherche n° 5917 — Mai 2006 — 72 pages

**Résumé :** Pour transcrire des documents sonores, les systèmes de reconnaissance de la parole font appel à des méthodes statistiques, notamment aux chaînes de Markov cachées et aux modèles N-grammes. Même si ces techniques se sont révélées performantes, elles approchent du maximum de leurs possibilités avec la mise à disposition de corpus de taille suffisante et il semble nécessaire, pour tenter d'aller au-delà des résultats actuels, d'utiliser des informations supplémentaires, en particulier liées au langage. Intégrer de telles connaissances linguistiques doit toutefois se faire en tenant compte des spécificités de l'oral (présence d'hésitations par exemple) et en étant robuste à d'éventuelles erreurs de reconnaissance de certains mots. Ce document présente un état de l'art des recherches de ce type, en évaluant l'impact de l'insertion des informations linguistiques sur la qualité de la transcription.

**Mots-clés :** reconnaissance de la parole, langue parlée, corpus oral, traitement automatique des langues, modèle de langage, connaissances linguistiques, disfluences

\* {shuet, sebillot, ggravier}@irisa.fr

# Using Linguistics in Speech Recognition: A State of the Art

**Abstract:** To transcribe speech, automatic speech recognition systems use statistical methods, particularly hidden Markov model and N-gram models. Although these techniques perform well and lead to efficient systems, they approach their maximum possibilities. It seems thus necessary, in order to outperform current results, to use additional information, especially bound to language. However, introducing such knowledge must be realized taking into account specificities of spoken language (hesitations for example) and being robust to possible misrecognized words. This document presents a state of the art of these researches, evaluating the impact of the insertion of linguistic information on the quality of the transcription.

**Key-words:** speech recognition, spoken language, spoken corpus, natural language processing, language model, linguistics knowledge, disfluencies

## Table des matières

<b>1</b>	<b>Introduction</b>	<b>4</b>
<b>2</b>	<b>Principes de la reconnaissance de la parole</b>	<b>6</b>
2.1	Difficultés de la transcription . . . . .	6
2.2	Modélisation statistique de la reconnaissance de la parole . . . . .	8
2.2.1	Extraction de caractéristiques . . . . .	9
2.2.2	Modèle acoustique . . . . .	11
2.2.3	Modèle de langage . . . . .	13
2.3	Sorties des systèmes de transcription . . . . .	18
2.4	Méthodes d'évaluation . . . . .	20
<b>3</b>	<b>Caractéristiques de la langue parlée</b>	<b>23</b>
3.1	Langue parlée et langue écrite . . . . .	23
3.2	Vocabulaire et syntaxe . . . . .	25
3.3	Phénomènes d'hésitation . . . . .	26
3.4	Corpus oraux . . . . .	28
3.4.1	Formes des corpus oraux . . . . .	29
3.4.2	Étiquetage des corpus oraux . . . . .	30
3.5	Transcription automatique des dialogues spontanés . . . . .	32
<b>4</b>	<b>La linguistique pour la reconnaissance de la parole</b>	<b>34</b>
4.1	Adaptation du processus de transcription . . . . .	35
4.1.1	Intégration du modèle acoustique avec le modèle de langage . . . . .	35
4.1.2	Linguistique et modèle de langage . . . . .	37
4.2	Quelles connaissances linguistiques? . . . . .	38
4.2.1	Phonologie et phonétique . . . . .	39
4.2.2	Morphologie . . . . .	40
4.2.3	Syntaxe . . . . .	41
4.2.4	Sémantique . . . . .	51
4.2.5	Pragmatique . . . . .	52
<b>5</b>	<b>Conclusion</b>	<b>56</b>
	<b>Références</b>	<b>58</b>

## 1 Introduction

De nombreux documents sonores contiennent de la parole. Rares sont les émissions radiophoniques ou audiovisuelles sans discours, dialogues ou encore commentaires, et il existe maintenant certaines bases de données sonores très volumineuses. En France, l'Institut national de l'audiovisuel collecte ainsi chaque année 80 000 heures de programmes radiophoniques ou télévisuels. Une des plus grandes archives digitales au monde, contenant les témoignages des survivants de la Shoah, contient quant à elle plus de 115 000 heures de discours non contraints, provenant de 52 000 locuteurs s'exprimant dans 32 langues différentes [FRWP03]. Face à la taille des données à traiter, le recours à des méthodes automatiques, notamment à celles de reconnaissance de la parole, facilite la manipulation des documents et apporte une aide pour certaines tâches, telles que l'indexation.

L'objectif d'un système de reconnaissance automatique de la parole (RAP) est de transcrire automatiquement un signal sonore en texte. Un tel système cherche dans un premier temps à reconnaître des mots, en se basant uniquement sur des critères d'ordre acoustique, sans essayer d'interpréter le « contenu » transmis par l'ensemble de ces mots. L'analyse du signal conduit alors à un ensemble d'hypothèses sur la succession des mots prononcés, auquel sont adjoints des scores dits acoustiques. Une seconde étape choisit la meilleure hypothèse en ne considérant plus le signal comme une suite de sons mais plutôt comme une succession de mots porteurs d'information. L'utilisation de la linguistique, en tant que science du langage, s'inscrit naturellement dans ce contexte puisque des connaissances sur la morphologie, la syntaxe ou le sens semblent pouvoir guider la sélection de la meilleure hypothèse. Néanmoins, bien souvent les systèmes de RAP ne s'appuient dans leur choix que sur des critères principalement statistiques, et ce, essentiellement pour des raisons historiques.

Pendant une longue période, les linguistes et les adeptes du traitement automatique des langues (TAL) ont en effet délaissé les méthodes empiriques permettant d'estimer la probabilité d'observation d'un phénomène à partir d'une grande collection de documents ou corpus. Noam Chomsky dit ainsi en 1969 : *But it must be recognized that the notion "probability of a sentence" is an entirely useless one, under any known interpretation of this term*<sup>1</sup>. Le point de vue des linguistes s'intéressant à l'écrit s'est alors éloigné des préoccupations des chercheurs en reconnaissance de la parole, à tel point que Frederick Jelinek, en 1988, alors à IBM, déclara : *Anytime a linguist leaves the group the recognition rate goes up*<sup>2</sup>. Les méthodes statistiques, à base de N-grammes, se sont en effet révélées beaucoup plus efficaces que les solutions proposées par les linguistes pour choisir la meilleure hypothèse de mots.

Depuis la fin des années 80, les méthodes statistiques ont toutefois commencé à atteindre leurs limites dans l'amélioration des systèmes de RAP, notamment avec la mise à disposition de corpus de tailles satisfaisantes. Dans le même temps, le domaine du TAL faisait de plus en plus appel aux modèles probabilistes. Une des pistes envisageables pour améliorer les performances de la reconnaissance de la parole peut donc consister à employer davantage de linguistique dans les systèmes de RAP.

<sup>1</sup>Mais il doit être reconnu que la notion de « probabilité de phrase » est absolument inutile, et ce, quelle que soit l'interprétation de ce terme.

<sup>2</sup>À chaque fois qu'un linguiste quitte le groupe, le taux de reconnaissance augmente.

---

Ce document se propose de faire une synthèse de l'introduction de connaissances linguistiques au cours de la reconnaissance de la parole. Il expose des tentatives effectuées pour utiliser certaines informations telles que la morphologie, la syntaxe ou la sémantique dans les différentes étapes du processus de transcription. Une première partie décrit le fonctionnement général d'un système de RAP, en insistant sur ses limitations à traiter certains phénomènes. Elle présente la succession des étapes nécessaires, à savoir l'extraction d'informations numériques pertinentes à partir du son, la conversion de ces valeurs en plusieurs hypothèses possibles de succession de mots et enfin le choix de la meilleure hypothèse. La deuxième partie examine les propriétés de la langue parlée. Les méthodes de TAL sont en effet souvent appliquées à des documents qui restent dans le domaine de l'écrit, comme des ouvrages ou des articles de journaux ; or, les documents analysés par les systèmes de RAP sont d'une autre nature. La dernière section expose à quel niveau du processus de transcription les connaissances linguistiques peuvent être mobilisées. Elle se focalise particulièrement sur l'introduction de connaissances linguistiques au sein des modèles de langages, un des deux constituants, avec le modèle acoustique, d'un système de RAP. Si le modèle acoustique utilise des ressources purement acoustiques, se limitant donc, sur le plan de la linguistique, à la phonétique et à la phonologie, le modèle de langage peut au contraire prendre en compte des connaissances linguistiques plus variées puisque son rôle est justement d'examiner les informations véhiculées par les hypothèses de mots. Ceci peut donc conduire à l'exploitation de morphologie, de syntaxe, de sémantique ou encore de pragmatique.



## 2 Principes de la reconnaissance de la parole

L'objectif d'un système de RAP est d'extraire les mots prononcés à partir du signal acoustique. Le résultat produit représente la finalité d'une application de dictée vocale mais peut également être utilisé par d'autres dispositifs. La transcription constitue ainsi une source d'informations pour indexer des documents audio ou audiovisuels. Les interfaces vocales homme-machine intègrent quant à elles un module de compréhension de la parole en sus d'un système de RAP.

La reconnaissance de la parole est un problème complexe, notamment du fait de la grande variabilité des signaux à traiter. Après avoir présenté les principales difficultés de la transcription, nous exposons la modélisation qui est employée pour décoder le signal acoustique. Nous évoquons ensuite sous quelles formes les systèmes de RAP produisent leurs résultats et nous terminons cette section par une description des techniques utilisées pour évaluer ces systèmes.

### 2.1 Difficultés de la transcription

Les difficultés de la transcription de la parole sont dues pour une grande part à la diversité des signaux à traiter. La parole produite pour une même phrase prononcée peut ainsi varier d'un individu à un autre. Outre le fait que chaque individu possède une voix qui lui est propre, on rencontre d'importantes différences telles que les variations homme/femme, le régionalisme ou encore les difficultés de prononciation rencontrées par des locuteurs non natifs. Cette variabilité est qualifiée d'*inter-locuteurs*. Il existe également une variabilité, dite *intra-locuteur*, correspondant à une modification de la parole produite par un même individu. Cette variabilité peut concerner aussi bien les caractéristiques de la voix, dans les cas d'un rhume ou d'un état émotionnel, ou bien la qualité d'élocution, selon que la parole intervient lors d'un discours formel ou d'un dialogue spontané. En sus des variabilités au niveau de la parole prononcée par le locuteur, les conditions d'enregistrement peuvent dégrader le signal qui sera traité par le système de RAP. Il peut par exemple y avoir une modification de la qualité du signal, notamment si celui-ci doit transiter par un canal de communication qui a une bande passante limitée, comme une ligne téléphonique. De même, l'environnement acoustique peut être disparate. Le bruit de fond peut être plus ou moins important et de natures diverses (musique, paroles d'autres locuteurs, parasites du micro, bruits de bouche...).

Par ailleurs, le lexique des documents à transcrire est un autre facteur-clé influençant la qualité des résultats et dépendant de chaque application. La taille du vocabulaire peut être très réduite (moins de 100 mots) dans le cas d'un système de navigation dans un menu, moyenne (quelques milliers de mots) pour des recherches d'information dans une base de données dans un domaine précis, ou large (plusieurs dizaines de milliers de mots) pour faire de la dictée vocale [Sto97].

La parole humaine, on le voit donc, est très variable. La modélisation d'une telle variation étant difficile à faire de manière compacte et la compréhension des mécanismes cognitifs intervenant dans la reconnaissance de la parole étant limitée, les systèmes de RAP utilisent

essentiellement des méthodes statistiques. Ces méthodes extraient automatiquement les informations sur le langage et la relation entre le son et les mots prononcés à partir de corpus dans lesquels les textes sont alignés avec les signaux acoustiques. Elles utilisent ainsi des modèles dont les paramètres sont appris sur des corpus d'apprentissage [DGP99].

Pour réduire les difficultés impliquées par la variabilité inter-locuteurs, certains systèmes de RAP requièrent de la part de chaque locuteur une prononciation préalable d'un certain nombre de mots. Toutefois, ce procédé étant contraignant, les systèmes sont généralement indépendants du locuteur. Afin d'augmenter la *robustesse* vis-à-vis du locuteur et de l'environnement, ils combinent plusieurs modèles statistiques [Jou96, GAAD<sup>+</sup>05]. Ainsi, dans le cas de la variabilité inter-locuteurs, un modèle spécifique peut être créé pour les hommes, un autre pour les femmes. De même, des mécanismes d'adaptation permettent de modifier le traitement acoustique en fonction des caractéristiques de la voix du locuteur. Dans le cas de changements de conditions d'enregistrement, les systèmes de RAP peuvent avoir un modèle particulier pour les entretiens téléphoniques et avoir des détecteurs « bruit/parole » ou « musique/parole » pour ne chercher à transcrire que les segments du signal contenant de la parole. En outre, un filtrage adaptatif peut être mené pour éliminer le bruit du signal.

En ce qui concerne le lexique, afin d'avoir un espace de recherche raisonnable lors du décodage du signal acoustique, les systèmes de RAP ont un vocabulaire fermé. Ils utilisent actuellement un lexique beaucoup plus important que ceux employés auparavant, qui pouvaient se limiter à quelques dizaines de mots. Les systèmes de RAP les plus perfectionnés, dits à très grand vocabulaire, ont ainsi un lexique dont le volume avoisine les 65 000, voire 200 000 mots. L'ensemble des mots reconnaissables est choisi en fonction de l'application. Dans le cas de reconnaissance d'un dialogue spontané, il pourra ainsi être utile d'inclure dans le lexique des mots qui n'en sont pas vraiment mais qui ont pourtant une fréquence élevée, tels que le « *eah* » marquant une hésitation (*cf.* section 3.3). Si l'on prend l'exemple d'une transcription d'émissions d'actualité diffusées à la radio, le vocabulaire pourra inclure les mots rencontrés récemment. Ce domaine d'application se heurte à des difficultés particulières dues à l'apparition des noms propres (noms de personnes ou de lieu) au gré de l'actualité. De manière à pouvoir modifier le vocabulaire pris en compte par le système de RAP, la reconnaissance doit être *flexible*, ce qui signifie qu'elle doit autoriser l'introduction de mots dans le dictionnaire qui n'ont pas été utilisés au cours de l'apprentissage des paramètres des modèles.

Notons au passage l'ambiguïté du terme *mot* [Pol03]. Il peut par exemple désigner un sens précis ou bien un signe linguistique. En reconnaissance de la parole, *mot* désigne un *mot-forme* défini par son orthographe. Ainsi, deux flexions ou dérivations d'un même lemme, *e.g.* « *mange* » et « *manges* », seront considérées comme deux mots différents. De même, deux homographes appartenant à deux catégories différentes (*e.g.* « *mérite* [VERBE] » et « *mérite* [NOM] ») ou deux sens différents (*e.g.* « *avocat* [AUXILIAIRE DE JUSTICE] » et « *avocat* [FRUIT COMESTIBLE] ») ne représenteront pas le même mot [Jel97]. Par la suite, chaque emploi du terme *mot* désignera en réalité un *mot-forme*.

Malgré les difficultés rencontrées, les systèmes de RAP parviennent à décoder le signal acoustique avec d'assez bonnes performances. La transcription de mots prononcés de

manière isolée, par un locuteur unique, est une technologie ancienne et bien maîtrisée. Son champ d'application est néanmoins très restreint puisque le locuteur doit marquer une pause brève entre chaque mot. Les chercheurs et ingénieurs en parole ont par la suite développé des systèmes de transcription de parole continue, ce qui a augmenté considérablement la complexité du problème. Les systèmes actuels les plus performants reconnaissent toutefois sans erreur plus de 90 % des mots d'une émission d'actualité en anglais [Pal03] et plus de 88 % des mots pour une émission du même type en français [GAAD<sup>+</sup>05]. Dans des situations plus complexes à analyser, comme des conversations téléphoniques, où chacun des locuteurs n'a pas préparé son discours et est donc sujet à de nombreuses hésitations, les systèmes de RAP peuvent reconnaître sans erreur jusqu'à 80 % des mots [Pal03, GAL<sup>+</sup>04].

## 2.2 Modélisation statistique de la reconnaissance de la parole

Le signal sonore à étudier peut être interprété comme une version de la phrase prononcée qui serait passée par un canal de communication. Ce canal introduit du « bruit » dans la version originale. L'objectif d'un système de RAP est de modéliser le canal de manière à retrouver la phrase prononcée après décodage (Fig. 1). Ceci revient à chercher parmi un très grand nombre de phrases sources potentielles celle qui a la plus grande probabilité de générer la phrase « bruitée » [JM00]. Autrement dit, dans cette métaphore du *canal bruité*, un système de RAP cherche à trouver la séquence de mots la plus probable  $W$  parmi toutes les séquences d'un langage  $\mathcal{L}$ , étant donné le signal acoustique  $A$ .

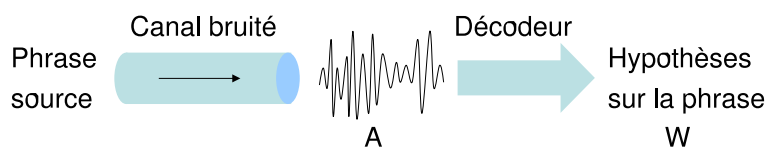


FIG. 1 – Modèle du canal bruité

L'entrée acoustique  $A$  représente une séquence d'observations  $a_1 \dots a_t$ , obtenue en découpant l'entrée par exemple toutes les 10 millisecondes et en associant à chaque morceau les coefficients représentant l'enveloppe spectrale du signal. La sortie  $W$  du système de RAP est une chaîne de mots  $w_1 \dots w_n$  appartenant à un vocabulaire fini.

En utilisant une formalisation statistique issue de la théorie de l'information, le problème de la reconnaissance de la parole se ramène alors à chercher :

$$\hat{W} = \arg \max_{W \in \mathcal{L}} P(W|A) \quad (1)$$

Cette équation peut être réécrite sous la forme suivante à l'aide de la formule de Bayes [Jel97] :

$$\hat{W} = \arg \max_{W \in \mathcal{L}} \frac{P(W)P(A|W)}{P(A)} \quad (2)$$

où  $P(W)$  est la probabilité que  $W$  soit prononcée,  $P(A|W)$  est la probabilité que le locuteur émette les sons  $A$  en souhaitant prononcer les mots  $W$  et  $P(A)$  est la probabilité moyenne que  $A$  soit produit.  $\hat{W}$  étant estimé en fixant  $A$ ,  $P(A)$  n'intervient pas et l'équation (2) devient :

$$\hat{W} = \arg \max_{W \in \mathcal{L}} P(W)P(A|W) \quad (3)$$

Le problème de la reconnaissance de la parole se ramène ainsi à l'extraction des indices acoustiques  $A$ , au calcul de la vraisemblance d'observation  $P(A|W)$  ainsi que de la probabilité *a priori*  $P(W)$ , et à la recherche de la séquence de mots  $W$  la plus probable. Pour ce faire, la transcription se décompose en plusieurs modules (Fig. 2) :

- l'extraction de caractéristiques produisant  $A$ ,
- l'utilisation du *modèle acoustique* (MA) calculant  $P(A|W)$  et cherchant les hypothèses  $W$  qui sont vraisemblablement associées à  $A$ ,
- l'utilisation du *modèle de langage* (ML) calculant  $P(W)$  pour choisir une ou plusieurs hypothèses sur  $W$  en fonction de connaissances sur la langue.

Pour évaluer  $P(W)$ , le ML doit disposer au préalable des hypothèses  $W$  établies par le MA sur les mots prononcés. Néanmoins, les systèmes de RAP actuels ne se limitent pas à une juxtaposition séquentielle des deux modules ; de manière à utiliser le plus tôt possible les informations sur la langue, ils font appel au ML dès qu'une hypothèse de mot est proposée par le MA et non à la fin du traitement de la totalité du signal.

Les sections suivantes décrivent les principes de fonctionnement de chacun de ces modules.

### 2.2.1 Extraction de caractéristiques

Le signal sonore à analyser se présente sous la forme d'une onde dont l'intensité varie au cours du temps. La première étape du processus de transcription consiste à extraire une succession de valeurs numériques suffisamment informatives sur le plan acoustique pour décoder le signal par la suite.

Le signal est susceptible de contenir des zones de silence, de bruit ou de musique. Ces zones sont tout d'abord éliminées afin de n'avoir que des portions du signal utiles à la transcription, *i.e.*, celles qui correspondent à de la parole. Le signal sonore est ensuite segmenté en ce que l'on qualifie de *groupes de souffle*, en utilisant comme délimiteurs des pauses silencieuses suffisamment longues (de l'ordre de 0,3 s). L'intérêt de cette segmentation est d'avoir un signal sonore continu de taille raisonnable par rapport aux capacités de calculs des modèles du système de RAP ; dans la suite du processus de transcription, l'analyse se fera séparément pour chaque groupe de souffle.

Pour repérer les fluctuations du signal sonore, qui varie généralement rapidement au cours du temps, le groupe de souffle est lui-même découpé en fenêtres d'étude de quelques millisecondes (habituellement de 20 ou 30 ms). De manière à ne pas perdre d'informations importantes se trouvant en début ou fin de fenêtres, on fait en sorte que celles-ci se chevauchent, ce qui conduit à extraire des caractéristiques toutes les 10 ms environ.

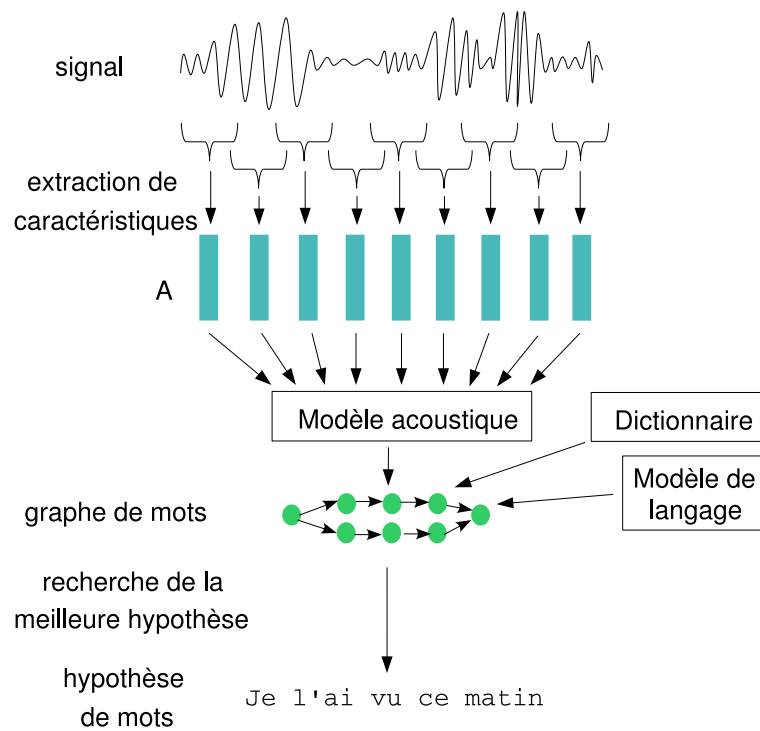


FIG. 2 – Constituants d'un système de transcription

À partir du signal contenu dans chaque fenêtre d'analyse sont calculées des valeurs numériques caractérisant la voix humaine. À l'issue de cette étape, le signal devient alors une succession de vecteurs dits acoustiques, de dimension souvent supérieure ou égale à 39.

### 2.2.2 Modèle acoustique

Une étape suivante consiste à associer aux vecteurs acoustiques, qui sont, comme nous venons de le voir, des vecteurs numériques, un ensemble d'hypothèses de mots, *i.e.*, des symboles. En se référant à l'équation (3) de la modélisation statistique, cela revient à estimer  $P(A|W)$ . Les techniques qui permettent de calculer cette valeur forment ce qu'on appelle le modèle acoustique.

L'outil le plus utilisé pour la modélisation du MA est la chaîne de Markov cachée (désignée aussi sous le terme de HMM pour *Hidden Markov Model*). Les HMM ont en effet montré dans la pratique leur efficacité pour reconnaître la parole. Même s'ils présentent quelques limitations pour modéliser certaines caractéristiques du signal, comme la durée ou la dépendance des observations acoustiques successives, les HMM offrent un cadre mathématique bien défini pour calculer les probabilités  $P(A|W)$  [Rab89]. Les MA font intervenir trois niveaux de HMM (Fig. 3).

Ils cherchent dans un premier temps à reconnaître les types de son, autrement dit à identifier des *phones*<sup>3</sup>. Pour ce faire, ils modélisent un phone par un HMM, généralement à trois états représentant ses début, milieu et fin. La variable cachée est alors un *sous-phone* et les observations sont des vecteurs acoustiques, *i.e.*, des vecteurs continus. Pour calculer les probabilités d'observation dans chaque état, deux approches sont souvent envisagées, l'une basée sur la représentation des densités de probabilité par des gaussiennes et l'autre reposant sur des réseaux de neurones. Ces différentes méthodes établissent des hypothèses sur la probabilité des phones prononcés. Or, l'objectif des MA est de déterminer une succession de mots. Les MA utilisent à cette fin un dictionnaire de prononciations, qui effectue la correspondance entre un mot et ses prononciations. Comme un mot est susceptible d'être prononcé de différentes manières, selon son prédécesseur et son successeur, ou tout simplement selon les habitudes du locuteur, il peut y avoir plusieurs entrées dans ce lexique pour un même mot. Les indications sont données au moyen des *phonèmes*<sup>4</sup> caractéristiques de la prononciation. Sur la figure 4, les phonèmes sont transcrits dans le système de représentation SAMPA.

Le deuxième niveau de HMM modélise les mots à partir des HMM représentant des phones et du lexique de prononciations. Il se présente sous la forme d'un arbre lexical contenant initialement tous les mots du vocabulaire, progressivement élagué au fur et à mesure que sont reconnus des phones. Puisque les HMM de premier niveau modélisent des phones et non des phonèmes, les phonèmes disponibles dans le dictionnaire de prononciations sont

<sup>3</sup>Sons prononcés par un locuteur et définis par des caractéristiques précises.

<sup>4</sup>Unité linguistique associée à un type de prononciation d'une langue donnée. Le phonème final /p/ pourra par exemple être prononcé en français par le phone [b] dans l'expression « *grippe aviaire* » et par le phone [p] dans « *grippe du poulet* ».

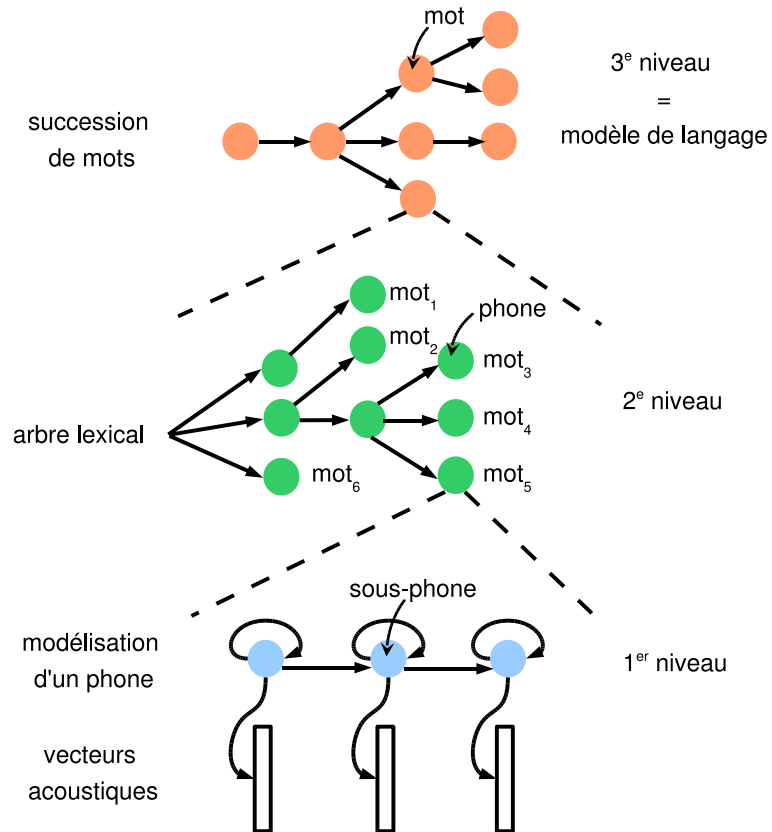


FIG. 3 – Niveaux de modélisation du modèle acoustique

adorateurs	a d O R a t 9 R z
adorateurs	a d O R a t 9 R
adoration	a d O R a s j o~
adore	a d O R @
adore	a d O R

FIG. 4 – Extrait d'un dictionnaire de prononciations

convertis en phones afin de reconnaître des mots. Des règles de transformation dépendant du contexte d'apparition du phonème sont alors utilisées.

Le troisième niveau modélise enfin la succession des mots  $W$  au sein d'un groupe de souffle et peut alors incorporer les connaissances apportées par le ML sur  $W$ . Pour établir ce HMM équivalent à un graphe de mots, le HMM correspondant à l'arbre lexical est dupliqué à chaque fois que le MA effectue l'hypothèse qu'un nouveau mot a été reconnu [ONA97].

Le fonctionnement du MA que nous venons de décrire se heurte à un problème majeur : l'espace de recherche du HMM de plus haut niveau devient fréquemment considérable, surtout si le vocabulaire est important et si le groupe de souffle à analyser contient plusieurs mots. Des algorithmes issus de la programmation dynamique permettent de calculer efficacement les probabilités ; il s'agit principalement de l'algorithme de Viterbi et le décodage par pile, appelé aussi décodage A\*. De plus, il est fait recours très régulièrement à l'élagage pour ne conserver que les hypothèses susceptibles d'être les plus intéressantes [DGP99].

Le rôle du MA consiste ainsi à aligner le signal sonore avec des hypothèses de mots en utilisant uniquement des indices d'ordre acoustique. Il inclut dans son dernier niveau de modélisation les informations sur les mots apportées par le ML.

### 2.2.3 Modèle de langage

Le ML a pour objectif de trouver les séquences de mots les plus probables, autrement dit celles qui maximisent la valeur  $P(W)$  de l'équation (3). Si l'on se réfère au HMM de plus haut niveau du MA (Fig. 3), les valeurs  $P(W)$  correspondent aux probabilités de succession de mots.

#### Fonctionnement d'un modèle de langage

En posant  $W = w_1^n = w_1 \dots w_n$ , où  $w_i$  est le mot de rang  $i$  de la séquence  $W$ , la probabilité  $P(W)$  se décompose de la manière suivante :

$$P(w_1^n) = P(w_1) \prod_{i=2}^n P(w_i | w_1 \dots w_{i-1}) \quad (4)$$

L'évaluation de  $P(W)$  se ramène alors au calcul des valeurs  $P(w_i)$  et  $P(w_i | w_1^{i-1})$  qui s'obtiennent respectivement à l'aide des égalités :

$$P(w_i) = \frac{C(w_i)}{\sum_{w \in \mathcal{V}} C(w)} \quad (5)$$

$$P(w_i | w_1^{i-1}) = \frac{C(w_1^i)}{\sum_{w_i} C(w_1^i)} \quad (6)$$

où  $\mathcal{V}$  est le vocabulaire utilisé par le système de RAP, et  $C(w_i)$  et  $C(w_1^i)$  représentent les nombres d'occurrences respectifs du mot  $w_i$  et de la séquence de mots  $w_1^i$  dans le corpus d'apprentissage.



Malheureusement, pour prédire la suite de mots  $w_1^n$ , le nombre des paramètres  $P(w_i)$  et  $P(w_i|w_1^{i-1})$  du ML à estimer augmente de manière exponentielle avec  $n$ . Dans le but de réduire ce nombre,  $P(w_i|w_1^{i-1})$  est modélisé par un *modèle N-gramme*, i.e., une chaîne de Markov d'ordre  $N - 1$  (avec  $N > 1$ ), à l'aide de l'équation suivante :

$$P(w_i|w_1^{i-1}) \approx P(w_i|w_{i-N+1}^{i-1}) \quad (7)$$

Cette équation indique que chaque mot  $w_i$  peut être prédit à partir des  $N-1$  mots précédents. Pour  $N = 2, 3$  ou  $4$ , on parle respectivement de modèle *bigramme*, *trigramme* ou *quadri-gramme*. Pour  $N = 1$ , le modèle est dit *unigramme* et revient à estimer  $P(w_i)$ . Généralement, ce sont les modèles bigrammes, trigrammes et quadrigrammes qui sont utilisés dans les ML des systèmes de RAP.

Cette approche rencontre des limites du fait de l'absence de nombreuses séquences de mots de taille  $N$ , appelées des *N-grammes*, dans les corpus d'apprentissage, bien que ceux-ci puissent être de taille conséquente. En effet, même en ayant une valeur de  $N$  réduite, de nombreux mots seront rares, voire absents du corpus. On dit à ce sujet que les mots suivent une loi de Zipf, stipulant que la fréquence d'apparition d'un mot décroît rapidement avec son rang d'apparition. Pour pallier cette difficulté, il est fait appel à des méthodes statistiques de *lissage*.

Un premier procédé de lissage, connu sous le nom de *discounting*, consiste à retrancher au comptage des N-grammes une certaine valeur qui sera en suite redistribuée vers le comptage des N-grammes absents du corpus d'apprentissage. Il existe de nombreuses méthodes de *discounting*. Un procédé très simple consiste par exemple à ajouter un à l'ensemble des comptages, y compris ceux associés aux N-grammes absents.

Un autre moyen d'effectuer le lissage est d'utiliser les fréquences d'apparition des N-grammes d'ordres inférieurs si le N-gramme étudié est peu présent dans le dictionnaire. On distingue alors la technique de l'*interpolation linéaire* de celle du *repli* (appelée aussi *backoff*).

L'interpolation linéaire consiste à évaluer la probabilité  $\hat{P}(w_i|w_{i-N+1}^{i-1})$  en faisant une combinaison linéaire des probabilités calculées pour les N-grammes d'ordre inférieur. Dans le cas d'un modèle trigramme par exemple, le calcul s'effectue à partir des probabilités unigramme, bigramme et trigramme de la manière suivante :

$$\hat{P}(w_i|w_{i-2}w_{i-1}) = \lambda_1 P(w_i|w_{i-2}w_{i-1}) + \lambda_2 P(w_i|w_{i-1}) + \lambda_3 P(w_i) \quad (8)$$

avec :

$$\sum_{k=1}^3 \lambda_k = 1 \quad (9)$$

de manière à ce que  $\hat{P}$  demeure une probabilité.  $P$  est calculée au moyen de l'équation (6) utilisant les comptages. Les valeurs  $\lambda_k$  sont estimées de façon à maximiser la vraisemblance de  $\hat{P}$  sur un corpus de test, différent du corpus d'apprentissage.

Le repli effectue lui aussi une combinaison linéaire avec les probabilités de N-grammes d'ordre inférieur mais, à la différence de l'interpolation linéaire, le recours aux N-grammes

de tailles plus réduites n'est pas systématique [Kat87]. Le calcul des probabilités s'exprime de la manière suivante :

$$\hat{P}(w_i|w_{i-N+1}^{i-1}) = \begin{cases} \tilde{P}(w_i|w_{i-N+1}^{i-1}) & \text{si } C(w_{i-N+1}^i) > 0 \\ \alpha(w_{i-N+1}^{i-1}) \times \hat{P}(w_i|w_{i-N+2}^{i-1}) & \text{sinon} \end{cases} \quad (10)$$

De même que pour les  $\lambda_i$  de l'équation (8), les coefficients  $\alpha$  sont calculés pour que  $\hat{P}$  soit une probabilité. Le symbole  $\tilde{\phantom{P}}$  sur le  $P$  sert à indiquer que  $\tilde{P}$  est souvent obtenu à partir d'un procédé de *discounting*.

Au-delà de ces techniques de lissage, il existe de nombreuses variantes basées sur leurs principes, parmi lesquelles figure la version modifiée du lissage de Kneser-Ney qui, d'après des études empiriques, donne de bons résultats [CG98].

Les modèles N-grammes, qu'ils utilisent ou non des techniques de lissage, sont des méthodes statistiques dont le nombre de paramètres à estimer est très grand, ce qui nécessite de disposer d'une grande quantité de données d'apprentissage. Depuis les débuts d'utilisation des N-grammes, de nombreux textes de natures différentes ont été collectés, ce qui a largement profité à l'amélioration de ces techniques. Toutefois, cette évolution suit depuis quelque temps une asymptote. Selon une estimation informelle d'IBM, les performances des modèles bigrammes n'enregistrent plus de gain important au-delà de quelques centaines de millions de mots, tandis que les modèles trigrammes semblent saturer à partir de quelques milliards de mots. Or, dans plusieurs domaines d'application de systèmes de RAP, de tels volumes de données ont déjà été collectés [Ros00b].

Pour tenter d'améliorer les performances des ML, des évolutions du mode de calcul des probabilités par les modèles N-grammes ont été envisagées. Le principal reproche qui est fait aux modèles N-grammes est l'hypothèse qui a permis l'élaboration des premiers ML performants et utilisables, à savoir la prise en compte d'un historique  $h_i$  de taille limitée. Le calcul des probabilités est en effet réalisé au moyen d'une égalité du type :

$$P(w_1^n) = P(w_1) \prod_{i=2}^n P(w_i|h_i) \quad (11)$$

Plusieurs études ont été menées pour examiner un historique plus étendu que  $w_{i-N+1}^{i-1}$  [Goo01].

### Modifications de l'historique dans le calcul des probabilités

La prise en compte d'un historique de très grande taille est susceptible d'améliorer les performances du ML. Toutefois, les modifications envisagées doivent tenir compte du fait que parmi l'ensemble des historiques possibles, beaucoup deviennent rares voire inexistants dans le corpus d'apprentissage quand on augmente le nombre de mots pris en compte.

Un moyen très simple pour étendre les N-grammes est d'utiliser l'interpolation linéaire (cf. équation (8)) pour combiner des modèles N-grammes d'ordres différents. Ceci permet d'avoir un historique étendu mais aussi de faire face à la rareté des données dans le corpus d'apprentissage. Ces modèles sont nommés modèles *polygrammes* [KNST94, GSTN96].

Cette méthode présente deux inconvénients. Elle nécessite tout d'abord l'évaluation de nombreux paramètres puisque les calculs de probabilité sont effectués à partir de plusieurs tailles d'historique. De plus, il n'est pas possible d'augmenter indéfiniment la taille de l'historique, même en utilisant des méthodes de lissage perfectionnées. On considère souvent à ce sujet que les modèles N-grammes d'ordre supérieur à 5 n'apportent pas de gain par rapport aux modèles d'ordres inférieurs [Goo01].

Une autre variation des modèles N-grammes consiste à ignorer certaines positions dans l'historique ; il s'agit alors de modèles *skipping*. L'interpolation linéaire de modèles 5-grammes *skipping* peut ainsi avoir la forme :

$$\lambda_1 P(w_i | w_{i-4} w_{i-3} w_{i-2} w_{i-1}) + \lambda_2 P(w_i | w_{i-4} w_{i-3} w_{i-1}) + \lambda_3 P(w_i | w_{i-4} w_{i-2} w_{i-1}) \quad (12)$$

Ces types de modèles permettent de prendre en compte des historiques qui sont proches mais non strictement identiques au contexte courant. Cette propriété est particulièrement importante quand on augmente la taille des N-grammes car il devient alors de plus en plus rare de trouver deux historiques identiques [Goo01].

Les modèles dits *permugrammes* adoptent une approche similaire en ce qui concerne les historiques, en effectuant une permutation des mots pris dans le contexte. Ils redéfinissent le calcul des probabilités par l'équation :

$$P(W) = P(w_{\pi(1)}) \prod_{i=2}^n P(w_{\pi(i)} | w_{\pi(i-N+1)} \dots w_{\pi(i-1)}) \quad (13)$$

où  $\pi$  est une permutation qui réordonne la succession des mots. Cette équation conduit ainsi à des historiques de la forme  $h_i = w_{i-3}, w_{i+1}$  [STHKN95]. L'utilisation de ces modèles peut parfois apporter une légère amélioration des performances des N-grammes. Ils augmentent toutefois de manière importante l'espace de recherche ainsi que le nombre de paramètres à calculer.

D'autres études ont considéré l'historique non pas comme une succession de mots mais comme une suite de groupes de mots. Cette approche s'appuie sur une segmentation en séquences de mots, en fixant une taille maximum  $M$  pour ces séquences. Le calcul de la probabilité  $P(w_1 w_2 w_3)$  devient par exemple, avec  $M = 3$  et en prenant un historique de longueur maximale 2 :

$$P(w_1 w_2 w_3) = \max \left\{ \begin{array}{l} P([w_1 w_2 w_3]) \\ P([w_1])P([w_2 w_3] | [w_1]) \\ P([w_1 w_2])P([w_3] | [w_1 w_2]) \\ P([w_1])P([w_2] | [w_1])P([w_3] | [w_1][w_2]) \end{array} \right\} \quad (14)$$

Il existe une variante de ce calcul, remplaçant *max* par la somme de toutes les segmentations possibles. Ces modèles, qualifiés de *multigrammes*, n'ont pas permis pour l'instant d'améliorer véritablement les performances par rapport aux ML classiques [BPLA95, DB95, AB05].

Un autre type d'extension des N-grammes repose sur un historique à longueur variable, ce qui permet de faire des distinctions supplémentaires dans certains cas ambigus, tout en conservant un nombre raisonnable de paramètres à estimer. Les N-grammes pris en compte dans le calcul des probabilités sont alors présentés sous la forme d'un arbre (Fig. 5). Dans le cas des modèles N-grammes classiques, où la taille de l'historique est fixe, cet arbre a une profondeur fixe égale à  $N - 1$ . Dans le cas des modèles *varigrammes* au contraire, la profondeur varie selon les branches puisque des nœuds proches dans l'arbre et associés à des probabilités conditionnelles similaires sont fusionnés [SO00, Kne96, NW96b]. La conception

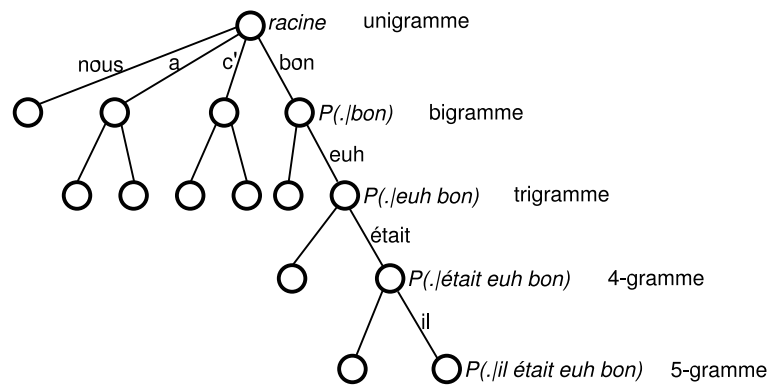


FIG. 5 – Représentation sous forme d'arbre d'un modèle varigramme

des modèles varigrammes conduit à une diminution de la taille des ML, sans modifier leurs performances. Les modèles dits *x-grammes* poursuivent également le même objectif, *i.e.*, la fusion des historiques très proches. Leur particularité est de modéliser les ML par des automates à états finis.

Enfin, on peut citer l'utilisation des modèles à base de *cache*. Ces modèles conduisent à l'utilisation d'historiques de l'ordre de plusieurs centaines de mots, *i.e.*, de taille beaucoup plus importante que les modèles N-grammes. Ils sont basés sur le principe que si un locuteur utilise un mot, la probabilité qu'il l'utilisera à nouveau dans un futur très proche augmente considérablement. Ils redéfinissent alors le calcul des probabilités  $P(W)$  par l'équation :

$$P(W) = \tilde{P}(w_1) \prod_{i=2}^n \tilde{P}(w_i | w_{i-N+1} \dots w_{i-1}) \quad (15)$$

où les probabilités  $\tilde{P}$  sont évaluées à partir des comptages des mots dans le *cache* et non plus dans un corpus d'apprentissage comme cela est le cas pour les modèles N-grammes. Les modèles à base de *cache* sont systématiquement combinés avec des modèles N-grammes. Ils peuvent contribuer à améliorer les performances des ML, mais ce, au détriment de la vitesse de calcul des probabilités lors de la transcription [Goo01, KDM90, CR97].

---

REF:	il AURA ALORS face à lui une fronde syndicale *** UNIE	
HYP:	il **** VALEUR face à lui une fronde syndicale EST PUNI	
	D      S	I      S

---

FIG. 6 – Alignement de la transcription automatique (HYP) et de la transcription de référence (REF)

On le voit donc, de nombreuses tentatives ont été faites pour essayer de saisir les dépendances à longue distance, *i.e.*, les relations entre les mots qui sont séparés par plus de  $N - 1$  mots. Les techniques qui viennent d’être présentées permettent souvent d’améliorer légèrement les performances par rapport à des modèles trigrammes. Toutefois, les progrès sont bien souvent peu significatifs, alors que leur utilisation engendre un accroissement de la complexité des calculs lors de la reconnaissance ou une augmentation importante du nombre de paramètres à estimer [Goo01].

Cette section a présenté les principes des ML, la dernière étape intervenant dans l’association du signal acoustique à une succession de mots. Il nous reste à voir sous quelles formes se présentent les sorties du processus de transcription de la parole.

### 2.3 Sorties des systèmes de transcription

Comme dit précédemment, le rôle d’un système de RAP est de produire une transcription d’un signal sonore. Le résultat pourra donc être naturellement un texte. Toutefois, selon le cadre d’utilisation d’un tel système, il existe d’autres types de sortie envisageables.

**Texte brut** Lorsqu’un système de RAP produit un texte, celui-ci correspond à la succession de mots qui a obtenu la plus haute probabilité de la part du système. Ce texte est organisé sous la forme de groupes de souffle, le décodage de la parole se basant sur la détection de pauses silencieuses. Cette forme purement textuelle laisse envisager la possibilité d’appliquer des techniques de TAL pour améliorer les résultats produits.

Dans le cas où l’on dispose d’une transcription de référence, obtenue généralement manuellement, celle-ci peut être alignée avec le texte produit par le système de RAP (Fig. 6) pour effectuer des calculs de performance (*cf.* section 2.4). L’alignement est obtenu en faisant correspondre un maximum de mots des deux transcriptions au moyen d’un algorithme de type Viterbi. Trois types d’erreur sont alors distingués :

- les insertions, repérées par des I, correspondant à un ajout d’un mot de la transcription automatique par rapport à la transcription de référence,
- les suppressions, marquées par des D (pour *deletion*), associées à un mot manquant dans la transcription automatique,
- les substitutions, représentées par des S, indiquant un remplacement d’un mot de la transcription de référence par un autre mot présent dans la transcription automatique.

Lorsque des traitements doivent être opérés après le processus de transcription, on choisit parfois de conserver non pas la meilleure hypothèse, mais plutôt la liste des  $N$  meilleures hypothèses. Ceci laisse la possibilité de réordonner les hypothèses, si on dispose de connaissances supplémentaires par la suite.

**Graphes de mots** La liste des  $N$  meilleures hypothèses présente de nombreuses redondances, les éléments de la liste différant bien souvent par un seul mot. Une sortie plus compacte qui lui est souvent préférée est le graphe de mots (Fig. 7). Ce graphe peut être une variante du HMM de plus haut niveau du MA (*cf.* section 2.2.2). Les arcs sont valués par les probabilités établies par le MA et le ML, et représentent des hypothèses ou sur les mots prononcés ou sur la présence de pauses silencieuses (notées par « *sil* » sur la figure 7). Les nœuds sont quant à eux associés aux instants possibles où un mot se termine et un autre débute. Puisque les informations concernant l'instant de prononciation des mots sont peu employées, le graphe est souvent compacté en supprimant les arcs et les nœuds qui représentent les mêmes hypothèses de succession de mots [DGP99]. De surcroît, la taille du graphe de mots construit pouvant être considérable, il peut être utile de l'élaguer en fixant par exemple un nombre maximum de nœuds. Un autre critère possible d'élagage est le nombre  $k$  d'hypothèses retenues sur les successions de mots.

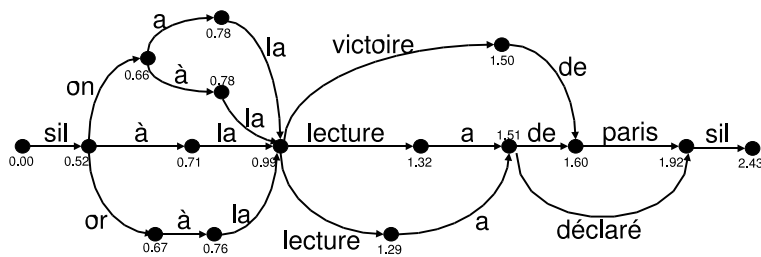


FIG. 7 – Exemple de graphe de mots (non valué)

**Réseaux de confusion** Cet autre type de sortie correspond à un compactage des graphes de mots, en conservant davantage d'informations que les  $k$  meilleures hypothèses. La construction des réseaux de confusion, également appelés « saussisses » (Fig. 8), consiste à aligner les hypothèses de succession de mots, un peu comme on le ferait si l'on souhaitait aligner la transcription automatique avec la transcription de référence (Fig. 6). L'extraction des meilleures hypothèses se fait ici en tentant de minimiser les erreurs d'insertion, de suppression ou de substitution, et non pas en déterminant la succession de mots associée aux plus faibles probabilités *a posteriori*, comme dans le cas des graphes de mots. Ceci est une propriété intéressante des réseaux de confusion dans la mesure où les systèmes de RAP

cherchent plutôt à obtenir une transcription ayant le minimum d'erreurs qu'un texte ayant un bon « score » de probabilité [MBS00].

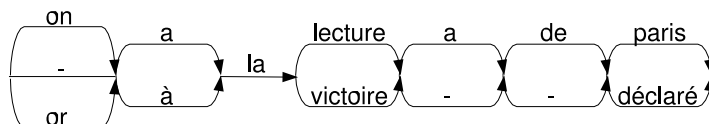


FIG. 8 – Exemple de réseau de confusion

À l'issue de la présentation des différents types de sorties de la transcription, la section suivante s'intéresse à la description de méthodes de mesure de la qualité des résultats produits.

## 2.4 Méthodes d'évaluation

La mesure qui est communément employée pour mesurer la qualité d'une transcription est le *taux d'erreur sur les mots* défini par :

$$\text{taux d'erreur} = \frac{\text{nb d'insertions} + \text{nb de suppressions} + \text{nb de substitutions}}{\text{nb de mots dans la transcription de référence}} \quad (16)$$

Dans ce calcul revenant à déterminer la distance d'édition entre la transcription automatique et la transcription de référence, les coûts d'insertion, de suppression ou de substitution ont une valeur fixe qui ne dépend pas du nombre de caractères erronés dans la transcription automatique. La transcription étant envisagée comme une succession de mots-formes (*cf.* section 2.1), ceci implique qu'une simple erreur d'accord en genre et en nombre aura le même coût qu'une substitution par un mot qui n'a aucun rapport sur le plan acoustique ou sémantique avec le mot correct. On voit ainsi que le taux d'erreur, même s'il est très utilisé, n'indique pas directement si le résultat produit est correct du point de vue des informations transmises. Bien entendu, s'il devient faible, la probabilité qu'il s'est produit une erreur sur les mots informatifs diminue mais ce n'est pas systématique.

Le taux d'erreur sur les mots permet d'évaluer globalement le système de RAP. Il existe d'autres mesures, telles que l'*entropie croisée* ou la *perplexité*, pour mesurer la qualité du ML seul [JM00]. Pour définir la première, précisons pour débiter la notion d'*entropie*.

L'entropie, notée  $H$ , est une mesure d'information qui se calcule pour une séquence de mots  $w_1^n$  de la manière suivante :

$$H(w_1^n) = - \sum_{w_1^n \in \mathcal{L}} P(w_1^n) \log_2 P(w_1^n) \quad (17)$$

Pour un langage  $\mathcal{L}$ , elle s'obtient par :

$$H(\mathcal{L}) = \lim_{n \rightarrow \infty} \frac{1}{n} H(w_1^n) \quad (18)$$

$$= \lim_{n \rightarrow \infty} -\frac{1}{n} \sum_{w_1^n \in \mathcal{L}} P(w_1^n) \log_2 P(w_1^n) \quad (19)$$

La mesure d'*entropie croisée*, directement dérivée de l'entropie et souvent abusivement appelée de manière identique, permet d'évaluer la performance d'un ML. Si l'estimation de la probabilité par un ML est notée  $\hat{P}$  et si  $P$  est la distribution réelle de probabilité, l'entropie croisée de ce modèle se calcule par :

$$H(P, \hat{P}) = \lim_{n \rightarrow \infty} -\frac{1}{n} \sum_{w_1^n \in \mathcal{L}} P(w_1^n) \log_2 \hat{P}(w_1^n) \quad (20)$$

En supposant que le langage  $\mathcal{L}$  possède de bonnes propriétés, le théorème de Shannon-McMillan-Breiman permet d'écrire :

$$H(P, \hat{P}) = \lim_{n \rightarrow \infty} -\frac{1}{n} \log_2 \hat{P}(w_1^n) \quad (21)$$

En pratique, pour comparer la performance de deux ML, on compare les entropies croisées calculées pour ces deux modèles sur un corpus de test  $T$  aussi grand que possible, de manière à approximer au mieux la limite vers l'infini :

$$H_T(P, \hat{P}) = -\frac{1}{t} \log_2 \hat{P}(T) \quad (22)$$

où  $t$  représente le nombre de mots du corpus de test. Cette équation est très proche de l'objectif de la modélisation statistique du ML (*cf.* équation (3)), *i.e.*, trouver  $\hat{P}$  maximisant  $\hat{P}(\mathcal{L})$ . Plus l'entropie croisée est faible, meilleur est ainsi le modèle.

Dans la plupart des études évaluant les ML, la *perplexité*  $PP$  est souvent préférée à l'entropie croisée. Elle se calcule sur un ensemble de test  $T$  à partir de :

$$PP_T(\hat{P}) = 2^{H_T(P, \hat{P})} \quad (23)$$

Bien que très utilisé pour la comparaison des ML, ce critère possède des limites. Une baisse importante de la perplexité ne correspond pas forcément à une baisse de même ordre du taux d'erreur sur les mots, si on intègre le ML dans un système de RAP ; seules des réductions de la perplexité d'au moins 10-20% semblent être vraiment significatives [Ros00b]. Son principal inconvénient est qu'il favorise les modèles accordant une plus grande probabilité aux mots présents dans le corpus de test, et ignore la manière avec laquelle sont distribuées les probabilités aux autres mots. Or, ces mots peuvent conduire, lors du décodage par le système de RAP, à des probabilités plus élevées que celles obtenues pour les mots corrects [CR99]. Dans certaines situations, on peut ainsi observer une augmentation du taux d'erreur en même temps qu'une diminution de la perplexité. Certains auteurs préconisent donc d'utiliser plutôt l'entropie croisée, qui semble plus corrélée avec le taux d'erreur [Goo01].



Cette section 2 dédiée aux principes de fonctionnement des systèmes de RAP nous a donné l'occasion d'évoquer les limitations des ML les plus populaires, *i.e.*, les modèles N-grammes. Avant de voir en détail comment il est possible d'intégrer davantage de connaissances linguistiques pour améliorer les performances de la transcription, la section suivante se focalise sur les spécificités de la langue parlée.

### 3 Caractéristiques de la langue parlée

Les documents traités par les systèmes de RAP sont par définition prononcés par des locuteurs et relèvent donc du domaine de la langue parlée. Ces types de documents ont malheureusement été peu étudiés par les linguistes et les adeptes du TAL, comparativement aux textes écrits. Cette constatation trouve sa source dans le mode de pensée de la culture occidentale, qui établit très souvent que l'écrit, considéré comme prestigieux, est seul digne d'intérêt. En outre, peu de corpus oraux existent et ceux qui sont disponibles sont souvent de taille peu importante. Les plus volumineux d'entre eux pour le français ne possèdent ainsi que quelques millions d'occurrences. Ceci a pour résultat que les systèmes de RAP actuels effectuent l'apprentissage des ML sur un corpus oral de taille réduite, accompagné d'un corpus écrit beaucoup plus volumineux, au risque de paramétrer les modèles avec la langue utilisée dans le journal *Le Monde* pour le français ou celui du *Wall Street Journal* pour l'anglais [Vér04].

Deux questions se posent alors. Peut-on modéliser la langue parlée à l'aide de corpus de la langue écrite, disponibles en plus grands volumes que les corpus oraux ? Quel est le comportement du TAL et du processus de transcription par rapport aux phénomènes spécifiques de la langue parlée ? Pour répondre à ces questions, nous précisons tout d'abord ce qui définit la langue parlée par rapport à la langue écrite. Nous présentons ensuite les principales caractéristiques de cette langue parlée en ce qui concerne son vocabulaire et sa syntaxe, avant de décrire les perturbations provoquées par les phénomènes d'hésitation très fréquents qu'elle contient. Nous abordons ensuite la constitution des corpus oraux, en montrant brièvement sous quelles formes ils se présentent et comment ils peuvent être annotés. Nous terminons cette section en étudiant le comportement des systèmes de RAP vis-à-vis des dialogues spontanés, qui sont la forme de l'oral se distinguant le plus des formes conventionnelles de l'écrit.

#### 3.1 Langue parlée et langue écrite

La *langue parlée* est définie comme étant ce qui est prononcé par des locuteurs à l'oral et fait donc appel à la voix et à l'oreille. La *langue écrite* représente quant à elle ce qui s'écrit et implique donc l'usage de la main et des yeux. La langue parlée et la langue écrite exploitent ainsi des canaux différents qui ont des contraintes importantes sur leur mode de production [Mel00].

Le canal oral impose une certaine linéarité lors son émission et de son écoute, même si on peut parfois avoir recours à des dispositifs d'enregistrement permettant de faire des retours en arrière. L'oral, produit en continu, exclut en général toute forme de préparation préalable, de planification. La configuration typique de l'oral est caractérisée par une réception immédiate par un ou plusieurs interlocuteurs, qui ont la possibilité de réagir au cours même de la production. Ceci peut obliger le locuteur à adapter son discours en fonction de ces réactions.

Le canal écrit au contraire exploite la page et fait intervenir deux dimensions : la largeur et la hauteur, facilitant ainsi les retours à des éléments antérieurs. Le texte est généralement produit en différé et peut donc être formulé et reformulé avant d'être livré aux destinataires.

Les récepteurs ne sont pas au contact du scripteur et celui-ci se trouve en général contraint de proposer un texte désambiguïsé, informatif et structuré.

Il existe en outre une opposition entre l'oral et l'écrit au niveau de la segmentation des productions. La *prosodie*, qui recouvre des informations sur l'intonation, l'accent, les pauses et même le débit, joue un rôle important dans la segmentation de la langue orale, sans toutefois permettre une structuration aussi importante que ce qu'on peut attendre pour les textes écrits. Il n'existe en effet pas à l'oral de démarcation aussi nette que celle permise par les signes de ponctuation ; on ne peut vraiment d'ailleurs y parler de phrase (*cf.* section 3.4.1). Les procédés de mise en page de l'écrit, à savoir les alinéas, les paragraphes ou encore les sections, autorisent au contraire une structuration de documents de taille importante.

Les modes de production des langues parlée et écrite, que nous venons d'évoquer, font que l'écrit est stable et la parole instable. Cette caractéristique a ainsi conduit la culture occidentale à valoriser l'écrit sur l'oral, même si finalement l'écrit peut être vu uniquement comme un simple système de codage de la langue parlée au moyen de signes visibles. Alors que les lois, les contrats, les textes religieux fondamentaux sont des documents écrits, l'oral ne nécessite pas un apprentissage à l'école aussi formel que l'écrit. Avec l'invention de l'imprimerie, qui a grandement favorisé la diffusion des textes, le code écrit a pris de plus en plus d'importance. L'apparition des ordinateurs et d'Internet a fait croître de manière exponentielle la production et la diffusion de documents écrits, mais favorise également l'apparition de nouvelles formes de documents, comme les courriels ou les blogs. Dans le même temps, les médias ont permis la diffusion massive de discours télévisés ou bien encore d'émissions radiophoniques. L'opposition entre langue parlée et langue écrite est ainsi loin d'être aussi claire que celle présentée au début de cette section et on assiste à une hybridation des codes [Pol03]. Si l'on peut facilement opposer les dialogues aux romans littéraires, on ne trouve pas toujours de contrastes aussi forts entre la langue parlée et la langue écrite, notamment en ce qui concerne l'interactivité et la préparation des documents. Le discours télévisé du chef de l'état à la nation sera par exemple préparé et n'autorise pas l'intervention des récepteurs. Le courrier électronique autorise quant à lui davantage d'interactions avec les destinataires et se rapproche plus du dialogue oral.

La langue qui est employée dans les documents écrits ou oraux peut subir d'importantes variations selon les régionalismes, les périodes considérées, les groupes sociaux et culturels ou encore les registres [Mel00]. Dans le domaine de l'oral, la linguiste Blanche-Benveniste définit six genres majeurs : les conversations face à face, les conversations par téléphone, les débats et entrevues en public, les émissions de radio ou de télévision, les discours non préparés et enfin les discours préparés [BB97].

Les langues parlée et écrite correspondent donc à deux codes différents puisqu'elles n'utilisent pas le même canal de communication. On ne peut ainsi considérer l'oral comme une simple forme dégradée de l'écrit. De même, bien que certaines écritures soient phonographiques et codent en partie les sons, comme cela est le cas du français, l'écrit est loin d'être une seule transcription de l'oral, ne serait-ce que parce que le scripteur choisit la manière de présenter les informations [GT04]. Les différences entre l'oral et l'écrit ne sont

toutefois pas suffisantes pour qu'il faille les étudier de manière totalement indépendante [BB90].

Les sections suivantes se proposent de dégager des phénomènes revenant plus fréquemment à l'oral qu'à l'écrit.

### 3.2 Vocabulaire et syntaxe

Même s'il est incorrect de dire que l'oral est une forme relâchée de l'écrit, le registre de la langue parlée est généralement moins soutenu dans la mesure où elle est souvent employée dans des cadres moins formels. Du point de vue lexical, la langue parlée se fonde sur un ensemble familier de mots, plus restreint que celui de la langue écrite. Une étude comparant deux corpus de 20 millions d'occurrences de mots chacun, le premier correspondant à des journaux écrits et le second à des émissions journalistiques radiophoniques et télévisées transcrites manuellement, a ainsi obtenu 127 000 mots distincts pour l'oral et 215 000 pour l'écrit, et ce, en considérant une même année de production. Toujours selon cette même étude, l'utilisation des catégories grammaticales suit des répartitions différentes dans les deux corpus. Un taux plus important de pronoms et un taux moins important de noms sont ainsi observés à l'oral [GAD02]. On y observe aussi l'introduction de « petits » mots, notamment appelés ligateurs, marqueurs de discours ou encore inserts, comme par exemple « *quoi* », « *bon* », « *donc* », « *enfin* » ou « *genre* ».

Au niveau de la syntaxe, il existe des tournures propres à l'oral. Si en français le « *pas* » est facultatif à l'écrit dans la négation « *ne ... pas* », c'est le « *ne* » qui le devient à l'oral. Le « *il* » prend également à l'oral un caractère facultatif dans les formules « *il faut* » et « *il y a* ». On peut aussi citer l'invariabilité du « *c'est* » dans des expressions telles que « *c'est les voisins qui sonnaient* ». Ces expressions sont dues au contexte d'utilisation de la langue orale qui favorise l'interactivité et impose donc un temps limité de formulation des idées. Pour la même raison, et même si cela est loin d'être systématique, les accords du participe passé seront moins respectés à l'oral qu'à l'écrit [Mel00].

Il semble également que l'ordre des mots soit un peu plus souple dans la langue parlée que dans l'écrite, bien que les langues rigides sur l'ordre des mots comme le français ou l'anglais le soient encore à l'oral [AG01]. Les éléments régis par le verbe peuvent, à l'oral, se placer avant le groupe sujet-verbe, comme dans l'exemple « *les haricots j'aime pas* ». On observe aussi davantage de clivées (« *c'est le coiffeur qui est content* »), de pseudos-clivées (« *ce qui l'intéresse c'est le pognon* »), de doubles marquages (« *moi j'en ai jamais vu en Suisse des immeubles* ») ou encore de dislocations (« *j'ai choisi la bleue de robe* ») [BB90].

Ces phénomènes ne justifient pas toutefois de proposer une grammaire spéciale pour la langue parlée. Pour le français, Blanche-Benveniste indique ainsi que la syntaxe de l'oral ne diffère en rien de celle de l'écrit, sauf en termes de proportions [BB90, BCD<sup>+</sup>04]. Bien que cela soit plus rare qu'à l'oral, on observe également à l'écrit des clivées ou même des erreurs d'accord de participe passé.

Cette section a montré les principales différences que l'on pouvait observer entre l'oral et l'écrit du point de vue du lexique et de la syntaxe. Toutefois, elle a occulté les perturbations

de la syntaxe par les phénomènes d'hésitation, plus nombreux dans la langue parlée. La section suivante étudie l'influence de ces phénomènes.

### 3.3 Phénomènes d'hésitation

Les marques d'hésitation ne sont pas propres à l'oral puisqu'elles se manifestent également dans les nouvelles formes de communication écrite comme les courriels, les forums, les chats ou encore les SMS [Ben04]. Elles sont également présentes dans les brouillons des textes écrits sous la forme de ratures. Néanmoins, elles sont avant tout caractéristiques de l'oral spontané. Il est ainsi estimé qu'elles représentent environ 5 % des mots dans les corpus spontanés et moins de 0,5 % des mots dans les corpus de parole lue ou préparée (*cf.* section 3.4). Elles ne peuvent donc être ignorées lors de la transcription.

Les phénomènes d'hésitation correspondent généralement à un travail de formulation de la part du locuteur, dans la mesure où le discours est composé au fur et à mesure de sa production. Ils peuvent cependant avoir d'autres rôles tels que la manifestation d'un doute, l'envoi d'un signal au récepteur pour garder la parole ou au contraire la lui céder, le marquage de frontières dans le discours, ou bien encore l'expression d'un stress ou d'une décontraction.

Comme les marques d'hésitation introduisent une rupture dans la continuité du discours, ou plus exactement dans le déroulement syntagmatique, elles sont également appelées *disfluences*. Il existe bien d'autres termes utilisés dans la littérature pour les nommer, certains ne désignant pas exactement les mêmes phénomènes : *turbulences*, *faux départs*, *lapsus*, *inattendus structurels*, *spontanéités*, *modes de production de la langue parlée*, *marques de réparation*, *marques du travail de formulation*, *extragrammaticalités*, *etc.* Une telle profusion de termes témoigne de la discordance entre certains linguistes sur ce qui relève ou non de l'hésitation à l'oral. Toutefois, depuis quelques années, plusieurs études se sont intéressées aux phénomènes d'hésitations et une nomenclature commence à se dégager.

Ces phénomènes peuvent être de natures différentes. Ils incluent par exemple les *pauses silencieuses*, mais pas toutes, certaines jouant d'autres rôles que l'hésitation comme la hiérarchisation et la structuration des constituants, ou encore la mise en valeur stylistique de certains syntagmes. Les *pauses remplies*, correspondant au « *um* » ou au « *uh* » en anglais et au « *euh* » en français, ou encore les *allongements vocaliques* en fin de mots sont des marques d'hésitation très employées [Hen02a, Can00]. On peut aussi citer [Shr01] :

- les répétitions : « *tous les - les jours*<sup>5</sup> »,
- les suppressions : « *c'est - il est arrivé lundi* »,
- les substitutions : « *tous les jours - toutes les semaines* »,
- les insertions : « *je suis convaincu - je suis intimement convaincu* »,
- les erreurs d'articulation : « *en jouin - en juin* ».

Ces quatre derniers types de disfluence sont parfois désignés sous le terme d'autocorrections [Hen02a].

Il est possible de faire une autre classification des phénomènes d'hésitation, en distinguant [PH04] :

<sup>5</sup>Le « - » représentant une pause silencieuse.

- les bribes, correspondant à une reprise à partir de syntagmes inachevés : « *il a quand même un - une fibre pédagogique* »,
- les amorces, associées aux mots inachevés : « *c'est pas malho - c'est pas malhonnête* ».

Les différents phénomènes d'hésitation sont rarement produits seuls mais plutôt en combinaison. Il peut ainsi y avoir des répétitions de fragments de mots (« *on le re- re- revendique encore une fois* ») [HP03]. De même, pour assurer un rôle d'hésitation, les pauses silencieuses sont généralement associées à des allongements vocaliques ou à des pauses remplies [CV04].

La fréquence d'apparition des marques d'hésitation est très variable selon le locuteur et le contexte d'élocution. Dans les dialogues homme-machine, on constate généralement moins de disfluences que dans les dialogues entre deux humains. Il semble également que les hommes produisent en moyenne davantage de disfluences que les femmes. La fréquence des marques d'hésitation est de plus dépendante de la position dans la phrase. Les disfluences ont ainsi tendance à se produire plutôt en début de phrase, lorsque l'effort de planification est le plus important. Pour la même raison, leur taux augmente souvent avec la longueur de la phrase [Shr01]. Les phénomènes d'hésitation affectent enfin différemment les mots. En distinguant deux types de mots, à savoir les *lexicaux*, représentant ceux qui ont une charge lexicale pleine<sup>6</sup> et les *grammaticaux*, participant à la structuration de la langue<sup>7</sup>, on constate que les répétitions concernent plus souvent les mots grammaticaux et que les amorces affectent plus fréquemment les lexicaux [PH04, Hen02b].

En analysant les disfluences, des psycholinguistes ont mis en évidence plusieurs régions (Fig. 9)[Shr94, Shr01] :

- le *reperandum*, désignant une partie ou la totalité de la séquence qui sera abandonnée au profit de la réparation,
- le *point d'interruption*, marquant une rupture dans la fluidité du discours,
- l'*interregnum*, pouvant être des pauses remplies, mais aussi des termes d'édition, comme « *ben* », « *hein* », « *tu vois* », « *disons pour simplifier* », « *je sais pas moi* » ou encore « *je me rappelle plus du nom* »,
- la *réparation*, représentant la partie corrigée du *reperandum*.

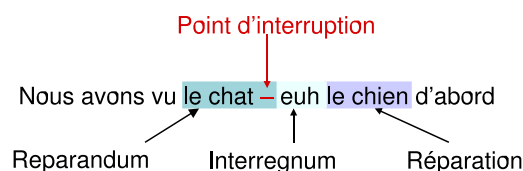


FIG. 9 – Régions dans une disfluence (d'après [Shr01])

La segmentation des disfluences en régions est liée à des modifications des propriétés acoustiques et phonétiques. On observe ainsi généralement un allongement des syllabes avant le

<sup>6</sup>En général les noms, les adjectifs, les verbes et les adverbes.

<sup>7</sup>Les pronoms, les déterminants, les prépositions, les conjonctions et les verbes auxiliaires.

point d'interruption, ou au contraire, dans le cas d'une détection d'erreur de la part du locuteur, un raccourcissement des syllabes. Cette structuration permettrait de réduire les énoncés à un oral « propre », proche de l'écrit, en éliminant le *reparandum* et l'*interregnum*. Malheureusement, cette approche ne rencontre pas l'unanimité des linguistes étudiant l'oral.

### 3.4 Corpus oraux

On dispose à l'heure actuelle de corpus écrits de taille volumineuse. Le *British National Corpus* (BNC) contient ainsi 100 millions d'occurrences pour l'anglais, tandis que la base *Frantext* comporte 210 millions d'occurrences pour le français. Les pages disponibles sur Internet, bien que souvent bruitées, peuvent être également considérées comme une source gigantesque d'informations [Vér04]. Il n'existe malheureusement pas pour l'oral de corpus de taille comparable. Le plus grand, à savoir la partie orale du BNC, contient 10 millions d'occurrences en anglais. On peut citer pour le français de Belgique le corpus Valibel comptant près de 4 millions d'occurrences et pour le français de l'hexagone, le corpus Corpaix comptant environ 2,5 millions d'occurrences [CVD05], et le *Corpus de référence du français parlé* comptant 440 000 occurrences [ÉD04].

Les corpus oraux, constitués au moyen de transcriptions manuelles, concernent aussi bien de la parole lue, préparée ou encore spontanée. La parole lue est associée à la lecture à haute voix de textes, qui bien souvent appartiennent au domaine de la langue écrite. La parole spontanée correspond à des situations de la vie quotidienne où le locuteur parle sans avoir préparé au préalable ce qu'il allait dire et en s'adaptant constamment à ses interlocuteurs. La parole préparée est une forme intermédiaire par rapport aux deux précédentes et est plus délicate à définir. Elle désigne des prises de parole où le locuteur a réfléchi auparavant aux idées qu'il souhaite transmettre; il peut même disposer d'un texte contenant une partie de son discours, qu'il n'ira pas toutefois jusqu'à lire mot à mot pour pouvoir interagir avec les interlocuteurs. Les questions posées par un journaliste et les entretiens donnés par un homme politique sont deux exemples de production de parole préparée.

Selon les conventions adoptées, les corpus peuvent avoir plusieurs formats; ils peuvent contenir des informations précises sur la prosodie, être ponctués ou bien encore avoir des annotations sur des événements autres que la prononciation des mots. Accessoirement, ils sont souvent étiquetés pour faciliter leur exploitation. Ils sont alors segmentés en mots ou groupes de mots, auxquels sont associées des informations (*e.g.* syntaxiques ou sémantiques) sur le rôle de ces constituants dans le groupe de souffle courant. La forme choisie du corpus dépend de l'utilisation prévue mais aussi des ressources humaines disponibles, compte tenu de la difficulté et du temps nécessaire pour effectuer une transcription manuelle. L'indication de la prosodie dans les corpus oraux, bien que souhaitable, demande notamment un temps considérable et se limite souvent à quelques informations. Même si des logiciels tels que *Transcriber* simplifient l'alignement du son avec la transcription [BGWL01], il a été estimé qu'environ 40 heures de travail sont nécessaires pour faire une transcription orthographique précise d'une heure de parole, et ce, avec de bonnes conditions d'enregistrement et sans alternance de locuteurs [Cam01, Vér04]. La suite de cette section décrit sous quelles formes

peuvent se présenter les corpus oraux, avant de présenter la manière dont ils peuvent être étiquetés.

### 3.4.1 Formes des corpus oraux

La transposition du signal sonore sous la forme d'un autre canal de communication, à savoir le texte, nécessite certaines conventions. Elle se heurte à des difficultés qui sont principalement : la présence d'événements non linguistiques dans le flux audio, l'existence de phénomènes propres à la langue parlée et la segmentation du discours.

La perception des mots à partir du signal sonore peut tout d'abord être perturbée par des événements non linguistiques qu'il est intéressant d'annoter dans la mesure où ils peuvent rendre la parole inaudible. Ces phénomènes peuvent être liés aux conditions d'enregistrement, tels les parasites ou les bruits de fond, ou bien encore être produits par les locuteurs, comme les inspirations, les rires ou les toux.

Des événements propres à la langue parlée nécessitent également des conventions. Les chevauchements entre les discours prononcés en même temps par plusieurs locuteurs sont ainsi problématiques, dans la mesure où il faut indiquer la simultanéité dans le texte, qui a un format séquentiel. Les phénomènes d'hésitation nécessitent quant à eux une attention particulière de la part du transcripteur car ils sont naturellement ignorés lors de la compréhension des dialogues. Ils présentent pourtant un intérêt pour les linguistes, qui en les observant peuvent voir la production du langage en train de se réaliser. Il faut donc adopter des règles pour indiquer les pauses silencieuses ou encore les amorces. Il existe aussi plusieurs types de prononciation pour certains mots, ce qu'il est utile d'indiquer. Ceci peut conduire à faire une transcription *phonétique* plutôt qu'une transcription plus classique, dite *orthographique*, sous forme de mots. On choisit généralement d'avoir une transcription orthographique, avec une convention pour distinguer les prononciations de mots problématiques, tels que les sigles.

Une des difficultés majeures de la constitution des corpus oraux réside dans leur segmentation. Pour faciliter leur lisibilité, ils ne peuvent en effet se présenter sous la forme d'une suite ininterrompue de mots. Un premier niveau de segmentation consiste à distinguer les  *tours de parole* , *i.e.*, les suites de mots prononcées par un locuteur donné avant qu'il ne cède la parole. Il n'existe pas d'unité qui soit aussi nettement définie que la phrase pour la langue écrite. L'unité de la langue parlée, que l'on qualifie d'*énoncé*, peut ainsi représenter l'ensemble du tour de parole. Certains linguistes et psycholinguistes caractérisent aussi un énoncé comme étant délimité par la prosodie. D'autres le définissent comme une succession de mots représentant une idée cohérente; dans les directives d'annotation *metadata* adoptées dans les évaluations *Rich Transcription* conduites par le NIST, on parle ainsi de *SU*, associée à plusieurs significations : *Sentential Unit*, *Syntactic Unit*, *Semantic Unit* ou *Slash Unit* [Str03].

Malgré la subjectivité que cela entraîne de la part du transcripteur, certaines recommandations de transcription demandent d'indiquer les signes de ponctuation dans les corpus oraux, ce qui les rapproche de textes plus conventionnels. La prosodie et les pauses silencieuses apportent alors des informations précieuses quant à la ponctuation. La figure 10 montre ainsi un extrait de corpus transcrit avec des ponctuations. D'autres conventions de



---

« *Le gouvernement américain va présenter, [i] euh, au conseil de sécurité de l'ONU, [i] un projet d(e) résolution, [mic] un texte qui permettrait, [i] la levée d(e) l'embargo Noëlle ^Véli.* »

où « [i] » marque une inspiration, « !ONU » représente un acronyme, « d(e) » indique une absence de prononciation du « e », « [mic] » témoigne de la présence d'un bruit de micro et « ^Véli » précise que le mot est inconnu par le transcripteur.

---

FIG. 10 – Extrait d'un corpus transcrit selon les conventions de Transcriber

---

« ils ont des ouvriers euh	payés	
	spécialisés	sup-
		sur les chantiers de fouille »

---

FIG. 11 – Mise en grille d'un énoncé (d'après [Gué05])

transcription recommandent au contraire de ne pas indiquer de signes de ponctuation car le découpage effectué par le transcripteur en syntagmes ou en phrases préjuge de l'analyse à faire [BB90]. Les pauses silencieuses se révèlent de plus imprécises pour détecter ces ponctuations, y compris quand le corpus transcrit correspond à un texte lu [CV02]. Pour améliorer la lisibilité des transcriptions, Blanche-Benveniste propose une segmentation plus fine sous forme de grille. Lors de la production de la parole, il existe en effet un travail de formulation important qui vient perturber la succession des mots et le discours s'apparente souvent à une énumération. Ce type de présentation conduit à positionner l'un en dessous de l'autre des éléments prononcés successivement et qui occupent une même place syntaxique (Fig. 11) [BB90].

### 3.4.2 Étiquetage des corpus oraux

Pour faciliter l'utilisation des corpus, qu'ils soient oraux ou écrits, un étiquetage morpho-syntaxique est souvent réalisé. Cet étiquetage consiste à attribuer à chaque mot, voire à des groupes de mots dans le cas de locutions, une étiquette que l'on appelle *partie du discours* ou *PoS* (pour *Part of Speech*). Une PoS correspond à une propriété grammaticale dans une phrase, telle que nom, verbe, adjectif, préposition, *etc.*, que l'on peut préciser par le genre, le nombre, la personne, *etc.* Le choix des PoS à considérer n'est pas fixe pour une langue donnée. Les étiqueteurs peuvent ainsi avoir des jeux d'étiquettes divers pour des mots très employés [CVD05].

Devant la taille des corpus à analyser, l'étiquetage automatique se révèle indispensable. Les premiers programmes d'étiquetage automatique sont apparus pour l'écrit dès les années 50. Ils étaient basés sur la production manuelle de règles. Depuis les années 80, certains étiqueteurs utilisent des méthodes statistiques, telles que les HMM [Bra00] ou les arbres

de décision [Sch94, Sch95], pour prédire la probabilité d'attribution des étiquettes. Mais ce n'est que récemment que des études ont été menées sur l'étiquetage de corpus oraux.

Il n'existe pas à l'heure actuelle, tout du moins à notre connaissance, d'étiqueteurs conçus spécifiquement pour l'oral. Des systèmes prévus initialement pour étiqueter des documents écrits ont été utilisés sur des corpus oraux, que ce soit pour l'anglais [Gar95], le français [VV99, CVD05], le suédois [NG01], le néerlandais [VEZD00], l'espagnol [MG03], le portugais [MAB03], l'italien [PPM04] ou encore le japonais [UNY<sup>+</sup>02]. Ces étiqueteurs sont souvent adaptés à l'oral en modifiant légèrement leur comportement pour certains mots et en appliquant au corpus un traitement préalable pour éliminer certaines caractéristiques de l'oral, telles que les amorces. Dans le cas où le corpus ne contient pas de signe de ponctuation, les marques de pauses peuvent être remplacées par des points de suspension, *i.e.*, par la ponctuation la plus neutre possible par rapport au fonctionnement de l'étiqueteur. Les amorces de mots sont en général ignorées par l'étiqueteur car il est parfois difficile de deviner le mot prononcé. Il peut être également utile d'éliminer des événements non lexicaux, tels que « *hein* » et « *eah* », ou encore des événements non linguistiques comme le signalement d'applaudissements. De plus, certains mots apparaissent beaucoup plus fréquemment à l'oral qu'à l'écrit, notamment les contractions de mots et les interjections ; les mots caractéristiques de l'oral sont ainsi insérés dans le lexique de l'étiqueteur. De même, on ajoute manuellement des règles ou on modifie manuellement les probabilités d'assignation des étiquettes pour les mots qui ont un comportement différent à l'écrit et à l'oral.

Pour mesurer les performances de l'étiquetage automatique, un extrait de corpus étiqueté automatiquement est comparé avec le même extrait étiqueté manuellement. L'évaluation est une opération délicate à effectuer dans la mesure où l'étiquetage de référence, *i.e.*, l'étiquetage manuel, peut différer selon les annotateurs pour certains mots problématiques. Il est le plus souvent choisi de considérer comme *acceptable* une étiquette dès qu'elle relève d'un point de discussion entre linguistes. Le critère de performance est alors le pourcentage d'étiquettes acceptables. La comparaison des performances entre étiqueteurs n'est pas non plus chose aisée car le jeu d'étiquettes ou la segmentation en unités à étiqueter peuvent différer [AMP<sup>+</sup>99]. En ce qui concerne l'étiquetage de l'écrit, les performances sont supérieures à 95 % d'étiquettes correctes. En utilisant des étiqueteurs conçus initialement pour l'écrit et adaptés pour l'oral, des expériences ont permis d'atteindre 98,75 % pour le français [VV99], entre 95 % et 97 % pour le suédois [NG01], 94,3 % pour le néerlandais [VEZD00] ou encore 98,3 % pour l'espagnol [MG03].

Les performances sont donc très proches de ce qui est observé pour l'écrit, contrairement à ce que l'on aurait pu penser. L'explication qui est souvent donnée pour justifier ces résultats est que les étiqueteurs se basent sur des phénomènes locaux et se révèlent donc peu sensibles aux phénomènes propres à la langue parlée. Une analyse qualitative des erreurs montre toutefois que l'étiquetage est mis en erreur par certaines particularités des corpus oraux comme l'absence de ponctuation, les répétitions ou les chevauchements entre des mots prononcés simultanément [VV99].

Outre l'étiquetage par des PoS, il existe d'autres niveaux d'annotation envisageables pour les corpus oraux. On peut ainsi citer :

- l’annotation pragmatique caractérisant les actes de discours, tels que la question, le conseil, la confirmation, les remerciements... [LMW97],
- l’annotation stylistique caractérisant la présentation du discours et de la pensée, *e.g.* par narration, discours direct, discours indirect, discours indirect libre... [LMW97]
- l’annotation syntaxique indiquant les dépendances entre mots ou groupes de mots [BCD<sup>+</sup>04, BV05].

Toutes ces annotations sont beaucoup moins courantes que l’étiquetage morphosyntaxique et ne peuvent être obtenues convenablement par des méthodes automatiques. L’analyse syntaxique de l’oral par des grammaires pose ainsi de nombreux problèmes. Aux difficultés rencontrées pour la langue écrite et qui sont également présentes dans la langue parlée, comme l’ambiguïté des analyses possibles, s’ajoutent des problèmes bien spécifiques à l’oral, tels que la présence de disfluences, le respect plus lâche des règles de la langue et l’absence de segmentation bien claire en phrases. Les distorsions de la langue parlée requièrent une grande robustesse de la part des analyseurs syntaxiques. Pour concevoir ces systèmes, une approche similaire à la démarche adoptée pour l’étiquetage morphosyntaxique consiste à utiliser un analyseur développé pour l’écrit et à lui adjoindre des procédures traitant les extragrammaticalités de l’oral [BDM98].

### 3.5 Transcription automatique des dialogues spontanés

La communauté de la reconnaissance de la parole s’intéresse de plus en plus aux dialogues spontanés, pour lesquels les locuteurs s’expriment sans préparer leur discours, ce qui entraîne une production accrue de disfluences. Les campagnes d’évaluation conduites par le NIST depuis 1987 témoignent de cette évolution [Pal03]. Les premières évaluations ont ainsi porté sur des applications au vocabulaire très réduit et avec une grammaire spécifiée, puis sur la lecture de journaux. C’est avec la transcription d’émissions d’actualité que les systèmes de RAP se sont heurtés à des documents plus variés. Ces émissions comportent en effet des discours préparés mais aussi des interviews ou des reportages, dans lesquels on retrouve fréquemment des dialogues spontanés, dans des conditions d’enregistrement bruitées. D’autres évaluations portent désormais sur la transcription de la base de données *Switchboard*, constituée de dialogues téléphoniques spontanés sur des sujets particuliers, ou plus récemment sur les échanges enregistrés entre les participants d’une réunion.

Les phénomènes caractéristiques de la langue parlée gênent considérablement la transcription automatique. L’ordre plus flexible des mots perturbe par exemple les modèles de langages qui basent leur calcul sur les séquences de mots (*cf.* section 2.2.3). Les disfluences compliquent quant à elles l’analyse du signal et viennent perturber le calcul des probabilités par les modèles N-grammes. Les amorces et les contractions de mots sont également difficiles à modéliser, le vocabulaire ayant une taille limitée. Enfin, les marqueurs de discours tels que « *en fait* », « *enfin* », « *voilà* » sont particulièrement problématiques car ils sont beaucoup plus fréquents à l’oral et sont souvent très mal articulés. Une étude menée sur 10 heures d’interviews télévisées a montré que les disfluences, représentant 8 % du corpus, expliquent 12 % du taux d’erreur (*cf.* section 2.4) global sur la transcription. Ce taux suggère que les

disfluences accroissent la difficulté de la transcription mais n'exercent pas un effet majeur sur les segments voisins [ADHB<sup>+</sup>04].

La segmentation du signal acoustique en groupes de souffle n'est pas non plus sans poser de problème. Le groupe de souffle étant en effet justifié par le mode de fonctionnement des systèmes de RAP et uniquement défini d'après des indices acoustiques, à savoir la détection de pauses, les unités sur lesquelles ont à opérer les ML ne sont pas aussi cohérentes que les phrases de la langue écrite. Une segmentation plus linguistique peut ainsi être envisagée avant d'appliquer des ML. Les meilleures perplexités obtenues sur un corpus de test non segmenté, avec un ML utilisant une segmentation linguistique plutôt qu'une segmentation acoustique, illustrent l'utilité de cette adaptation [MI96]. Toutefois, la détection des segments linguistiques, effectuée notamment à l'aide d'indices sur les classes grammaticales et la prosodie, restent une tâche difficile puisque les meilleurs systèmes détectent les fins de segments avec des taux d'erreurs variant entre 30 % et 50 %, selon le type de discours à analyser [LSHS04].

Pour traiter les dialogues spontanés, il est nécessaire d'apporter des modifications importantes au processus de transcription. Un système de RAP initialement conçu pour être appliqué sur des émissions d'information transcrivait ainsi des conversations téléphoniques avec un taux d'erreur de l'ordre de 50 %. Après plusieurs raffinements, ce taux a finalement pu être ramené à 21 % [GAL<sup>+</sup>04]. Parmi les adaptations à effectuer pour prendre en compte la parole spontanée, on peut citer :

- au niveau du modèle acoustique, l'ajout de trois phones associés aux pauses, aux hésitations et aux respirations ;
- au niveau du dictionnaire de prononciations, l'ajout de répétitions et de contractions pour certains mots et l'intégration des interjections ;
- au niveau du modèle de langage, l'inclusion des inspirations et des pauses remplies dans l'historique des modèles N-grammes.

De manière à limiter les perturbations engendrées par les disfluences, certains systèmes de RAP cherchent à les détecter et à les éliminer. Les disfluences qui sont alors principalement corrigées sont celles où le *reparandum* et la réparation (Fig. 9, page 27) sont très proches au niveau lexical [DGA<sup>+</sup>93, JC04]. Certaines sont toutefois plus problématiques à détecter et peuvent de surcroît apporter des informations pour la prédiction des mots. Il a été ainsi constaté qu'une suppression des pauses remplies dans les corpus d'apprentissage et de test entraînait une augmentation de la perplexité [SS96]. Une alternative est donc d'utiliser les disfluences comme source d'informations pour la transcription [SO96].

La section 3 a exposé les caractéristiques de la langue parlée, en montrant brièvement les formes sous lesquelles se présentent les corpus oraux. La section suivante décrit comment des connaissances linguistiques peuvent être introduites dans des systèmes de RAP. Les techniques utilisées à cette fin doivent pouvoir s'intégrer dans le schéma de fonctionnement que nous avons décrit en section 2 et être suffisamment flexibles pour ne pas être perturbées par les phénomènes de l'oral.

## 4 La linguistique pour la reconnaissance de la parole

Comme nous l'avons vu en section 2.2, les systèmes de RAP adoptent une modélisation hiérarchique pour décoder la parole. Le MA identifie à partir des caractéristiques extraites du signal des phones puis des mots ; le ML est quant à lui chargé de reconnaître les successions de mots les plus probables pour un groupe de souffle. Ce mode de fonctionnement présente des similarités avec celui de l'être humain qui, pour reconnaître le sens véhiculé par la parole, identifie successivement des phones, des syllabes, des mots et des « phrases » [All94]. Mais à la différence de l'être humain qui utilise différents niveaux de connaissances linguistiques, les informations sur le langage dans les systèmes de RAP se limitent essentiellement à la connaissance des phones et des phonèmes d'une langue, à la réalisation d'un dictionnaire de prononciations et à l'apprentissage de modèles de langages sur des corpus oraux. Il est notamment souvent reproché aux ML les plus employés, *i.e.*, aux modèles N-grammes, de ne considérer de manière arbitraire qu'un historique de  $N - 1$  mots. Il existe en effet bien des configurations où leurs hypothèses de conception sont mises en défaut. Ainsi, dans l'exemple suivant : « *les oiseaux sur l'arbre chantaient* », un modèle quadrigramme n'aura pas les éléments utiles, dans son historique de trois mots, pour prédire l'accord du verbe « *chantaient* », du fait de l'insertion du complément de lieu entre le sujet et le verbe.

Une expérience menée par Brill *et al.* corrobore l'hypothèse d'un apport de la linguistique à l'amélioration de la reconnaissance de la parole. Elle a consisté à analyser les ressources qu'un être humain utiliserait pour corriger des transcriptions automatiques [BFHM98]. Trois corpus de parole ont été à cette fin traités par un système de RAP pour produire les listes des dix meilleures hypothèses (*cf.* section 2.3) associées à chaque groupe de souffle. Il a alors été demandé à des sujets humains de sélectionner parmi chacune de ces listes l'hypothèse qui leur semblait la plus juste. Les choix effectués par les sujets se sont souvent révélés judicieux puisqu'ils ont permis une nette diminution du taux d'erreur par rapport à la meilleure hypothèse désignée par le système de RAP pour chacun des groupes de souffle. Un questionnaire sur les connaissances utilisées pour faire leur sélection a montré que les humains se basaient principalement sur des informations linguistiques telles que l'emploi correct des prépositions et des déterminants, les accords en genre et en nombre, l'examen du temps pour les verbes, la connaissance de syntagmes idiomatiques ou encore l'analyse de la structure des hypothèses proposées.

Si l'introduction de connaissances linguistiques supplémentaires est donc souhaitable, elle peut s'effectuer de diverses manières. D'une part, elle implique des choix qui influent sur le processus de transcription. Si l'on souhaite utiliser au plus tôt des informations sur le langage, certains modes de couplage MA-ML sont ainsi à privilégier par rapport à d'autres. Il existe également plusieurs modes possibles de calcul de probabilités de successions de mots. D'autre part, les techniques diffèrent suivant le type de connaissances envisagé. Les méthodes introduisant des informations morphologiques ne sont ainsi pas les mêmes que celles prenant en compte de la sémantique. Cette section présente dans un premier temps les techniques adaptant le processus de transcription à l'introduction de la linguistique, avant de rentrer plus en détail dans les méthodes spécifiques à chaque type de connaissances.

## 4.1 Adaptation du processus de transcription

Les ML classiques ne manipulent que des mots et se limitent à un historique restreint dans leurs calculs de probabilités. L'intégration de connaissances linguistiques conduit à traiter des informations supplémentaires au cours de la transcription, telles que la décomposition morphologique des mots à reconnaître, ou bien encore à effectuer des calculs de probabilités prenant en compte l'intégralité du groupe du souffle courant, comme cela est le cas de certaines méthodes d'analyse syntaxique. Le mode de fonctionnement du système de RAP doit alors être modifié afin d'introduire ces connaissances sans trop dégrader la rapidité de la transcription. Les adaptations nécessaires sont effectuées au niveau de l'interface MA-ML et du calcul des probabilités par le ML, ce qui fait l'objet des deux sections suivantes.

### 4.1.1 Intégration du modèle acoustique avec le modèle de langage

L'intervention du MA et du ML dans le processus de transcription se révèle souvent plus complexe que la description qui en a été faite en section 2.2. Pour des raisons de rapidité, le décodage de la parole est généralement réalisé en plusieurs passes, chaque passe correspondant à une application d'un MA et d'un ML et conduisant à la création d'un graphe de mots. Au fur et à mesure des passes, les graphes de mots produits sont de plus en plus réduits, tandis que les MA et les ML sont de plus en plus informatifs et, par conséquent, lents. On utilisera par exemple plutôt un ML trigramme dans une première passe et un ML quadrigramme dans une seconde passe. L'intérêt de ces systèmes multipasses consiste ainsi à pouvoir utiliser des modèles complexes au niveau des dernières passes, sans engendrer une augmentation trop importante du temps de calcul. En sus du mode d'organisation général du processus de transcription, le système de RAP peut être adapté au niveau de l'interface entre un MA et un ML donnés. Il existe ainsi trois types d'intégration des modules : le couplage étroit, le couplage lâche et le couplage modéré [HJM<sup>+</sup>94].

Un système à *couplage étroit* intègre toutes les connaissances dans un ensemble de processus interdépendants et non séparables. Dans un tel système, le ML est directement intégré dans le MA, ce qui est permis par certaines structures de ML. Celui-ci doit ainsi être particulièrement rapide puisque le nombre d'hypothèses examinées par le MA est très grand et n'a pas encore été réduit. Les ML qui prennent uniquement en compte des contraintes locales sont ainsi bien adaptés à ce type de couplage car ils facilitent la mise en œuvre d'algorithmes dynamiques. Un modèle N-grammes, à condition que  $N$  ne soit pas trop grand, ou encore une grammaire à états finis [Moo99], peuvent par exemple convenir à ce type de couplage.

Les systèmes à couplage étroit présentent l'avantage d'utiliser au plus tôt les connaissances linguistiques du ML pour réduire l'espace de recherche du MA. Ils sont néanmoins difficiles à faire évoluer puisque ce couplage impose des contraintes fortes sur le fonctionnement du ML. En outre, la forte imbrication des composants complique grandement l'évaluation de l'impact de chaque connaissance sur les performances du système.

Un système à *couplage lâche* isole le MA et le ML en modules relativement indépendants et communiquant entre eux. Le rôle du ML est alors de filtrer ou réordonner les hypothèses fournies par le MA. Trois possibilités sont généralement envisagées pour faire l'interface entre

le MA et le ML : utiliser un graphe de mots [CR89, HJJ<sup>+</sup>99], une liste des N meilleures hypothèses [CS89], ou bien encore un *graphe d'homophones* [GAAD<sup>+</sup>05, BNSd99]. Dans ce dernier type d'interface, chaque mot de la meilleure hypothèse est remplacé par tous ses homophones possibles, *i.e.*, par tous les mots du langage qui ne sont pas discriminables au moyen d'indices acoustiques. Cette approche convient parfaitement à la détection des erreurs dues à des problèmes d'accord en genre et en nombre.

Généralement, les systèmes à couplage lâche sont plus faciles à faire évoluer que ceux à couplage étroit et on ne constate pas d'explosion combinatoire quand on augmente le nombre de connaissances prises en compte. De plus, ils n'imposent pas de contraintes importantes sur la structure du ML et permettent d'évaluer facilement les performances des MA et des ML, qui se présentent sous la forme de modules bien différenciés. Ils ont toutefois le désavantage de ne pas réduire l'espace de recherche lors du décodage par le MA.

Un système à *couplage modéré* a un comportement intermédiaire entre les deux précédents. Il utilise le ML pour guider le MA, sans y être intégré. Il se différencie du couplage lâche par le fait que le ML ne peut pas être supprimé sans modifier le processus de recherche conduit par le MA. Deux approches ont été envisagées pour le couplage modéré : l'approche descendante et l'approche ascendante [HW94]. Dans l'*approche descendante*, le ML est invoqué à des points de décision où il prédit des hypothèses. Le MA est ensuite chargé de sélectionner la meilleure hypothèse. Un analyseur syntaxique LR peut par exemple être utilisé pour prédire des phones qui sont ensuite vérifiés par un HMM de phones [KKS89]. Les phones qui constituent un mot sont alors spécifiés par des règles de grammaire. Dans l'*approche ascendante*, les scores acoustiques sont calculés en premier et le ML est appliqué pour vérifier les hypothèses, en réduisant éventuellement le nombre de candidats acoustiques. Cette organisation est très proche de celle des systèmes à couplage lâche mais il est fait appel au ML à chaque point de décision et non pas à la fin de l'analyse du groupe de souffle. L'approche ascendante a été par exemple mise en œuvre par un analyseur syntaxique tabulaire par îlots [BP03]. Les hypothèses de mots sont alors fournies les unes après les autres par le MA au ML. Quand celui-ci échoue dans son analyse syntaxique, il demande au MA de produire des hypothèses supplémentaires à des endroits précis du groupe de souffle. Une autre possibilité envisagée a consisté à construire dynamiquement le ML par un réseau représentant une grammaire, de manière à n'inclure que les transitions nécessaires à l'analyse du groupe de souffle courant [MPM89].

Les systèmes à couplage modéré constituent donc un compromis original avec les deux autres types. Malheureusement, ils posent de nombreux problèmes d'ingénierie pour faire interagir les modules et imposent des contraintes sur la structure des ML telles que la nécessité d'une analyse gauche-droite.

Au terme de la présentation des trois modes de couplage du MA et du ML, les deux alternatives qui paraissent les plus intéressantes sont donc le couplage étroit et le couplage lâche, le premier autorisant l'intégration de connaissances linguistiques au plus tôt, le second imposant moins de contraintes sur la structure et la rapidité du ML. Les systèmes à couplage modéré semblent quant à eux être une approche difficile à mettre en œuvre. Il existe d'ailleurs peu de systèmes de RAP à ce jour qui soient basés sur cette organisation. Cette section a

présenté les manières de réaliser l'interface MA-ML en vue d'introduire des informations linguistiques; la suivante décrit comment des ML prenant en compte des connaissances diverses peuvent être agencés pour calculer des probabilités de succession de mots.

#### 4.1.2 Linguistique et modèle de langage

Dans les systèmes de RAP actuels, le ML contient généralement un modèle N-grammes. Il existe plusieurs schémas d'intégration de connaissances linguistiques supplémentaires dans ce modèle [CRAR99].

Une première solution consiste à construire un nouveau ML, plus linguistique, que l'on utilise en remplacement d'un modèle N-grammes. Ceci implique qu'il possède des propriétés communes avec le modèle N-grammes; il doit être notamment capable de fournir un score et de faire une analyse gauche-droite en un temps raisonnable. Ce type de solution est généralement envisagé quand on dispose de connaissances *a priori* sur un domaine particulier, mais pas d'un corpus d'apprentissage suffisant, ou quand le nouveau ML intègre simultanément plusieurs types de connaissances, y compris des connaissances lexicales qui seraient redondantes avec celles apportées par un modèle N-grammes [WH02].

Une deuxième possibilité est d'utiliser les connaissances linguistiques pour lisser le calcul des probabilités des modèles N-grammes [JWS<sup>+</sup>95]. Ceci se fait par exemple en utilisant une grammaire pour générer des phrases qui viennent ensuite enrichir le corpus d'apprentissage du modèle N-grammes [WMH00]. Les probabilités d'un modèle N-grammes peuvent également être évaluées directement à partir du nouveau ML [SS94].

Dans une troisième solution, les connaissances sont utilisées séquentiellement. L'intérêt est que si chaque source d'informations est susceptible d'améliorer les performances globales du système de RAP, chacune a une influence et une complexité de calcul potentiellement très différentes. Par exemple, un modèle bigramme pourra réduire la perplexité d'un facteur de dix, avec peu de calculs, tandis qu'un ML apportant des informations plus complètes sur la syntaxe et la sémantique pourra demander beaucoup de calculs pour une réduction de perplexité moindre [SA90]. Dans cette approche séquentielle, les modèles les plus rapides sont utilisés en premier lieu pour produire l'ensemble des hypothèses les plus probables. Cet ensemble, qui se présente sous la forme d'une liste des meilleures hypothèses [CS89] ou d'un graphe de mots [SCL92], est ensuite filtré et réordonné au moyen des sources de connaissance restantes.

Une dernière solution consiste à combiner plusieurs modèles, apportant chacun des connaissances, pour constituer un seul ML. Cette combinaison peut se faire à l'aide de l'interpolation linéaire ou du repli, dont les principes ont déjà été exposés (*cf.* section 2.2.3). Dans le cas de l'interpolation linéaire, technique la plus utilisée du fait de sa simplicité, la combinaison de  $M$  modèles différents, associés aux distributions de probabilité  $P_k$  avec  $k = 1 \dots M$ , s'effectue de la manière suivante :

$$P(w_i|w_1^{i-1}) = \sum_{k=1}^M \lambda_k P_k(w_i|w_1^{i-1}) \quad (24)$$



où  $\sum_{k=1}^M = 1$ . En ce qui concerne le repli, le calcul pour combiner M modèles, avec M fixé ici à 2 pour simplifier l'équation, s'écrit [NW96a] :

$$P(w_i|w_1^{i-1}) = \begin{cases} P_1(w_i|w_1^{i-1}) & \text{si } w_i \in \Phi_1(w_1^{i-1}) \\ \alpha(w_1^{i-1}) \times P_2(w_i|w_1^{i-1}) & \text{sinon} \end{cases} \quad (25)$$

où  $\Phi_k(w_1^{i-1})$  représente l'ensemble des mots dans le contexte  $w_1^{i-1}$  pour lequel le  $k^{\text{ième}}$  modèle est à utiliser en priorité, et  $\alpha$  est un coefficient de normalisation. L'intérêt de cette méthode est d'utiliser d'abord les modèles les plus informatifs quand on dispose de suffisamment d'informations dans le contexte courant. Une autre possibilité pour combiner des modèles est d'utiliser des *modèles exponentiels*, appelés encore à *entropie maximale* [Ros96, Goo01]. Ceux-ci évaluent les probabilités sous la forme :

$$P(w_i|w_1^{i-1}) = \frac{\exp(\sum_k \lambda_k f_k(w_1^i))}{z(w_1^{i-1})} \quad (26)$$

où  $z$  est une fonction de normalisation telle que :

$$\sum_{w_i} P(w_i|w_1^{i-1}) = 1 \quad (27)$$

Les  $\lambda_k$  sont des coefficients obtenus grâce à un algorithme d'apprentissage, tandis que les  $f_k$  sont des fonctions de contraintes retournant typiquement 0 ou 1. Le principal intérêt de ce type de modèle est que les  $f_k$  permettent de représenter des modèles N-grammes, des modèles à base de cache (*cf.* section 2.2.3), des modèles N-classes (*cf.* section 4.2.2) ou encore des modèles *triggers* (*cf.* section 4.2.3). Dans le cas d'une fonction pour le trigramme  $w_a w_b w_c$ , on aura ainsi :

$$f_{w_a w_b w_c}(w_1^i) = \begin{cases} 1 & \text{si } w_{i-2} = w_a, w_{i-1} = w_b \text{ et } w_i = w_c \\ 0 & \text{sinon} \end{cases} \quad (28)$$

Les modèles exponentiels permettent ainsi d'intégrer plusieurs sources d'informations de manière élégante. Leur temps d'apprentissage est cependant extrêmement long et ils se révèlent assez lents lors de leur utilisation. En outre, mis à part avec les modèles *triggers*, il semble qu'ils n'aient pas encore permis de réduire la perplexité de manière significative [Goo01].

Après avoir décrit les adaptations possibles du processus de transcription pour intégrer de nouvelles informations, nous nous tournons désormais vers la présentation des sortes de connaissances linguistiques pouvant compléter celles généralement apportées par les MA et les ML, des méthodes qui ont déjà été envisagées pour les utiliser et de leur influence sur la qualité de la transcription produite.

## 4.2 Quelles connaissances linguistiques ?

Cette section fait un tour d'horizon des types d'informations linguistiques qui ont été pris en compte en RAP. Les MA ne pouvant inclure que des connaissances concernant

l'acoustique, nous nous consacrons plus particulièrement ici aux expériences menées dans le cadre de l'intégration de ces connaissances au sein de ML. La présentation est structurée selon une typologie à cinq niveaux généralement reconnue en linguistique :

- la *phonologie* et la *phonétique* qui étudient les sons,
- la *morphologie* qui étudie la structure des mots,
- la *syntaxe* qui étudie la structure des syntagmes et des phrases,
- la *sémantique* qui étudie les sens des mots, des locutions, des phrases ou des textes,
- la *pragmatique* qui étudie la relation entre le langage et son contexte d'utilisation.

Les techniques d'insertion se heurtent à plusieurs difficultés ; elles doivent être robustes aux distorsions de l'oral, suffisamment rapides selon le type de couplage MA-ML envisagé et doivent pouvoir s'intégrer dans un ML statistique pour sélectionner les hypothèses les plus probables. Ce dernier point complique notamment la prise en compte de connaissances symboliques au cours de la transcription [AG99].

#### 4.2.1 Phonologie et phonétique

Les ensembles de phonèmes et de phones d'une langue sont finis et il est possible d'établir des règles quant à la succession de leurs éléments. En français, l'emploi consécutif des deux phones [d] et [s] est par exemple illicite. Les connaissances issues de la phonologie et de la phonétique présentent la particularité de pouvoir être prises en compte à la fois dans le MA, le dictionnaire de prononciations et le ML.

Les MA les plus performants utilisent par exemple davantage d'informations sur le contexte de prononciation des phones que ce qui a été présenté en section 2.2.2 ; les états des HMM peuvent être des *triphones*, prenant en considération les influences des phones précédant et suivant le phone courant, plutôt que des phones. La prédiction des phénomènes de liaison entre les mots, dans des séquences telles que « *les enfants* » [BDMADG03], peut également être envisagée pour affiner l'utilisation du dictionnaire de prononciations par le MA. Cette section se consacrant plutôt à l'amélioration du ML, nous ne décrivons toutefois pas davantage les modifications possibles du MA et du dictionnaire de prononciations.

Peu d'informations sur la phonologie et la phonétique ont été intégrées dans le ML. La seule étude que nous ayons rencontrée porte sur la détermination des événements impossibles [LBSH03]. Ce sont des successions de mots qui ne peuvent se produire du fait de contraintes sur la langue. En français, l'expression « *je aime* » est par exemple interdite. La détermination de ces événements permet d'affiner le calcul des probabilités des modèles N-grammes, qui s'effectue au moyen de l'équation (6) (*cf.* page 13), en dénombrant l'ensemble des événements  $S$  (pour *Seen*) présents dans le corpus d'apprentissage. Les techniques de lissage prennent en compte les séquences de mots absentes de ce corpus mais elles ne distinguent pas les événements  $U$  (pour *Unseen*) non rencontrés mais pourtant possibles des événements  $I$  impossibles. Il paraît intéressant d'intégrer uniquement  $S$  et  $U$  dans l'évaluation des probabilités. Une des contraintes examinées dans l'étude, d'ordre phonologique, repose sur des règles d'élision, imposant que certains mots terminés par une voyelle ne peuvent être suivis par une voyelle (*e.g.* « *le arbre* »). Cette contrainte a permis de considérer comme impossible

une fraction très modeste (0,1 %) de l'ensemble des bigrammes et n'apporte donc pas de gain significatif au niveau de la qualité de la transcription.

#### 4.2.2 Morphologie

La structuration des mots constitue une source d'informations intéressante pour des langues morphologiquement très riches comme le turc, le russe ou l'arabe ou même pour des langues à haut taux de flexion comme le français. Dans le cas des langues agglutinantes notamment, pour avoir une couverture lexicale similaire à celle obtenue pour l'anglais, il est nécessaire d'envisager un nombre considérable de mots ; le fait de décomposer les mots en plusieurs constituants élémentaires permet de réduire le nombre d'événements à envisager lors du calcul des probabilités. Dans le cas des langues flexionnelles, l'analyse morphologique permet de rassembler dans une même classe des mots qui jouent le même rôle dans la phrase, comme « *mangeait* » et « *mangerons* » en français.

Les informations morphologiques peuvent être introduites dans un ML en utilisant des modèles N-grammes basés sur des classes (on parle alors de *modèles N-classes*), et non plus des modèles basés sur des mots. Les classes contiennent alors l'ensemble des mots possédant le même lemme<sup>8</sup>. Généralement, un mot est associé à un seul lemme mais il peut arriver en cas d'ambiguïté qu'il corresponde à plusieurs. Il en existe ainsi deux pour le participe passé « *plu* » : « *pleuvoir* » et « *plaire* ».

Dans un modèle N-classes, si  $\mathcal{C}_i$  représente l'ensemble des classes  $c_i$  auxquelles peut appartenir un mot  $w_i$ , le calcul des probabilités se fait de la manière suivante :

$$P(w_1^n) = \sum_{c_1 \in \mathcal{C}_1 \dots c_n \in \mathcal{C}_n} P(w_1^n c_1^n) \quad (29)$$

$$= \sum_{c_1 \in \mathcal{C}_1 \dots c_n \in \mathcal{C}_n} \prod_{i=1}^n P(w_i | w_1^{i-1} c_1^i) P(c_i | w_1^{i-1} c_1^{i-1}) \quad (30)$$

où  $P(w_i | w_1^{i-1} c_1^i)$  est appelée la *probabilité lexicale* et  $P(c_i | w_1^{i-1} c_1^{i-1})$  la *probabilité contextuelle*. En supposant que la probabilité de  $w_i$  dépend uniquement des classes  $c_i$  auxquelles il peut appartenir dans le groupe de souffle, on obtient :

$$P(w_1^n) \approx \sum_{c_1 \in \mathcal{C}_1 \dots c_n \in \mathcal{C}_n} \prod_{i=1}^n P(w_i | c_i) P(c_i | w_1^{i-1} c_1^{i-1}) \quad (31)$$

De manière similaire aux modèles N-grammes de mots qui ne prennent en compte que les  $N - 1$  mots précédents dans le calcul,  $P(c_i | w_1^{i-1} c_1^{i-1})$  est approximé en ne considérant que

<sup>8</sup>Forme « simple » d'un mot obtenue par un processus de *lemmatisation*. Le lemme associé à un verbe conjugué sera par exemple sa forme à l'infinitif ; pour un adjectif ou un nom, ce sera sa forme au masculin singulier. Cette notion permet d'associer à un même lemme l'ensemble de mots qui ne se distinguent que par la flexion.

les  $N - 1$  classes attribuées précédemment :

$$P(w_1^n) \approx \sum_{c_1 \in \mathcal{C}_1 \dots c_n \in \mathcal{C}_n} \prod_{i=1}^n P(w_i | c_i) P(c_i | c_{i-N+1}^{i-1}) \quad (32)$$

L'intérêt principal des modèles N-classes est de réduire considérablement le nombre d'événements possibles par rapport aux modèles N-grammes puisque le nombre total de classes est généralement beaucoup plus petit que la taille du vocabulaire du système de RAP. Ceci permet de limiter le recours aux techniques de lissage, qui introduisent des approximations lors du calcul des probabilités.

Les modèles N-classes utilisant les lemmes sont généralement combinés, par interpolation linéaire, avec des modèles N-classes à base de PoS (*cf.* section 4.2.3) [EBD90, MM92]. Ceci permet d'améliorer légèrement la perplexité par rapport à un simple modèle N-classes de PoS. Le comportement n'a toutefois pas été testé dans des systèmes de RAP.

Une alternative aux modèles N-classes est le *ML factorisé*. Ce ML décompose chaque mot  $w_i$  en  $k$  caractéristiques  $f_i^{1:k}$ , aussi appelées facteurs. Elles représentent des informations morphologiques, syntaxiques ou sémantiques sur le mot, en plus du mot lui-même. Les ML factorisés probabilistes utilisant une approximation trigramme calculent les probabilités grâce à l'expression :

$$P(f_1^{1:k}, f_2^{1:k} \dots f_n^{1:k}) = \prod_{i=3}^n P(f_i^{1:k} | f_{i-2}^{1:k}, f_{i-1}^{1:k}) \quad (33)$$

Ces ML factorisés ont notamment été appliqués à l'arabe, avec comme facteurs le radical<sup>9</sup>, la racine<sup>10</sup> et la classe morphologique. Leur utilisation au sein d'un système de RAP pour transcrire des émissions d'actualité a permis de réduire le taux d'erreur sur les mots de 57,6 % à 56,1 % [VKDS04]. Un inconvénient principal est qu'ils nécessitent des adaptations importantes pour être intégrés dans le processus de transcription, celui-ci étant généralement conçu pour opérer sur des mots, plutôt que sur des facteurs.

### 4.2.3 Syntaxe

L'information syntaxique peut être utilisée sous plusieurs formes : les ML peuvent ainsi tenir compte des classes grammaticales attribuées aux mots, d'une nouvelle segmentation du groupe de souffle en réunissant plusieurs mots au sein d'un même constituant, ou encore d'une analyse syntaxique du groupe de souffle. Nous présentons successivement un état de l'art des tentatives faites pour intégrer ces divers types de connaissances.

<sup>9</sup>Support morphologique d'un mot. C'est la partie qui contient le sens d'un mot, après avoir supprimé tout ce qui relevait de la flexion. Par exemple, les radicaux respectifs de « *déstabiliser* » et « *nationaliser* » sont « *déstabilis-* » et « *nationalis-* ».

<sup>10</sup>Constituant d'un mot qui porte la partie principale de son sens. À la différence du radical, la racine ne peut pas être décomposée en d'autres éléments porteurs de sens ou morphologiquement simples. Ainsi, les racines respectives de « *déstabiliser* » et « *nationaliser* » sont « *stabil-* » et « *nation* ».

### Parties du discours et classes statistiques

Il existe des successions de parties de discours (PoS) (*cf.* section 3.4.2) qui se produisent de manière rarissime dans une langue donnée. Il est ainsi très peu fréquent que deux noms communs se suivent en français. Des ML basés sur ce principe ont été conçus en considérant soit des classes syntaxiques déterminées *a priori* (en général des PoS souvent accompagnées d'informations morphologiques sur le genre, le nombre, le temps ou encore le mode), soit des classes produites automatiquement par des méthodes statistiques.

Les PoS sont généralement intégrées au système de RAP au moyen de modèles N-classes, où chaque PoS correspond à une classe [MM92, Goo01]. Les modèles N-classes sont toutefois moins performants que les modèles N-grammes, en considérant des historiques de même longueur. On observe en revanche dans certains cas une amélioration de la perplexité par rapport aux modèles N-grammes quand on combine des modèles N-classes avec des modèles N-grammes. Cette baisse de la perplexité reste cependant peu importante même quand on dispose de suffisamment de données pour apprendre les paramètres du ML. Diverses améliorations des modèles N-classes ont donc été envisagées.

Dans une première solution, le mode d'intégration dans le système de RAP des modèles N-classes utilisant les PoS est révisé. Au lieu d'être employé sur le graphe de mots produit aux cours des passes précédentes, le ML est appliqué de manière plus sélective sur le graphe d'homophones établi à partir de la meilleure hypothèse trouvée (*cf.* section 4.1.1). Le rôle du modèle N-classes est alors de sélectionner un homophone possible parmi ceux générés pour chaque hypothèse de mots [BNSd99]. Ce mode est particulièrement adapté pour corriger des fautes d'accord en genre et en nombre, notamment en français où les formes d'un même mot au singulier et au pluriel sont souvent homophones. L'application d'un modèle N-classes adoptant cette méthode a permis de réduire le taux d'erreur de 10,7% à 10,5% sur la transcription d'émissions d'actualité en français [GAAD<sup>+</sup>05].

Une deuxième solution consiste à reconsidérer le mode de calcul des probabilités. En modifiant l'approximation faite sur la probabilité lexicale dans l'équation (31) par :

$$P(w_i | w_1^{i-1} c_1^i) \approx P(w_i | c_{i-N+1}^i) \quad (34)$$

une légère amélioration de l'entropie croisée a été observée [Goo01]. On peut même aller jusqu'à supprimer les approximations faites à la fois sur les probabilités lexicale et contextuelle. Une étude [Hee99] propose ainsi de redéfinir l'objectif d'un système de RAP (*cf.* section 2.2), de manière à ce que les PoS  $C$  associées aux mots  $W$  à reconnaître soient considérées comme partie intégrante de la sortie de la transcription et non plus comme des objets intermédiaires. La finalité de la RAP revient alors à estimer :

$$\hat{W}, \hat{C} = \arg \max_{W, C} P(W, C | A) \quad (35)$$

En éliminant les approximations, le nombre d'événements à examiner pour évaluer les probabilités augmente considérablement ; une méthode basée sur des arbres de décision a donc été élaborée. Cette technique a donné des résultats satisfaisants pour transcrire des dialogues portant sur des sujets précis puisque le ML triclassé modifié a conduit à une réduction du

taux d'erreur de 26,0 % à 24,9 % par rapport à des ML trigrammes, tandis que le ML triclasse conventionnel faisait quant à lui augmenter le taux d'erreur.

Une troisième solution pour améliorer les modèles N-classes est de modifier la construction des classes. Au lieu d'avoir une classe par PoS, une possibilité est de regrouper au sein d'une même classe l'ensemble des mots ayant les mêmes PoS possibles avec le même ordre de vraisemblance. Ceci supprime l'ambiguïté des classes que l'on peut associer aux mots. Bien que l'attribution d'une PoS à un mot soit une technique relativement maîtrisée (cf. section 3.4.2), il demeure toujours des erreurs qui peuvent venir perturber le calcul des probabilités. Un ML triclasse utilisant ce type de classes, combiné avec un ML trigramme, a permis de réduire le taux d'erreur sur les mots pour traiter de la parole lue [SR99].

On peut également envisager un système à nombre très restreint de classes, en poussant à l'extrême la propriété principale des modèles N-classes, qui est de réduire le total des événements à envisager pour le calcul des probabilités. Normalement, il est possible de considérer d'une vingtaine à une centaine de PoS différentes ; dans un nouveau système, on ne discerne que deux types de PoS : les *classes ouvertes*, correspondant aux mots lexicaux, et les *classes fermées*, correspondant aux mots grammaticaux (cf. section 3.3). Cette distinction est faite avec l'idée que la séquence des mots grammaticaux reflète les contraintes syntaxiques du groupe de souffle, alors que la séquence de mots lexicaux est contrôlée par des relations sémantiques entre les mots. Les classes ouvertes et fermées sont généralement utilisées en adaptant les modèles N-classes [IM94, Geu96]. Dans le cas où on considère un historique de longueur deux, seuls les derniers mots lexical et grammatical rencontrés sont considérés. Si  $o_{i-1}$  est le dernier mot de classe ouverte et  $f_{i-1}$  le dernier mot de classe fermée dans  $w_1^{i-1}$ , le calcul des probabilités se fait comme suit :

$$P(w_i|w_1^{i-1}) \approx \begin{cases} P(w_i|w_{i-1}, o_{i-1}) & \text{si } w_{i-1} \text{ est un mot de classe fermée} \\ P(w_i|w_{i-1}, f_{i-1}) & \text{si } w_{i-1} \text{ est un mot de classe ouverte} \end{cases} \quad (36)$$

La combinaison d'un tel modèle avec un ML trigramme a permis de réduire le taux d'erreur de 29,4 % à 29,0 % pour transcrire un corpus de parole spontané en allemand [Geu96]. La prise en compte de ces deux types de classes peut être faite différemment en utilisant un modèle N-grammes classique pour les classes ouvertes et un modèle spécialement conçu pour les mots de classes fermées [PS01]. Le modèle spécifique prédit alors les mots grammaticaux à partir des  $M - 1$  mots grammaticaux précédents. La conception de ce modèle est justifiée par le fait qu'en anglais par exemple, les mots des classes fermées représentent 30 % du langage écrit et sont en moyenne distants de 1,9 mots. L'utilisation de ce modèle a permis une légère amélioration de la perplexité par rapport à un ML trigramme.

Outre les modèles N-classes, la connaissance sur les PoS peut être introduite grâce aux modèles de cache (cf. section 2.2.3). Un cache, limité à 200 mots, peut être construit pour chaque PoS [KDM90]. Le cache d'une PoS donnée contiendra alors les derniers mots rencontrés, étiquetés par cette PoS. Ces modèles sont conçus avec l'idée que chaque PoS a une répartition particulière d'occurrences. Les mots lexicaux ont ainsi tendance à apparaître par vagues, au gré des sujets traités, tandis que les mots grammaticaux sont répartis plus uniformément. Ce type de modèle possède des propriétés intéressantes puisque la combinaison

d'un modèle de cache avec un ML trigramme a conduit à une réduction de perplexité d'un facteur supérieur à trois par rapport à un modèle trigramme [KDM90]. Un modèle assez similaire a été également conçu pour discriminer les formes singulier et pluriel homophones [BNSd99].

Une alternative aux PoS déterminées *a priori* consiste à construire des classes statistiques [BDPd<sup>+</sup>92, KN93, TK95, FIO96, Jar96]. L'objectif de ces classes est de regrouper les mots qui ont le même « comportement ». Elles peuvent ainsi réunir ceux qui possèdent les mêmes PoS et qui sont liés sur le plan sémantique, comme par exemple « *gens* », « *hommes* », « *femmes* » et « *enfants* ». Elles sont obtenues par des méthodes de classification automatique, en cherchant notamment à maximiser la perplexité sur un ensemble de test ou à maximiser l'information mutuelle moyenne entre les classes. Lors du processus de construction des classes statistiques, la position du mot peut être prise en compte. Des classes dites *prédictives* peuvent ainsi être créées quand le mot est situé en position  $w_i$  dans le calcul de  $P(w_i|w_1^{i-N+1})$ , tandis que d'autres dites *conditionnelles* correspondent aux mots positionnés dans l'historique  $w_1^{i-N+1}$  [YS99]. Si on considère par exemple en anglais « *a* » et « *an* », ils peuvent suivre les mêmes mots puisque ce sont deux articles indéfinis ; ils appartiennent par conséquent à la même classe prédictive. En revanche, du fait qu'il existe très peu de mots qui peuvent à la fois se positionner après « *a* » et « *an* », ils sont associés à deux classes conditionnelles différentes. L'utilisation de deux types de classes permet ainsi d'introduire des connaissances de granularité plus fine.

Généralement, on attribue une seule classe statistique à chaque mot, contrairement aux ML à base de PoS où plusieurs PoS sont associables à un mot. Cette propriété permet de simplifier le calcul des probabilités de l'équation (32) :

$$P(w_1^n) \approx \prod_{i=1}^n P(w_i|c_i)P(c_i|c_{i-N+1}^{i-1}) \quad (37)$$

Des modèles triclassés utilisant des classes statistiques ont permis de réduire la perplexité par rapport à des modèles trigrammes, que ce soit pour l'anglais, le français ou l'allemand [Jar96]. Les gains en perplexité sont plus importants pour des langues à haut taux de flexion telles que le russe [WW01]. L'intégration des classes statistiques dans des modèles varigrammes (*cf.* section 2.2.3) a conduit quant à elle à une baisse du taux d'erreur relative de 7% par rapport à des modèles varigrammes sans classe pour transcrire de la parole lue en anglais [Bla99]. Il semble que ces classes statistiques conduisent à une réduction plus grande de perplexité que les PoS [NWW98]. Toutefois, le fait que la perplexité diminue davantage quand on combine les modèles N-classes à base de classes statistiques avec ceux à base de PoS [KN93, PMVGL03] laisse penser que les deux types de classes portent des informations complémentaires. Les bons résultats des classes statistiques pourraient être expliqués par le fait que celles-ci dépassent généralement une taille critique, ce qui permet aux modèles N-classes de ne pas faire intervenir des nombres d'occurrences faibles dans le calcul des probabilités, contrairement aux modèles N-grammes. La baisse du taux d'erreur, constatée en classifiant uniquement les mots peu fréquents pour transcrire de la parole spontanée dans

le domaine du trafic aérien, va dans ce sens [FIO96]. Il apparaît toutefois que l'association des mots rares à leur classe par des méthodes automatiques est peu fiable, ce qui diminue l'intérêt des modèles N-classes par rapport aux techniques de lissage pour évaluer la probabilité des événements absents de l'ensemble d'apprentissage [Ros00a].

En sus des modèles N-classes, les réseaux de neurones constituent un autre mode de construction des ML à base de classes statistiques [BDVJ03]. Leur particularité est de représenter chaque mot non plus par une classe discrète, mais par un vecteur de caractéristiques à  $m$  dimensions (e.g.  $m = 30$ ), où  $m$  est beaucoup plus petit que la taille du vocabulaire. L'objectif du nouvel espace est que deux mots qui jouent des rôles syntaxiques et sémantiques similaires, comme « *chat* » et « *chien* » par exemple, aient des vecteurs de caractéristiques proches. La projection de chaque mot dans le nouvel espace et les valeurs  $P(w_i|w_{i-N+1}^{i-1})$  sont déterminées simultanément lors de la phase d'apprentissage d'un réseau de neurones, où la première couche cachée représente les vecteurs de caractéristiques des mots  $w_{i-N+1}^{i-1}$  et les sorties les probabilités  $P(w_i|w_{i-N+1}^{i-1})$  pour chaque mot  $w_i$  du vocabulaire. Les ML à base de neurones se révèlent performants puisqu'ils ont permis de réduire le taux d'erreur de 22,6 % à 21,8 % pour transcrire de la parole conversationnelle en anglais, par rapport à un ML quadrigramme [SG04].

Enfin, une autre application envisageable de la connaissance des classes statistiques est de réduire le nombre d'événements impossibles (cf. section 4.2.1). Les mots sont regroupés en classes selon leurs ressemblances syntaxiques et sémantiques et reçoivent l'étiquette de leur classe. L'ensemble des suites possibles de deux classes consécutives étant connu, les séquences impossibles sont celles qui sont absentes d'un corpus représentatif de la langue, voire présentes de manière peu significative si on tient compte des erreurs de typographie de ce corpus. Cette technique a également été étendue aux PoS, en remarquant que des règles de certaines langues prohibent des successions telles que « [ARTICLE DÉFINI PLURIEL] [NOM MASCULIN SINGULIER] ». Dans l'expérience décrite dans [LBSH03], la prise en compte de 200 classes statistiques et de 233 PoS a ainsi permis de réduire de près de 15 % le nombre de bigrammes possibles.

### Multimots

Une autre utilisation de connaissances syntaxiques consiste à regrouper plusieurs mots au sein d'unités d'ordre supérieur, comme des locutions ou des syntagmes, et à les ajouter au vocabulaire du système de RAP. Ces groupes de mots, que nous désignons par la suite par le terme *multimots*, peuvent être de natures très diverses. Il peut s'agir de mots qui cooccurrent fréquemment dans un corpus, tels que « *demain matin* » ou encore « *millions de dollars* », de mots composés comme « *New York* » ou « *vice président* », ou d'entités nommées concernant des dates ou des noms de personne. Dans le cas d'applications ciblées, ce peut être des expressions propres à un domaine telles que « *vous êtes la bienvenue* » ou « *pouvez-vous s'il vous plaît me mettre en contact avec* ».

Les multimots sont généralement sélectionnés par des méthodes automatiques parmi l'ensemble des combinaisons possibles de mots du vocabulaire du système de RAP. Il existe deux approches principales pour les obtenir : les méthodes purement statistiques et les



systèmes à base de règles, même si cette distinction est parfois un peu arbitraire, certaines techniques combinant les deux approches [BNSd99]. Parmi les méthodes purement statistiques figurent les modèles multigrammes, déjà évoqués en section 2.2.3, qui choisissent pour chaque groupe de souffle la meilleure segmentation parmi plusieurs possibles, en maximisant la probabilité d'observation. Ce type de modèle détermine les multimots en ligne pour chaque nouveau groupe de souffle, en fixant le nombre maximal de mots qu'ils peuvent contenir. La plupart des méthodes statistiques prennent cependant en compte les multimots en les sélectionnant hors ligne suivant un critère donné, tel que les fréquences de cooccurrences, ou encore l'évaluation de la perplexité par validation croisée, puis en les intégrant au vocabulaire du système de RAP au même titre que les mots [SW94, RBW96, KR99]. Parmi les méthodes à base de règles, on peut citer celles utilisant les automates à états finis [NEB<sup>+</sup>99, BNSd99] ou les grammaires non contextuelles probabilistes [GW98, WMH00, MSZ02, SWH03].

L'intérêt principal des multimots par rapport aux mots est d'autoriser la prise en compte de phénomènes entre mots distants, tels que les accords en genre et en nombre, sans avoir à augmenter la taille de l'historique des ML [BNSd99]. Ces unités peuvent également être utilisées avec profit pour améliorer la modélisation acoustique des liaisons entre les mots. Elles permettent enfin d'introduire des connaissances spécifiques à un domaine en indiquant des expressions à reconnaître [KR99] ou encore des règles pour identifier certaines entités nommées [MSZ02].

La qualité des multimots détectés diffère selon le domaine étudié [KR99]. Dans des applications ciblées, telles que le routage d'appels téléphoniques, les constructions stéréotypées telles que « *can you please get me* » sont très naturellement modélisées. Le choix des multimots se révèle plus délicat dans le cadre de la parole lue, où les constructions et le vocabulaire sont riches et variés. L'apport des multimots dans des modèles trigrammes est d'ailleurs plus important en domaines spécialisés.

L'ajout de multimots au vocabulaire ayant l'inconvénient d'augmenter le nombre d'événements possibles, des modèles N-classes sont souvent envisagés [DS98, RBW96, NEB<sup>+</sup>99]. Les classes peuvent indifféremment contenir des mots et des multimots [DS98, RBW96], ce qui peut notamment permettre de regrouper des expressions différant uniquement par la présence de disfluences (*cf.* section 3.3), comme « *eh je veux* » et « *je veux* ». Ce type de modèle a conduit à une réduction du taux d'erreur de 29,5 % à 27,9 % par rapport à un ML trigramme pour transcrire de la parole spontanée en allemand [RBW96]. Les classes peuvent également ne concerner que les multimots, ce qui est surtout le cas quand on utilise des automates d'états finis ou des grammaires [NEB<sup>+</sup>99, GW98, WMH00, MSZ02, SWH03]. Les multimots reconnus par une même règle et correspondant à un concept précis, tels que le type de précipitation ou le nom d'une ville, se retrouvent ainsi au sein d'une même classe. Cette méthode a permis de faire baisser le taux d'erreur par rapport à des modèles N-classes standard pour des requêtes sur la météo (de 18,3 % à 18,0 %) ainsi que dans le domaine du trafic aérien (de 15,6 % à 15,0 %) [SWH03].

### Analyse syntaxique

Nous venons de voir, au travers de l'extraction de multimots, un emploi un peu particulier

de l'analyse syntaxique. Les techniques présentées ici l'utilisent dans un cadre plus formel, en prenant en compte la structure du groupe de souffle. Les *grammaires non contextuelles*, notamment leurs versions probabilistes, ont été le premier type d'analyse envisagé dans les systèmes de RAP [SS94, JWS<sup>+</sup>95, Sen92, SMZ95, LBS04]. Le mode d'attribution des probabilités pour ces grammaires est souvent jugé limité puisqu'il ne dépend que du non-terminal de la partie gauche de chacune des règles; d'autres formes d'analyse sont souvent préférées. Les *grammaires lexicalisées* [JM00] étendent ainsi les grammaires non contextuelles en choisissant pour chaque constituant détecté un mot jouant le rôle de tête. Le calcul des probabilités est alors conditionné par la tête [CJ00] et parfois par d'autres non-terminaux rencontrés auparavant dans l'analyse [Roa01, Cha01]. Les *grammaires de dépendance* offrent elles aussi un plus haut degré de lexicalisation que les formes non contextuelles, en déterminant les liens qui s'établissent entre les mots (par exemple à l'aide de *grammaires de liens* [LST92, BP98]) ou en exprimant les dépendances par des contraintes syntaxiques et sémantiques (on parle alors de *grammaires de dépendance par contraintes* [HH95, WH02]). Si les règles de toutes ces grammaires sont généralement déterminées *a priori*, l'évaluation des probabilités se fait automatiquement soit à partir de corpus annotés syntaxiquement tels que le Penn Treebank, soit à partir de corpus analysés par des méthodes automatiques.

Le principal avantage de ce type de connaissance est de prendre en compte les dépendances syntaxiques entre les constituants d'un même groupe de souffle, et ce, même si ces constituants se trouvent à des positions assez éloignées. Le calcul des probabilités intègre alors des informations plus précises que l'influence des  $N - 1$  mots précédents sur le mot courant, comme cela est le cas des modèles N-grammes.

La difficulté la plus importante à laquelle se trouvent confrontées les méthodes d'analyse syntaxique concerne la conception de grammaires suffisamment robustes. Il est déjà difficile d'élaborer pour l'écrit des analyseurs syntaxiques ayant une large couverture, même si certaines grammaires lexicalisées parviennent à une précision<sup>11</sup> et un rappel<sup>12</sup> proches de 90 % pour analyser une partie du corpus du *Wall Street Journal* [Cha00]. Les difficultés intrinsèques de l'oral (*cf.* section 3), notamment le manque de ponctuation, la présence de disfluences et l'absence éventuelle de majuscules dans la transcription, ajoutées aux erreurs de reconnaissance des systèmes de RAP, compliquent encore la réalisation d'analyseurs. L'utilisation de l'analyse syntaxique par les systèmes de RAP a, pour ces raisons, longtemps été confinée à des applications homme-machine où les tournures de phrase et le vocabulaire étaient très spécifiques [Sen92]. Les méthodes d'analyse partielle, qui ne nécessitent pas de construire un arbre syntaxique décrivant la structure détaillée du groupe de souffle entier, sont particulièrement adaptées pour concevoir des solutions robustes. Un ML a ainsi été défini en segmentant les suites de mots à analyser en constituants non récursifs, que l'on appelle *chunks* [ZW98].

La prise en compte de l'analyse syntaxique au niveau du ML s'envisage différemment selon que l'on associe ou non des probabilités aux règles de la grammaire. Dans le cas des

<sup>11</sup>Définie par  $\frac{\text{nb de constituants communs à GOLD et à TEST}}{\text{nb de constituants dans TEST}}$ , où GOLD et TEST sont les arbres d'analyse obtenus respectivement manuellement et automatiquement.

<sup>12</sup>Défini par  $\frac{\text{nb de constituants communs à GOLD et à TEST}}{\text{nb de constituants dans GOLD}}$ .

versions non probabilistes, on distingue trois types d'intégration pour guider le choix des hypothèses de transcription. Une solution simple consiste à utiliser l'analyse pour filtrer les  $N$  meilleures hypothèses fournies les unes après les autres par le système de RAP, en arrêtant dès qu'une hypothèse a pu être analysée entièrement. Une deuxième solution repose sur le calcul du nombre de mots du groupe de souffle couverts par l'analyse. Plus ce nombre est grand, meilleure est considérée l'hypothèse [ZW98]. Dans une dernière solution, un modèle  $N$ -grammes est construit en examinant les séquences de constituants de l'analyse et non plus celles de mots. Dans le cas où on utilise des *chunks*, l'examen du groupe de souffle « [NP *you*] [VC *weren't born*] [ADVP *just*] [NP *to soak up*] [NP *sun*]<sup>13</sup> » conduira ainsi à l'étude de la séquence « NP VC ADVP VC NP » [ZW98].

L'intégration des versions probabilistes dans les systèmes de RAP dépend quant à elle du type de grammaire envisagé. En ce qui concerne les grammaires non contextuelles probabilistes, il existe des algorithmes rapides [JL91, Sto95] permettant d'estimer avec exactitude les probabilités  $P(w_1^i) = P(S \xrightarrow{*} w_1 w_2 \dots w_i \dots)$  des chaînes préfixes  $w_1^i$ . Ces techniques peuvent être utilisées pour définir des probabilités de modèles bigrammes à partir de grammaires [SS94, JWS<sup>+</sup>95], en considérant que :

$$P(w_i | w_1^{i-1}) = \frac{P(w_1^i)}{P(w_1^{i-1})} \quad (38)$$

Un autre procédé introduisant une grammaire non contextuelle probabiliste dans un ML consiste à la convertir en un réseau stochastique. Dans le système TINA, les règles  $X \rightarrow ABC$  et  $X \rightarrow BCD$  sont par exemple transformées en un graphe (Fig. 12) dont les arcs sont valués par des probabilités apprises sur des exemples et dépendant de la partie gauche  $X$  de la règle utilisée ainsi que du mot qui vient d'être rencontré. Cette méthode présente l'avantage de fournir des probabilités explicites d'un mot, étant donné une séquence de mots [Sen92, SMZ95].

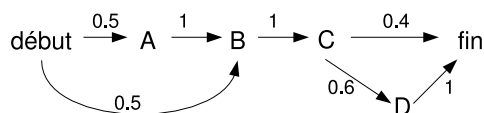


FIG. 12 – Réseau stochastique obtenu à partir de 2 règles ayant la même partie gauche

Une autre technique, employée aussi bien par les grammaires non contextuelles que par les lexicalisées, repose sur les probabilités associées aux arbres de dérivation [LBS04, Roa01, Cha01]. Si  $D_{w_1^i}$  représente l'ensemble des dérivations associées à  $w_1^i$ , on peut estimer  $P(w_1^i)$  par :

$$P(w_1^i) = \sum_{d \in D_{w_1^i}} P(d) \quad (39)$$

<sup>13</sup>Où NP = groupe nominal, VC = verbe complexe et ADVP = groupe adverbial.

où  $P(d)$  est obtenue en faisant le produit des probabilités de toutes les règles utilisées dans la dérivation  $d$ . Il en résulte que :

$$P(w_i|w_1^{i-1}) = \frac{P(w_1^i)}{P(w_1^{i-1})} = \frac{\sum_{d \in D_{w_1^i}} P(d)}{\sum_{d \in D_{w_1^{i-1}}} P(d)} \quad (40)$$

En pratique, l'ensemble  $D_{w_1^i}$  peut être très grand, ce qui conduit à utiliser une pile conservant uniquement les dérivations les plus probables.

Les modèles de langages structurés [CJ00], à base d'une grammaire lexicalisée, utilisent un mode de calcul légèrement différent de l'équation (40) :

$$P(w_i|w_1^{i-1}) = \frac{\sum_{d \in D_{w_1^{i-1}}} P(w_i, h_{-1,d}, h_{0,d})P(d)}{\sum_{d \in D_{w_1^{i-1}}} P(d)} \quad (41)$$

où  $h_{-1,d}$  et  $h_{0,d}$  représentent les têtes des deux derniers constituants de la dérivation  $d$ . Ce calcul a l'avantage de prendre explicitement en compte des dépendances à longue distance. Dans l'exemple de la figure 13, le mot « *after* » est ainsi prédit à partir de « *contract* » et « *ended* », et non à partir de « *7* » et « *cents* » comme cela serait le cas avec un ML trigramme classique.

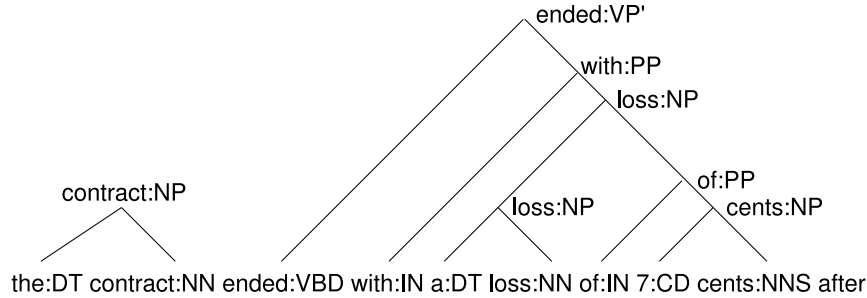


FIG. 13 – Arbre de dérivation partielle dont les feuilles sont des mots accompagnés de leurs PoS et dont les nœuds sont annotés par la tête et le type du constituant auquel ils sont associés (d'après [CJ00])

Il existe, selon les méthodes d'analyse syntaxique employées, des façons encore différentes d'intégrer les connaissances syntaxiques au ML. C'est par exemple le cas pour l'analyse par grammaire de liens qui utilisent d'autres méthodes probabilistes comme les modèles exponentiels (*cf.* section 4.1.2) [BP98]. La fonction de contraintes déterminée sur les liens du groupe de souffle s'exprime ici sous la forme :

$$f_{u \sim v}(w_1^i) = \begin{cases} 1 & \text{si } w = v \text{ et si } \exists u \in w_1^{i-2} \text{ tel que } u \text{ et } v \text{ sont liés} \\ 0 & \text{sinon} \end{cases} \quad (42)$$

Category: Verb
Features: {verdtype=past, voice=active, inverted=yes, gapp=yes, mood=whquestion, agr=all}
Role=G, Label=VP, (PX > MX) ( <i>Gouverné par un mot à sa gauche</i> )
Role=Need1, Label=S, (PX < MX) ( <i>Nécessite un modifieur à sa droite</i> )
Role=Need2, Label=S, (PX < MX) ( <i>Nécessite un modifieur à sa droite</i> )
Role=Need3, Label=S, (PX = MX) ( <i>Pas de modifieur</i> )
Dependent Positional Constraints: MX[G] < PX = MX[Need3] < MX[Need1] < MX[Need2]

FIG. 14 – Super-étiquette de « *did* » dans le groupe de souffle « *what did you learn* », où  $G$  représente le rôle de gouverneur,  $Need1$ ,  $Need2$  et  $Need3$  représentent des contraintes grammaticales sur le mot « *did* »,  $PX$  et  $MX$  représentent respectivement la position de « *did* » et d'un de ses modifieurs (d'après [WLH02])

La prédiction est réalisée à partir des liens établis avec le contexte gauche, mais elle peut également être faite en prenant en compte à la fois les contextes gauche et droit [LST92]. Les modèles construits à partir de ce type de grammaire sont très similaires aux ML à base de mots *triggers* [TN97]. Les paires de mots *triggers* sont deux mots qui apparaissent souvent dans le même contexte, tels que « *demande* » et « *répond* », et qui peuvent être assimilés à des bigrammes longue distance. Les relations entre mots *triggers* s'apparentent à des liens établis par des grammaires de liens mais elles présentent la particularité d'être déterminées *a priori* et non en fonction de l'analyse du groupe de souffle courant.

Enfin, les grammaires de dépendance par contraintes font quant à elles appel à des méthodes proches de celles rencontrées pour intégrer les PoS dans un ML. Le principe de ce type d'analyse est d'attribuer à chaque mot un rôle en fonction de sa position dans le groupe de souffle. Ce rôle, qui contient des informations lexicales, syntaxiques mais aussi sémantiques, est assimilable à une super-étiquette (Fig. 14), plus informative qu'une PoS, ce qui permet de prendre en compte ces grammaires en utilisant des modèles N-classes, où chaque classe est associée à une super-étiquette [WH02].

Les grammaires ont été principalement intégrées dans les systèmes de RAP pour des applications ciblées, du fait de leur manque de robustesse et de leur lenteur. Les non contextuelles probabilistes ont ainsi permis de diminuer le taux d'erreur dans les domaines de réservations pour le trafic aérien (de 34,6 % à 29,6 % par rapport à un modèle bigramme) [JWS<sup>+</sup>95] et de la restauration (de 6,9 % à 6,7 % par rapport à un modèle quadrigramme) [SMZ95]. Les grammaires lexicalisées ont quant à elles conduit à une baisse de la perplexité par rapport aux modèles trigrammes sur un corpus du *Wall Street Journal* [CJ00, Roa01, Cha01]. Les progrès récents de l'analyse syntaxique, obtenus en prenant en compte un contexte de plus en plus grand lors de l'attribution des probabilités à chacune des règles, ont ainsi autorisé son utilisation sur de la parole lue. En outre, leur combinaison avec des modèles N-grammes

a fait baisser davantage encore la perplexité, ce qui semble indiquer qu'elles apportent des connaissances supplémentaires à ces modèles. Les modèles conçus avec ces grammaires sont cependant trop coûteux en calculs pour être utilisés au sein de systèmes de RAP dans des applications à grande échelle, même si l'un de ces modèles a pu être intégré dans un ML multipasses en rescorant des graphes de mots [HJ04]. L'intégration des grammaires de dépendance par contraintes semble quant à elle plus aisée puisqu'elle a permis d'obtenir des ML de complexité raisonnable, et ce, en constatant une diminution du taux d'erreur pour transcrire des émissions d'actualité (de 14,7 % à 14,3 % par rapport à un modèle trigramme) [WHS03] mais aussi des conversations téléphoniques (avec une baisse relative de 6,2 % par rapport à des ML 4-grammes et 4-classes interpolés) [WSH04].

#### 4.2.4 Sémantique

L'introduction de connaissances sémantiques vise à favoriser les hypothèses qui possèdent plusieurs mots proches au niveau de leur sens, en supposant que les groupes de souffle à reconnaître ont une certaine cohérence sémantique. On attribuera ainsi un meilleur score à des hypothèses de décodage contenant les mots « *action* » « *obligation* » et « *bourse* » qu'à une hypothèse n'ayant pas de champ sémantique bien défini.

Une première possibilité pour introduire ce type d'information repose sur l'utilisation de connaissances *a priori*. Un dictionnaire, contenant les domaines d'emplois et les définitions des sens de 36 000 lemmes en anglais, a ainsi permis d'établir des associations sémantiques  $S(x, y)$  entre paires de mots [DAS97]. Pour ne pas avoir à désambiguïser les mots, les descriptions sémantiques de tous les sens associés à un même mot sont fusionnées.  $S(x, y)$  est alors fonction du nombre de mots en commun que possèdent les descriptions sémantiques de  $x$  et  $y$ . Un ML a été construit en calculant pour chaque groupe de souffle  $W$  un score  $Score(W)$  :

$$Score(W) = \frac{1}{k} \sum_{i \neq j} S(w_i, w_j) \quad (43)$$

où  $k$  est un facteur de normalisation dépendant de la taille de  $W$ . Ce ML n'a pas été comparé à des modèles N-grammes mais des expériences ont montré qu'il était informatif lorsque le contexte pris en compte était important, *i.e.*, que le groupe de souffle à analyser était long.

Une autre technique étudiée pour introduire des connaissances sémantiques est l'*analyse sémantique latente* [Bel98]. Le principe de cette méthode, utilisée en recherche d'information, est de trouver les relations sémantiques qui s'établissent entre les mots d'un document. Pour ce faire, elle suppose que deux mots sont proches s'ils ont tendance à apparaître dans les mêmes documents.  $W$ , la matrice d'occurrences de mots dans des documents d'un corpus d'apprentissage, est calculée de façon à ce que  $w_{ij}$  représente le nombre d'occurrences pondéré du mot  $w_i$  dans le document  $d_j$ . La pondération utilise un coefficient dépendant de  $w_i$ , qui vise à traduire que deux mots apparaissant avec les mêmes nombres d'occurrences dans  $d_j$  ne portent pas nécessairement autant d'information ; cela dépend également de leur distribution dans la collection entière de documents, ceux apparaissant dans beaucoup de documents étant considérés comme moins informatifs. Afin de réduire ses dimensions,  $W$  est

transformée dans un espace  $\mathcal{S}$  de dimension  $R \times R$  (avec  $100 \leq R \leq 200$ ), par des techniques proches de l'analyse en composantes principales :

$$W = USV^T \quad (44)$$

où  $S$  est une matrice diagonale de dimension  $R \times R$ ,  $U$  est une matrice dont les lignes  $u_i$  sont les représentations de chaque mot  $w_i$  dans  $\mathcal{S}$  et  $V$  est une matrice dont les lignes  $v_j$  sont les projections de chaque document  $d_j$  dans  $\mathcal{S}$ . D'après la définition de  $W$ , la manière dont  $w_i$  et  $d_j$  sont associés est déterminée par  $w_{ij}$ . Ceci peut être également caractérisé selon l'équation (44) par le produit scalaire de  $u_i S^{1/2}$  et  $v_j S^{1/2}$ , et une distance mesurant la proximité entre  $u_i$  et  $v_j$ , *i.e.*, entre  $w_i$  et  $d_j$ , est :

$$K(u_i, v_j) = \cos(u_i S^{1/2}, v_j S^{1/2}) \quad (45)$$

L'analyse sémantique latente permet de définir sur ce principe un nouveau type de ML [Bel98, Bel00]. Le calcul des probabilités  $P(w_i|h_i)$  est établi en considérant que  $h_i$  représente le document courant  $\tilde{d}_{i-1}$  jusqu'au mot  $w_i$ .  $\tilde{d}_{i-1}$  pouvant être vu comme une colonne supplémentaire de  $W$ , on calcule sa représentation  $\tilde{v}_{i-1}$  dans  $\mathcal{S}$ .  $P(w_i, \tilde{d}_{i-1})$  est alors déterminée à partir de la mesure de proximité  $K(w_i, \tilde{d}_{i-1})$  entre  $w_i$  et l'historique  $\tilde{d}_{i-1}$ , en réalisant une normalisation pour obtenir une probabilité bien définie. Grâce à ce mode de calcul,  $P(w_i, \tilde{d}_{i-1})$  est plus grande pour les mots dont le sens se rapproche le plus de celui des mots de l'historique et plus petite dans le cas contraire. Pour avoir des informations sur l'ordre des mots, le calcul des probabilités prend en outre en considération les N-grammes. L'analyse sémantique latente est une technique intéressante puisque ce type de ML a permis une baisse relative du taux d'erreur de 16 % par rapport à des ML trigrammes pour transcrire de la parole lue en anglais [Bel00].

#### 4.2.5 Pragmatique

L'introduction de la pragmatique dans la RAP recouvre un ensemble de techniques qui visent à adapter le processus de transcription au contexte d'utilisation. L'adaptation peut se faire de deux façons :

- soit en concevant des modèles adaptatifs, qui se spécialisent automatiquement au cours du processus de RAP, en fonction de ce que le système a déjà reconnu,
- soit en utilisant des corpus d'adaptation, obtenus par des systèmes de recherche d'information et plus proches du texte à transcrire que le corpus d'apprentissage.

#### Modèles adaptatifs

Il existe plusieurs méthodes d'adaptation des ML, parmi lesquelles on retrouve les modèles à base de cache (*cf.* section 2.2.3) et les modèles à base de mots *triggers* (*cf.* section 4.2.3). Dans le cas des modèles à base de cache, les probabilités  $P(w_i|h_i)$  sont réestimées en fonction des mots contenus dans le cache, en supposant que les mots qui sont apparus récemment voient leur probabilité d'apparition augmenter. Les modèles classiques à base de cache n'apportant pas un gain significatif, il en existe plusieurs variantes. Il a été par exemple envisagé

de ne considérer dans le cache que les mots rares [Ros94], ou bien encore de faire décroître les probabilités du cache de manière exponentielle en fonction de la distance entre le mot courant et les apparitions précédentes de ce mot [CR97].

Si ces modèles modifient uniquement les probabilités des mots présents dans le cache, ceux à base de mots *triggers* font l'hypothèse que certains mots présents dans l'historique ont une influence sur l'apparition des mots auxquels ils sont corrélés. Le mot « avion » pourra par exemple favoriser la prédiction du mot « vol ». Ce type de modèle peut ainsi prendre en compte des dépendances entre mots distants dans le document à transcrire [Ros96].

Un autre type de ML adaptatif repose sur la détection du thème traité par le document à transcrire. Il est en effet constaté que pour des documents sonores abordant plusieurs thèmes, notamment les émissions d'actualité, les tournures de phrase et les termes employés diffèrent selon le sujet traité. Dans cette approche, le corpus d'apprentissage est divisé en plusieurs ensembles correspondant chacun à un sous-langage ou un thème; un ML thématique est alors construit pour chacun de ces ensembles. Pour partitionner le corpus d'apprentissage, il est généralement supposé que chaque document du corpus est associé à un unique thème. Ces documents sont le plus souvent regroupés en utilisant des méthodes automatiques non supervisées [IO99, CR97, MLN97, FY99, GH99] mais peuvent être aussi classés grâce à des méthodes automatiques supervisées [LKMN05, KW99] ou de manière manuelle en annotant chaque document [KS93, Bru03]. Le nombre de groupes obtenu à partir du corpus d'apprentissage est très variable. Il peut par exemple être inférieur à 10 si les thèmes détectés sont généraux, comme l'économie, l'histoire ou la politique [Bru03, BDMS00] ou être de plus de 5 000, permettant ainsi une distinction plus fine du sujet abordé [SR97].

Il existe plusieurs procédés pour utiliser les modèles adaptatifs basés sur la détection de thèmes. Tout d'abord, la reconnaissance du sujet abordé par le document à transcrire peut se faire à différents niveaux. Elle est souvent réalisée à partir du document entier pour limiter l'influence des erreurs sur les mots reconnus, puisque ce sont des hypothèses de transcription, susceptibles d'être partiellement erronées, qui sont utilisées pour la détection du thème. Elle peut néanmoins être effectuée au niveau de chaque groupe de souffle. Malgré une difficulté plus importante, ce type de détection conduit dans certains cas à une plus grande baisse du taux d'erreur de la transcription [KW99]. Pour limiter les erreurs de détection du thème sur un groupe de souffle, les sujets peuvent être hiérarchisés; on utilise alors des modèles correspondant à des thèmes plus ou moins précis selon la confiance attribuée au groupe de souffle [LKMN05]. Un autre niveau de segmentation suppose de découper le document à transcrire pour obtenir des ensembles consécutifs de groupes de souffle associés à un seul thème. Les segments obtenus étant plus longs, cette solution présente l'avantage de limiter l'influence des erreurs de transcription, tout en détectant des thèmes plus précis qu'en prenant en compte l'ensemble du document [CGL<sup>+</sup>01].

Les modèles thématiques diffèrent en outre selon la manière dont ils sont combinés pour calculer les probabilités. Le ML correspondant au thème détecté à partir d'une première hypothèse de transcription peut être directement utilisé dans la deuxième passe du système de RAP [LKMN05]. Cependant, les modèles spécifiques à un thème ayant été appris sur des corpus de taille réduite, ils sont généralement combinés avec le modèle général, construit



à partir de l'ensemble du corpus d'apprentissage. Cette combinaison peut être réalisée au moyen d'une interpolation linéaire [Bru03, IO99, SR97] ou d'un modèle exponentiel (cf. section 4.1.2) [KW99]. L'avantage du dernier est qu'il modifie, pour chaque thème, les probabilités du modèle général pour un nombre limité de mots, ce qui permet de ne pas accroître beaucoup la taille du ML. Des méthodes propres à l'adaptation ont été de plus conçues pour combiner plusieurs modèles spécifiques à un thème. Si  $t_k$  est un thème présent dans le corpus d'apprentissage, les ML dits à *mélange de modèles* calculent les probabilités de la manière suivante :

$$P(w_i|w_1^{i-1}) = \sum_k \lambda_k(w_1^{i-1})P(w_i|w_1^{i-1}, t_k) \quad (46)$$

Ce procédé diffère de l'interpolation linéaire dans la mesure où les paramètres  $\lambda_k$  sont déterminés dynamiquement en fonction des mots rencontrés précédemment [KS93, MLN97, CR97]. Il existe une autre variante des mélanges de modèles, consistant à évaluer les paramètres  $\lambda_k$  au niveau des phrases (ou groupes de souffles) et non pas au niveau des N-grammes [IO99]. Notons que les coefficients  $\lambda_k$  peuvent être assimilés aux probabilités  $P(t_k|w_1^{i-1})$  [FY99, GH99]. L'intérêt de ce type de modèle est de pouvoir associer plusieurs thèmes au document à transcrire.

D'une manière générale, les modèles adaptatifs ont permis une réduction significative de la perplexité, mais cette réduction s'est traduite par une baisse limitée du taux d'erreur de transcription. Une propriété intéressante des modèles thématiques est qu'ils semblent apporter des informations complémentaires à celles fournies par une analyse syntaxique. En combinant ces deux sources d'informations, il a été ainsi constaté que les gains de chaque source sont presque additifs en ce qui concerne la baisse du taux d'erreur [WK99].

Une technique alternative d'adaptation consiste à utiliser des corpus spécialisés en fonction du texte à transcrire.

### Utilisation de corpus d'adaptation

Les ML statistiques sont particulièrement sensibles à l'adéquation du corpus d'apprentissage vis-à-vis du type de document à transcrire. Dans le cas d'émissions d'actualité, le poids attribué aux transcriptions manuelles de programmes radio est ainsi beaucoup plus important, eu égard à son volume de données, que celui de journaux écrits. On voit ainsi l'intérêt d'avoir un corpus adapté au document à transcrire.

Le corpus d'adaptation peut être déjà disponible, ce qui correspond à une adaptation *statique*. Par rapport au corpus général, ce corpus se rapporte davantage au domaine des textes à transcrire mais sa taille réduite, souvent de l'ordre de quelques milliers de mots, ne lui permet pas d'être directement utilisable comme corpus d'apprentissage du ML. Les paramètres du ML sont alors calculés principalement à partir du corpus général et, dans les situations où le volume de données est suffisant, à partir du corpus spécialisé [GLL00, Fed99].

Le corpus d'adaptation peut être aussi régulièrement mis à jour ; on parle alors d'adaptation *dynamique*. Cette mise à jour est effectuée soit à partir de documents qui ont été conçus durant la même période que ce qui est à transcrire [AG03, KW98], soit à partir de documents dont le contenu est similaire à ce que l'on souhaite décoder [BM98, CGLA03, BHD04].

Dans le deuxième cas, il est fait appel à un système de recherche d'information (SRI) ; un ML général établit alors une première hypothèse de transcription utilisée par le SRI pour retourner les documents les plus proches. L'ensemble des résultats de la recherche forme un corpus d'adaptation, permettant de rectifier le ML, et le ML modifié est utilisé dans une deuxième passe du processus de transcription.

La recherche d'information peut se faire à partir de sources de types différents. Une recherche parmi l'ensemble de documents du corpus d'apprentissage permet d'avoir un corpus spécialisé, ce qui est intéressant quand le corpus d'apprentissage se rapporte à plusieurs thèmes [CGLA03]. Une autre possibilité repose sur l'utilisation d'Internet [BHD04]. Cette vaste source d'informations évolutive présente également, au sein de certains sites tels que les blogs ou les serveurs de news, des caractéristiques proches de la langue parlée (*cf.* section 3.1) [VAR99].

Le système de RAP prend en compte les corpus d'adaptation de deux manières différentes :

- en modifiant son vocabulaire,
- en adaptant le calcul des probabilités du ML au corpus d'adaptation.

L'adaptation du vocabulaire est particulièrement pertinente pour transcrire des émissions d'actualité, des entités nommées apparaissant au gré des événements [KW98, AG03, BHD04]. La collecte régulière d'informations sur des sites de dépêches d'agences de presse ou de quotidiens nationaux permet de réduire le taux de mots hors vocabulaire. Pour ne pas augmenter la taille du vocabulaire, des mots présents initialement sont supprimés pour laisser place à des mots apparaissant souvent dans l'actualité récente. L'ajout d'un mot au système de RAP nécessite à la fois d'associer une transcription phonétique à ce mot et de l'inclure dans les distributions N-grammes du ML. La modification du vocabulaire présente donc l'inconvénient de devoir effectuer un réapprentissage fréquent des MA et des ML. Toutefois, une méthode a permis de construire un système de RAP à vocabulaire ouvert, avec une diminution constatée du taux d'erreur de 25,5 % à 24,9 % pour la transcription d'émissions d'actualité en français [AG05].

Le calcul des probabilités du ML peut être modifié selon plusieurs procédés pour prendre en compte le corpus d'adaptation. Une première possibilité consiste à utiliser les mélanges de modèles thématiques (*cf.* section 4.2.5). Les paramètres  $\lambda_k$  de l'équation (46) sont alors appris non plus à partir des mots transcrits précédemment comme dans le cas des modèles adaptatifs, mais à partir du corpus d'adaptation. Le mélange de modèles peut se faire en outre de manière dynamique en adaptant le calcul des probabilités  $P(w_i|w_1^{i-1}, t_k)$  au corpus d'adaptation [CGLA04]. D'autres techniques consistent à spécialiser un ML général sur le corpus d'adaptation en utilisant un critère de maximum *a posteriori* (MAP) [Fed96, BM98, CGLA04] ou un critère de minimum d'information discriminante (MDI pour *Minimum Discrimination Information*) [Fed99, CGLA04]. L'utilisation de corpus d'adaptation permet une réduction non négligeable du taux d'erreur de la transcription. Il a ainsi été constaté une baisse de 17,1 % à 16,3 % dans la transcription d'émissions d'actualité en anglais. Les adaptations au moyen de mélanges dynamiques de modèles et du critère MDI semblent être les procédés les plus performants [CGLA04].

## 5 Conclusion

La description du principe de fonctionnement de la transcription met en évidence que les systèmes de RAP actuels s'appuient sur une modélisation statistique. Ce type de conception a conduit à la construction de MA à base de HMM et de ML N-grammes. Il a permis de réaliser différents systèmes capables de transcrire des émissions d'actualité mais aussi des dialogues spontanés, même si le traitement de ce dernier type de documents est rendu problématique par la présence non négligeable de phénomènes de la langue parlée tels que les disfluences. Les connaissances linguistiques prises en compte se limitent bien souvent à la conception d'un lexique de prononciations et à l'apprentissage de ML N-grammes, ce qui laisse penser que des améliorations de la qualité de la transcription sont possibles en exploitant davantage d'informations sur le langage.

La présentation que nous avons faite des ML N-grammes et leurs variantes a montré que deux limitations sont souvent mises en avant quant à leur capacité à donner une valeur correcte aux probabilités de séquences de mots. D'une part, les ML N-grammes sont basés sur l'hypothèse un peu simpliste et arbitraire, mais conduisant à des méthodes de calculs rapides, qui est d'examiner uniquement les  $N-1$  mots précédents pour prédire le mot courant. Afin de remédier à ce premier point, des essais ont été conduits de façon à intégrer des connaissances supplémentaires telles que les dépendances syntaxiques ou les similarités sémantiques entre les mots d'un groupe de souffle. Des méthodes d'adaptation tentent également de prendre en compte des relations entre différents groupes de souffle. D'autre part, malgré des techniques de lissage perfectionnées, les calculs pour prédire des événements rares voire même absents du corpus sont imprécis. Pour tenter de corriger ce problème, des études ont réuni des mots possédant la même partie du discours ou le même lemme au sein d'une même classe, de manière à réduire le nombre d'événements possibles.

Au final, les résultats obtenus en intégrant des connaissances linguistiques supplémentaires montrent que les améliorations en terme de taux d'erreur sur les mots reconnus sont généralement assez peu significatives, d'autant plus que les nouveaux ML proposés ne sont presque jamais comparés avec des ML 5-grammes utilisant un lissage performant [Goo01]. Parmi les raisons qui expliquent cet état de fait, les nouvelles informations introduites par ces méthodes sont souvent redondantes avec les connaissances déjà apportées par les ML N-grammes de mots. En outre, les particularités des transcriptions produites automatiquement, notamment la flexibilité de la langue parlée, la segmentation en groupes de souffle ou encore les erreurs de reconnaissance, viennent compliquer la conception de méthodes extrayant automatiquement des connaissances linguistiques. De plus, les techniques employées ont souvent le défaut d'augmenter considérablement le temps de décodage du signal acoustique, ce qui fait qu'elles sont utilisées principalement au niveau de la dernière passe du processus de transcription.

Quelques méthodes apparaissent toutefois prometteuses pour corriger certaines erreurs de transcription : les ML utilisant des grammaires lexicalisées probabilistes, bien qu'ils soient encore trop coûteux au niveau des calculs, l'introduction de connaissances sémantiques, les modèles thématiques ou encore l'utilisation de corpus d'adaptation. Les modèles N-classes, qui sont rapides lors de leur utilisation, peuvent également réduire le taux d'erreur. La

combinaison de plusieurs types de connaissances semble de plus souhaitable pour apporter des informations complémentaires.

Notons enfin que cette synthèse s'est limitée à l'amélioration de la qualité de la transcription mais que le couplage TAL-RAP peut également s'exprimer à l'issue de la reconnaissance. De nombreuses recherches s'intéressent ainsi actuellement au repérage d'entités nommées dans les textes produits, à la détection de thèmes ou encore à la réalisation de résumés.

## Références

- [AB05] P. ALAIN et O. BOËFFARD. « Évaluation des modèles de langage N-gramme et N/M-multigramme ». Dans *Actes de la 12ème conférence sur le Traitement Automatique des Langues Naturelles (TALN)*, volume 1, Dourdan, France, 2005.
- [ADHB<sup>+</sup>04] M. ADDA-DECKER, B. HABERT, C. BARRAS, G. ADDA, P. BOULA DE MAREÛIL et P. PAROUBEK. « Une étude des disfluences pour la transcription automatique de la parole spontanée et l'amélioration des modèles de langage ». Dans *Actes des 25èmes Journées d'Études sur la Parole (JEP)*, Fès, Maroc, 2004.
- [AG99] J.-Y. ANTOINE et D. GENTHIAL. « Méthodes hybrides issues du TALN et du TAL Parlé : état des lieux et perspectives ». Dans *Actes de la 6ème conférence sur le Traitement Automatique des Langues Naturelles (TALN)*, Cargèse, France, 1999.
- [AG01] J.-Y. ANTOINE et J. GOULIAN. « Word Order Variations and Spoken Man-Machine Dialogue in French: a Corpus Analysis on the ATIS Domain ». Dans *Proc. of Corpus Linguistics*, Lancaster, Royaume-Uni, 2001.
- [AG03] A. ALLAUZEN et J.-L. GAUVAIN. « Adaptation automatique du modèle de langage d'un système de transcription de journaux parlés ». *Traitement Automatique des Langues (TAL)*, 44(1):11–31, 2003.
- [AG05] A. ALLAUZEN et J.-L. GAUVAIN. « Open Vocabulary ASR for Audiovisual Document Indexation ». Dans *Proc. of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, volume 1, Philadelphie, Pennsylvanie, États-Unis, 2005.
- [All94] J. B. ALLEN. « How do Humans Process and Recognize Speech? ». *IEEE Transactions on Speech and Audio Processing*, 2(4):567–577, 1994.
- [AMP<sup>+</sup>99] G. ADDA, J. MARIANI, P. PAROUBEK, M. RAJMAN et J. LECOMTE. « Métrique et premiers résultats de l'évaluation GRACE des étiqueteurs morphosyntaxiques pour le français ». Dans *Actes de la 6ème conférence sur le Traitement Automatique des Langues Naturelles (TALN)*, Cargèse, France, 1999.
- [BB90] C. BLANCHE-BENVENISTE. *Le français parlé : études grammaticales*. Paris : Éditions du CNRS, 1990.
- [BB97] C. BLANCHE-BENVENISTE. *Approches de la langue parlée en français*. Gap - Paris : Ophrys, 1997.
- [BCD<sup>+</sup>04] C. BENZITOUN, E. CAMPIONE, J. DEULOFEU, S. HENRY, F. SABIO, S. TESTON, A. VALLI et J. VÉRONIS. « L'analyse syntaxique de l'oral : problèmes et méthode ». Dans *Actes de la journée d'étude de l'ATALA sur l'annotation syntaxique de corpus*, Paris, France, 2004.

- [BDM98] N. BOUFADEN, S. DELISLE et B. MOULIN. « Analyse syntaxique robuste de dialogue retranscrits : peut-on vraiment traiter l'oral à partir de l'écrit ? ». Dans *Actes de la 5ème conférence sur le Traitement Automatique des Langues Naturelles (TALN)*, Paris, France, 1998.
- [BDMADG03] P. BOULA DE MAREÛIL, M. ADDA-DECKER et V. GENDNER. « Liaisons in French: a Corpus-Based Study Using Morpho-Syntactic Information ». Dans *Proc. of the 15th International Congress of Phonetic Sciences*, Barcelone, Espagne, 2003.
- [BDMS00] B. BIGI, R. DE MORI et T. SPRIET. « Reconnaissance thématique à partir de textes dictés et adaptation dynamique de modèles de langage thématiques ». Dans *Actes des 23èmes Journées d'Études sur la Parole (JEP)*, Aussois, France, 2000.
- [BDPd<sup>+</sup>92] P. F. BROWN, V. J. DELLA PIETRA, P. V. DESOUZA, J. C. LAI et R. L. MERCER. « Class-Based N-Gram Models of Natural Language ». *Computational Linguistics*, 18(4):467–480, 1992.
- [BDVJ03] Y. BENGIO, R. DUCHARME, P. VINCENT et C. JAUVIN. « A Neural Probabilistic Language Model ». *Journal of Machine Learning Research*, 3(2):1137–1155, 2003.
- [Bel98] J. R. BELLEGARDA. « A Multispan Language Modeling Framework for Large Vocabulary Speech Recognition ». *IEEE Transactions on Speech and Audio Processing*, 6(5):456–467, 1998.
- [Bel00] J. R. BELLEGARDA. « Large Vocabulary Speech Recognition with Multispan Statistical Language Models ». *IEEE Transactions on Speech and Audio Processing*, 8(1):76–84, 2000.
- [Ben04] C. BENZITOUN. « L'annotation syntaxique de corpus oraux constitue-t-elle un problème spécifique ? ». Dans *Actes de la 8ème Rencontre des Étudiants Chercheurs en Informatique pour le Traitement Automatique des Langues (RECITAL)*, Fès, Maroc, 2004.
- [BFHM98] E. BRILL, R. FLORIAN, J. C. HENDERSON et L. MANGU. « Beyond N-Grams: Can Linguistic Sophistication Improve Language Modeling? ». Dans *Proc. of the 36th Annual Meeting of the Association for Computational Linguistics and the 17th International Conference on Computational Linguistics (COLING-ACL)*, volume 1, Montréal, Canada, 1998.
- [BGWL01] C. BARRAS, E. GEOFFROIS, Z. WU et M. LIBERMAN. « Transcriber: Development and Use of a Tool for Assisting Speech Corpora Production ». *Speech Communication*, 33(1-2):5–22, 2001.
- [BHDM04] B. BIGI, Y. HUANG et R. DE MORI. « Vocabulary and Language Model Adaptation Using Information Retrieval ». Dans *Proc. of the 8th International Conference on Spoken Language Processing (ICSLP)*, volume 2, île de Jeju, Corée du Sud, 2004.

- [Bla99] R. BLASIG. « Combination of Words and Word Categories in Varigram Histories ». Dans *Proc. of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, volume 1, Phoenix, Arizona, États-Unis, 1999.
- [BM98] A. BERGER et R. MILLER. « Just-in-Time Language Modelling ». Dans *Proc. of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, volume 2, Seattle, Washington, États-Unis, 1998.
- [BNSd99] F. BÉCHET, A. NASR, T. SPRIET et R. DE MORI. « Modèles de langage à portée variable : application au traitement des homophones ». Dans *Actes de la 6ème conférence sur le Traitement Automatique des Langues Naturelles (TALN)*, Cargèse, France, 1999.
- [BP98] A. BERGER et H. PRINTZ. « Recognition Performance of a Large-Scale Dependency-Grammar Language Model ». Dans *Proc. of the 5th International Conference on Spoken Language Processing (ICSLP)*, volume 6, Sydney, Australie, 1998.
- [BP03] R. BEUTLER et B. PFISTER. « Integrating Statistical and Rule-Based Knowledge for Continuous German Speech Recognition ». Dans *Proc. of the 8th European Conference on Speech Communication and Technology (Eurospeech)*, Genève, Suisse, 2003.
- [BPLA95] F. BIMBOT, R. PIERACCINI, E. LEVIN et B. ATAL. « Variable-Length Sequence Modeling: Multigrams ». *Signal Processing Letters, IEEE*, 2(6):111–113, 1995.
- [Bra00] T. BRANTS. « TnT - A Statistical Part-of-Speech Tagger ». Dans *Proc. of the Sixth Applied NLP*, Seattle, Washington, États-Unis, 2000.
- [Bru03] A. BRUN. « Détection de thème et adaptation des modèles de langage pour la reconnaissance automatique de la parole ». Thèse de doctorat, Université Henri Poincaré - Nancy 1, France, 2003.
- [BV05] C. BENZITOUN et J. VÉRONIS. « Problèmes d’annotation d’un corpus oral dans le cadre de la campagne EASY ». Dans *Actes de la 12ème conférence sur le Traitement Automatique des Langues Naturelles (TALN)*, volume 2, Dourdan, France, 2005.
- [Cam01] E. CAMPIONE. « Étiquetage prosodique semi-automatique de corpus oraux : algorithmes et méthodologie ». Thèse de doctorat, Université de Provence, Aix-en-Provence, France, 2001.
- [Can00] M. CANDEA. « Contribution à l’étude des pauses silencieuses et des phénomènes dits d’«hésitation» en français oral spontané. Étude sur un corpus de récits en classe de français ». Thèse de doctorat, Université Paris III, France, 2000.
- [CG98] S. F. CHEN et J. GOODMAN. « An Empirical Study of Smoothing Techniques for Language Modeling ». Rapport technique, Harvard University, Cambridge, Massachusetts, États-Unis, 1998.

- [CGL<sup>+</sup>01] L. CHEN, J.-L. GAUVAIN, L. LAMEL, G. ADDA et M. ADDA-DECKER. « Using Information Retrieval Methods for Language Model Adaptation ». Dans *Proc. of the 7th European Conference on Speech Communication and Technology (Eurospeech)*, Aalborg, Danemark, 2001.
- [CGLA03] L. CHEN, J.-L. GAUVAIN, L. LAMEL et G. ADDA. « Unsupervised Language Model Adaptation for Broadcast News ». Dans *Proc. of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, volume 1, Hong Kong, Chine, 2003.
- [CGLA04] L. CHEN, J.-L. GAUVAIN, L. LAMEL et G. ADDA. « Dynamic Language Modeling for Broadcast News ». Dans *Proc. of the 8th International Conference on Spoken Language Processing (ICSLP)*, île de Jeju, Corée du Sud, 2004.
- [Cha00] E. CHARNIAK. « A Maximum-Entropy-Inspired Parser ». Dans *Proc. of the 1st Conference of the North American Chapter of the Association for Computational Linguistics*, Seattle, Washington, États-Unis, 2000.
- [Cha01] E. CHARNIAK. « Immediate-Head Parsing for Language Models ». Dans *Proc. of the 39th Annual Meeting of the Association for Computational Linguistics (ACL)*, Toulouse, France, 2001.
- [CJ00] C. CHELBA et F. JELINEK. « Structured Language Modeling ». *Computer Speech and Language*, 14(4):283–332, 2000.
- [CR89] Y.-L. CHOW et S. ROUKOS. « Speech Understanding Using a Unification Grammar ». Dans *Proc. of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, volume 2, Glasgow, Royaume-Uni, 1989.
- [CR97] P. R. CLARKSON et A. J. ROBINSON. « Language Model Adaptation Using Mixtures and an Exponentially Decaying Cache ». Dans *Proc. of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, volume 2, Munich, Allemagne, 1997.
- [CR99] P. CLARKSON et T. ROBINSON. « Towards Improved Language Model Evaluation Measures ». Dans *Proc. of the 6th European Conference on Speech Communication and Technology (Eurospeech)*, volume 5, Budapest, Hongrie, 1999.
- [CRAR99] J.-C. CHAPPELIER, M. RAJMAN, R. ARAGÜÉS et A. ROZENKNOP. « Lattice Parsing for Speech Recognition ». Dans *Actes de la 6ème conférence sur le Traitement Automatique des Langues Naturelles (TALN)*, Cargèse, France, 1999.
- [CS89] Y.-L. CHOW et R. SCHWARTZ. « The N-Best Algorithm: An efficient Procedure for Finding Top N Sentence Hypotheses ». Dans *Proc. of the DARPA Speech and Natural Language Workshop*, Philadelphie, Pennsylvanie, États-Unis, 1989.



- [CV02] E. CAMPIONE et J. VÉRONIS. « Étude des relations entre pauses et ponctuations pour la synthèse de la parole à partir de texte ». Dans *Actes de la 9ème conférence sur le Traitement Automatique des Langues Naturelles (TALN)*, Nancy, France, 2002.
- [CV04] E. CAMPIONE et J. VÉRONIS. « Pauses et hésitations en français spontané ». Dans *Actes des 25èmes Journées d'Études sur la Parole (JEP)*, Fès, Maroc, 2004.
- [CVD05] E. CAMPIONE, J. VÉRONIS et J. DEULOFEU. « *C-ORAL-ROM, Integrated Reference Corpora for Spoken Romance Languages* », Chapitre 3. The French corpus, pages 111–133. Amsterdam: John Benjamins, 2005.
- [DAS97] G. DEMETRIOU, E. ATWELL et C. SOUTER. « Large-Scale Lexical Semantics for Speech Recognition Support ». Dans *Proc. of the 5th European Conference on Speech, Communication, Technology (Eurospeech)*, Rhodes, Grèce, 1997.
- [DB95] S. DELIGNE et F. BIMBOT. « Language Modeling by Variable Length Sequences: Theoretical Formulation and Evaluation of Multigrams ». Dans *Proc. of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, Detroit, Michigan, États-Unis, 1995.
- [DGA<sup>+</sup>93] J. DOWDING, J. M. GAWRON, D. APPELT, J. BEAR, L. CHERNY, R. MOORE et D. MORAN. « Gemini: A Natural Language System for Spoken Language Understanding ». Dans *Proc. of the 31st Annual Meeting of the Association for Computational Linguistics (ACL)*, Columbus, Ohio, États-Unis, 1993.
- [DGP99] N. DESHMUKH, A. GANAPATHIRAJU et J. PICONE. « Hierarchical Search for Large Vocabulary Conversational Speech Recognition ». *IEEE Signal Processing Magazine*, 16(5):84–107, 1999.
- [DS98] S. DELIGNE et Y. SAKISAGA. « Learning a Syntagmatic and Paradigmatic Structure from Language Data with a Bi-Multigram Model ». Dans *Proc. of the 36th Annual Meeting of the Association for Computational Linguistics and the 17th International Conference on Computational Linguistics (COLING-ACL)*, volume 1, Montréal, Canada, 1998.
- [EBD90] M. EL-BÈZE et A.-M. DEROUAULT. « A Morphological Model for Large Vocabulary Speech Recognition ». Dans *Proc. of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, volume 1, Albuquerque, Nouveau Mexique, États-Unis, 1990.
- [Fed96] M. FEDERICO. « Bayesian Estimation Methods for N-Gram Language Model Adaptation ». Dans *Proc. of the 4th International Conference on Spoken Language Processing (ICSLP)*, volume 1, Philadelphie, Pennsylvanie, États-Unis, 1996.
- [Fed99] M. FEDERICO. « Efficient Language Model Adaptation through MDI Estimation ». Dans *Proc. of the 6th European Conference on Speech, Communication, Technology (Eurospeech)*, volume 4, Budapest, Hongrie, 1999.

- [FIO96] A. FARHAT, J.-F. ISABELLE et D. O'SHAUGHNESSY. « Clustering Words for Statistical Language Models Based on Contextual Word Similarity ». Dans *Proc. of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, volume 1, Atlanta, Géorgie, États-Unis, 1996.
- [FRWP03] M. FRANZ, B. RAMABHADRAN, T. WARD et M. PICHENY. « Information Access in Large Spoken Archives ». Dans *Proc. of the ISCA Multilingual Spoken Document Retrieval Workshop*, Macao/Hong Kong, Chine, 2003.
- [FY99] R. FLORIAN et D. YAROWSKY. « Dynamic Nonlocal Language Modeling via Hierarchical Topic-Based Adaptation ». Dans *Proc. of 37th Annual Meeting of the Association for Computational Linguistics (ACL)*, College Park, Maryland, États-Unis, 1999.
- [GAAD+05] J.-L. GAUVAIN, G. ADDA, M. ADDA-DECKER, A. ALLAUZEN, V. GENDNER, L. LAMEL et H. SCHWENK. « Where are we in Transcribing French Broadcast News? ». Dans *Proc. of the 9th European Conference on Speech Communication and Technology (Eurospeech)*, Lisbonne, Portugal, 2005.
- [GAD02] V. GENDNER et M. ADDA-DECKER. « Analyse comparative de corpus oraux et écrits français : mots, lemmes et classes morpho-syntaxiques ». Dans *Actes des 24èmes Journées d'Études sur la Parole (JEP)*, Nancy, France, 2002.
- [GAL+04] J.-L. GAUVAIN, G. ADDA, L. LAMEL, F. LEFÈVRE et H. SCHWENK. « Transcription de la parole conversationnelle ». Dans *Actes des 25èmes Journées d'Études sur la Parole (JEP)*, Fès, Maroc, 2004.
- [Gar95] R. GARSIDE. « *Spoken English on Computer: Transcription, Mark-up and Application* », Chapitre Grammatical Tagging of the Spoken Part of the British National Corpus: A Progress Report, pages 161–167. Harlow: Longman, 1995.
- [Geu96] P. GEUTNER. « Introducing Linguistic Constraints into Statistical Language Modeling ». Dans *Proc. of the 4th International Conference on Spoken Language Processing (ICSLP)*, Philadelphie, Pennsylvanie, États-Unis, 1996.
- [GH99] D. GILDEA et T. HOFMANN. « Topic-Based Language Models Using EM ». Dans *Proc. of the 6th European Conference on Speech Communication and Technology (Eurospeech)*, Budapest, Hongrie, 1999.
- [GLL00] J. GAO, M. LI et K.-F. LEE. « N-Gram Distribution Based Language Model Adaptation ». Dans *Proc. of the 6th International Conference on Spoken Language Processing (ICSLP)*, volume 1, Pékin, Chine, 2000.
- [Goo01] J. T. GOODMAN. « A Bit of Progress in Language Modeling, Extended Version ». Rapport technique, Microsoft Research, Redmond, Washington, États-Unis, 2001.
- [GSTN96] F. GALLWITZ, E. G. SCHUKAT-TALAMAZZINI et H. NIEMANN. « Integrating Large Context Language Models into a Real Time Word Recognizer ». Dans *Proc. of the 3rd Slovenian-German and the 2nd SDRV Workshop*, Ljubljana, Slovénie, 1996.

- [GT04] J. GARDES TAMINE. *Pour une grammaire de l'écrit*. Paris : Belin, 2004.
- [Gué05] M.-L. GUÉNOT. « Parsing de l'oral : traiter les disfluences ». Dans *Actes de la 12ème conférence sur le Traitement Automatique des Langues Naturelles (TALN)*, volume 1, Dourdan, France, 2005.
- [GW98] J. GILLET et W. WARD. « A Language Model Combining Trigrams and Stochastic Context-Free Grammars ». Dans *Proc. of the 5th International Conference on Spoken Language Processing (ICSLP)*, volume 6, Sydney, Australie, 1998.
- [Hee99] P. A. HEEMAN. « POS Tags and Decision Trees for Language Modeling ». Dans *Proc. of the Joint SIGDAT Conference on Empirical Methods in Natural Language Processing and Very Large Corpora*, College Park, Maryland, États-Unis, 1999.
- [Hen02a] S. HENRY. « Quelles répétitions à l'oral ? Esquisse d'une typologie ». Dans *Actes des 2èmes Journées de Linguistique de Corpus*, Lorient, France, 2002.
- [Hen02b] S. HENRY. « Étude des répétitions en français parlé spontané pour les technologies de la parole ». Dans *Actes de la 6ème Rencontre des Étudiants Chercheurs en Informatique pour le Traitement Automatique des Langues (RECITAL)*, Nancy, France, 2002.
- [HH95] M. P. HARPER et R. A. HELZERMEN. « Extensions to Constraint Dependency Parsing for Spoken Language Processing ». *Computer Speech and Language*, pages 187–234, 1995.
- [HJ04] K. HALL et M. JOHNSON. « Attention Shifting for Parsing Speech ». Dans *Proc. of the 42nd Meeting of the Association for Computational Linguistics (ACL)*, Barcelone, Espagne, 2004.
- [HJJ+99] M. P. HARPER, M. T. JOHNSON, L. H. JAMIESON, S. A. HOCKEMA et C. M. WHITE. « Interfacing a CDG Parser with an HMM Word Recognizer Using Word Graphs ». Dans *Proc. of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, volume 2, Phoenix, Arizona, États-Unis, 1999.
- [HJM+94] M. P. HARPER, L. H. JAMIESON, C. D. MITCHELL, G. YING, S. POTISUK, P. N. SRINIVASAN, R. CHEN, C. B. ZOLTOWSKI, L. L. MCPHETERS, B. PELLOM et R. A. HELZERMEN. « Integrating Language Models with Speech Recognition ». Dans *Proc. of the AAAI94 Workshop on the Integration of Natural Language and Speech Processing*, Seattle, Washington, États-Unis, 1994.
- [HP03] S. HENRY et B. PALLAUD. « Word Fragments and Repeats in Spontaneous Spoken French ». Dans *Proceedings of Disfluency in Spontaneous Speech Workshop (DISS)*, Göteborg, Suède, 2003.
- [HW94] A. HAUNSTEIN et H. WEBER. « An Investigation of Tightly-Coupled Time-Synchronous Speech Language Understanding Using a Unification Grammar ». Dans *Proc. of the 12th National Conference on Artificial Intelligence*

- Workshop on the Integration of Natural Language and Speech Processing*, Seattle, Washington, États-Unis, 1994.
- [IM94] R. ISOTANI et S. MATSUNAGA. « Speech Recognition Using a Stochastic Language Model Integrating Local and Global Constraints ». Dans *Proc. of the ARPA SLT Workshop*, 1994.
- [IO99] R. IYER et M. OSTENDORF. « Modeling Long Distance Dependence in Language: Topic Mixtures *versus* Dynamic Cache Models ». *IEEE Transactions on Speech and Audio Processing*, 7(1):30–39, 1999.
- [Jar96] M. JARDINO. « Multilingual Stochastic N-Gram Class Language Models ». Dans *Proc. of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, volume 1, Atlanta, Géorgie, États-Unis, 1996.
- [JC04] M. JOHNSON et E. CHARNIAK. « A TAG-Based Noisy Channel Model of Speech Repairs ». Dans *Proc. of the 42nd Annual Meeting of the Association for Computational Linguistics (ACL)*, Barcelone, Espagne, 2004.
- [Jel97] F. JELINEK. *Statistical Methods for Speech Recognition*. The MIT Press, 1997.
- [JL91] F. JELINEK et J. D. LAFFERTY. « Computation of the Probability of Initial Substring Generation by Stochastic Context-Free Grammars ». *Computation Linguistics*, 17(3):315–323, 1991.
- [JM00] D. JURAFSKY et J. H. MARTIN. *Speech and Natural Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition*. Prentice-Hall, 2000.
- [Jou96] D. JOUVET. « Robustesse et flexibilité en reconnaissance automatique de la parole ». *L'écho des recherches*, 165:25–38, 1996.
- [JWS+95] D. JURAFSKY, C. WOOTERS, J. SEGAL, A. STOLCKE, E. FOSLER, G. TAJCHMAN et N. MORGAN. « Using a Stochastic Context-Free Grammar as a Language Model for Speech Recognition ». Dans *Proc. of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, volume 1, Detroit, Michigan, États-Unis, 1995.
- [Kat87] S. M. KATZ. « Estimation of Probabilities from Sparse Data for the Language Model Component of a Speech Recognizer ». *IEEE Transactions on Acoustics, Speech, and Signal Processing*, 35(3):400–401, 1987.
- [KDM90] R. KUHN et R. DE MORI. « A Cache-Based Natural Language Model for Speech Recognition ». *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 12(6):570–583, 1990.
- [KKS89] K. KITA, T. KAWABATA et H. SAITO. « HMM Continuous Speech Recognition Using Predictive LR Parsing ». Dans *Proc. of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, Glasgow, Royaume-Uni, 1989.

- [KN93] R. KNESER et H. NEY. « Improved Clustering Techniques for Class-Based Statistical Language Modelling ». Dans *Proc. of the 3rd European Conference on Speech Communication and Technology (Eurospeech)*, volume 2, Berlin, Allemagne, 1993.
- [Kne96] R. KNESER. « Statistical Language Modeling Using a Variable Context Length ». Dans *Proc. of the 4th International Conference on Spoken Language Processing (ICSLP)*, volume 1, Philadelphie, Pennsylvanie, États-Unis, 1996.
- [KNST94] T. KUHN, H. NIEMANN et E. G. SCHUKAT-TALAMAZZINI. « Ergodic Hidden Markov Models and Polygrams for Language ». Dans *Proc. of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, volume 1, Adélaïde, Australie, 1994.
- [KR99] H.-K. J. KUO et W. REICHL. « Phrase-Based Language Models for Speech Recognition ». Dans *Proc. of the 6th European Conference on Speech Communication and Technology (Eurospeech)*, Budapest, Hongrie, 1999.
- [KS93] R. KNESER et V. STEINBISS. « On the Dynamic Adaptation of Stochastic Language Models ». Dans *Proc. of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, volume 2, Minneapolis, Minnesota, États-Unis, 1993.
- [KW98] T. KEMP et A. WAIBEL. « Reducing the OOV Rate in Broadcast News Speech Recognition ». Dans *Proc. of the 5th International Conference on Spoken Language Processing (ICSLP)*, Sydney, Australie, 1998.
- [KW99] S. KHUDANPUR et J. WU. « A Maximum Entropy Language Model to Integrate N-Grams and Topic Dependencies for Conversational Speech Recognition ». Dans *Proc. of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, Phoenix, Arizona, États-Unis, 1999.
- [LBS04] D. LINARES, J.-M. BENEDÍ et J.-A. SÁNCHEZ. « A Hybrid Language Model Based on a Combination of N-Grams and Stochastic Context-Free Grammars ». *ACM Transactions on Asian Language Information Processing (TALIP)*, 3(2):113–127, 2004.
- [LBSH03] D. LANGLOIS, A. BRUN, K. SMAÏLI et J.-P. HATON. « Événements impossibles en modélisation stochastique du langage ». *Traitement Automatique des Langues (TAL)*, 44(1):33–61, 2003.
- [LKMN05] I. R. LANE, T. KAWAHARA, T. MATSUI et S. NAKAMURA. « Dialogue Speech Recognition by Combining Hierarchical Topic Classification and Language Model Switching ». *IEICE Transactions on Information and Systems*, E88-D(3):446–454, 2005.
- [LMW97] G. LEECH, A. MCENERY et M. WYNNE. « *Corpus Annotation* », Chapitre Further Levels of Annotation, pages 85–101. London: Longman, 1997.

- [LSHS04] Y. LIU, A. STOLCKE, M. P. HARPER et E. SHRIBERG. « Comparing and Combining Generative and Posterior Probability Models: Some Advances in Sentence Boundary Detection in Speech ». Dans *Proc. of the Conference on Empirical Methods in Natural Language Processing (EMNLP)*, Barcelone, Espagne, 2004.
- [LST92] J. LAFFERTY, D. SLEATOR et D. TEMPERLEY. « Grammatical Trigrams: A Probabilistic Link Grammar ». Dans *Proc. of the AAAI Fall Symposium on Probabilistic Approaches to Natural Language*, Cambridge, Massachusetts, États-Unis, 1992.
- [MAB03] A. MENDES, R. AMARO et M. F. BACELAR DO NASCIMENTO. « Reusing Available Resources for Tagging a Spoken Portuguese Corpus ». Dans *Proc. of the Workshop on Tagging and Shallow Processing of Portuguese (TASHA)*, Lisbonne, Portugal, 2003.
- [MBS00] L. MANGU, E. BRILL et A. STOLCKE. « Finding Consensus in Speech Recognition: Word Error Minimization and Other Applications of Confusion Networks ». *Computer Speech and Language*, 14(4):373–400, 2000.
- [Mel00] L. MELIS. « Le français parlé et le français écrit, une opposition à géométrie variable ». *Romanesque*, 25(3):56–66, 2000.
- [MG03] A. MORENO et J. M. GUIRAO. « Tagging a Spontaneous Speech Corpus of Spanish ». Dans *Proc. of Recent Advances in Natural Language Processing (RANLP)*, Borovets, Bulgarie, 2003.
- [MI96] M. METEER et R. IYER. « Modeling Conversational Speech for Speech Recognition ». Dans *Proc. of the Conference on Empirical Methods in Natural Language Processing (EMNLP)*, Philadelphie, Pennsylvanie, États-Unis, 1996.
- [MLN97] S. C. MARTIN, J. LIERMANN et H. NEY. « Adaptive Topic Dependent Language Modelling Using Word-Based Varigrams ». Dans *Proc. of the 5th European Conference on Speech Communication and Technology (Eurospeech)*, Rhodes, Grèce, 1997.
- [MM92] G. MALTESE et F. MANCINI. « An Automatic Technique to Include Grammatical and Morphological Information in a Trigram-Based Statistical Language Model ». Dans *Proc. of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, volume 1, San Francisco, Californie, États-Unis, 1992.
- [Moo99] R. C. MOORE. « *Computational Models of Speech Pattern Processing* », Chapitre Using Natural-Language Knowledge Sources in Speech Recognition, pages 304–327. Springer-Verlag, 1999.
- [MPM89] R. MOORE, F. PEREIRA et H. MURVEIT. « Integrating Speech and Natural-Language Processing ». Dans *Proc. of the DARPA Speech and Natural Language Workshop*, Philadelphie, Pennsylvanie, États-Unis, 1989.

- [MSZ02] X. MOU, S. SENEFF et V. ZUE. « Integration of Supra-Lexical Linguistic Models with Speech Recognition Using Shallow Parsing and Finite State Transducers ». Dans *Proc. of the 7th International Conference on Spoken Language Processing (ICSLP)*, Denver, Colorado, États-Unis, 2002.
- [NEB<sup>+</sup>99] A. NASR, Y. ESTÈVE, F. BÉCHET, T. SPRIET et R. DE MORI. « A Language Model Combining N-Grams and Stochastic Finite State Automata ». Dans *Proc. of the 6th European Conference on Speech Communication and Technology (Eurospeech)*, volume 5, Budapest, Hongrie, 1999.
- [NG01] J. NIVRE et L. GRÖNQVIST. « Tagging a Corpus of Spoken Swedish ». *International Journal of Corpus Linguistics*, 6(1):47–78, 2001.
- [NW96a] T. R. NIESLER et P. C. WOODLAND. « Combination of Word-Based and Category-Based Language Models ». Dans *Proc. of the 4th International Conference on Spoken Language Processing (ICSLP)*, volume 1, Philadelphie, Pennsylvanie, États-Unis, 1996.
- [NW96b] T. R. NIESLER et P. C. WOODLAND. « A Variable-Length Category-Based N-Gram Language Model ». Dans *Proc. of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, volume 1, Atlanta, Géorgie, États-Unis, 1996.
- [NWW98] T. R. NIESLER, E. W. D. WHITTAKER et P. C. WOODLAND. « Comparison of Part-of-Speech and Automatically Derived Category-Based Language Models for Speech Recognition ». Dans *Proc. of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, Seattle, Washington, États-Unis, 1998.
- [ONA97] S. ORTMANN, H. NEY et X. AUBERT. « A Word Graph Algorithm for Large Vocabulary Continuous Speech Recognition ». *Computer, Speech and Language*, 11(1):43–72, 1997.
- [Pal03] D. S. PALLETT. « A Look at NIST’s Benchmark ASR Tests: Past, Present, and Future ». Dans *Proc. of the IEEE Workshop Automatic Speech Recognition and Understanding*, St. Thomas, îles Vierges, États-Unis, 2003.
- [PH04] B. PALLAUD et S. HENRY. « Amorces de mots et répétitions : des hésitations plus que des erreurs en français parlé ». Dans *Actes des 7èmes Journées internationales d’Analyse statistique des Données Textuelles (JADT)*, Louvain-la-Neuve, Belgique, 2004.
- [PMVGL03] F. PERRAUD, E. MORIN, C. VIARD-GAUDIN et P.-M. LALLICAN. « Modèles N-grammes et N-classes pour la reconnaissance de l’écriture manuscrite en ligne ». *Traitement Automatique des Langues (TAL)*, 44(1):63–92, 2003.
- [Pol03] A. POLGUÈRE. *Lexicologie et sémantique lexicale : notions fondamentales*. Les Presses de l’Université de Montréal, 2003.
- [PPM04] A. PANUNZI, E. PICCHI et M. MONEGLIA. « Using PiTagger for Lemmatization and PoS Tagging of a Spontaneous Speech Corpus: C-Oral-Rom Ital-

- ian ». Dans *Proc. of the 4th International Conference on Language Resources and Evaluation (LREC)*, volume 2, Lisbonne, Portugal, 2004.
- [PS01] F. PENG et D. SCHUURMANS. « A Simple Closed-Class/Open-Class Factorization for Language Modeling ». Dans *Proc. of the 6th Natural Language Processing Pacific Rim Symposium (NLPRS)*, Tokyo, Japon, 2001.
- [Rab89] L. RABINER. « A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition ». *Proc. of the IEEE*, 77(2):257–285, 1989.
- [RBW96] K. RIES, F. D. BUØ et A. WAIBEL. « Class Phrase Models for Language Modeling ». Dans *Proc. of the 4th International Conference on Spoken Language Processing (ICSLP)*, volume 1, Philadelphie, Pennsylvanie, États-Unis, 1996.
- [Roa01] B. ROARK. « Probabilistic Top-Down Parsing and Language Modelling ». *Computational Linguistics*, 27(2):249–276, 2001.
- [Ros94] R. ROSENFELD. « A Hybrid Approach to Adaptive Statistical Language Modeling ». Dans *Proc. of the ARPA Workshop on Human Language Technology*, Plainsboro, New Jersey, États-Unis, 1994.
- [Ros96] R. ROSENFELD. « A Maximum Entropy Approach to Adaptive Statistical Language Modeling ». *Computer, Speech and Language*, 10:187–228, 1996.
- [Ros00a] R. ROSENFELD. « Incorporating Linguistic Structure into Statistical Language Models ». *Philosophical Transactions: Mathematical, Physical and Engineering Sciences*, 358:1311–1324, 2000.
- [Ros00b] R. ROSENFELD. « Two Decades of Statistical Language Modeling: Where do we Go from Here? ». *Proc. of the IEEE*, 88(8):1270–1278, 2000.
- [SA90] R. SCHWARTZ et S. AUSTIN. « Efficient, High-Performance Algorithms for N-Best Search ». Dans *Proc. of the DARPA Speech and Natural Language Workshop*, Hidden Valley, Pennsylvanie, États-Unis, 1990.
- [Sch94] H. SCHMID. « Probabilistic Part-of-Speech Tagging Using Decision Trees ». Dans *Proc. of the International Conference on New Methods in Language Processing*, Manchester, Royaume-Uni, 1994.
- [Sch95] H. SCHMID. « Improvements in Part-of-Speech Tagging with an Application to German ». Dans *Proc. of the ACL SIGDAT Workshop*, Dublin, Irlande, 1995.
- [SCL92] K.-Y. SU, T.-H. CHIANG et Y.-C. LIN. « A Unified Framework to Incorporate Speech and Language Information in Spoken Language Processing ». Dans *Proc. of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, volume 1, San Francisco, Californie, États-Unis, 1992.
- [Sen92] S. SENEFF. « TINA: a Natural Language System for Spoken Language Applications ». *Computational Linguistics*, 18(1):61–86, 1992.
- [SG04] H. SCHWENK et J.-L. GAUVAIN. « Neural Network Language Models for Conversational Speech Recognition ». Dans *Proc. of the 8th International*



- Conference on Spoken Language Processing (ICSLP)*, île de Jeju, Corée du Sud, 2004.
- [Shr94] E. SHRIBERG. « *Preliminaries to a Theory of Speech Disfluencies* ». Thèse de doctorat, University of California, Berkeley, Californie, États-Unis, 1994.
- [Shr01] E. SHRIBERG. « To “Errrr” is Human: Ecology and Acoustics of Speech Disfluencies ». *Journal of the International Phonetic Association*, 31(1):153–169, 2001.
- [SMZ95] S. SENEFF, M. MCCANDLESS et V. ZUE. « Integrating Natural Language into the Word Graph Search for Simultaneous Speech Recognition and Understanding ». Dans *Proc. of the 4th European Conference on Speech Communication and Technology (Eurospeech)*, Madrid, Espagne, 1995.
- [SO96] M.-H. SIU et M. OSTENDORF. « Modeling Disfluencies in Conversational Speech ». Dans *Proc. of the 4th International Conference on Spoken Language Processing (ICSLP)*, volume 1, Philadelphie, Pennsylvanie, États-Unis, 1996.
- [SO00] M.-H. SIU et M. OSTENDORF. « Variable N-Grams and Extensions for Conversational Speech Language Modeling ». *IEEE Transactions on Speech and Audio Processing*, 8(1):63–75, 2000.
- [SR97] K. SEYMORE et R. ROSENFELD. « Using Story Topics for Language Model Adaptation ». Dans *Proc. of the 5th European Conference on Speech Communication and Technology (Eurospeech)*, Rhodes, Grèce, 1997.
- [SR99] C. SAMUELSSON et W. REICHL. « Class-Based Language Model for Large-Vocabulary Speech Recognition Extracted from Part-of-Speech Statistics ». Dans *Proc. of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, volume 1, Phoenix, Arizona, États-Unis, 1999.
- [SS94] A. STOLCKE et J. SEGAL. « Precise N-Gram Probabilities from Stochastic Context-Free Grammars ». Dans *Proc. of the 32nd Annual Meeting of the Association for Computational Linguistics (ACL)*, Las Cruces, Nouveau Mexique, États-Unis, 1994.
- [SS96] A. STOLCKE et E. SHRIBERG. « Statistical Language Modeling for Speech Disfluencies ». Dans *Proc. of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, volume 1, Atlanta, Géorgie, États-Unis, 1996.
- [STHKN95] E. G. SCHUKAT-TALAMAZZINI, R. HENDRYCH, R. KOMPE et H. NIEMANN. « Permugram Language Models ». Dans *Proc. of the 4th European Conference on Speech Communication and Technology (Eurospeech)*, volume 3, Madrid, Espagne, 1995.
- [Sto95] A. STOLCKE. « An Efficient Probabilistic Context-Free Parsing Algorithm that Computes Prefix Probabilities ». *Computational Linguistics*, 21(2):165–202, 1995.

- [Sto97] A. STOLCKE. « Linguistic Knowledge and Empirical Methods in Speech Recognition - Natural Language Processing ». *AI Magazine*, 18(4), 1997.
- [Str03] S. STRASSEL. « *Simple Metadata Annotation Specification. Version 5.0* ». Linguistic Data Consortium, 2003.
- [SW94] B. SUHM et A. WAIBEL. « Towards Better Language Models for Spontaneous Speech ». Dans *Proc. of the 3rd International Conference on Spoken Language Processing (ICSLP)*, volume 2, Yokohama, Japon, 1994.
- [SWH03] S. SENEFF, C. WANG et T. J. HAZEN. « Automatic Induction of N-Gram Language Models from a Natural Language Grammar ». Dans *Proc. of the 8th European Conference on Speech Communication and Technology (Eurospeech)*, Genève, Suisse, 2003.
- [TK95] M. TAMOTO et T. KAWABATA. « Clustering Word Category Based on Binomial Posteriori Cooccurrence Distribution ». Dans *Proc. of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, Detroit, Michigan, États-Unis, 1995.
- [TN97] C. TILLMANN et H. NEY. « Word Triggers and the EM Algorithm ». Dans *Proc. of the Workshop Computational Natural Language Learning (CoNLL 97)*, Madrid, Espagne, 1997.
- [UNY<sup>+</sup>02] K. UCHIMOTO, C. NOBATA, A. YAMADA, S. SEKINE et H. ISAHARA. « Morphological Analysis of the Spontaneous Speech Corpus ». Dans *Proc. of the 19th International Conference on Computational Linguistics (COLING)*, volume 2, Taipei, Taiwan, 2002.
- [VAR99] D. VAUFREYDAZ, M. AKBAR et J. ROUILLARD. « Internet Documents: a Rich Source for Spoken Language Modeling ». Dans *Proc. of the IEEE Workshop Automatic Speech Recognition and Understanding (ASRU)*, Keystone, Colorado, États-Unis, 1999.
- [VEZD00] F. VAN EYNDE, J. ZAVREL et W. DAELEMANS. « Part of Speech Tagging and Lemmatisation for the Spoken Dutch Corpus ». Dans *Proc. of the Conference on Language Resources and Evaluation (LREC)*, Athènes, Grèce, 2000.
- [VKDS04] D. VERGYRI, K. KIRCHHOFF, K. DUH et A. STOLCKE. « Morphology-Based Language Modeling for Arabic Speech Recognition ». Dans *Proc. of the 8th International Conference on Spoken Language Processing (ICSLP)*, île de Jeju, Corée du Sud, 2004.
- [VV99] A. VALLI et J. VÉRONIS. « Étiquetage grammatical de corpus oraux : problèmes et perspectives ». *Revue française de linguistique appliquée*, 4(2):113–133, 1999.
- [Vér04] J. VÉRONIS. « Le traitement automatique des corpus oraux ». *TAL*, 45(2):7–14, 2004.
- [WH02] W. WANG et M. P. HARPER. « The SuperARV Language Model: Investigating the Effectiveness of Tightly Integrating Multiple Knowledge Sources ».

- Dans *Proc. of the Empirical Methods in Natural Language Processing Conference (EMNLP)*, Philadelphie, Pennsylvanie, États-Unis, 2002.
- [WHS03] W. WANG, M. P. HARPER et A. STOLCKE. « The Robustness of an Almost-Parsing Language Model Given Errorful Training Data ». Dans *Proc. of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, volume 1, Hong Kong, Chine, 2003.
- [WK99] J. WU et S. KHUDANPUR. « Combining Nonlocal, Syntactic and N-Gram Dependencies in Language Modeling ». Dans *Proc. of the 5th European Conference on Speech Communication and Technology (Eurospeech)*, Budapest, Hongrie, 1999.
- [WLH02] W. WANG, Y. LIU et M. P. HARPER. « Rescoring Effectiveness of Language Models Using Different Levels of Knowledge and their Integration ». Dans *Proc. of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, volume 1, Orlando, Florida, États-Unis, 2002.
- [WMH00] Y.-Y. WANG, M. MAHAJAN et X. HUANG. « A Unified Context-Free Grammar and N-Gram Model for Spoken Language Processing ». Dans *Proc. of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, volume 3, Istanbul, Turquie, 2000.
- [WSH04] W. WANG, A. STOLCKE et M. P. HARPER. « The Use of a Linguistically Motivated Language Model in Conversational Speech Recognition ». Dans *Proc. of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, volume 1, Montréal, Canada, 2004.
- [WW01] E. W. D. WHITTAKER et P. C. WOODLAND. « Efficient Class-Based Language Modelling for Very Large Vocabularies ». Dans *Proc. of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, volume 1, Salt Lake City, Utah, États-Unis, 2001.
- [YS99] H. YAMAMOTO et Y. SAGISAKA. « Multi-Class Composite N-Gram Based on Connection Direction ». Dans *Proc. of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, volume 1, Phoenix, Arizona, États-Unis, 1999.
- [ZW98] K. ZECHNER et A. WAIBEL. « Using Chunk Based Partial Parsing of Spontaneous Speech in Unrestricted Domains for Reducing Word Error Rate in Speech Recognition ». Dans *Proc. of the 36th Annual Meeting of the Association for Computational Linguistics and the 17th International Conference on Computational Linguistics (COLING-ACL)*, Montréal, Canada, 1998.
- [ÉD04] Équipe DELIC. « Présentation du corpus de référence du français parlé ». *Recherches sur le français parlé*, 18, 2004.



---

Unité de recherche INRIA Rennes  
IRISA, Campus universitaire de Beaulieu - 35042 Rennes Cedex (France)

Unité de recherche INRIA Futurs : Parc Club Orsay Université - ZAC des Vignes  
4, rue Jacques Monod - 91893 ORSAY Cedex (France)

Unité de recherche INRIA Lorraine : LORIA, Technopôle de Nancy-Brabois - Campus scientifique  
615, rue du Jardin Botanique - BP 101 - 54602 Villers-lès-Nancy Cedex (France)

Unité de recherche INRIA Rhône-Alpes : 655, avenue de l'Europe - 38334 Montbonnot Saint-Ismier (France)

Unité de recherche INRIA Rocquencourt : Domaine de Voluceau - Rocquencourt - BP 105 - 78153 Le Chesnay Cedex (France)

Unité de recherche INRIA Sophia Antipolis : 2004, route des Lucioles - BP 93 - 06902 Sophia Antipolis Cedex (France)

---

Éditeur  
INRIA - Domaine de Voluceau - Rocquencourt, BP 105 - 78153 Le Chesnay Cedex (France)

<http://www.inria.fr>

ISSN 0249-6399