

[MIG] Systèmes Embarqués

Julien CAILLARD, Adrien DE LA VAISSIÈRE, Thomas DEBARRE,
Matthieu DENOUX, Maxime ERNOULT, Axel GOERING,
Clément JOUDET, Nathanaël KASRIEL, Anis KHLIF,
Sofiane MAHIOU, Paul MUSTIÈRE, Clément ROIG, David VITOUX

18/11/13 - 6/12/13

Table des matières

0.1	Présentation des enjeux	3
0.2	Objectifs du projet	3
0.3	Approches de la reconnaissance vocale	4
0.3.1	Acoustique-phonétique	4
0.3.2	Reconnaissance de motifs	5
I	Démarche Technique	6
1	Principe général du traitement du signal	6
1.1	Objectifs	6
1.2	Schéma global	7
2	Analyse, formatage du signal	8
2.1	Introduction	8
2.2	Prérequis	8
2.2.1	Qu'est-ce que le son ?	8
2.2.2	Comment le son est-il représenté dans l'ordinateur ?	8
2.2.3	Pré accentuation des aigus	9
2.3	Enregistrement, recadrage, filtrage HF	9
2.3.1	Synchronisation	9
2.3.2	Filtrage passe-haut	10
2.4	Echantillonnage, fenêtrage	10
2.5	Transformée de Fourier	11
2.6	Mel	12
2.7	Transformée inverse	12
3	Modèles de Markov cachés (MMC)	14
3.1	Prérequis et principe	14
3.2	Principaux algorithmes sur les modèles de Markov	15
3.3	Application à notre objectif	15
3.4	Phase d'apprentissage	15
3.5	Phase de reconnaissance	16
II	Approche commerciale	18
1	Approche développement web	18
1.0.1	Choix d'une architecture optimale pour notre projet	18
1.0.2	Réalisation du SpeechServer	19
1.0.3	Système de Gestion de Base de Données (SGBD)	20
1.0.4	SpeechApp	21

1.0.5	SpeechRecorder	21
2	Applications	22
3	Budget, modèle économique	24
3.1	Introduction	24
3.2	Les salaires	24
3.3	Le compte de résultat prévisionnel	25
3.4	Le bilan	25
3.5	Les impôts	26
3.6	Conclusion et vue sur le long terme	26
III	Le Code	27
1	Python	27
2	C/C++	29

Introduction

0.1 Présentation des enjeux

La reconnaissance vocale automatisée est l’objet d’intenses recherches depuis plus de 50 ans. Malgré son caractère d’abord futuriste, comme cela peut se retrouver dans de nombreuses oeuvres de science-fiction, elle a pris sa place dans nos quotidiens avec la prolifération de systèmes qui embarquent une telle technologie, par exemple avec le logiciel Siri dans les téléphones d’Apple[1]. Les perspectives économiques s’ouvrant au détenteur d’un système de reconnaissance fiable, robuste, et portable sont innombrables et l’on ne saurait surestimer son importance, (systèmes embarqués, commandes vocales, aide aux sourds/muets, ...) et s’inscrit dans le domaine prolifique beaucoup plus large du traitement du signal. Les derniers systèmes les plus aboutis offrent des performances remarquables, mais le problème reste toujours ouvert et suscite plus d’engouement que jamais en raison de la croissante puissance de calcul disponible et les dernières avancées et applications découvertes. La complexité de ce problème s’explique notamment par la grande diversité des thèmes qui lui sont connexes et que tout système se voulant performant se doit d’incorporer (traitement du signal, théorie de l’information, acoustique, linguistique, intelligence artificielle, physiologie, psychologie, ...) et de part l’approche unidirectionnelle de certaines recherches. La reconnaissance vocale requiert des connaissances trop diverses pour être maîtrisées par un seul individu et la capacité à savoir exploiter des ressources dont on est pas expert devient un atout capital. Elle ne se réduit pas à la seule détermination d’une suite de mots prononcés, mais peut s’étendre à divers autres applications telles que la reconnaissance de langage, d’accent, déterminer le sexe et l’âge du locuteur, si il est stressé ou calme, dans quel environnement est-il, tant ces paramètres influent de manière capitale sur l’analyse.

0.2 Objectifs du projet

Ce MIG s’est placé dans une perspective résolument plus humble en raison du temps limité imparti. Il ne s’agissait pas de réaliser un programme prétendant rivaliser avec les actuels systèmes de reconnaissance, fruits de nombreuses années de recherches et de développement ; mais plutôt, à l’instar de l’ingénieur généraliste, de prendre connaissance d’un sujet et d’une problématique et tâcher, en équipe, d’y apporter une solution qui soit la plus optimale possible compte tenu des exigences temporelles et matérielles. La complexité de la discipline fut un des principaux obstacles, et une phase d’appropriation des techniques requises, de part la lecture de livres dédiés, d’articles de recherches ainsi que de thèses a été le poumon du projet. Le caractère abscon de certains articles a rajouté à la difficulté. Le projet des MIG ne réduisant pas non plus à une réalisation technique il s’agissait de garder en vue les perspectives économiques et les composantes juridiques, indissociables d’un tel projet, comme garde fou de toute pérégrination informatique.

Les rôles ont été attribués dès le début selon les goûts et compétences de chacun mais la pertinente répartition des tâches, la diversité intrinsèque au projet et l’angle avec lequel nous

l'avons abordé a permis à chacun d'exploiter un panel très diversifié de ses compétences tout en apportant la valeur ajoutée de sa spécialité. Chaque fonction dépendant très fortement de ce qui précède et de ce qui suit, une bonne communication interne était indispensable pour un développement juste et efficace. Si la coordination spontanée d'une équipe de treize personnes a été au début délicate, une indéniable rigueur et discipline adjointe à l'exploitation de ressources adaptées ont vite imposé une organisation naturelle. Par exemple l'utilisation de la plateforme github[2] pour l'échange de fichiers et de mises à jour s'est révélée particulièrement efficace et permettait à chacun d'incorporer en temps réels les dernières modifications. Ceci a permis à chacun d'être en permanence en totale connaissance des rôles de chacun, et de la ligne directrice de chaque sous-fonction sans pour autant en connaître tous les détails. Une efficacité dans la coordination des tâches et dans leur réalisation en a indiscutablement découlé.

0.3 Approches de la reconnaissance vocale

Avant de rentrer dans des considérations techniques, il est nécessaire de définir un principe d'étude, une stratégie de résolution qui dictera l'orientation générale du projet en plus de rendre les objectifs et les enjeux plus clairs. Cette partie a pour but de donner un aperçu des différents angles d'attaques du problème donné pouvant être considérés, ainsi que de présenter celui que nous avons choisi, avec quelles motivations.

Dans son livre *Fundamentals of speech recognition*, Lawrence Rabiner[3] dégage des travaux de ces prédécesseurs trois approches conceptuelles du problème. Ces approches sont les suivantes : l'approche acoustique-phonétique, l'approche par reconnaissance de motifs et l'approche par intelligence artificielle. Cette dernière n'étant, d'après Rabiner, qu'un avatar de la première ; nous ne présenterons que l'acoustique phonétique et la reconnaissance de motifs que nous avons choisi pour notre projet.

0.3.1 Acoustique-phonétique

L'approche acoustique-phonétique est indubitablement celle qui paraît la plus naturelle et directe pour faire de la reconnaissance vocale et est celle qui s'impose a priori à l'esprit. Le principe est le suivant : l'ordinateur tâche de découper l'échantillon sonore de manière séquentielle en se basant sur les caractéristiques acoustiques observées et sur les relations connues entre caractéristiques acoustiques et phonème. Ceci dans le but d'identifier une suite de phonèmes et d'ainsi reconnaître un mot.

Définition Wikipédia d'un phonème : En phonologie, domaine de la linguistique, un phonème est la plus petite unité discrète ou distinctive (c'est-à-dire permettant de distinguer des mots les uns des autres) que l'on puisse isoler par segmentation dans la chaîne parlée. Un phonème est en réalité une entité abstraite, qui peut correspondre à plusieurs sons. Il est en effet susceptible d'être prononcé de façon différente selon les locuteurs ou selon sa position et son environnement au sein du mot.

Cette approche suppose qu'il existe un ensemble fini de phonèmes différenciables et que leurs propriétés sont suffisamment manifestes pour être extraites d'un signal ou de la donnée de son spectre (tableau des fréquences et de leur amplitude associée, composant un signal à un instant donné) au cours du temps. Même si il est évident que ces caractéristiques dépendent très largement du sujet parlant, on part du principe que les règles régissant la modification des paramètres peuvent être apprises et appliquées.

Bien qu'elle ait été vastement étudiée et soit viable on lui préférera l'approche par reconnaissance de motifs qui, pour plusieurs raisons, l'a supplantée dans les systèmes appliqués. C'est celle que nous avons choisi et que nous présentons dans le prochain paragraphe.

0.3.2 Reconnaissance de motifs

Elle diffère de la méthode précédente par le fait qu'elle ne cherche pas à exhiber des caractéristiques explicites. Elle se compose de deux étapes : « l'entraînement » des motifs, et la reconnaissance via la comparaison de ces motifs. L'idée sous-jacente au concept d'entraînement repose sur le principe selon lequel si l'on dispose d'un ensemble suffisamment grand de version d'un motif à reconnaître, on doit être capable de caractériser pertinemment les propriétés acoustiques du motif. Notons que les motifs en question peuvent être de nature très diverses, comme des sons, des mots, des phrases ; ce qui sous-tend l'idée d'un grand nombre d'applications théoriques comme présenté en introduction. La machine apprend alors quelles propriétés acoustiques sont fiables et pertinentes. On effectue ensuite une comparaison entre le signal à reconnaître et les motifs tampons, afin de le classifier en fonction du degré de concordance.

Sans plus entrer dans les détails, les avantages de cette approche qui nous on poussé à l'adopter sont les suivants :

- Elle est simple à appréhender, et est très largement comprise et utilisée
- Elle est robuste, c'est-à-dire qu'elle dépend peu du locuteur et de l'environnement
- Elle donne lieu à de très bons résultats

Première partie

Démarche Technique

1. Principe général du traitement du signal

1.1 Objectifs

Bien que la reconnaissance vocale telle qu'elle est aujourd'hui mise-en-place dans les différents matériels semble immédiate, le travail à effectuer pour reconnaître un mot est complexe. La première étape pour faire de la reconnaissance vocale est de parvenir à trouver un moyen de caractériser efficacement et uniformément un mot. Cela désigne un mot par un certain motif puis permet par le même procédé appliqué sur un enregistrement quelconque, de parvenir à identifier deux motifs proches qui correspondraient alors au même mot. Il s'agit donc tout d'abord de traiter le signal pour en découvrir certaines caractéristiques. En effet, une même personne ne prononce pas toujours les mots de la même façon, au même débit, avec les mêmes hauteurs de son, ce qui rend ardue une simple identification par comparaisons temporelles.

1.2 Schéma global

Afin de gérer ces difficultés, nous avons mis en place plusieurs étapes de traitement supplémentaires afin d'obtenir cette fameuse « trace » qui caractériserait un enregistrement, c'est-à-dire un mot. Nous avons pour cela utilisé plusieurs techniques de traitement du signal communément connues (échantillonnage, fenêtrage, transformée de Fourier directe et inverse). Cette figure explique globalement le traitement que nous avons choisi de mettre-en-place afin de reconnaître le mot prononcé. Il y a donc plusieurs étapes qui s'enchaînent pour parvenir à un objet que nous pourrions manipuler en le sachant représentatif et caractéristique du son.

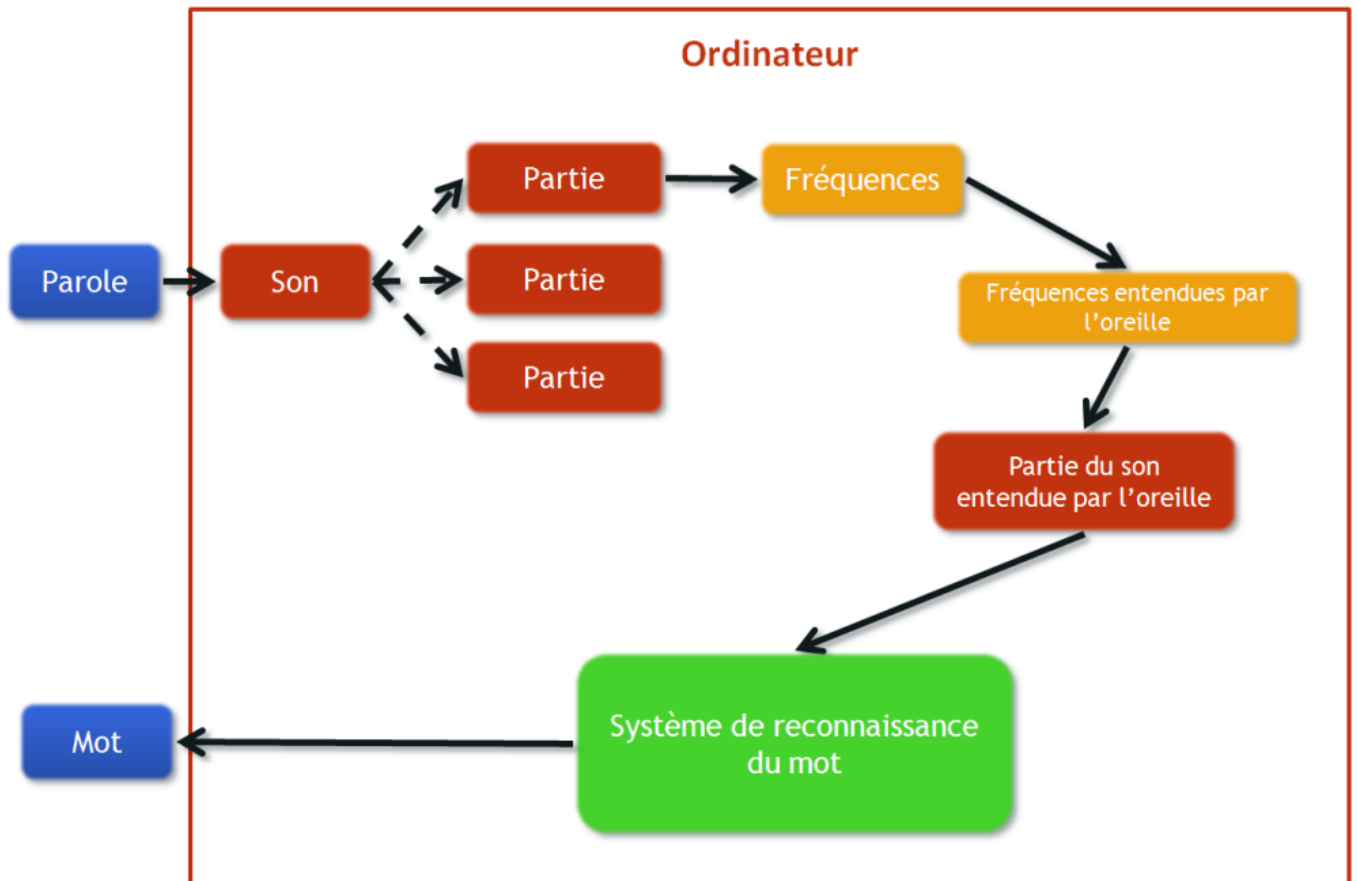


FIGURE 1.1 – Traitement du son pour le reconnaître

Enregistrement du son La première étape consiste à simplement enregistrer le son sur le disque dur de l'ordinateur. Nous utilisons pour cela un module intégré à Python appelé PyAudio[4]. Cela permet d'enregistrer avec une certaine fréquence (donc un certain nombre de captures par seconde) les amplitudes du son captées par le micro.

Découpage en fenêtre Le son est découpé ensuite en petites fenêtres de quelques dizaines de millisecondes ce qui permet d'isoler les événements sonores qui pourraient avoir une importance. Il s'agit d'un *fenêtrage*.

Passage en fréquence Jusque là, le son étudié se trouvait représenter à peu près temporellement ce qui avait été entendu. Néanmoins, il est difficile d'étudier un son tel quel et on utilise alors le lien entre les fréquences et le signal temporel. Il est ensuite plus facile d'étudier

et de transformer un ensemble de fréquences pour appliquer par un exemple des filtres qui rapprochent le programme du fonctionnement de l'oreille.

Utilisation de l'échelle de Mel Puisque le programme doit savoir *faire la différence entre des mots*, c'est-à-dire des sons identifiés tels quels par une oreille *humaine*, il faut donner au programme un comportement similaire à celui d'une oreille humaine. On utilise pour cela une échelle qui accentue certaines fréquences. En effet, il a été montré[5] (et ensuite appliqué [6]) que l'oreille ne perçoit pas toutes les fréquences de la même façon.

2. Analyse, formatage du signal

2.1 Introduction

Comme nous l'avons mentionné, même le plus élémentaire des systèmes de reconnaissance vocale utilise des algorithmes au carrefour d'une grande diversité de disciplines : reconnaissance de motifs statistiques, théorie de l'information, traitement du signal, analyse combinatoire, linguistique entre autres, le dénominateur commun étant le traitement du signal qui transforme l'onde acoustique de la parole en une représentation paramétrique plus apte à l'analyse automatisée. Le principe est simple : garder les traits distinctifs du signal et s'absoudre au maximum de tout ce qui pourra en parasiter l'étude. Cette conversion ne se fait donc pas sans perte d'information, et la délicatesse de la discipline tient en la sélection judicieuse des outils les plus adaptés afin de trouver le meilleur compromis entre perte d'information et représentation fidèle du signal.

2.2 Prérequis

2.2.1 Qu'est-ce que le son ?

Le son est une onde mécanique se traduisant par une variation de la pression au cours du temps. Cette onde est caractérisée par différents facteurs comme son amplitude à chaque instant, qui est en d'autres termes la valeur de la dépression à cet instant, et par les fréquences qui la composent et qui changent au cours du temps.

2.2.2 Comment le son est-il représenté dans l'ordinateur ?

En se propageant, l'onde mécanique qu'est le son fait vibrer la membrane du micro. L'amplitude de la vibration dépend directement de l'amplitude du son. La position de la membrane est enregistrée à intervalles de temps réguliers définis par l'échantillonnage. L'échantillonnage correspond au nombre de valeurs prélevées en une seconde (principe [7]). Par exemple un échantillonnage à 44100 Hz correspond à relever la position de la membrane 44100 fois par secondes. La valeur de la position de la membrane est alors enregistrée sous la forme d'un entier signé codé sur n bits (n valant généralement 8,16,32 ou 64). Plus n est grand, plus la position de la membrane sera stockée de manière précise, et donc plus la qualité du son sera bonne. Grâce à l'échantillonnage et à, on définit aisément le bitrate, qui correspond au débit d'information

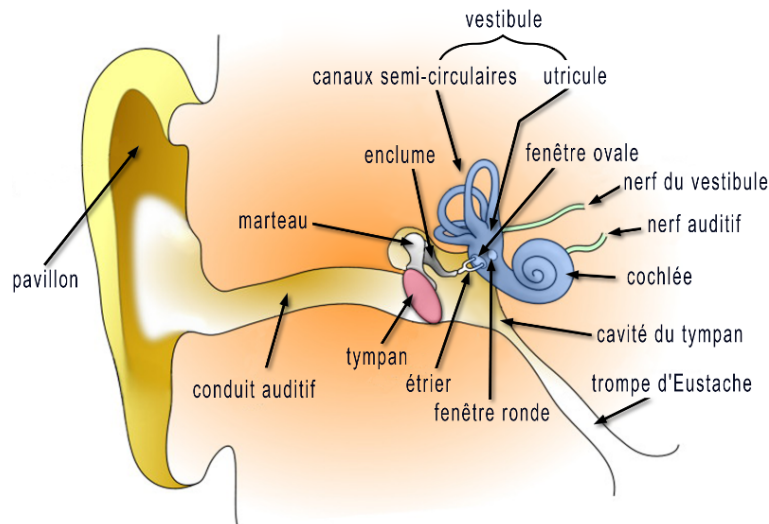


FIGURE 2.1 – Oreille humaine

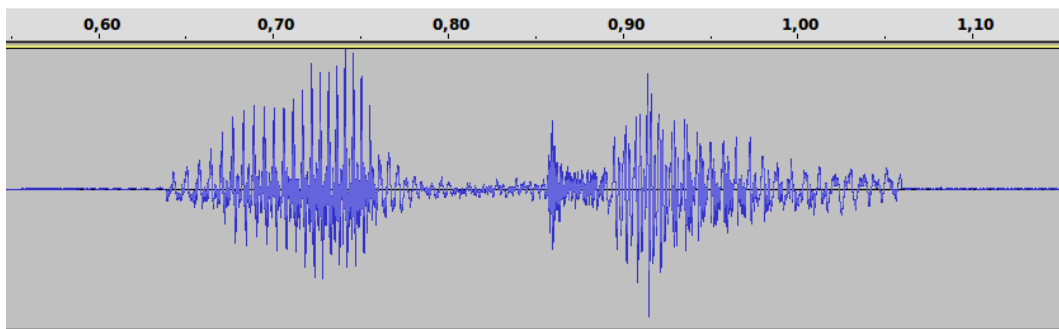


FIGURE 2.2 – Exemple audiogramme prononciation du mot "VICA"

par seconde, de la façon suivante : $bitrate = n * \text{échantillonnage}$. Ce dont nous disposons donc pour analyser un signal, est la donnée de l'amplitude en fonction du temps la caractérisant.

2.2.3 Pré accentuation des aigus

2.3 Enregistrement, recadrage, filtrage HF

2.3.1 Synchronisation

Afin de synchroniser le début des enregistrements d'un mot, et de leur donner la même durée, il a été nécessaire de détecter les silences avant et après le mot pour les couper. Le signal est lissé à l'aide d'une moyenne sur plusieurs échantillons pour que les fluctuations inhérentes à l'enregistrement ne gênent pas notre fonction. On détecte alors le moment où le signal (en valeur absolue) dépasse pour la première fois une valeur seuil et celui à partir duquel le signal ne dépasse plus celle-ci. On sait alors où couper le signal d'origine, en élargissant légèrement la coupe afin d'éviter de supprimer des consonnes peu sonores. Cela permet en plus d'afficher un message d'erreur suspectant un enregistrement ayant commencé trop tard ou fini trop tôt.

Deux problèmes se posent : en pratique, un bruit trop important perturbe le signal et le mot n'est plus détectable par l'amplitude des oscillations. Toutefois, pour l'enregistrement de notre base de données, une pièce calme et un micro de bonne qualité nous ont permis un

découpage satisfaisant ce qui ne résoud pas définitivement le problème : l'utilisateur ne pouvant pas toujours se placer dans ces conditions, le signal est traité par un filtre anti-bruit.

Ce filtre consiste en l'utilisation de bibliothèques, SoX et ffmpeg([8] et [9]), qui permettent par l'étude d'un court laps de temps de bruit de soustraire le bruit de l'enregistrement. Nous n'avons pas cherché à traiter nous-même le bruit car il s'agit d'un problème complètement à part et qui ne demande pas les mêmes compétences que le traitement du signal effectué jusque là.

De plus, il a fallu déterminer la valeur de nos constantes de découpe (coefficient de lissage, coefficient de coupe, intervalle de temps de sécurité), qui dépendent bien sûr les unes des autres. Ceci a été fait de manière empirique sur plusieurs enregistrements de mots différents, permettant une découpe automatique la plus satisfaisante possible pour l'ensemble des mots.

2.3.2 Filtrage passe-haut

Les performances de tout système de reconnaissances dépendent fortement de la variabilité des données (locuteur, environnement, bruit, réverbération, ...). Plus ces données sont variables, plus le taux d'erreur sera grand et un système de reconnaissance qui se veut être utilisable dans la vie de tous les jours : voiture, endroits bruyants ; se doit d'y remédier. Ces effets se font particulièrement sentir dans les basses fréquences, c'est pourquoi le conditionnement du signal en vue de son étude comprend inmanquablement un filtre passe haut c'est-à-dire une accentuation de l'amplitude associée aux hautes fréquences et une diminution des basses fréquences. C'est le même principe qui est utilisé dans les égaliseurs des lecteurs de musique d'aujourd'hui qui propose d'augmenter les basses ou les aigus. Les filtres passe-haut améliorent significativement les résultats de reconnaissance comme en témoignent les expériences de H.G. Hirsch P. Meyer et H.W. Ruehl dans leur papier . Utiliser un filtre passe-haut présente comme avantage de ne pas nécessiter de procéder au préalable à une reconnaissance de silence contrairement aux techniques de réduction du bruit et de soustraction spectrale.

Nous effectuons ce filtrage selon une méthode résolution basée sur les coefficients spectraux de mel[6]. Le signal étant caractérisé par une suite (x_n) d'amplitudes, comme présenté dans les prérequis, où n représente un instant de la musique déterminé par l'échantillonnage ; on opère linéairement la transformation suivante sur le signal : $y_0 = x_0$ et $y_n = x_n - 0.95 * x_{(n-1)}$ pour $n > 0$, où y représente le signal de sortie après transformation.

Cette opération consiste effectivement en un filtre passe-haut, en effet une telle formule part du principe que 95% d'un échantillon a pour origine l'échantillon précédent. Ce constat étant plus pertinent pour les hautes fréquences (car les pics de l'onde associée sont plus rapprochés et engendrent donc un pic d'amplitude plus régulièrement), l'influence des basses fréquences est donc discriminée.

2.4 Echantillonnage, fenêtrage

L'analyse du signal, pour accéder au domaine fréquentiel, s'affranchit de la dépendance temporelle. Le spectre obtenu ne correspond plus à une perception physique, mais à une moyenne temporelle du spectre perçu. Le procédé que nous avons mis en place pour pallier à ce problème est celui le plus couramment utilisé dans ce domaine : l'échantillonnage. Nous avons découpé le signal à traiter en petites séquences, qui, juxtaposées, approximent une échelle temporelle continue.

La taille des échantillons est un paramètre déterminant sur la qualité et la précision de l'analyse combinée finale. Une fois calculé, le spectre ne reflète plus du tout de dépendance temporelle. La durée d'un échantillon correspond ainsi à la durée minimale d'un événement

sonore détectable. Il faut donc réduire cette durée autant que possible, pour obtenir une discrétisation temporelle le plus proche possible de la continuité. Il est en revanche nécessaire de conserver un certain nombre de points par échantillons. En effet, le spectre obtenu par l'analyse sera plus précis et proche de la réalité fréquentielle si le nombre de point du signal analysé est important. La meilleure technique pour contourner ce compromis est d'augmenter la fréquence d'échantillonnage. On obtient alors un nombre important de points qui s'étirent peu dans le temps.

Le théorème de Nyquist-Shannon[10] assure qu'un signal reproduit fidèlement toutes les fréquences inférieures à la moitié de sa fréquence d'échantillonnage. Une fréquence d'échantillonnage de 44100Hz (parfois 48000Hz) est donc suffisante pour couvrir la totalité d'une oreille humaine en bonne santé. L'utilisation la plus courante de l'enregistrement audio étant (à notre niveau) la restitution, le matériel et logiciel à notre disposition se cantonnait à ces fréquences d'échantillonnage. Nous avons ainsi dû trouver un compromis entre résolution fréquentielle et précision temporelle. L'hypothèse principale a été que les événements sonores et variations s'étalant sur une durée inférieure à 20 millisecondes n'étaient pas significatifs pour notre analyse. Le nombre de points a été par cette donnée, couplée à notre fréquence d'échantillonnage lors des enregistrements, à 44100Hz.

L'échantillonnage introduit par ailleurs des discontinuités aux bornes des morceaux, qui ne sont pas présentes dans le signal original. Le fenêtrage permet de réduire l'effet de ces discontinuités virtuelles. On découpe le signal en plus de morceaux, tout en conservant la même durée pour chaque échantillon. On obtient des "fenêtres", qui se recoupent les unes les autres. Pour que la même partie du signal ne soit pas retraitée à l'identique, on applique une fonction - dite fonction de fenêtrage, ou dans notre cas, fonction de Hann - qui diminue l'importance des valeurs situées aux extrémités de la fenêtre. Ce procédé a le désavantage de démultiplier le temps de calcul des étapes suivantes de l'algorithme (le nombre d'échantillons est bien plus important pour un signal de même longueur). Certaines applications (notamment pour les téléphones portables) devant réduire la complexité au maximum en font donc abstraction. Notre reconnaissance privilégiant plutôt la précision, et disposant d'une puissance de calcul largement suffisante pour conserver un rendu de l'ordre de la seconde, nous avons opté pour un fenêtrage important (recouvrement total d'un échantillon à l'autre), au prix d'une multiplication du temps de calcul par deux.

2.5 Transformée de Fourier

Le domaine temporel est parfait pour l'acquisition et la restitution de l'audio, car il représente fidèlement la vibration de la membrane d'un micro ou d'une enceinte. L'oreille humaine base sa perception et sa reconnaissance sur le domaine fréquentiel. Il faut donc passer de l'un à l'autre, et ce grâce à l'utilisation de la transformation de Fourier. L'algorithme "intuitif" de calcul ayant, pour trouver le spectre d'un unique échantillon, une complexité en $O(N^2)$ (avec N le nombre de points par échantillons), il est nécessaire de trouver d'autres méthodes si l'on envisage des applications proches du temps réel. Heureusement, plusieurs approches se sont ouvertes à nous pour l'optimisation du temps de calcul.

Le calcul de la transformée de Fourier est incontournable en analyse du signal, et il a donné lieu à de nombreuses études. Des algorithmes optimisés pour diverses utilisations sont disponibles, et notre travail a surtout été d'identifier lequel s'adapterait à notre projet. Notre fonction se base sur l'algorithme de Cooley-Tukey, qui permet de réduire la complexité à $O(N \log_2(N))$, et qui repose sur le fonctionnement diviser pour régner. Le principe est dans un premier temps de diviser le signal à analyser en sous-tableaux de mêmes tailles, de manière croisée (par exemple deux sous-tableaux, pour les indices pairs et impairs). On calcule ensuite les transformées de Fourier de ce sous-tableaux, en opérant récursivement, jusqu'à obtenir des sous-tableaux dont

la taille est un entier. On calcule leur transformée de Fourier, et on recombine les résultats obtenus. Cette méthode a l'avantage de pouvoir être couplée à d'autres algorithmes pour calculer les spectres des sous-tableaux dont la taille n'est pas un produit d'entier. Le meilleur cas est alors instinctivement un signal initial dont la longueur est une puissance de deux. Il est même intéressant d'utiliser la technique du bourrage de zéros (zero padding), qui consiste à rajouter des zéros à la suite du signal pour atteindre la puissance de deux la plus proche. Cela ne change pas le spectre obtenu et augmente les performances. Dans notre cas, nous avons eu la possibilité d'ajuster la taille des échantillons. Nous avons ainsi choisi des échantillons de 1024 points, ce qui correspond, avec notre fréquence d'échantillonnage de 44100Hz, à une durée d'environ 23ms. Seul le dernier échantillon du signal est complété par des zéros.

De plus, comme les données sur lesquelles nous travaillons sont réelles, et que les calculs de la Transformée de Fourier Rapide (Fast Fourier Transform, ou FFT) s'effectuent avec des complexes, la première idée d'optimisation que nous avons eu est de calculer le spectre de deux échantillons à la fois, en créant des complexes à partir des deux signaux réels (l'un représente la partie réelle, l'autre imaginaire). On obtient rapidement les coefficients respectifs des deux échantillons par une simple opération sur le spectre résultant. Cependant, cette méthode ne divise le temps de calcul que par deux, et notre FFT demeure trop lente (plusieurs secondes pour un signal d'environ une seconde), surtout au regard du temps de calcul total de la reconnaissance en elle-même. La deuxième optimisation que nous avons donc appliqué est de passer le code de Python à C++, langage compilé beaucoup plus rapide. De plus, nous avons repensé les fonctions, de façon à éviter les appels récursifs. En effet, le travail sur des tableaux force une recopie à chaque appel de fonction, ce qui démultiplie la complexité du calcul. Le résultat est un algorithme qui s'effectue en moins d'une seconde, et qui peut s'inscrire dans un contexte d'exploitation en temps réel.

2.6 Mel

Des études de psycho acoustique ont montré que l'oreille humaine ne percevait pas les fréquences selon une échelle linéaire. Il a donc été utile de définir une nouvelle échelle plus subjective : à chaque fréquence f , exprimée en Hertz, on fait correspondre une nouvelle fréquence selon une fonction censée représenter le comportement de l'oreille humaine. Par convention, la fréquence de 1000 Hz correspond à 1000 mel. Les autres fréquences mel sont ajustées de façon à ce qu'une augmentation de la fréquence mel corresponde à la même augmentation de la tonalité perçue. Cela conduit à la fonction *mel* suivante :

$$mel(f) = 2595 * \log(1 + f/700)$$

On remarque que le poids des hautes fréquences (supérieures à 1000 Hz) est diminué tandis que le poids des basses fréquences (inférieure à 1000 Hz) est augmenté.

Il est préférable d'employer cette échelle de fréquence dans l'algorithme de reconnaissance : ce dernier doit en effet différencier plusieurs mots selon la perception humaine, c'est-à-dire en simulant le comportement de l'oreille humaine.

2.7 Transformée inverse

Nam eu sollicitudin massa. Duis sagittis velit mi. Nunc dictum risus ac interdum lacinia[3].

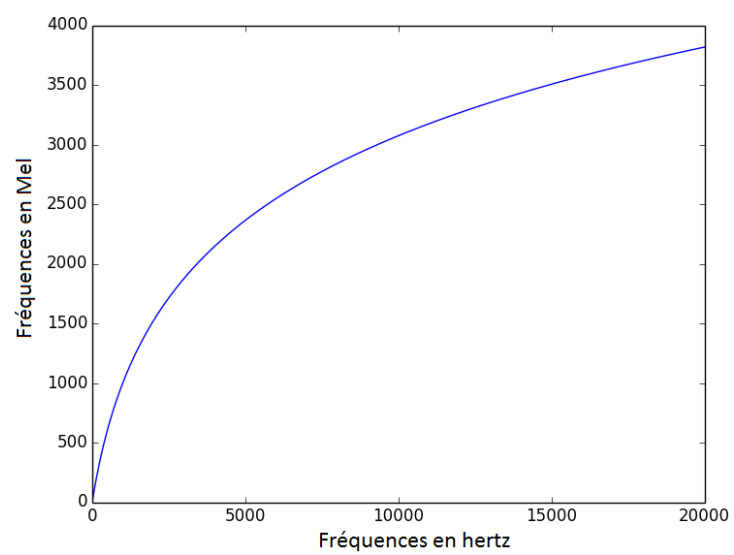


FIGURE 2.3 – Graphe de conversion

3. Modèles de Markov cachés (MMC)

3.1 Prérequis et principe

Un modèle de Markov caché est un modèle statistique qui peut modéliser des processus physiques. Il fait appel aux structures d'automates[11]. Un automate représente un système physique. Il est composé d'états (les cercles sur la figure), qui correspondent aux états du système réel, et de transitions (les flèches sur la figure), pour passer d'un état à l'autre. Il existe aussi la notion de chemin : par exemple pour passer de 0 à 3 sur la figure, il faut passer par 1 puis 2 : le chemin de 0 à 3 est 0,1,2,3.

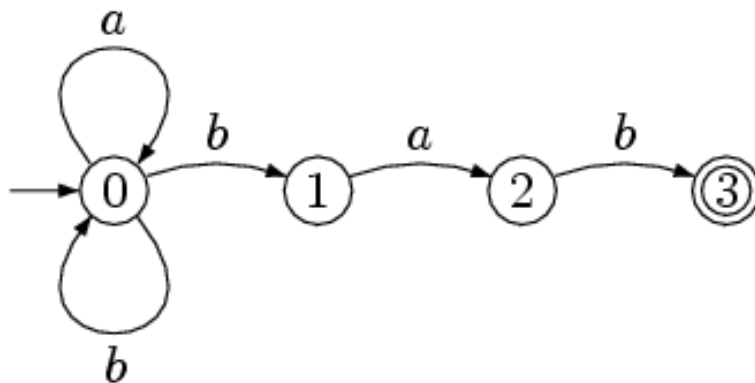


FIGURE 3.1 – Exemple d'automate « classique »

Les modèles de Markov cachés sont largement répandus dans la reconnaissance vocale([12], [3] et [13]). Entre un modèle discret et un modèle continu, nous avons choisi ce dernier car les données en entrée ne font pas partie d'un ensemble fini : il existe une infinité de sons possibles pour un même phonème. Les modèles de Markov cachés sont particulièrement adaptés pour la reconnaissance vocale car ils permettent un apprentissage constant de la part du programme : celui-ci est capable d'apprendre de nouveaux mots de manière autonome, et de s'améliorer au-fur-et-à-mesure que la base de données de mots grandit.

Nous avons modélisé chaque mot par un automate, dont les états sont les différents phonèmes du mot. Lorsque l'on prononce un mot, on se dirige dans l'automate grâce aux phonèmes prononcés, jusqu'à rencontrer l'état final. Ceci permet de reconnaître le mot même si une syllabe dure plusieurs secondes : dans ce cas, on se contente de tourner en rond (en restant sur l'état 0 de la figure par exemple) dans l'automate jusqu'à rencontrer un nouveau phonème. Dans l'automate, la transition de l'état i à k représente la probabilité de passer de l'état i à k , c'est-à-dire la probabilité que le phonème $n^o k$ vienne tout de suite après le phonème $n^o i$.

Exemple:

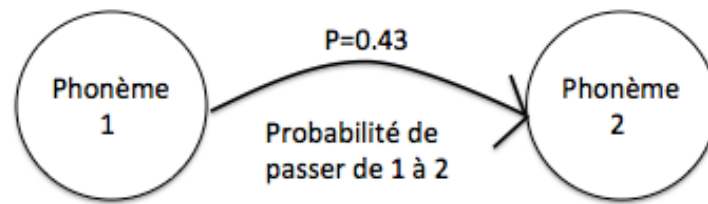


FIGURE 3.2 – Exemple de deux phonèmes et de la probabilité de passer du phonème 1 au phonème 2

3.2 Principaux algorithmes sur les modèles de Markov

Lorsque l'on fait passer un mot dans un automate, ie. qu'on s'oriente dans l'automate à l'aide des phonèmes, on peut calculer la probabilité que le mot corresponde à cet automate : on multiplie toutes les probabilités rencontrées pendant le parcours. Elles dépendent bien sûr du chemin parcouru (i-e des transitions rencontrées). C'est le principe de l'algorithme *forward*.

L'algorithme de Baum-Welch permet d'optimiser un automate. En se plaçant dans l'ensemble des modèles de Markov, on cherche à faire converger une suite d'automates définis à l'aide de plusieurs versions d'un même mot vers un automate optimisé qui corresponde au mieux au mot.

3.3 Application à notre objectif

Résumons la situation lorsque l'on lance notre programme : d'un côté une base de données de mots, représentés chacun par un automate ; de l'autre, un fichier audio : le mot prononcé par l'utilisateur. Le programme se déplace dans chaque automate grâce au fichier audio, il s'oriente en fonction des phonèmes prononcés. Nous appellerons cette opération "faire passer un mot dans un automate".

L'algorithme *forward* permet donc de calculer la probabilité qu'un automate corresponde au mot prononcé : en comparant les probabilités dans chacun des automates, on sélectionne la plus grande et on a l'automate qui correspond le mieux au mot sélectionné.

L'algorithme de Baum-Welch permet l'apprentissage de nouveaux mots : pour chaque nouveau mot il crée un nouvel automate, et le rend le plus optimisé possible en s'appuyant sur la bibliothèque existante. C'est ce que fait la partie logicielle de notre programme, pour que les programmeurs puissent agrandir la base de données.

3.4 Phase d'apprentissage

Une fois l'algorithme de reconnaissance vocale implémenté, il nous a fallu l'améliorer. Deux aspects demandent un apprentissage de la part du programme. Il doit d'abord faire grossir l'ensemble des mots reconnus, de manière à pouvoir en reconnaître le plus possible. Mais il est aussi intéressant de lui faire apprendre un mot par des locuteurs différents. Plus le nombre de locuteurs est grand, plus l'algorithme peut être précis.

Enregistrer plusieurs personnes permet d'obtenir une diversité de spectres qui accroît la précision du programme.

Une fois un mot appris, il est également très utile qu'un même locuteur enregistre de nombreuses versions du mot. Nous avons fait pour notre locuteur 10 versions de chaque mot.

Pour mettre en place un apprentissage, nous avons des besoins matériels (stocker l'ensemble des mots reconnus) mais aussi des besoins humains, et en l'occurrence une diversité de voix.

3.5 Phase de reconnaissance

La phase de reconnaissance constitue le cœur du programme. Comme dit précédemment, le programme effectue l'algorithme *forward* sur chacun des automates et renvoie le mot le plus probable, après avoir comparé toutes les probabilités.

A l'origine, la phase de reconnaissance a été codée en Python. Cependant le temps d'exécution était trop long, nous l'avons donc codé en C++, ce qui a permis de diviser le temps d'exécution par 50 000. Grâce à ce travail laborieux, le programme s'effectue en un temps proche de la seconde. Tout a été mis en place, notamment en amont avec le codage en C++ de la transformée de Fourier rapide, pour privilégier la rapidité de l'exécution.

Au départ nous n'avions qu'un seul locuteur pour faire la base de donnée des mots reconnus, ce qui ne permettait de faire fonctionner le programme que pour un seul utilisateur : celui qui avait enregistré les mots. Cependant nous avons réussi à enregistrer plusieurs locuteurs, ce qui a permis au programme de reconnaître plusieurs utilisateurs, même un utilisateur qui n'avait enregistré aucun mot.

Récapitulatif

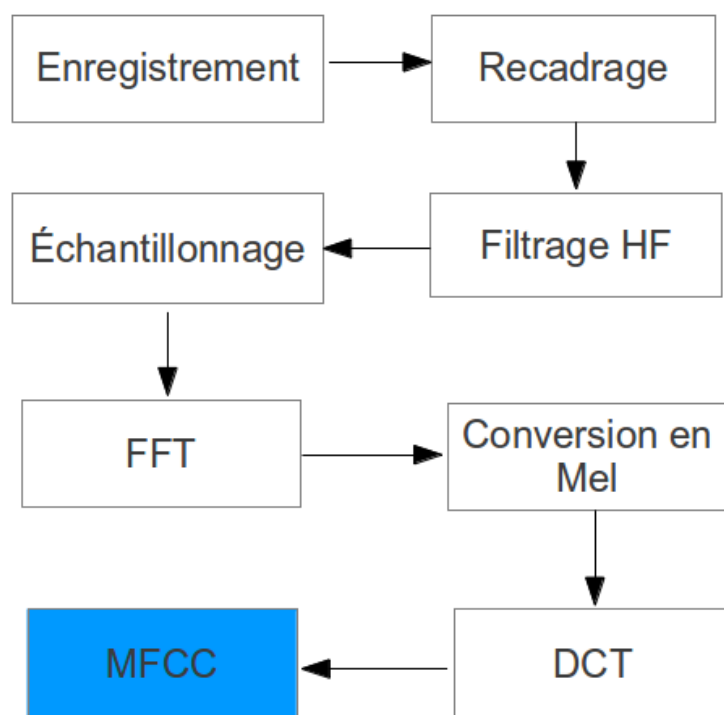


FIGURE 3.3 – exemple

Deuxième partie

Approche commerciale

1. Approche développement web

Pour assurer une rentabilité à notre projet, il nous faut le penser, le structurer en vue d'une large distribution sous de multiples formes.

Nous sommes dotés d'une identité porteuse de ce projet. Le groupe de travail est baptisé The SpeechApp Company. Visuellement, elle se constitue en premier lieu d'un logo : Un micro, élément central du projet, dont la tête est le logo de Mines ParisTech, signe de notre appartenance à l'école et de l'aide qu'elle nous a apporté dans le projet.

1.0.1 Choix d'une architecture optimale pour notre projet

Distribuer notre projet tel quel présenterait à ce stade de nombreux défauts : - le coeur de notre technologie de reconnaissance vocale est directement accessible à tous. - une interface unique en ligne de commande constitue un blocage majeur pour la majorité des utilisateurs finaux et empêche une intégration large à des applications tierces.

Étudions l'opportunité d'adopter une architecture client/serveur pour ce projet.

Dans ce scénario, divers clients logiciels, potentiellement indépendant de The SpeechApp Company pourraient communiquer par requêtes/réponses (spécifiées par une API) avec les serveurs de The SpeechApp Company. Ces derniers seuls auraient accès au coeur algorithmique du projet, qui resterait ainsi exclusivement entre nos mains. Par leurs requêtes, les clients demanderaient l'analyse automatique de mots, l'ajout de nouveaux mots ainsi que toute autre opération pertinentes relative à l'analyse et la gestion d'une base de données de mots. L'accès à notre API serait monétisable forfaitairement ou à l'utilisation.

Les mots enregistrés par les clients seraient conservés dans des bases de données chez The SpeechApp Company. La location de ses bases de données hébergées serait monétisable. Alors, The SpeechApp Company pourrait prioritairement développer deux applications connectables au serveur : la première, SpeechCreator, permettrait l'enregistrement aisé de nouveaux mots dans les bases de données clients. La seconde, SpeechApp, permettrait, au travers d'une application Web riche, de tester la reconnaissance vocale en ligne.

Cette configuration permettrait aussi à une multitudes d'applications tierces d'utiliser notre technologie en ne voyant de l'extérieur qu'une API définissant le format des requêtes et réponses dans la communication entre clients logiciel et serveur.

Nous aboutirions alors à l'architecture représentée par le schéma suivant :

Plus précisément dans le cadre des échanges entre le SpeechServer Les requêtes pourraient être traitées de la façon suivante : Le client (au sens logiciel toujours) envoie au SpeechServer une requête HTTP POST contenant un formulaire avec en particulier son identifiant, son mot de

passer, la base de données qu'il veut utiliser, l'action qu'il veut faire effectuer au SpeechServer, et les données d'entrée qui lui sont associées. La requête analysée par le SpeechServer, les opérations adéquates ayant été réalisées par le coeur algorithmique, le SpeechServer répond au client par une réponse HTTP POST contenant des données au format XML. Le client peut alors lire et interpréter la réponse donnée par le SpeechServer.

Avec ses spécifications, nous obtiendrions le cycle suivant pour la reconnaissance d'un mot par SpeechApp :

L'architecture client/serveur proposée présenterai pour nous l'avantage de

- permettre la création d'un écosystème varié d'applications basées sur le coeur algorithmique de The SpeechApp Company via l'API de son SpeechServer, et générant ainsi des revenus
- conserver le coeur de notre travail entre nos mains et même de nous donner le contrôle sur toute la chaîne

L'architecture client/serveur proposée présenterai pour nos clients l'avantage de

- ne pas se soucier du coeur algorithmique de la reconnaissance vocale, en n'y voyant que l'API de SpeechServer. Cette API peut offrir par ailleurs une grande liberté d'action
- n'avoir pas ou peu d'investissement initial de développement à effectuer, nos applications propriétaires SpeechApp et SpeechRecorder pouvant être intégrées sous forme de widgets aux applications tierces
- ne pas avoir à faire de lourds calculs eux-mêmes, ceux-ci étant réalisés par les machines de The SpeechApp Company
- les cycles de mises à jour seraient en majorité invisibles chez les clients, l'API restant immuable sur des cycles plus long (Long Term Support)

Au vu des nombreux avantages qu'elle présente, *Nous avons donc opté pour une architecture modulaire client/serveur pour notre projet.*

1.0.2 Réalisation du SpeechServer

Le SpeechServer a été codé en Python. Python a une librairie standard suffisamment riche pour n'avoir à traiter ce problème qu'à un haut niveau (en réception de requêtes selon leurs méthodes). De plus, ce choix facilite les interactions avec le coeur algorithmique : des imports et appels de fonctions depuis le SpeechServer suffisent.

Finalement, le SpeechServer prend seulement la forme d'un programme python à lancer sur un ordinateur.

Il écoute alors les requêtes sur le port 8010 (par défaut) de l'ordinateur. Lorsqu'il en reçoit, il interagit avec le coeur algorithmique et le système de gestion de bases de données mis en place.

1.0.3 Système de Gestion de Base de Données (SGBD)

Le SGBD doit permettre de stocker et gérer les fichiers audio associés aux mots (au moins une dizaine d'enregistrement par mot), les modèles de markov cachés qui leur sont associés ainsi que les données d'authentification des applications clientes.

Le standard actuel de gestion de bases de données est le modèle relationnel basé sur le langage SQL. Néanmoins, dans le cas précis de stockage de fichiers relativement lourds (> 0.1 Mo), la lecture/écriture des données directement sur le disque dur s'avère plus performante.

Nous avons donc fait le choix de stocker nos données sur le disque dur du serveur, en enregistrant les fichiers audio en format brut, et les autres données (modèles de Markov, données d'authentification) comme des objets Python, avec le module pickle de la librairie standard.