

[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	9
2.4	Echantillonnage, fenêtrage	10
2.5	Transformée de Fourier	11
2.6	Mel	11
2.7	Transformée inverse	11
3	Modèles de Markov cachés (MMC)	13
3.1	Prérequis et principe	13
3.2	Principaux algorithmes sur les modèles de Markov	14
3.3	Application à notre objectif	14
3.4	Phase d'apprentissage	14
3.5	Phase de reconnaissance	15
II	Approche commerciale	17
1	Côté web	17
2	Applications	18

3	Budget, modèle économique	20
3.1	Introduction	20
3.2	Les salaires	20
3.3	Le compte de résultat prévisionnel	21
3.4	Le bilan	21
3.5	Les impôts	22
3.6	Conclusion et vue sur le long terme	22
III	Le Code	23
1	Python	23
2	C/C++	25

Introduction

0.1 Présentation des enjeux

La reconnaissance vocale automatisée est l’objet d’intenses recherches depuis plus de 50 ans. Malgré son caractère tout 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.

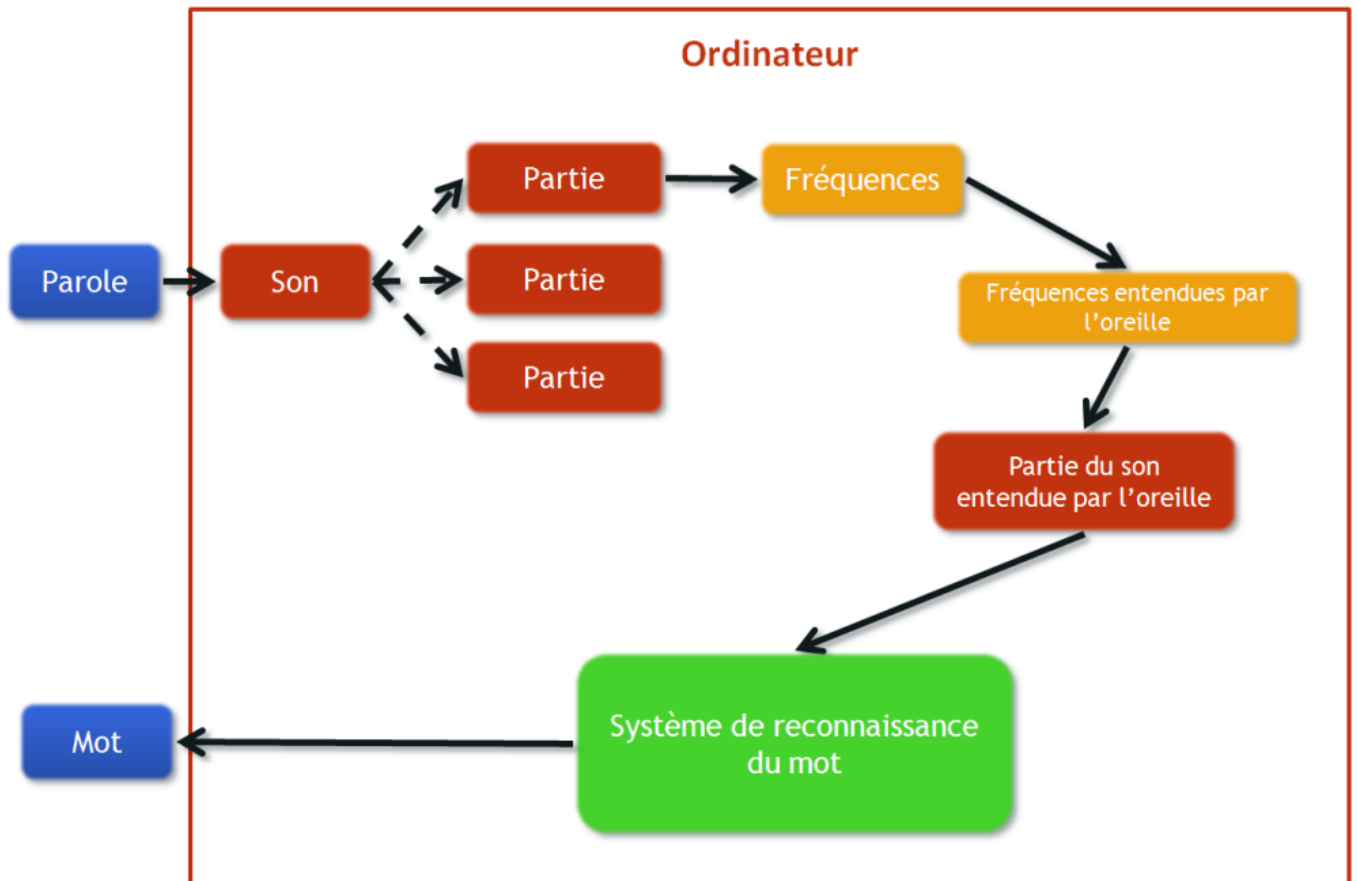


FIGURE 1.1 – Traitement du son pour le reconnaître

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

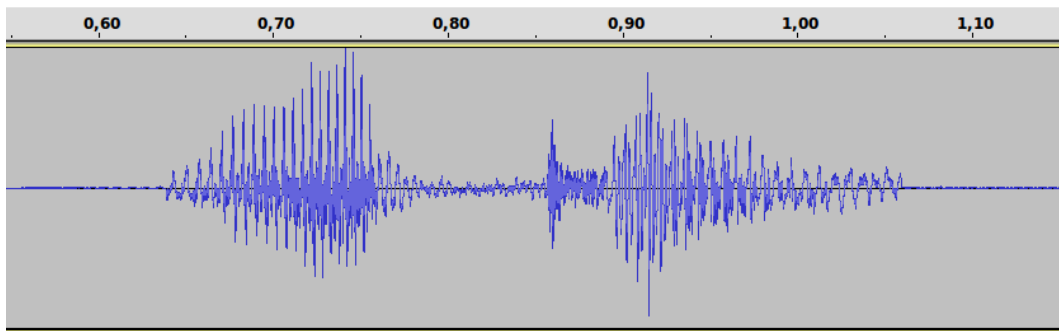


FIGURE 2.1 – Exemple audiogramme prononciation du mot "VICA"

par seconde, de la façon suivante : $bitrate = n \times \text{échantillonnage}$. 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 la 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

Nam eu sollicitudin massa.[3]. Aliquam a fermentum lectus. Praesent dapibus molestie mauris sed vestibulum. Curabitur sodales egestas est a pellentesque. Quisque id vulputate erat, a sagittis turpis.

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].

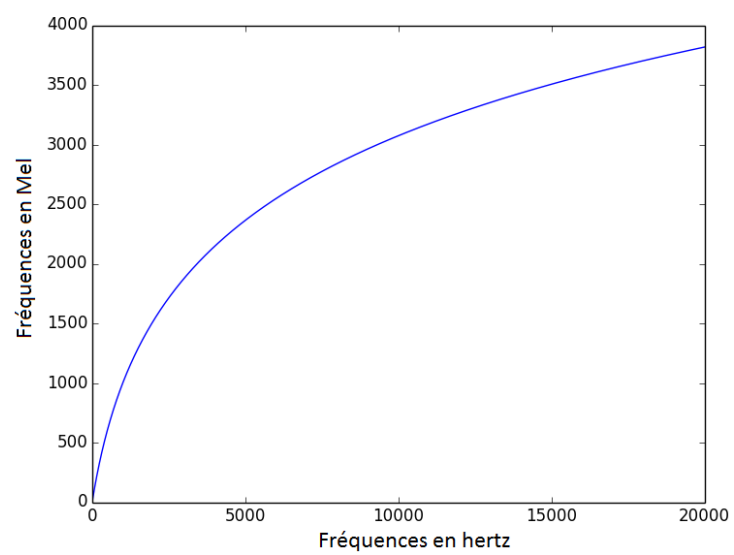


FIGURE 2.2 – Graphe de conversion

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.

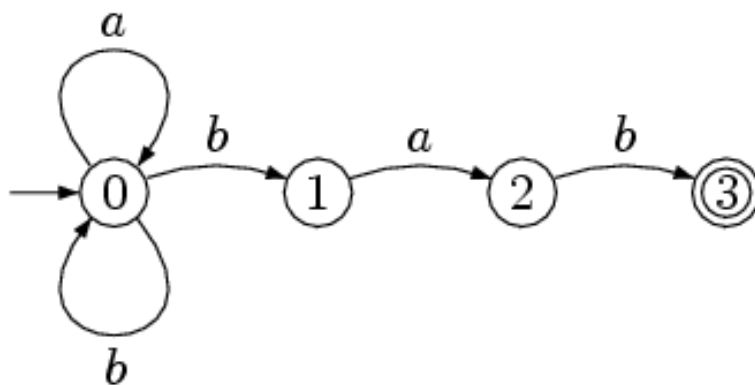


FIGURE 3.1 – Exemple d'automate « classique »

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:

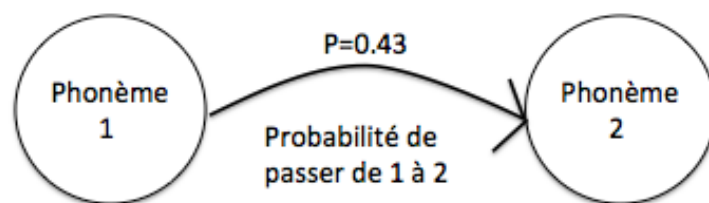


FIGURE 3.2 – Exemple de deux phonèmes et de la probabilité de passer du phonème 1 au phonème 2

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

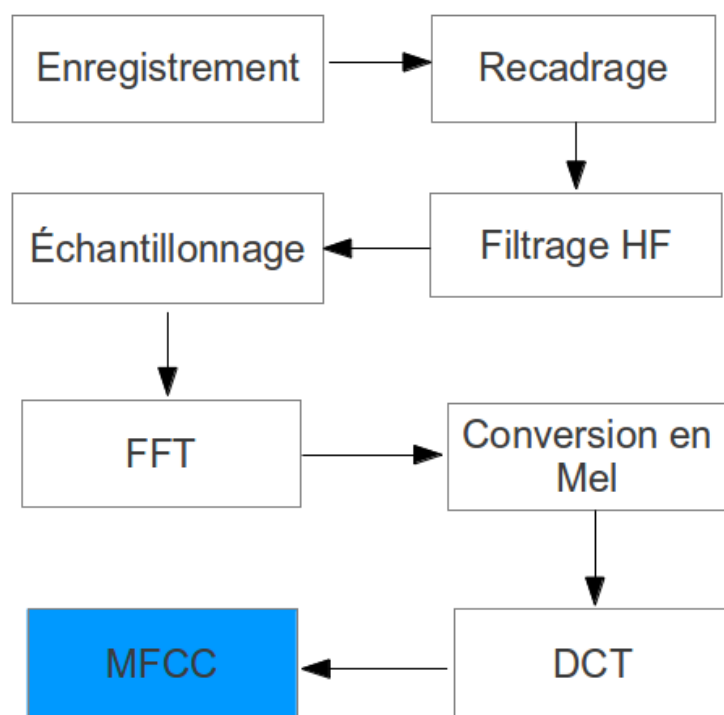


FIGURE 3.3 – exemple

Deuxième partie

Approche commerciale

1. Côté web

2. Applications

La reconnaissance vocale est une technologie promise à un futur radieux ; les plus grands noms de l'informatique, dont Bill Gates, annonçaient il y a quelques années qu'elle allait remplacer les claviers d'ici peu. Il s'avère aujourd'hui que leurs prédictions ne sont pas encore réalisées, il est tout à fait possible qu'elle se réalise plus tard que prévu. Le principal obstacle à l'explosion de cette technologie étant le manque de fiabilité totale, mais avec les progrès à venir, la technologie deviendra de plus en plus sûre.

L'armée étatsunienne a bien compris le potentiel de cette technologie : elle investit massivement depuis des années dans la recherche pour la développer[14]. Elle est d'ailleurs déjà utilisée sur certains avions de chasse, et pas seulement aux Etats-Unis : en France, en Angleterre et en Suède aussi notamment. Vu les investissements massifs, il y a fort à penser que les armées de ces différents pays ont des techniques bien plus avancées que celles connues du grand public, qui sont déjà plutôt performantes. Pour le moment, les commandes vocales ne servent pas encore à des fonctions critiques comme lancer un missile, et elles demandent toujours la confirmation du pilote avant d'exécuter une action. Elles libèrent néanmoins considérablement le pilote de beaucoup de tâches secondaires, ce qui lui permet de se concentrer sur les fonctions critiques. La technologie est également utilisée sur certains hélicoptères, notamment le célèbre Puma de l'armée française. Dans les deux cas, elle demande une grande fiabilité dans des conditions de stress et de bruit ambiant énorme (en particulier pour les hélicoptères, dans lesquels les pilotes n'ont souvent pas de casque anti bruit). Dans ce domaine, les perspectives sont donc très intéressantes financièrement mais elles demandent un savoir-faire qui est totalement hors de notre portée.

La reconnaissance vocale est également utilisée dans le contrôle aérien, et pourrait à terme remplacer les contrôleurs aériens. En effet, les phrases utilisées dans ce contexte sont très typées, ce qui favorise la reconnaissance (phrases souvent identiques, syntaxe très simple, prononciation très articulée). La technologie est donc moins avancée que dans le domaine de l'armée, et elle est déjà utilisée aux Etats-Unis, en Australie, en Italie, au Brésil et au Canada. Notre produit pourrait servir à ce type d'application, en créant une base de données spécifique au contrôle aérien.

La reconnaissance vocale se développe dans de nombreux domaines professionnels où les tâches administratives prennent beaucoup de temps, notamment la médecine, le droit et la police. En médecine, elle permet de remplir des rapports médicaux automatiquement : une simple relecture est alors nécessaire. Elle est notamment déjà utilisée dans 95% des hôpitaux aux Pays-Bas. Pour le droit, elle pourrait remplacer le travail du greffier pour prendre des notes dans les tribunaux. Et pour la police, elle permet de rédiger des rapports environ trois fois plus vite qu'au clavier. Le besoin de fiabilité est bien moindre dans ces domaines que dans les domaines de l'armée ou du contrôle aérien, une relecture est souvent largement nécessaire. Dans le domaine du droit, il faut néanmoins prendre en compte les conditions particulières d'enregistrement (brouhaha ambiant, émotions dans la voix, volume variable...). Notre produit peut tout à fait servir à ce type d'applications, à condition de créer une base de données spécifique aux domaines concernés.

Une autre application possible de la reconnaissance vocale est l'aide aux handicapés, par

exemple des commandes vocales pour une chaise roulante. Les phrases utilisées sont très typées (avancer, reculer,...) donc la technologie n'a pas besoin d'être très avancée. De plus, avec la possibilité qu'offre notre produit d'ajouter ses propres mots à la base de données, l'utilisateur lui-même peut rentrer les commandes ce qui assure un taux de reconnaissance très élevé. Notre produit peut donc bien s'adapter à cette utilisation.

La technologie est également très utilisée pour un usage plus ludique : fonctions de recherche dans les téléphones mobiles, les ordinateurs, robotique, jeux vidéo, traduction automatique,... Notre produit, dans sa version pour les particuliers, peut servir à ces usages même si la concurrence ne manque pas.

Enfin, la reconnaissance vocale peut servir à des fins sécuritaires, pour des vérifications d'identité. Il s'agit alors de reconnaître le locuteur, ce que notre produit ne permet pas.

Pour conclure, les applications pour notre produit sont assez nombreuses, et la demande est de plus en plus forte, ce qui montre sa pertinence.

3. Budget, modèle économique

3.1 Introduction

Après les études techniques et théoriques, l'étude économique est une nécessité. Elle est au coeur des problématiques de l'ingénieur, car c'est elle qui permet de dire si le projet est viable ou non.

Dans le cas de la programmation d'un logiciel de reconnaissance vocale, divers facteurs sont à prendre en compte, comme les salaires des employés, la communication sur le produit ou les impôts à payer. Il s'agit également de trouver le meilleur moyen pour vendre le logiciel. Faut-il le vendre pour iPhone sur l'App Store ? Le réserver à un public restreint (majoritairement des entreprises) ou le proposer également à des particuliers ?

La concurrence importante nous oblige à être à la fois ambitieux et prudent. Nous avons donc décidé d'envisager à la fois la vente sur notre site internet d'un logiciel pour les particuliers, et de proposer des licences en parallèle, permettant notamment aux entreprises d'accéder à nos bases de données, les compléter et créer leurs propres dictionnaires.

3.2 Les salaires

SALAIRES					
Catégorie	Salaire brut	Charges salariales	Salaire net	Charges patronales	Budget
Personnel					
David Vitoux	2290,00 €	503,80 €	1786,20 €	1007,60 €	3 297,60 €
Axel Goering	2290,00 €	503,80 €	1786,20 €	1007,60 €	3 297,60 €
Sofiane Mahiou	2290,00 €	503,80 €	1786,20 €	1007,60 €	3 297,60 €
Maxime Ernoult	2290,00 €	503,80 €	1786,20 €	1007,60 €	3 297,60 €
Adrien De La Vaissière	2290,00 €	503,80 €	1786,20 €	1007,60 €	3 297,60 €
Clément Joudet	2290,00 €	503,80 €	1786,20 €	1007,60 €	3 297,60 €
Clément Roig	2290,00 €	503,80 €	1786,20 €	1007,60 €	3 297,60 €
Anis Khlif	2290,00 €	503,80 €	1786,20 €	1007,60 €	3 297,60 €
Paul Mustière	2290,00 €	503,80 €	1786,20 €	1007,60 €	3 297,60 €
Matthieu Denoux	2290,00 €	503,80 €	1786,20 €	1007,60 €	3 297,60 €
Julien Caillard	2290,00 €	503,80 €	1786,20 €	1007,60 €	3 297,60 €
Nathanaël Kasriel	2750,00 €	605,00 €	2145,00 €	1210,00 €	3 960,00 €
Thomas Debarre	2750,00 €	605,00 €	2145,00 €	1210,00 €	3 960,00 €
Total	30690,00 €	6751,80 €	23938,20 €	13503,60 €	44193,60 €

FIGURE 3.1 – Salaires

Treize employés travaillent sur le projet, pendant un temps effectif d'environ un mois. Deux d'entre eux s'occupent des ressources humaines, pour un salaire de 2750 € brut mensuel, les autres étant considérées comme des développeurs de moins de deux ans d'expérience, avec un salaire de 2290 € brut mensuel.

Sur ce salaire brut, l'employé paye environ 22% de charges salariales, et l'entreprise 44% de charges patronales.

3.3 Le compte de résultat prévisionnel

Le compte de résultat prévisionnel dresse l'ensemble des charges (fixes et variables) de l'entreprise, ainsi que ses produits (recettes).

Pour parvenir à un équilibre budgétaire, il nous faut, pour la première année, vendre 24 000 logiciels à un prix de 4,17 € hors taxes, et une dizaine de licences permettant d'accéder à nos bases de données pour un prix de 833,33 €. Au niveau des charges, l'ensemble des salaires cités plus haut est à prendre en compte, ainsi que le coût de notre campagne de publicité. Celle-ci peut être décrite en deux principaux pôles : des articles de journaux spécialisés, gratuits, et des annonces google. On peut estimer le prix d'une telle annonce à 10 centimes d'€ le clic. En estimant que 10% des visiteurs du site par l'intermédiaire de l'annonce vont acheter le produit, on peut évaluer le coût de la publicité à 24 000 €.

La différence des produits et des charges donne alors un chiffre de 56 315 €.

La différence des produits et des charges donne alors un chiffre de 164 423€.

COMPTE DE RÉSULTAT PRÉVISIONNEL							
	Produit					Charges	
	Vente	Prix unité (Hors taxes)	Prix unité TTC	Nombre	Total	Salaires	45 777,00 €
	Logiciel	4,17 €	5,00 €	24000	120 096,00 €	Frais pub (google)	24 000,00 €
	Licences	833,33 €	1 000,00 €	10	10 000,00 €		
Total	60319,00 €						

FIGURE 3.2 – Compte de résultat prévisionnel

3.4 Le bilan

Actifs		Passif	
Actifs incorporels	0,00 €	Fonds propres	0,00 €
Créances	0,00 €	Dettes long terme	100000,00 €
Actifs immobiliers	0,00 €	Compte de Résultat prévisionnel	60319,00 €
Créances clients	0,00 €		
Trésorerie	0,00 €		

FIGURE 3.3 – Bilan

Le bilan prend en compte l'actif et le passif de l'entreprise.

Cette année, celle-ci n'a pas d'actif réel. Pas de trésorerie, de créances ou d'actifs immobiliers et incorporels. Son passif ne contient pas de fonds propres, et le compte de résultat prévisionnel a été explicité plus haut. On peut en revanche considérer que nous avons effectué un prêt à long terme de 100 000 €, afin de financer les primes du projet.

3.5 Les impôts

S'agissant des impôts, nous devons dans un premier temps reverser à l'Etat la TVA sur les produits que nous vendons, à un taux de 20% à compter du 1er Janvier 2014. Le logiciel étant vendu 4,17 € et la licence 833,33 €, le total de la TVA à reverser sera de 26 019 €.

Ensuite, l'impôt sur les sociétés est à un taux de 33% sur les bénéfices. A partir du bilan et de la TVA, on peut estimer nos bénéfices à 34 300 €, et donc un impôt sur les bénéfices à hauteur de 11 319 €.

Impôts		
Sur les sociétés	33% des bénéfices	11 318,93 €
TVA	20% sur les ventes	26 019,20 €

FIGURE 3.4 – Impôts

3.6 Conclusion et vue sur le long terme

En considérant un prêt à un taux de 3% sur cinq ans, et le prêt de locaux et matériels par un incubateur (l'école des Mines par exemple) l'entreprise est viable la première année à partir de 23 000 téléchargements.

En utilisant les mêmes calculs, pour les quatre années qui suivent, sans faire de mise à jour, il faudrait en moyenne 11 000 ventes de logiciels par an, et 3 ventes de licences.

Troisième partie

Le Code

1. Python

```
1 # -*- coding: utf-8 -*-
2
3 import numpy as np
4 import scipy as sc
5 import math
6 from operator import add
7
8
9 TAILLE_TABLEAU_MEL_ENTREE = 24
10
11 NOMBRE_COMPOSANTES_GARDEES = 13
12
13 B = TAILLE_TABLEAU_MEL_ENTREE
14
15
16 def inverseDCTI(x): # x represente le tableau en mel donne par les
    fonctions precedentes
17     X = np.zeros(B)
18     for k in range(B):
19         X[k] = (0.5 * (x[0] + math.pow(-1, k) * x[B-1]) + \
20             reduce(add, [x[n] * math.cos(math.pi * n * k / (B -
21                 1)) \
22                     for n in range(1,B - 1)])) * math.sqrt(2. / (B -
23                 1))
24     return X
25
26 def inverseDCTII(x):
27     X = np.zeros(B)
28     for k in range(B):
29         X[k] = reduce(add, [x[n]*math.cos(math.pi*(n + 0.5)
30             * k / B) \
31                 for n in range(B)]) * math.sqrt(2. / B)
32     return X
33
34 def inverseDCTIII(x):
35     X = np.zeros(B)
36     for k in range(B):
```



```

34         X[k] = (0.5*x[0]+reduce(add, [x[n] *
math.cos(math.pi * (k + 0.5) / B)\
35         for n in range(1, B)])) * math.sqrt(2. / B)
36     return X

```

2. C/C++

```
1 cDouble* fftCT(cDouble *sig)
2 {
3     int N = sizeof(sig)/sizeof(cDouble);
4     //int iMax = (int)(log(N)/log(2));
5     int i,j,k,p=0,f=1;
6     cDouble ekN;
7     cDouble **tmp = (cDouble**) malloc(2*sizeof(cDouble*));
8
9     for (i=0;i<2;i++)
10         tmp[i] = (cDouble*) malloc(N*sizeof(cDouble));
11     for (i=0;i<N;i++)
12         tmp[0][i] = sig[i];
13     for (i=N/2;i!=1;i/=2)
14     {
15         for (j=0;j<i;j++)
16             for (k=0;k<N/(2*i);k++)
17             {
18                 ekN = e(k,N)*tmp[p][i*(2*k+1)+j];
19                 tmp[f][i*k+j] = tmp[p][i*(2*k)+j] + ekN;
20                 tmp[f][i*k+j+N/2] = tmp[p][i*(2*k)+j] - ekN;
21             }
22         p = f;
23         f = (p+1)%2;
24     }
25     return tmp[p];
26 }
```

Conclusion

Le marché de la reconnaissance vocale est pour le moment assez restreint, mais est appelé à grandir dans les prochaines années. Si les systèmes de reconnaissance vocale fleurissent sur les objets multimédias à usage personnel, comme les ordinateurs portables ou les téléphones mobiles, ils servent uniquement à simplifier un peu certaines tâches de l'utilisateur, et ne sont en pratique que très peu utilisés, ce qui s'explique par leurs performances moyennes. Le représentant le plus utilisé de ce type d'usage de la reconnaissance vocale est probablement Siri sur les téléphones mobiles iPhone d'Apple, mais il reste assez peu utilisé malgré la grande popularité de l'iPhone.

Dans le domaine des logiciels payants, pour un usage plus sérieux, le marché est dominé par les logiciels Dragon NaturallySpeaking de la firme américaine Nuance. Les prix, selon les modèles, varient entre environ 100\$ pour le modèle de base et environ 1000\$ pour les versions spécialisées dans un domaine professionnel. Le principe est que plus la base de données de mots est grande, plus les erreurs sont fréquentes ; Dragon NaturallySpeaking propose donc des versions adaptées à un domaine particulier. Par exemple, il existe une version "juriste" avec une base de données contenant surtout du vocabulaire technique de droit, et une version "médecine" avec des termes techniques médicaux. Ces versions visant une cible très précise donc plus restreinte, ils sont vendus considérablement plus cher que les versions plus classiques. Cependant, la demande étant en constante augmentation - un tiers des radiologues français utilisent cette technologie, tout comme 95% des hôpitaux aux Pays Bas -, le marché est assez prometteur. En effet, cette technologie réduit considérablement les tâches administratives de ces professions : une simple relecture au plus est nécessaire. La réussite est renforcée par des taux de réussite exceptionnels avec des bases de données adaptées, et par l'absence de concurrence très forte.

Cependant, nous avons choisi de concevoir un logiciel avec une base de données moins spécialisée pour un usage personnel : en effet, dans le temps qui nous est imparti, créer des bases de données étudiées spécialement pour un certain domaine (droit, médecine) nous paraissait très compliqué. Il aurait fallu faire une étude linguistique très poussée pour construire la base de données, alors que nous avons concentré l'essentiel de nos efforts sur l'algorithme de reconnaissance lui-même. Notre produit est donc destiné à un usage plus ludique, ou du moins personnel. Notre cible est donc légèrement différente, puisque les professionnels intéressés par notre produit doivent construire eux-mêmes leur base de données spécifique à leur domaine. L'inconvénient de cette approche est le désagrément de devoir ajouter soi-même les mots, l'avantage étant que la reconnaissance sera plus fiable puisque la voix de l'utilisateur elle-même sert de comparateur, et elle permet d'avoir une base de données réellement personnalisée (celles de Dragon, bien que dédiées à un domaine, ne sont pas totalement personnelles). Le prix envisagé de la licence pour cette utilisation de notre produit est comparable (de l'ordre de 1000€) à celui des versions personnalisées de Dragon.

Nous prévoyons également de mettre en vente une version à usage personnel, sans possibilité d'ajouts de mots, au prix de 5€. Il est difficile de prévoir le potentiel de cette version, puisque les concurrents sont très nombreux, de qualité et de prix très variables.

Bibliographie

- [1] Apple. Application siri. <http://www.apple.com/fr/ios/siri/>, 2013.
- [2] Tom Preston-Werner, Chris Wanstrath, and PJ Hyett. Github. <http://www.github.com/>, 2013.
- [3] Lawrence Rabiner. *Fundamentals of Speech Recognition*. Prentice Hall PTR, 1993.
- [4] MIT. Pyaudio. <http://people.csail.mit.edu/hubert/pyaudio/>, 2006.
- [5] Stanley Smith Stevens, John Volkman, and Edwin B. Newman. A scale for the measurement of the psychological magnitude pitch. *Journal of the Acoustical Society of America*, 1937.
- [6] Begam, Elamvazuthi, and Muda. Voice recognition algorithms using mel frequency cepstral coefficient (mfcc) and dynamic time warping (dtw). *Journal of Computing*, 2010.
- [7] Vincent Arsigny. Modélisation par un champ de markov du signal de parole et application à la reconnaissance vocale. Technical report, École Nationale Supérieure des Télécom de Paris, 2000.
- [8] Zohar Babin. How to do noise reduction using ffmpeg and sox. <http://www.zoharbabin.com/how-to-do-noise-reduction-using-ffmpeg-and-sox/>, 2011.
- [9] Chris Bagwell. Sox website. <http://sox.sourceforge.net/Docs/Documentation>, 2009.
- [10] Anonyme. Théorème de nyquist-shannon. http://fr.wikipedia.org/wiki/Th%C3%A9or%C3%A8me_d%27%C3%A9chantillonnage_de_Nyquist-Shannon, 2013.
- [11] Luc Maranget. *Introduction à la programmation*. École Polytechnique, 2008-2009.
- [12] Maurice Charbit. Reconnaissance de mots isolés (utilisation des modèles HMM). *Inconnu*, Oct. 2002.
- [13] Franck Bonnet, Benjamin Devèze, Mathieu Fouquin, and Julien Jeany. Reconnaissance automatique de la parole. *Epita*, 2004.
- [14] Aharon Etengoff. Nuance clinches speech-recognition deal with us army. <http://www.itexaminer.com/nuance-clinches-speech-recognition-deal-with-us-army.aspx>, 2009.