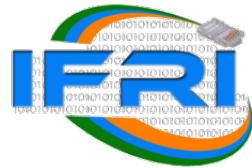




RÉPUBLIQUE DU BÉNIN
MINISTÈRE DE L'ENSEIGNEMENT SUPÉRIEUR
ET DE LA RECHERCHE SCIENTIFIQUE

UNIVERSITÉ D'ABOMEY-CALAVI

INSTITUT DE FORMATION ET DE
RECHERCHE EN INFORMATIQUE



BP 526 Cotonou Tel : +229 21 14 19 88
<http://www.ifri-uac.net> Courriel : contact@ifri.uac.bj

MÉMOIRE

pour l'obtention du

Diplôme de Master en Informatique

Option : Systèmes d'Information et Réseaux Informatiques

Présenté par :

Sero Abdel-Kader SAKA

Mise en place d'un Cloud Rest API pour l'identification et l'authentification de locuteur par empreinte vocale

Sous la supervision :

Dr John AOGA

Membres du jury :

Nom et prénoms du président	Grade	Entité	Président
Nom et prénoms de l'examineur	Grade	Entité	Examinateur
Nom et prénoms du rapporteur	Grade	Entité	Rapporteur

Sommaire

Dédicace	ii
Remerciements	iii
Résumé	iv
Abstract	v
List of Figures	vi
Liste des Algorithmes	viii
Introduction	1
1 ETAT DE L'ART DE IDENTIFICATION PAR BIOMETRIE VOCALE	3
2 MODELISATION ET LE DEVELOPPEMENT DE NOTRE MODEL D'IA	22
3 MATERIELS ET METHODES	27
4 Résultats et discussions	49
Conclusion	55
Bibliographie	56
Bibliographie	56
Table des matières	57

Dédicace

Je tiens à dédier ce travail à la mémoire de mon père disparu trop tôt. J'espère que, du monde qui est sien maintenant, il apprécie cet humble geste comme preuve de reconnaissance de la part d'un fils qui a toujours prié pour le salut de son âme. Puisse Dieu, le tout puissant, l'avoir en sa sainte miséricorde !

Remerciements

C'est avec un grand plaisir que je réserve ces quelques lignes en signe de gratitude et de profonde reconnaissance à tous ceux qui, de près ou de loin, ont contribué à la réalisation et l'aboutissement de ce travail.

Je tiens tout d'abord à remercier **Le Docteur John AOGA** pour son soutien, son accompagnement, et surtout pour son aide précieuse tout au long de l'élaboration de ce travail.

Notre parfaite gratitude au corps professoral de l'**IFRI** en général et plus particulièrement au **Prof. Eugène C. EZIN, Directeur de l'Institut de Formation et de Recherche en Informatique (IFRI)**.

Je m'acquitte, enfin, volontiers d'un devoir de gratitude et de remerciements à **tous mes enseignants** pour la qualité de l'enseignement qu'ils ont bien voulu me prodiguer durant mes études afin de me fournir une formation efficiente.

Résumé

La parole est le moyen principal et le plus pratique de communication entre les gens. Que ce soit par curiosité technologique de construire des machines qui imitent les humains ou par le désir d'automatiser le travail avec des machines, la recherche en reconnaissance de la parole et des locuteurs, comme première étape vers une communication naturelle homme-machine, a suscité beaucoup d'enthousiasme au cours des cinq dernières décennies.

La reconnaissance de locuteur (RL) est le processus d'identification d'un locuteur en fonction des caractéristiques vocales de la parole donnée. Cela diffère de la reconnaissance de la parole où le processus d'identification est limité au contenu plutôt qu'au locuteur. Le processus de SR repose sur l'identification et l'extraction des caractéristiques uniques de la parole du locuteur. Les caractéristiques de la voix de la personne sont également connues sous le nom de biométrie vocale.

Un système de Reconnaissance de Locuteur (RL) est utilisé pour identifier et distinguer les locuteurs et extraire des caractéristiques uniques qui peuvent être utilisées pour la vérification ou l'authentification de l'utilisateur. L'identification de locuteur (IL) est connue comme le processus d'identification du locuteur à partir d'une énonciation donnée en comparant la biométrie vocale de l'échantillon donné du locuteur. Lorsque la voix est utilisée pour l'autorisation, on parle de vérification du locuteur. Le domaine d'application clé de la RL est la sécurité et les sciences médico-légales.

Les systèmes de RL sont également utilisés comme remplacement des mots de passe et d'autres processus d'authentification des utilisateurs (mot de passe vocal). Les sciences médico-légales appliquent la RL pour comparer les échantillons vocaux de la personne prétendument impliquée avec d'autres preuves obtenues, telles que des conversations téléphoniques ou d'autres preuves enregistrées. Ce processus est également appelé détection de locuteur.

L'aspect le plus important de l'utilisation des systèmes d'IL est l'automatisation de processus tels que la redirection des courriers des clients vers la bonne boîte aux lettres, la reconnaissance des interlocuteurs dans une discussion, l'avertissement des systèmes de reconnaissance de discours en cas de changements de locuteur, la vérification si un client est enregistré dans le système, etc. Ces systèmes d'IL peuvent fonctionner sans la connaissance de l'échantillon vocal d'un client car ils se limitent uniquement à l'identification d'un locuteur d'entrée à partir de la base de données existante de locuteurs.

Par conséquent, l'objectif de cette étude est de réaliser une revue systématique de la littérature sur les diverses approches de reconnaissance de locuteur et d'implémenter l'approche la plus prometteuse.

Mots clés : Reconnaissance de la parole et des locuteurs, biométrie vocale, échantillon vocal, biométrie vocale, reconnaissance de la parole

Abstract

Speech is the primary, and the most convenient means of communication between people. Whether due to technological curiosity to build machines that mimic humans or desire to automate work with machines, research in speech and speaker recognition, as a first step toward natural human-machine communication, has attracted much enthusiasm over the past five decades.

Speaker Recognition (SR) is the process of identifying the speaker according to the vocal features of the given speech. This is different to speech recognition where the identification process is confined to the content rather than speaker. The process of SR is based on identifying and extracting unique characteristics of the speaker's speech. The characteristics of voices of the person is also known as voice biometrics.

A SR system is used to identify and distinguish speakers and extract unique characteristics that may be used for user verification or authentication. Speaker Identification (SI) is known as the process of identifying the speaker from a given utterance by comparing voice biometrics of the given sample of the speaker.

When voice is used for authorization, it is termed as Speaker Verification. The key application area of SR is security and forensic science. SR systems are also used as a replacement for password and other user authentication processes (voiced password). Forensic science applies SR to compare the voice samples of the person claimed to be with other evidences obtained like telephone conversation or other recorded evidence.

This process is also referred as speaker detection. The most important aspect of using SI systems is for automating processes like directing clients' mails to the right mailbox, recognizing talkers in discussion, cautioning discourse acknowledgment frameworks of speaker changes, checking if a client is enlisted in the framework as of, and so on. These SI systems may work without the knowledge of a client's voice sample since they rely only on identifying an input speaker from the existing database of speakers.

Therefore, the objective of this study is to conduct a systematic literature review on various speaker recognition approaches and implement the most promising approach.

Key words: research in speech and speaker recognition, Speaker Recognition (SR), speech recognition, voice biometrics, voice sample

List of Figures

1.1	Architecture d'une IA de reconnaissance par biométrie vocale.	4
1.2	Ajout d'un individu à une IA de reconnaissance parempreinte vocale.	6
1.3	Intelligence Artificielle, Machin learning et Deep learning.	10
1.4	Machin learning Nouveau paradigme de programmation.	11
1.5	Un Reseau de neurone pour la reconnaissance des caractères.	13
1.6	modèle de classification de caractères manuscrits.	14
1.7	Paramétrisation d'un réseau de neurones	15
1.8	Détermination de la fonction Loss	16
1.9	Le Loss Score utilisé pour optimiser le Weigth à l'entrée	17
2.1	paramètres d'exécution du Colab	23
2.2	Ressource de l'environnement d'exécution Colab	24
2.3	Structure d'un RN de type TDNN	26
3.1	Structure d'un RN de type TDNN	29
3.2	caractéristiques des enregistrements vocaux	30
3.3	Le dataset dansle bucket Cloud Storage	30
3.4	enregistrementsvocaux d'un seul utilisateur	31
3.5	enregistrementsvocaux d'un seul utilisateur	31
3.6	enregistrementsvocaux d'un seul utilisateur	32
3.7	récupération des données	33
3.8	Fichier manifestes d'entraînement	34
3.9	Fichier des paramètres d'entraînement	35
3.10	paramètres de base d'entraînement	36
3.11	paramètres des chemins des fichiers manifestes	36
3.12	paramètres des chemins des fichiers manifestes	37
3.13	Autres paramètres	37
3.14	Paramètres d'augmentation des données avec les bruits	38
3.15	Paramètres d'augmentation des données avec les bruits	39
3.16	Code d'entraînement	40
3.17	commande de separation des données	40
3.18	Code de separation des données	41
3.19	ligne du fichier manifeste	42
3.20	Classe Brain	42
3.21	compute forward	43
3.22	augmentation des données	43
3.23	definition de la fonction de perte	44

3.24 Logs de l'entraînement variation de la function loss	45
3.25 Logs de l'entraînement, passage de l'Epoch 14 à l'Epoch 15	46
3.26 Les deux models	46
3.27 Préparation du model	47
3.28 Inférence sur le modèle pour l'authentification	47
3.29 Inférence sur le modèle pour l'identification	48
4.1 Liste des POI	50
4.2 Meilleure modèle sauvegardé sur Cloud Storage	51
4.3 récupération du meilleur modèle	51
4.4 Identification de locuteur	52
4.5 Identification de locuteur	53
4.6 Identification de locuteur	54

Liste des Algorithmes

Introduction Générale

Les temps ont changé, l'heure est à la mobilité. Or les usages et les règles de sécurité qui fonctionnaient à domicile sur un PC sont plus difficilement applicables en situation de mobilité, où l'usage du mot de passe atteint ses limites. Et pourtant, plus que jamais, il est primordial de protéger nos contenus, nos transactions et notre identité.

La problématique de l'authentification par mot de passe ou code pin est liée à la sécurité des données et à la facilité d'utilisation. En effet, d'un côté, l'utilisation d'un mot de passe ou d'un code pin peut garantir la confidentialité des informations personnelles et des données sensibles, en empêchant l'accès aux tiers non autorisés. D'un autre côté, ces méthodes d'authentification peuvent être facilement compromises, par exemple, en utilisant des attaques de force brute ou de phishing, ou en découvrant le mot de passe ou le code pin par l'observation ou la divulgation.

De plus, la complexité et la longueur des mots de passe peuvent rendre leur utilisation difficile et fastidieuse pour les utilisateurs, tandis que les codes pin peuvent être plus faciles à retenir mais également plus faciles à deviner.

Ainsi, la question de l'authentification par mot de passe ou code pin est un compromis entre la sécurité et la convivialité, et doit être abordée avec prudence et en prenant en compte les risques potentiels et les besoins des utilisateurs.

La reconnaissance par biométrie vocale est une technologie qui utilise la voix comme moyen d'authentification et de reconnaissance d'identité. Elle peut être utilisée dans de nombreux domaines tels que la sécurité, l'accès à des services en ligne, la surveillance ou encore la gestion des appels téléphoniques. Cependant, cette technologie soulève des problématiques importantes en termes de protection de la vie privée, de fiabilité des systèmes et d'éthique. De plus, les données vocales collectées peuvent être utilisées à des fins malveillantes si elles tombent entre de mauvaises mains. Enfin, la biométrie vocale pose des questions quant à la protection de la vie privée et des libertés individuelles, notamment en ce qui concerne l'utilisation des données vocales pour la surveillance de masse ou la reconnaissance automatique de la voix dans les médias.

Comment garantir la sécurité et la confidentialité des données collectées ? Comment éviter les erreurs de reconnaissance et les usurpations d'identité ? Quelles sont les conséquences éthiques de l'utilisation de la biométrie vocale ? Autant de questions auxquelles il est nécessaire de répondre pour assurer une utilisation responsable et efficace de cette technologie.

Face à la fraude documentaire et au vol d'identité, aux menaces du terrorisme ou de la cyber-criminalité, et face à l'évolution logique des réglementations internationales, de nouvelles solutions technologiques sont progressivement mises en œuvre.

Parmi ces technologies, la biométrie s'est rapidement distinguée comme la plus pertinente pour identifier et authentifier les personnes de manière fiable et rapide, en fonction de caractéristiques biologiques uniques.

Aujourd'hui, de nombreuses applications font appel à cette technologie.

Ce qui était autrefois réservé à des applications sensibles telles que la sécurisation de sites militaires est devenue une application grand public en développement rapide.

Un simple geste de la main, une pression du doigt sur un capteur, une expression prononcée de la voix, ou un regard d'une seconde vers une caméra suffisent pour prouver son identité.

L'authentification biométrique vocale facilite la vie des consommateurs et citoyens qui sont de plus en plus mobiles et connectés en leur apportant une alternative simple aux mots de passe et aux codes PIN pour s'authentifier.

Ainsi nous souhaiterons mettre en place un Cloud API pour l'authentification et l'identification des individus par biométrie vocale.

L'objectif principal de notre étude est de concevoir développer et déployer un Cloud AI pouvant être appelé par REST API capable d'authentifier et d'identifier les individus à travers leur biometrie vocale. La matérialisation de cet objectif, passera par l'aboutissement des objectifs spécifiques que sont :

De façon spécifique, il s'agira :

- la mise en place la sauvegarde et l'organisation du dataset;
- le choix des outils, des technologies et de architecture à adopter;
- la mise en place de modèles d'IA permettant les opérations d'identification et d'authentification par biometrie vocale ;
- la création mise en place des REST APIs ;

Outre l'introduction et la conclusion, le présent mémoire comprend quatre (04) chapitres.

Le premier chapitre présente l'état de l'art de l'identification par biometrie vocale. Dans le deuxième chapitre nous explorons les matériels et méthodes. Le troisième chapitre nous permet d'exposer la modélisation et le développement de notre modèle d'IA. Le quatrième chapitre nous permet de présenter les résultats découlant de nos travaux. Nous faisons ensuite une discussion et présentons les difficultés rencontrées dans la réalisation de cette étude.

ETAT DE L'ART DE IDENTIFICATION PAR BIOMETRIE VOCALE

Introduction

L'identification par biométrie vocale est le processus de reconnaissance et de vérification de l'identité d'une personne à partir de leur enregistrement vocal. Elle a des applications importantes telles que la sécurité a criminalistique, l'IoT, les centres d'appel, et l'interaction homme-machine. Elle a été révolutionnée par les techniques d'apprentissage profond ces dernières années, permettant une identification plus précise, même dans des environnements acoustiques difficiles avec des parasites sonores. Ces techniques exploitent de grandes quantités de données vocales pour former des modèles de réseaux neuronaux qui peuvent apprendre à extraire automatiquement des caractéristiques discriminantes des signaux vocaux.

Dans ce chapitre, nous clarifierons d'abord les concepts clés liés à notre étude. Par la suite nous présenterons les quelques techniques de mise en place de ce type de modèles d'Intelligence Artificielle.

1.1 Définitions

1.1.1 Biométrie

La biométrie comprend tout un ensemble de technologies et procédés de reconnaissance, d'authentification et d'identification des personnes à partir de certaines de leurs caractéristiques physiques (empreintes digitales, la forme de la main, du doigt, le réseau veineux, l'œil), biologiques (ADN, le sang, la salive) ou comportementales (reconnaissance vocale, la dynamique des signatures). Ces caractéristiques doivent être à la fois [5] :

- Universelles : pour être utilisables par tous ;
- Uniques : pour distinguer les personnes sans équivoque ;
- Invariables : pour permettre une utilisation tout au long de la vie ;

- Mesurables : pour permettre la comparaison.

Les différentes techniques utilisées font l'objet de recherches régulières, de développements et bien entendu, d'améliorations constantes. Toutefois, les différentes sortes de mesures n'ont pas le même niveau de fiabilité. On estime que les mesures physiologiques ont l'avantage d'être plus stables dans la vie d'un individu. Par exemple, elles ne subissent pas autant les effets du stress, contrairement à l'identification par mesure comportementale [5].

1.1.2 Biométrie vocale

l'authentification par Biométrie vocale consiste à utiliser la voix d'une personne comme caractéristique biologique d'identification unique afin de l'authentifier. La Biométrie vocale est considérée comme la méthode d'authentification la plus solide car elle lie l'identité à un individu réel ("ce que vous êtes" plutôt que "ce que vous savez" ou "ce que vous avez"). Elle extrait un ensemble de caractéristiques audio qui garantissent qu'un locuteur est une personne réelle plutôt qu'un enregistrement ou un synthétique voix d'un individu [9].

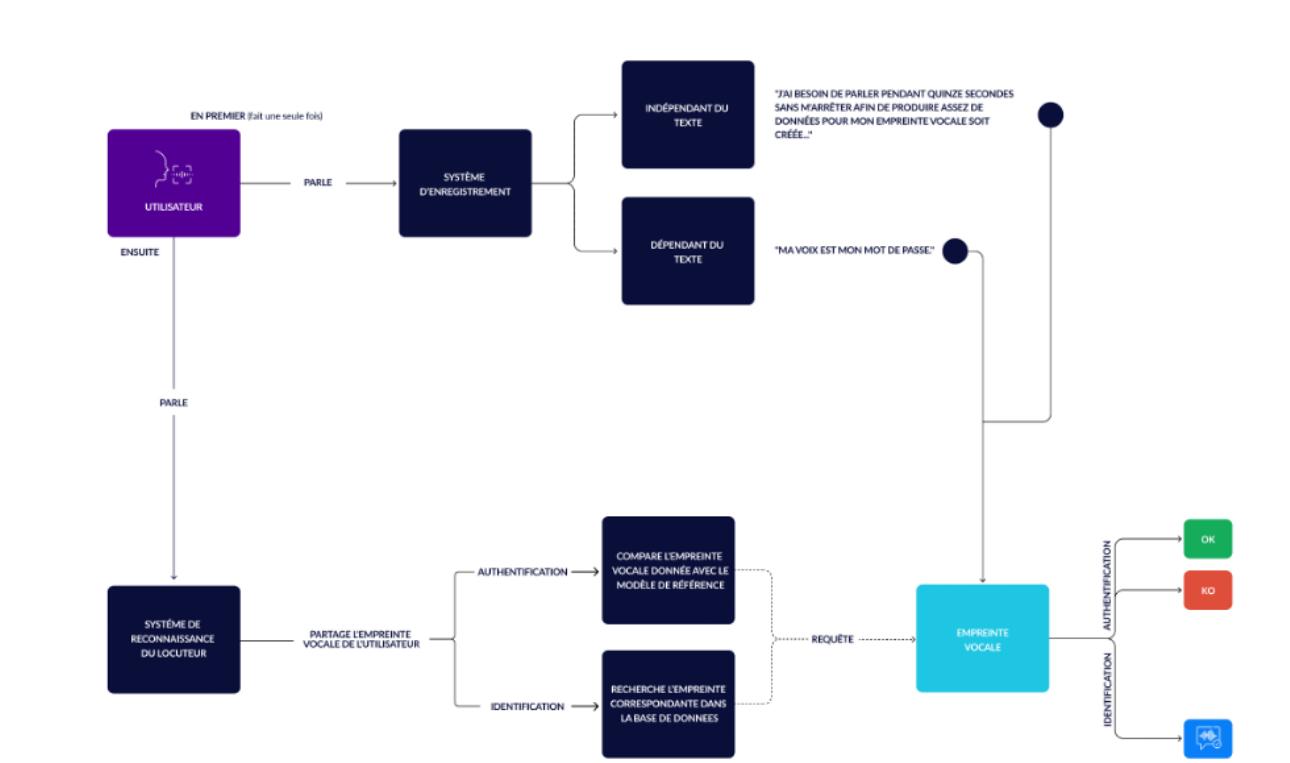


FIGURE 1.1 : Architecture d'une IA de reconnaissance par biométrie vocale.

Source: Source of the image.

1.1.3 Identification et authentification biométrique

La biométrie permet l'identification et l'authentification d'une personne à partir de données reconnaissables et vérifiables, qui lui sont propres et qui sont uniques.

L'identification consiste à déterminer l'identité d'une personne. Il s'agit de saisir une donnée biométrique de cette personne, en prenant par exemple une photo de son visage, en enregistrant sa voix, ou en captant l'image de son empreinte digitale. Ces données sont ensuite comparées aux données

biométriques de plusieurs autres personnes qui figurent dans une base; **ici, on essaye de répondre à la question : « qui êtes-vous ? ».**

L'authentification, appelée également vérification, est le processus qui consiste à comparer les données caractéristiques provenant d'une personne, au modèle de référence biométrique de cette dernière (« template »), afin de déterminer la ressemblance. Le modèle de référence est préalablement enregistré et stocké dans une base de données, dans un équipement ou objet personnel sécurisé. On vérifie ici que la personne présentée est bien la personne qu'elle prétend être. Dans ce cas on essaye de répondre à la question : « êtes-vous bien celui que vous prétendez être ? ».

1.1.4 Deep learning

Le deep learning ou apprentissage profond est une technique d'apprentissage automatique basée sur des réseaux de neurones artificiels. Dans le domaine de la reconnaissance vocale, l'apprentissage profond est largement utilisé pour améliorer la précision de la reconnaissance [10].

1.2 historique de la biométrie vocale

La biométrie répond à une préoccupation très ancienne de prouver son identité, de manière irréfutable, et en utilisant nos différences et particularités.

1952 : le premier système de reconnaissance vocale a été développé par Lawrence G. Roberts et ses collègues au laboratoire Lincoln du MIT. Le système utilisait l'analyse spectrographique de fichiers vocaux comparé à une base de données de d'empreintes vocales de personnes bien identifiées.

1962 : les laboratoires Bell ont créé le système "voiceprint", qui utilise un spectrogramme pour analyser les modèles de fréquence uniques de la voix d'un locuteur afin de les identifier.

Années 1970 : Plusieurs systèmes d'identification de d'individus par empreinte vocale commerciaux ont été développés, y compris le système VeriVox de Raytheon et le système VoiceID de Bolt Beranek et Newman (BBN).

Années 1990 : Les chercheurs ont commencé à explorer l'utilisation de réseaux neuronaux pour l'identification des individus par voix, ce qui a permis une identification plus précise et plus fiable.

Années 2000 : L'identification des individus par voix est devenue plus largement utilisée dans les applications d'application de la loi et de sécurité, et les chercheurs ont commencé à explorer l'utilisation d'algorithmes d'apprentissage automatique pour identifier automatiquement les locuteurs sans avoir besoin d'une base de données préexistante de locuteurs connus.

Aujourd'hui, l'identification des individus par voix continue d'être un domaine de recherche actif, avec des applications dans divers domaines, notamment l'application de la loi, la sécurité et la criminalistique.

1.3 Types de reconnaissances vocales

Les principes de la biométrie vocale sont généralement séparés en 2 types :

- Dépendant du texte :

Ce type de reconnaissance du locuteur exige que l'utilisateur dise exactement l'expression enregistré. Cela manque de flexibilité étant donné que l'utilisateur doit se souvenir de ladite expression, mais offre une précision et une rapidité intéressantes.

- Indépendant du texte :

Ce processus de vérification n'a pas la contrainte de contenu. L'utilisateur peut parler librement. Toutefois, l'entraînement et les tests d'énonciation prendront plus de temps pour atteindre la performance attendue [2].

1.4 Fonctionnement

La Biométrie vocale utilise des modèles de voix pour générer une identification unique pour chaque individu, en travaillant avec plus de 100 caractéristiques physiques et comportementales. Il comprend la vitesse de la parole, l'accent, la prononciation, l'emphase, en plus des facteurs physiques de vos voies vocales, de votre bouche et de vos voies nasales [6].

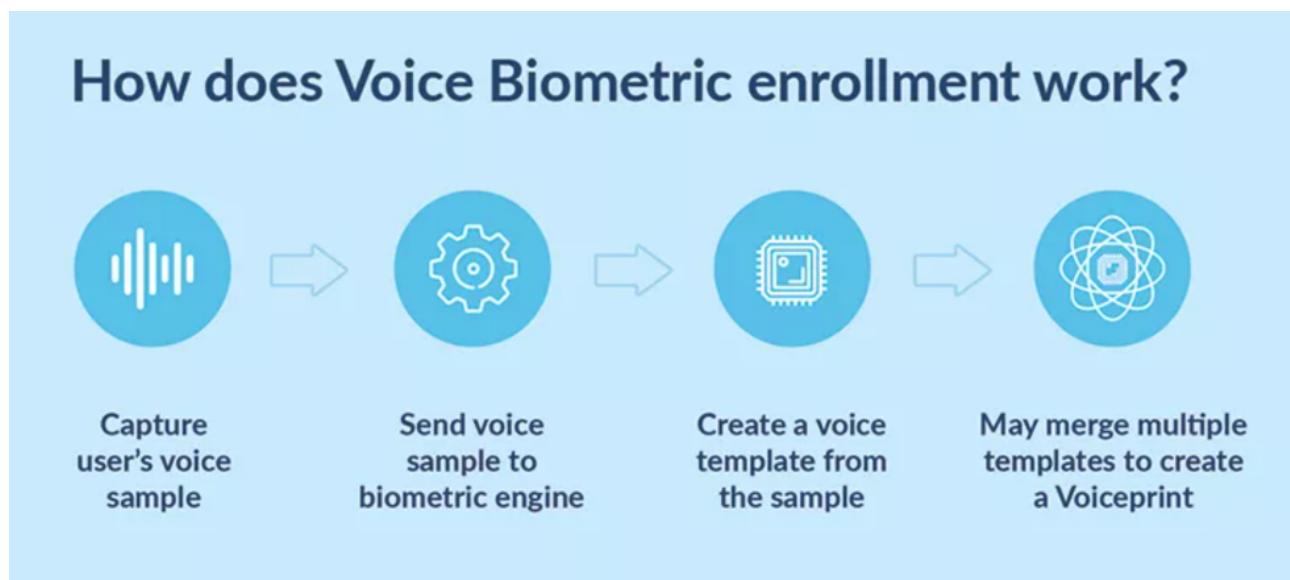


FIGURE 1.2 : Ajout d'un individu à une IA de reconnaissance par empreinte vocale.

Selon la méthode d'authentification (dépendante du texte ou indépendante du texte), une IA de biométrie vocale collecte le modèle vocal d'un utilisateur [12]. Cependant, il n'authentifie pas ce que l'utilisateur dit. Il vérifie seulement qui parle. Il extrait les caractéristiques qui distinguent le discours d'une personne des autres. Le résultat est une empreinte vocale ou un modèle vocal, tel qu'on peut avoir une empreinte digitale. Cela signifie-t-il qu'une personne avec une voix similaire peut contourner le système ? La voix d'une personne est extrêmement difficile à forger à des fins de

comparaison biométrique en raison de son caractère unique inhérent comme le dialecte, le style de parole etc.

Cela signifie simplement que même si une usurpation d'identité vocale est ressemblant pour une oreille humaine, une analyse détaillée de l'empreinte vocale effectuée à l'aide d'algorithmes informatiques peut aider à la distinguer de l'échantillon. Cela grâce a Plus de 100 caractéristiques physiques du corps, chacune avec une taille et une forme unique, qui contribuent à la façon dont une personne parle. La biométrie vocale repose sur des caractéristiques vocales fortement corrélées aux qualités physiologiques de la façon dont une personne crée la parole. Maintenant que nous avons étudié l'authentification vocale et son fonctionnement, jetons un coup d'œil à certaines de ses applications réelles.

1.5 cas d'utilisation de la biométrie vocale

Des centres d'appel (call center), et des applications mobiles aux applications de messagerie et aux appareils domotiques intelligents, la biométrie vocale peut fonctionner dans divers cas d'utilisation [15].

Voici un aperçu détaillé de ceux-ci :

- les applications mobiles :

Le principal cas d'utilisation de l'authentification vocale auprès des consommateurs est l'authentification mobile mains libres. Tout ce que vous avez à faire est de fournir une commande vocale pour vous connecter ou autoriser les achats, éliminant ainsi le besoin de mémoriser les identifiants et les mots de passe. Ceci est idéal pour les téléphones portables ou d'autres paramètres où la reconnaissance faciale et d'autres formes d'authentification biométrique peuvent être gênantes. De plus, l'authentification vocale peut également être utile pour les solutions d'assistant virtuel telles que Google Home, Alexa d'Amazon et Siri. On peut l'utiliser pour passer des commandes et effectuer d'autres fonctions qui nécessitent une certaine authentification.

- les Centres d'appels et systèmes IVR(Interactive Voice Response) :

Les méthodes de sécurité obsolètes comme les mots de passe traditionnels ou les questions ne sont plus suffisamment sécurisées. Les systèmes de biométrie vocale offrent une résistance contre l'imitation de la voix grâce à des algorithmes intrinsèques utilisés pour l'analyse biométrique et offrent une liste de blocage. Cela rend la technologie particulièrement utile dans le secteur de l'assistance téléphonique. On peut également utiliser la reconnaissance vocale comme authentificateur lors des appels d'assistance client. Les appelants peuvent trouver cela plus pratique et plus sûr que de partager des données personnelles telles que leur numéro de licence ou de carte de crédit à des fins de vérification d'identité.

- les Applications web :

On peut ajouter des systèmes de vérification vocale à des pages Web ou à des applications dans les secteurs de la banque et du commerce électronique. L'authentification vocale dans les applications Web peut être utile pour l'identification à distance des utilisateurs. De plus, l'inscription passive ou l'authentification indépendante du texte facilite l'intégration de nouveaux utilisateurs pour votre service sans aucune inscription. Les clients sont automatiquement vérifiés en temps réel lorsqu'ils interagissent avec un agent IVR ou un centre de contact.

- Les objetc connectés :

Les applications IoT offrent des moyens nouveaux et innovants de communication et d'interaction entre les humains et les machines. Une mise en œuvre appropriée de l'authentification vocale peut offrir une expérience utilisateur plus flexible que les méthodes traditionnelles telles que les écrans tactiles. Et comme l'authentification vocale peut fournir une couche de sécurité supplémentaire, vous pouvez facilement accéder à votre appareil domotique IoT sans aucun souci. De toute évidence, l'authentification biométrique vocale semble rendre les choses beaucoup plus faciles. Cependant, avant de décider d'utiliser l'authentification vocale pour votre entreprise, examinons ses avantages et ses défis [14].

1.6 Les avantages de la biométrie vocale

La Biométrie vocale peut être utilisé dans le cadre d'un processus d'authentification à deux facteurs pour augmenter la sécurité et les 3 principaux avantages sont :

- Améliorer l'expérience client avec une authentification rapide et fluide
- Améliorer la sécurité et minimiser les violations dues aux mots de passe compromis, au phishing, etc...
- Identifier instantanément les utilisateurs et personnalisez l'interaction. Grâce à la technologie de Biométrie vocale, les appelants n'ont plus besoin de fournir de mots de passe ou de codes PIN ni de répondre à des questions de sécurité pour vérifier leur identité. Cela rend la biométrie vocale idéale pour les déploiements multicanaux. Une fois qu'un utilisateur est inscrit, on peut utiliser son empreinte vocale sur tous les canaux d'assistance d'une organisation [2].

1.7 Les inconvenients de la biométrie vocale

- Authentification via des contrefaçons audio :

Les récents progrès de la technologie ont permis aux gens de créer des contrefaçons profondes. Il s'agit de fausses voix produites synthétiquement d'une personne, identiques à leur voix d'origine. Les contrefaçons profondes sont de plus en plus courantes et peuvent faire croire à un programme d'IA en son authenticité. Alors, comment empêcher les utilisateurs non autorisés d'entrer dans la base de données? Vous pouvez créer une liste d'autorisation d'empreintes vocales et les stocker dans un répertoire actif. Au cours de ce processus, le système de reconnaissance vocale inscrit l'utilisateur dans une liste de membres autorisés. Ainsi, chaque fois qu'un utilisateur essaie d'accéder au système, son empreinte vocale est comparée à la fois à la liste d'autorisation et à une liste de blocage des empreintes vocales des fraudeurs . Et pendant que l'authentification est en cours, la détection passive des fraudes peut envoyer des alertes si l'empreinte vocale correspond à la base de données de la liste noire.

- Manque de précision :

Le bruit de fond ou son parasite est l'un des principaux facteurs qui affectent la reconnaissance

automatique de la parole. Cela peut avoir un impact sur la qualité du modèle de voix de l'orateur et, à son tour, diminuer le niveau de précision du processus d'authentification. Un système d'authentification vocale peut ne pas être en mesure de faire la différence entre le discours de l'utilisateur, les autres personnes qui parlent et le bruit ambiant, ce qui entraîne des confusions et des erreurs [1].

Cela signifie qu'il peut être difficile d'utiliser l'authentification vocale dans des environnements bruyants tels que des bureaux très fréquentés ou des espaces publics.

Pour une authentification transparente, on peut utiliser des microphones proches ou des casques antibruit qui permettent au logiciel de se concentrer sur le discours de l'utilisateur.

L'identification du locuteur est le processus de détermination de l'identité d'une personne en fonction de sa voix. Il s'agit d'une tâche essentielle dans de nombreuses applications, y compris l'application de la loi, la sécurité et la reconnaissance vocale.

L'état de l'art en matière d'identification du locuteur implique l'utilisation de techniques d'apprentissage en profondeur, qui ont considérablement amélioré la précision et la fiabilité du processus. Les modèles d'apprentissage en profondeur les plus populaires pour l'identification du locuteur comprennent les réseaux de neurones convolutifs (CNN), les réseaux de neurones récurrents (RNN) et leurs variantes.

L'un des principaux défis de l'identification du locuteur est de gérer la variabilité des signaux vocaux. Des facteurs tels que le bruit de fond, l'accent et l'état émotionnel peuvent affecter de manière significative la qualité des signaux vocaux, ce qui rend difficile l'identification précise du locuteur. Pour relever ce défi, les chercheurs ont développé diverses techniques d'extraction de caractéristiques, telles que les coefficients cepstraux de fréquence Mel (MFCC), qui sont couramment utilisés dans le traitement de la parole.

Ces dernières années, les chercheurs ont également exploré l'utilisation d'autres modalités biométriques, telles que la reconnaissance faciale et la reconnaissance de l'iris, en conjonction avec l'identification du locuteur pour améliorer la précision globale du processus. Ces approches multimodales tirent parti des atouts de différentes modalités biométriques pour améliorer le processus d'identification.

Dans l'ensemble, l'état de l'art en matière d'identification des locuteurs continue d'évoluer rapidement, les recherches en cours se concentrant sur le développement de modèles plus robustes et plus précis pour identifier les locuteurs dans des environnements difficiles.

1.8 Apprentissage profond

1.8.1 Notion d'IA

L'intelligence artificielle est née dans les années 1950, lorsqu'une poignée de pionniers du domaine naissant de l'informatique ont commencé à se demander si les ordinateurs pouvaient être amenés à "penser" - une question dont nous explorons encore les ramifications aujourd'hui. Une définition concise du domaine serait la suivante : l'effort d'automatisation des tâches intellectuelles normalement exécutées par les humains. En tant que tel, l'IA est un domaine général qui englobe l'apprentissage automatique et l'apprentissage en profondeur, mais qui comprend également de nombreuses autres approches qui n'impliquent aucun apprentissage. Les premiers programmes d'échecs, par exemple, n'impliquaient que des règles codées en dur conçues par des programmeurs et n'étaient

pas considérées comme de l'apprentissage automatique.

Pendant assez longtemps, de nombreux experts ont cru que l'intelligence artificielle au niveau humain pouvait être obtenue en demandant aux programmeurs de créer à la main un ensemble suffisamment large de règles explicites pour manipuler les connaissances. Cette approche est connue sous le nom d'IA symbolique et a été le paradigme dominant de l'IA des années 1950 à la fin des années 1980. Il a atteint son apogée pendant le boom des systèmes experts des années 1980.

Bien que l'IA symbolique se soit avérée adaptée pour résoudre des problèmes logiques bien définis, comme jouer aux échecs, il s'est avéré impossible de trouver des règles explicites pour résoudre des problèmes plus complexes et flous, tels que la classification d'images, la reconnaissance vocale et le langage. Une nouvelle approche est apparue pour prendre la place symbolique de l'IA : l'apprentissage automatique.

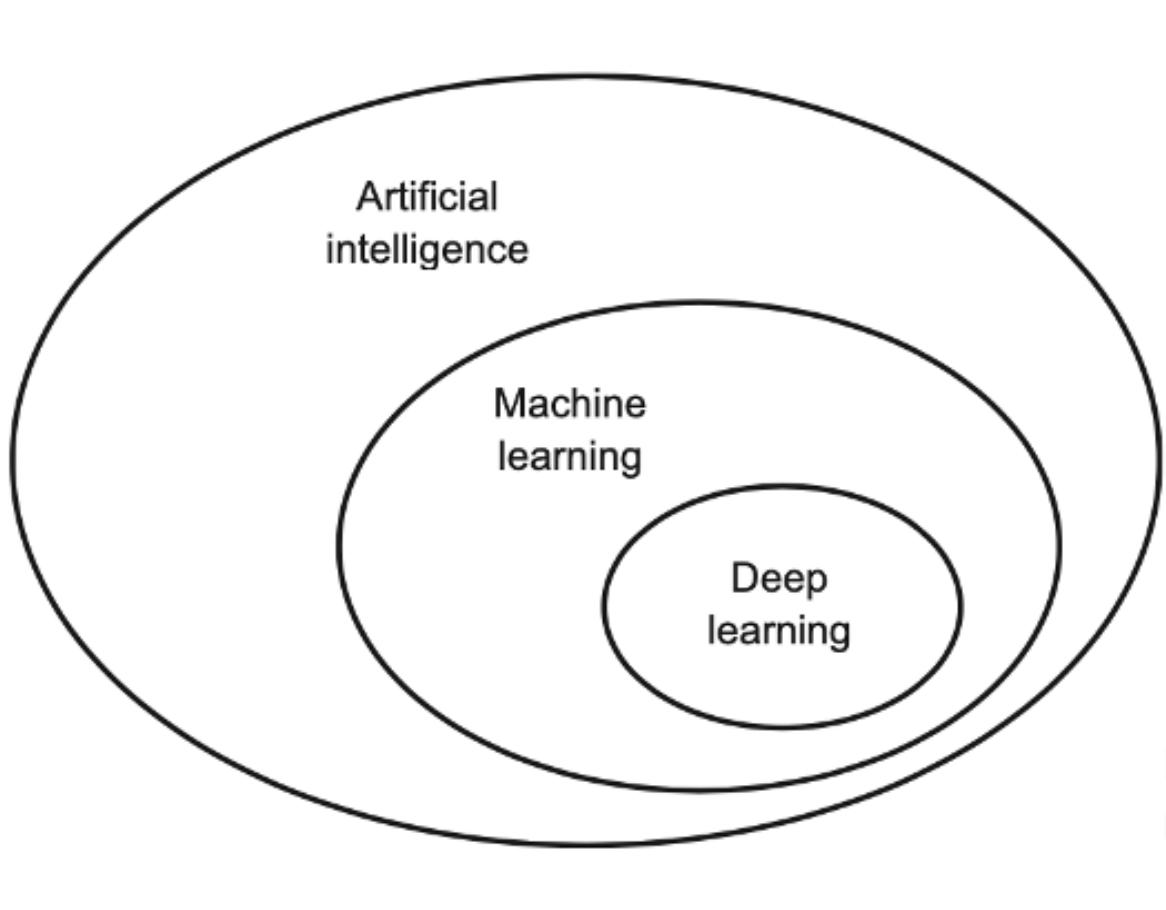


FIGURE 1.3 : Intelligence Artificielle, Machin learning et Deep learning.

1.8.2 Notion d'apprentissage automatique

Dans l'Angleterre victorienne, Lady Ada Lovelace était une amie et collaboratrice de Charles Babbage, l'inventeur de la machine analytique : le premier ordinateur mécanique polyvalent connu. Bien que visionnaire et très en avance sur son temps, l'Analytical Engine n'était pas conçu comme un ordinateur à usage général lorsqu'il a été conçu dans les années 1830 et 1840, car le concept de calcul à usage général n'était pas encore inventé. Il s'agissait simplement d'un moyen d'utiliser des opéra-

tions mécaniques pour automatiser certains calculs du domaine de l'analyse mathématique, d'où le nom de moteur analytique.

En 1843, Ada Lovelace a fait remarquer à propos de l'invention : « La machine analytique n'a aucune prétention à créer quoi que ce soit. Il peut faire tout ce que nous savons comment lui ordonner de fonctionner... Son rôle est de nous aider à mettre à disposition ce que nous connaissons déjà. Cette remarque a ensuite été citée par le pionnier de l'IA, Alan Turing, comme "l'objection de Lady Lovelace" dans son article historique de 1950 "Computing Machinery and Intelligence",¹ qui a présenté le test de Turing ainsi que les concepts clés qui allaient façonner l'IA. Turing citait Ada Lovelace alors qu'il se demandait si les ordinateurs à usage général pouvaient être capables d'apprentissage et d'originalité, et il en est venu à la conclusion qu'ils le pouvaient.

L'apprentissage automatique découle de cette question : un ordinateur pourrait-il aller au-delà de "ce que nous savons lui ordonner d'exécuter" et apprendre par lui-même comment effectuer une tâche spécifiée ? Un ordinateur pourrait-il nous surprendre ? Plutôt que des programmeurs élaborant des règles de traitement de données à la main, un ordinateur pourrait-il automatiquement apprendre ces règles en examinant les données ? Cette question ouvre la porte à un nouveau paradigme de programmation.

En programmation classique, le paradigme de l'IA symbolique, les humains saisissent des règles (un programme) et des données à traiter selon ces règles, et en sortent des réponses (voir figure 1.2).

Avec l'apprentissage automatique, les humains saisissent les données ainsi que les réponses attendues des données, et sortent les règles. Ces règles peuvent ensuite être appliquées à de nouvelles données pour produire des réponses originales.

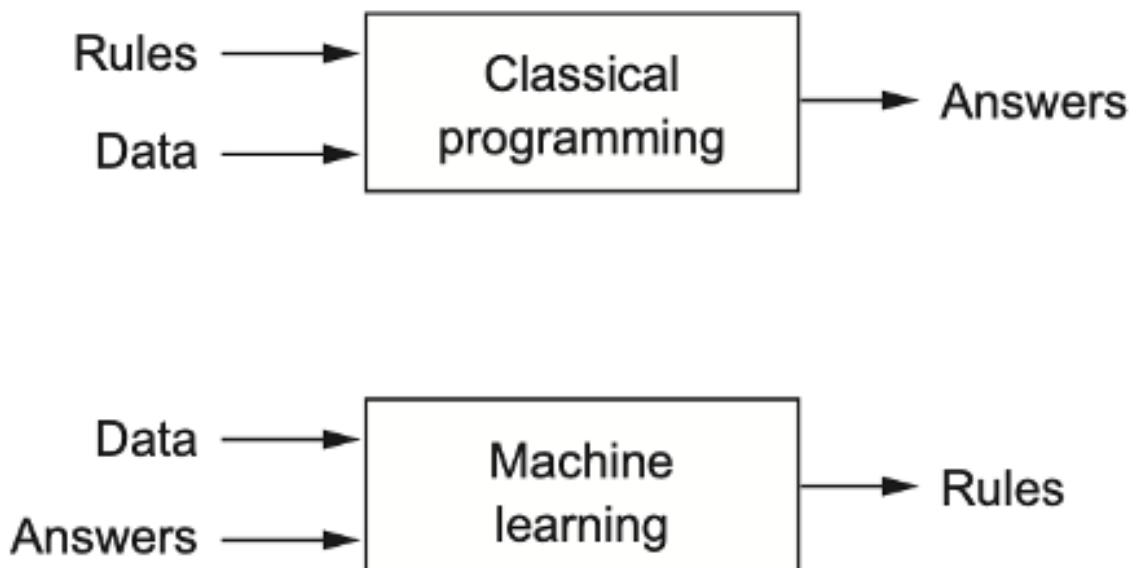


FIGURE 1.4 : Machin learning Nouveau paradigme de programmation.

Un système d'apprentissage automatique est entraîné plutôt qu'explicitement programmé. Il est présenté avec de nombreux exemples pertinents pour une tâche, et il trouve une structure statistique dans ces exemples qui permet finalement au système de proposer des règles pour automatiser la tâche.

Par exemple, si vous souhaitez automatiser la tâche d'étiquetage de vos photos de vacances, vous pouvez présenter un système d'apprentissage automatique avec de nombreux exemples d'images déjà étiquetées par des humains, et le système apprendrait des règles statistiques pour associer des images spécifiques à des étiquettes spécifiques.

Un système d'apprentissage automatique est entraîné plutôt qu'explicitement programmé. Il est présenté avec de nombreux exemples pertinents pour une tâche, et il trouve une structure statistique dans ces exemples qui permet finalement au système de proposer des règles pour automatiser la tâche. Par exemple, si vous souhaitez automatiser la tâche d'étiquetage de vos photos de vacances, vous pouvez présenter un système d'apprentissage automatique avec de nombreux exemples d'images déjà étiquetées par des humains, et le système apprendrait des règles statistiques pour associer des images spécifiques à des étiquettes spécifiques.

Bien que l'apprentissage automatique n'ait commencé à prospérer que dans les années 1990, il est rapidement devenu le sous-domaine de l'IA le plus populaire et le plus réussi, une tendance stimulée par la disponibilité d'un matériel plus rapide et d'ensembles de données plus volumineux.

L'apprentissage automatique est étroitement lié aux statistiques mathématiques, mais il diffère des statistiques de plusieurs manières importantes. Contrairement aux statistiques, l'apprentissage automatique a tendance à traiter de grands ensembles de données complexes (tels qu'un ensemble de données de millions d'images, chacune composée de dizaines de milliers de pixels) pour lesquels une analyse statistique classique telle que l'analyse bayésienne ne serait pas pratique. En conséquence, l'apprentissage automatique, et en particulier l'apprentissage en profondeur, présente relativement peu de théorie mathématique peut-être trop peu - et est orienté vers l'ingénierie. C'est une discipline pratique dans laquelle les idées sont prouvées empiriquement plus souvent que théoriquement.

Pour resumer, Le machine learning (apprentissage automatique) est une branche de l'intelligence artificielle qui permet à un système informatique d'apprendre à partir de données et d'expériences passées, sans être explicitement programmé. Il s'agit d'une technique d'analyse de données qui utilise des algorithmes pour identifier des modèles et des tendances dans les données, puis les utilise pour faire des prédictions ou des recommandations. Le machine learning est utilisé dans de nombreux domaines, tels que la reconnaissance vocale, la reconnaissance faciale, la détection de fraudes, la recommandation de produits, etc.

1.8.3 Notions de Deep learning

1.8.3.1 définition

Le deep learning, ou apprentissage profond en français, est une branche de l'intelligence artificielle qui utilise des réseaux de neurones artificiels pour apprendre à partir de données complexes

et massives. Il s'agit d'une technique d'apprentissage automatique qui permet aux machines d'apprendre à partir d'un grand nombre de données en les nourrissant à travers des couches de neurones interconnectés [3].

Le *deep* dans *deep learning* ne fait référence à aucune sorte de compréhension plus profonde obtenue par l'approche; il représente plutôt cette idée de couches successives de représentations. Le nombre de couches qui contribuent à un modèle de données est appelé la profondeur du modèle.

Dans l'apprentissage profond, ces représentations en couches sont (presque toujours) apprises via des modèles appelés réseaux de neurones, structurés en couches littérales empilées les unes sur les autres. Le terme réseau de neurones fait référence à la neurobiologie, mais bien que certains des concepts centraux de l'apprentissage en profondeur aient été développés en partie en s'inspirant de notre compréhension du cerveau, les modèles d'apprentissage en profondeur ne sont pas des modèles du cerveau. Il n'y a aucune preuve que le cerveau implémente quelque chose comme les mécanismes d'apprentissage utilisés dans les modèles modernes d'apprentissage en profondeur.

Le deep learning est utilisé dans de nombreux domaines tels que la reconnaissance d'image, la reconnaissance vocale, la traduction automatique, la prédiction de résultats, etc. Cette technologie est devenue très populaire ces dernières années grâce à ses performances exceptionnelles dans la reconnaissance de formes et dans l'analyse de données complexes.

Examinons comment un réseau de plusieurs couches de profondeur (voir la figure 1.5) transforme l'image d'un chiffre afin de reconnaître de quel chiffre il s'agit.

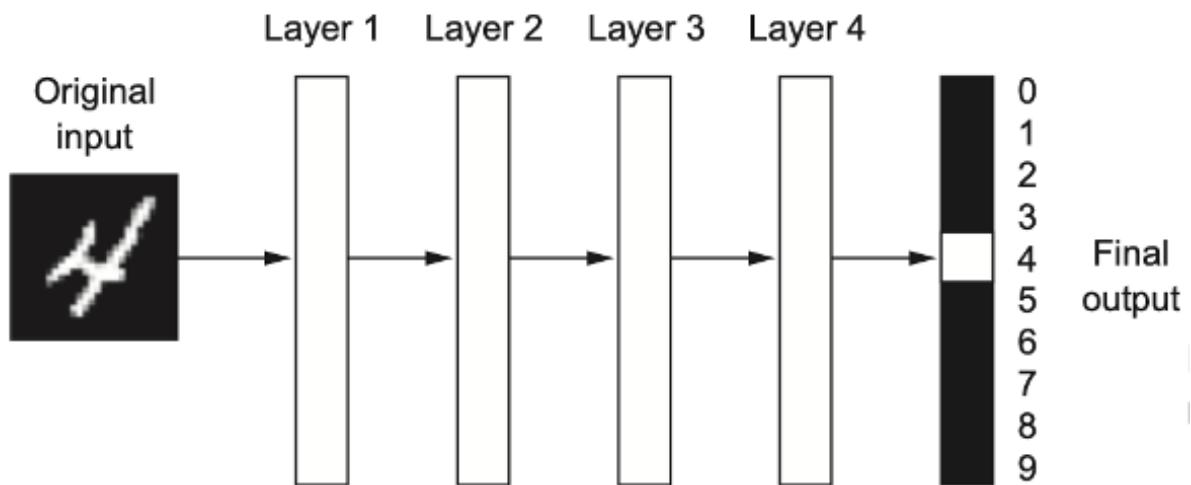


FIGURE 1.5 : Un Reseau de neurone pour la reconnaissance des caractères.

Comme on peut le voir sur la figure 1.6, le réseau transforme l'image numérique en représentations de plus en plus différentes de l'image originale et de plus en plus informatives sur le résultat final. On peut considérer un réseau profond comme une opération de distillation d'informations à plusieurs étapes, où les informations passent par des filtres successifs et sortent de plus en plus purifiées (c'est-à-dire utiles pour certaines tâches).

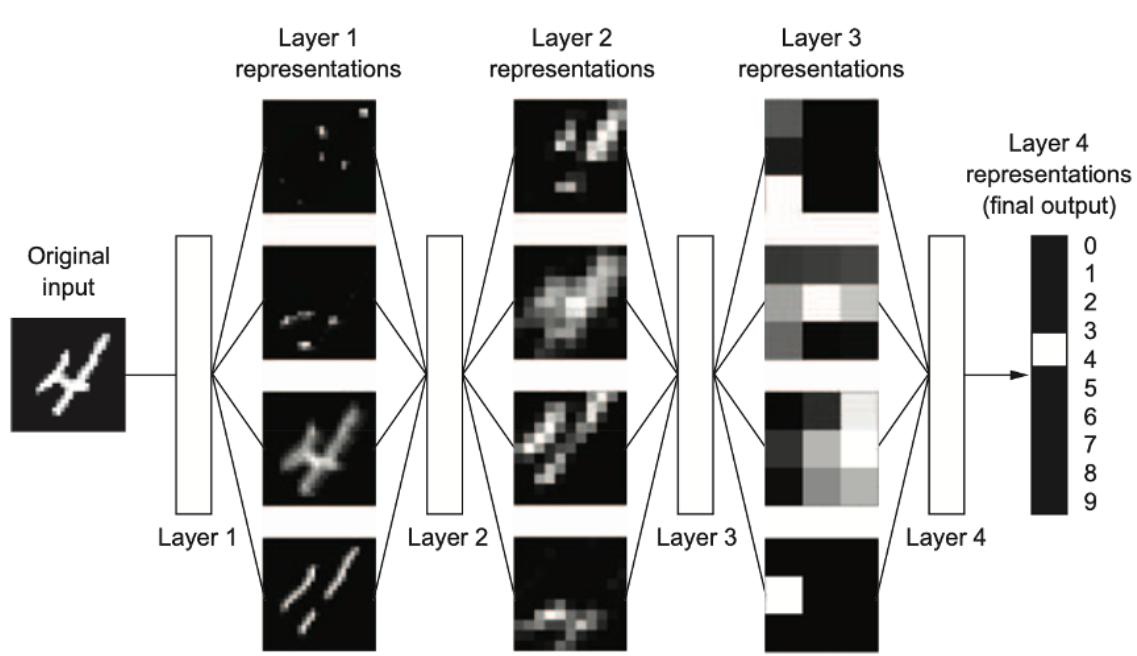


FIGURE 1.6 : modèle de classification de caractères manuscrits.

1.8.3.2 Fonctionnement

La spécification de ce qu'une couche fait à ses données d'entrée est stockée dans les poids de la couche, qui sont essentiellement un ensemble de nombres. En termes techniques, on dirait que la transformation mise en œuvre par une couche est paramétrée par ses poids (voir figure 1.7). Dans ce contexte, l'apprentissage signifie trouver un ensemble de valeurs pour les pondérations de toutes les couches d'un réseau, de sorte que le réseau mappe correctement les exemples d'entrées sur leurs cibles associées. Mais voici le problème : un réseau de neurones profonds peut contenir des dizaines de millions de paramètres. Trouver la valeur correcte pour chacun d'eux peut sembler une tâche ardue, d'autant plus que la modification de la valeur d'un paramètre affectera le comportement de tous les autres !

Pour contrôler la sortie d'un réseau de neurones, on doit être en mesure de mesurer à quelle distance cette sortie (la prédition, \hat{Y}) [3] est de ce que vous attendiez (le résultat attendu, Y). C'est le travail de la fonction de perte du réseau, également appelée fonction objectif. La fonction de perte prend les prédictions du réseau et la véritable cible (ce que vous voulez que le réseau produise) et calcule un score de distance, capturant la performance du réseau sur cet exemple spécifique (voir figure 1.7).

L'astuce fondamentale dans l'apprentissage profond consiste à utiliser ce score comme un signal de rétroaction pour ajuster un peu la valeur des poids, dans une direction qui abaissera le score de perte pour l'exemple actuel (voir figure 1.9). Cet ajustement est le travail de l'optimiseur, qui implémente ce qu'on appelle l'algorithme de rétropropagation : l'algorithme central de l'apprentissage en

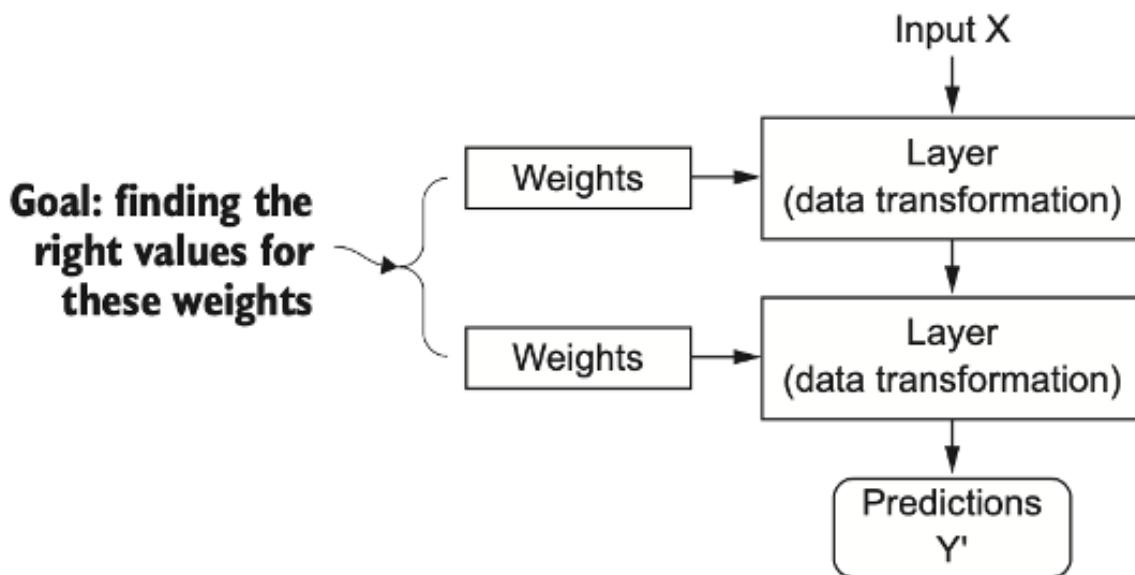


FIGURE 1.7 : Paramétrisation d'un réseau de neurones

profondeur.

Initialement, les poids du réseau reçoivent des valeurs aléatoires, de sorte que le réseau implémente simplement une série de transformations aléatoires. Naturellement, son rendement est loin de ce qu'il devrait être idéalement, et le score de perte est donc très élevé. Mais avec chaque exemple traité par le réseau, les pondérations sont légèrement ajustées dans la bonne direction et le score de perte diminue. Il s'agit de la boucle d'apprentissage qui, répétée un nombre suffisant de fois (généralement des dizaines d'itérations sur des milliers d'exemples), donne des valeurs de poids qui minimisent la fonction de perte. Un réseau avec une perte minimale est un réseau dont les sorties sont aussi proches que possible des cibles : un réseau entraîné.

En résumé, le deep learning implique la collecte et la préparation des données, [3] la construction et l'entraînement du réseau de neurones, l'évaluation du modèle et enfin la prédiction. Le processus de deep learning est souvent itératif, et il peut être nécessaire d'ajuster les paramètres du modèle ou de modifier l'architecture du réseau de neurones pour améliorer les performances du modèle.

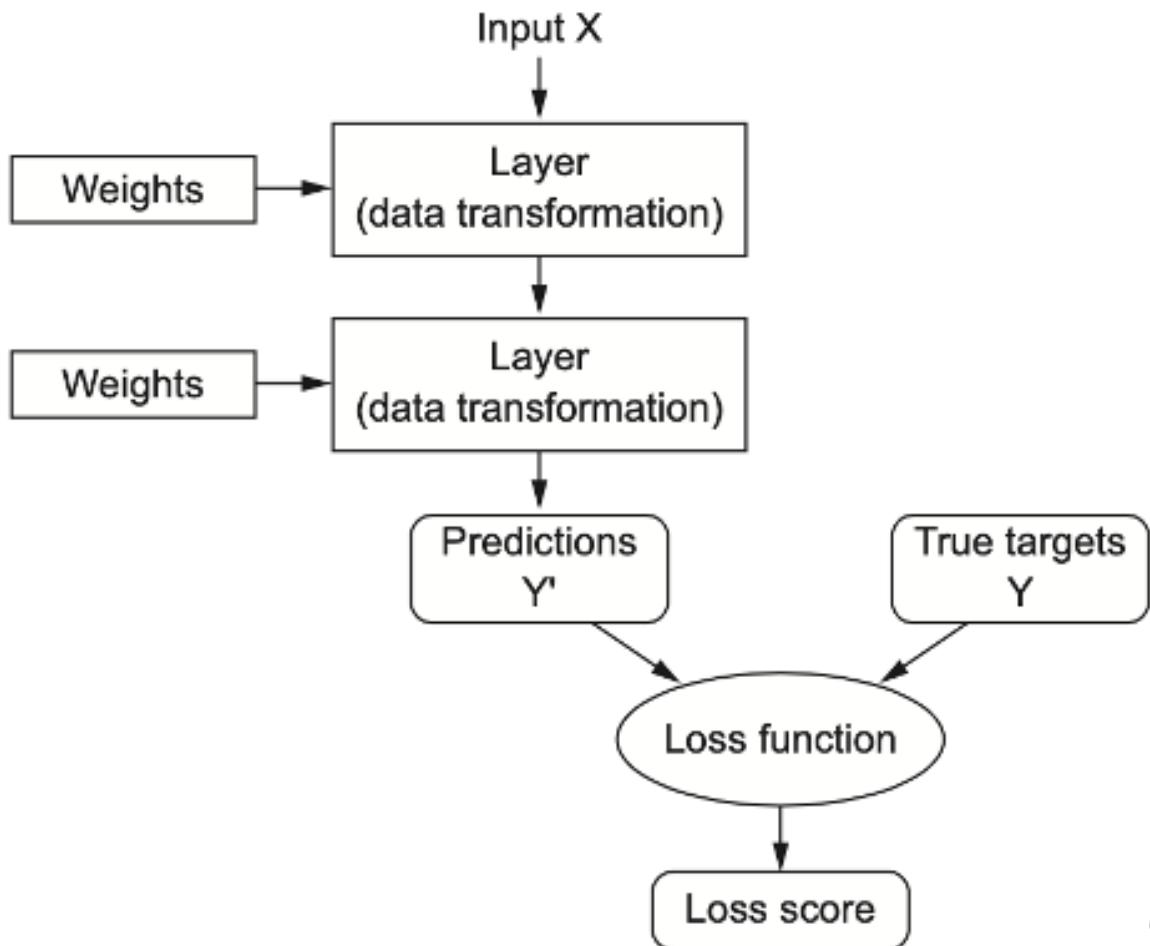


FIGURE 1.8 : Détermination de la fonction Loss

1.9 deep Learning pour la reconnaissance vocale

Dans le domaine de la reconnaissance de locuteur, le Deep Learning est utilisé pour améliorer la précision des systèmes de reconnaissance de la parole en analysant les caractéristiques du signal vocal. Plusieurs méthodes de Deep Learning sont utilisées pour la reconnaissance de locuteur :

- **Réseaux de neurones convolutionnels (CNN)** : Les CNN sont largement utilisés pour l'analyse d'images, mais ils ont également été appliqués à la reconnaissance de locuteur en utilisant des images spectrographiques du signal vocal. Les CNN peuvent extraire des caractéristiques significatives du signal vocal qui peuvent être utilisées pour identifier le locuteur.
- **Réseaux de neurones récurrents (RNN)** : Les RNN sont utilisés pour modéliser des séquences de données, telles que des séquences de phonèmes dans la parole. Les RNN peuvent être utilisés pour identifier le locuteur en analysant les modèles de parole spécifiques à chaque locuteur.
- **Réseaux de neurones profonds (DNN)** : Les DNN sont utilisés pour l'apprentissage en pro-

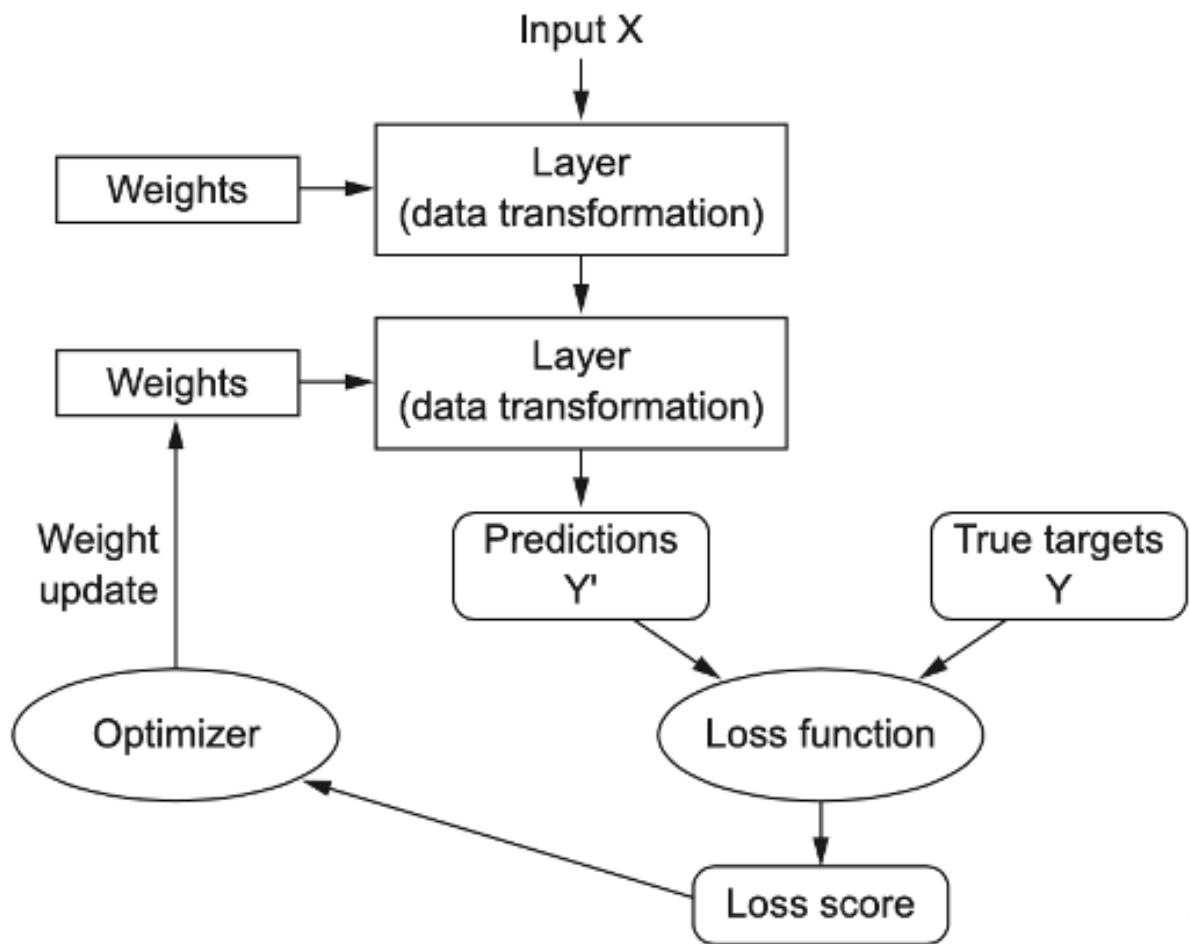


FIGURE 1.9 : Le Loss Score utilisé pour optimiser le Weighth à l'entrée

fondeur en utilisant des couches de neurones artificiels. Les DNN ont été utilisés pour la reconnaissance de locuteur en analysant les caractéristiques du signal vocal telles que la fréquence, l'intensité, la durée et les modèles de parole.

- **Réseaux adversaires génératifs (GAN)** : Les GAN sont utilisés pour générer des données synthétiques en imitant les données réelles. Les GAN ont été utilisés pour la reconnaissance de locuteur en générant des images spectrographiques du signal vocal qui peuvent être utilisées pour entraîner des modèles de reconnaissance de locuteur.

Les TDNN (Time-Delay Neural Networks) représentent une architecture de réseaux de neurones artificiels, qui utilisent des couches de neurones à délai temporel pour traiter des données séquentielles ou temporelles. Ces réseaux sont capables de modéliser des dépendances temporelles entre des événements, en prenant en compte l'ordre et le timing des entrées. Les TDNN sont souvent utilisés pour la reconnaissance de la parole, la classification de séries temporelles, la prédiction de séquences, etc.

En somme, ces techniques et méthodes de Deep Learning permettent d'améliorer la précision de la reconnaissance de locuteur en analysant les caractéristiques du signal vocal de manière plus fine et précise [16].

1.9.1 que représentent les x-vector dans la reconnaissance vocale ?

Les x-vector sont une représentation numérique de haut niveau d'un segment de signal vocal, qui est utilisée dans la reconnaissance vocale pour l'identification des individus. Les x-vector sont calculés à partir de caractéristiques acoustiques telles que les coefficients cepstraux de fréquence (MFCC) extraits d'un segment de parole [7].

Les x-vector sont généralement extraits en utilisant des réseaux de neurones profonds appelés réseaux de neurones convolutionnels (CNN) et des réseaux de neurones récurrents (RNN), qui sont entraînés à prédire l'identité du locuteur à partir des caractéristiques acoustiques. Le réseau de neurones est ensuite utilisé pour extraire les x-vector à partir de nouvelles données de parole, qui sont ensuite utilisées pour identifier l'individus.

Les x-vector ont plusieurs avantages par rapport aux méthodes traditionnelles de reconnaissance vocale basées sur les traits acoustiques. Ils peuvent être utilisés pour représenter des segments de parole de longueur variable, sont robustes aux variations de canal et peuvent être utilisés pour l'identification des individus dans des scénarios multi-locuteurs. Les x-vector ont été utilisés avec succès dans de nombreux systèmes de reconnaissance vocale, notamment dans les systèmes de sécurité et de surveillance.

1.9.2 que représentent les DNN

Les DNN (Deep Neural Networks) représentent une classe de modèles de réseaux de neurones artificiels qui sont capables de traiter des données complexes et d'effectuer des tâches de classification, de reconnaissance de motifs, de traitement de langage naturel, etc. Les DNN utilisent plusieurs couches de neurones pour extraire des caractéristiques de plus en plus abstraites à partir des données d'entrée, ce qui leur permet d'apprendre des représentations profondes des données et de réaliser des tâches de manière plus précise et efficace que les modèles de réseaux de neurones plus simples. Les DNN ont trouvé des applications dans de nombreux domaines, notamment la vision par ordinateur, le traitement du langage naturel, la reconnaissance de la parole, la biologie, la médecine, la finance, etc.

1.9.3 Les ECAPA-TDNN

Les ECAPA-TDNN (Ensemble Classifiers using Adaptive Partitioning Algorithm - Time-Delay Neural Network) sont un type de modèle de classification utilisé en intelligence artificielle et en reconnaissance de formes. Ils combinent les techniques d'ensemble learning (utilisation de multiples classificateurs pour améliorer la précision de la classification) avec les réseaux de neurones à retards temporels (TDNN) pour résoudre des problèmes de classification complexes. Les ECAPA-TDNN sont particulièrement adaptés pour la reconnaissance de formes dans des signaux temporels tels que les signaux audio, les séries chronologiques, les images vidéo, etc. Ils sont largement utilisés dans des domaines tels que la reconnaissance de la parole, la reconnaissance de gestes et la vision par ordinateur [10].

1.9.4 Les ECAPA-TDNN

Les PLDA (Probabilistic Linear Discriminant Analysis) sont des modèles d'apprentissage automatique utilisés dans la classification de données. Ils sont principalement utilisés pour l'analyse discriminante linéaire, qui est une méthode statistique qui permet de déterminer la catégorie à laquelle une observation appartient en se basant sur un ensemble de variables prédéfinies. Les PLDA sont souvent utilisés dans le domaine de la reconnaissance vocale pour l'identification des locuteurs, mais peuvent également être appliqués à d'autres types de données. Ils sont considérés comme un outil efficace pour la réduction de dimensionnalité et pour l'analyse des données à grande échelle.

1.10 Les bonnes pratiques en matière de reconnaissance vocale

Comme pour toute méthode de sécurité, il existe certaines meilleures pratiques de base que les organisations doivent suivre pour s'assurer que les informations de leurs clients ne sont pas en danger [8].

- Établir une méthode d'authentification principale avant toute autre méthode ; cela signifie choisir une forme de validation, qu'il s'agisse d'un code PIN, d'un mot de passe ou d'une biométrie, avant d'accepter d'utiliser plus d'authentifications.
- Obtenir le consentement explicite de l'utilisation prévue des utilisateurs pour la mesure de sécurité ; par exemple, accord d'utilisation de la reconnaissance faciale pour permettre les transactions de paiement.
- Exiger la méthode d'authentification principale toutes les 72 heures.
- Utiliser un pipeline entièrement sécurisé pour toutes les données biométriques et leur manipulation.
- Conserver toutes les données biométriques dans un environnement sécurisé et isolé pour empêcher leur acquisition par des fraudeurs.

L'utilisation des étapes ci-dessus fournit une excellente base à toute entreprise pour assurer la sécurité des informations de ses clients. Il fournit également une base pour commencer à ajouter plus de niveaux d'authentification afin de créer un système de vérification encore plus puissant.

1.11 À quel point la reconnaissance vocale est-elle sécurisée ?

La commodité qu'offre la reconnaissance vocale se fait-elle au détriment de la sécurité ? Il y a eu de nombreuses histoires dans les nouvelles sur les pirates capables d'infiltrer les maisons et les entreprises au détriment de leurs propriétaires en utilisant des problèmes de sécurité de reconnaissance vocale. Comme pour toute mesure de sécurité, croire qu'une seule mesure à elle seule est infaillible est imprudent [13].

Il reste encore du chemin à parcourir pour parvenir à une identification vocale absolument sécurisée. L'état actuel de la biométrie vocale est vulnérable. Des recherches ont montré que des échantillons de voix provenant de quelque chose comme une vidéo YouTube peuvent être acceptés comme modèles de parole approuvés. Les pirates ont pu enterrer des commandes malveillantes dans du bruit blanc pour contrôler les appareils à commande vocale.

Selon Lior Atzi, directeur de la gestion des produits chez NICE, deux technologies majeures sont utilisées pour lutter contre la fraude dans l'authentification vocale : la détection de la vivacité et l'authentification continue.

La détection de la vivacité, comme son nom l'indique, garantit que l'empreinte digitale ou l'échantillon vocal utilisé est réel. Certaines tentatives pour contrecarrer cette mesure de sécurité sont des voix synthétiques et d'autres implémentations d'intelligence artificielle, mais elles n'ont pas encore suffisamment progressé pour provoquer une menace réelle.

L'authentification continue vérifie à plusieurs reprises l'identité d'un individu tout au long d'une session, plutôt qu'une seule fois. Cela aide à surmonter les problèmes potentiels tels que les appelants qui changent au milieu d'un appel téléphonique ou d'autres astuces qu'un mauvais acteur pourrait utiliser pour accéder au compte de quelqu'un.

Conclusion

L'authentification vocale rend les connexions aussi simples que parler. Cela peut sembler futuriste, mais c'est la sécurité qui protège les entreprises des violations quotidiennes.

La reconnaissance vocale est-elle sécurisée ? Oui, la reconnaissance vocale est sécurisée, surtout par rapport aux connexions classiques qui nécessitent un nom d'utilisateur et un mot de passe. Semblable à d'autres biométries, la reconnaissance vocale est plus sécurisée car une personne doit interagir avec un identifiant plutôt que de simplement saisir un code.

De toute évidence, l'authentification vocale offre de nombreux avantages, tels que des économies de coûts, une sécurité accrue et une expérience client améliorée.

Cependant, cela vient avec son propre ensemble de défis. Par exemple, parce que la technologie peut être sensible aux bruits et aux signaux parasites.

Elle reste tout de même un excellent choix pour accroître la sécurité des systèmes utilisés comme une vérification supplémentaire, car elle ajoute des couches de protection supplémentaires que les codes d'accès manuels pourraient ne pas fournir. L'authentification vocale est bénéfique à la fois pour les clients et pour l'entreprise car elle élimine la frustration associée aux processus de connexion fastidieux.

La commodité et la sécurité qu'elle offre en font une méthode de vérification dont l'utilisation ne fera qu'augmenter dans tous les secteurs.

Dans les prochaines décennies, l'authentification par empreinte vocale est susceptible de gagner encore plus en popularité en raison des améliorations de la précision, causées en grande partie par les progrès de l'IA. Cette croissance est également due aux attentes accrues des clients pour un accès rapide et facile à l'information. La biométrie vocale permet un accès rapide, fluide et hautement sécurisé pour une gamme de cas d'utilisation allant des centres d'appels aux applications mobiles et en ligne, en passant par les chatbots, les appareils IoT et l'accès physique.

Chapitre 2

MODELISATION ET LE DEVELOPPEMENT DE NOTRE MODEL D'IA

Introduction

Nous abordons, dans ce chapitre, la présentation de la méthodologie et des choix techniques. Nous voulons mettre en place un modèle de reconnaissance vocal. Ce model sera des réseaux de neurones profonde (DNN). Nous avons vu précédemment qu'il existait plusieurs méthodes de Deep Learning pour la reconnaissance par empreinte vocale. Dans ce chapitre, nous allons effectuer un choix et approfondir la méthode choisie. Nous allons commencer par la méthode de Deep Learning choisie et l'architecture du système que nous voulons mettre en place. Nous aborderons également les matériels ainsi que les technologies utilisées dans la mise en œuvre du modèle.

2.1 Architecture du système

La finalité de ce travail est de mettre en place un modèle d'IA pour identification et l'authentification par empreinte vocale. L'IA sera ensuite déployée comme un cloud REST API pour permettre à divers client web ou mobile de pouvoir y accéder. Ainsi l'architecture sera composée de trois parties :

1. le code de paramétrage de l'IA et la gestion du dataset
2. Le notebook Google Colab pour l'entraînement
3. la partie backend pour servir les APIs d'identification et d'authentification.

Chacune de ces parties joue un rôle spécifique :

- la partie 1 permet de gérer les paramètres d'exécution et gère l'entraînement du modèle.
- Le notebook va nous permettre de préparer nos données d'entraîner le modèle, d'effectuer l'inférence et d'exporter le modèle optimal vers le backend pour servir le cloud API.

- la partie backend va mettre en œuvre un projet Django qui va offrir les différents API d'identification et d'authentification entre autre.

2.2 Matériels

Le code d'entraînement sera exécuté sur Colab qui nous offre une machine linux avec toutes les ressources CPU, GPU nécessaire à une telle opération.

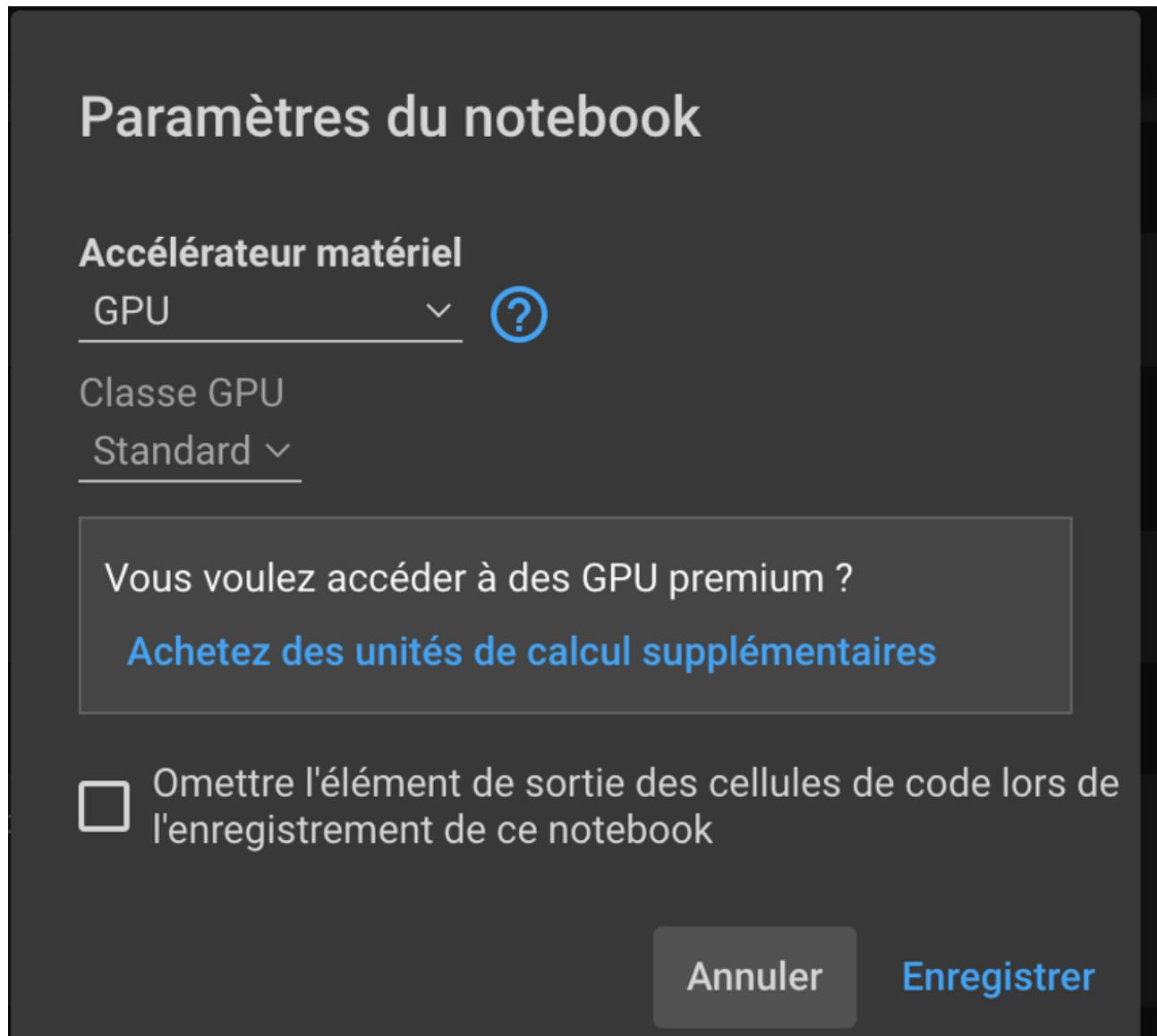


FIGURE 2.1 : paramètres d'exécution du Colab

Le data-set sera sauvegardé sur Google Cloud Storage. En ce qui concerne le backend il sera déployé sur un backend sur le cloud.

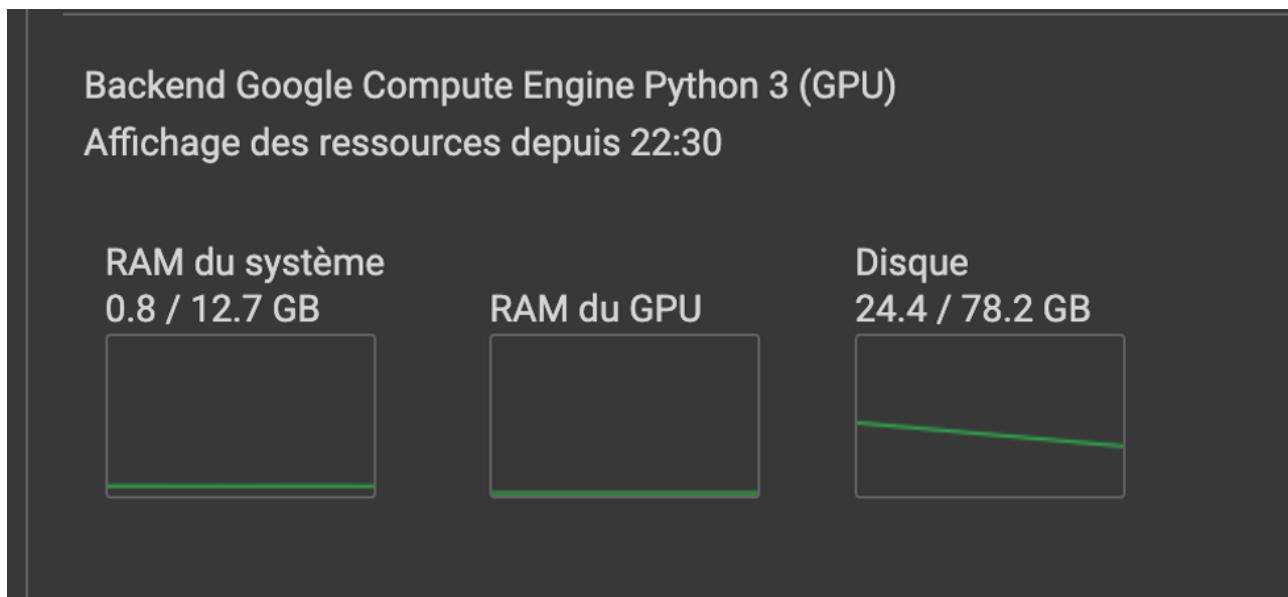


FIGURE 2.2 : Ressource de l'environnement d'exécution Colab

2.3 Choix technologique

2.3.1 Définitions

La reconnaissance de locuteurs est utilisée dans de nombreux domaines, tels que la sécurité, la biométrie, la communication et l'analyse de données. Pour réaliser cette reconnaissance, plusieurs méthodes peuvent être utilisées, notamment la reconnaissance automatique de la parole, l'analyse acoustique, la modélisation de la voix et l'apprentissage automatique.

Les x-vectors sont des vecteurs d'embedding de haut niveau extraits à partir de caractéristiques acoustiques et utilisés pour la reconnaissance automatique de la parole et la vérification du locuteur. Pour extraire les x-vectors, un modèle de réseau de neurones à décalage temporel (TDNN) est souvent utilisé. Le modèle TDNN pour les x-vectors est entraîné sur de grandes quantités de données de parole pour apprendre à extraire les caractéristiques acoustiques les plus discriminantes pour la reconnaissance du locuteur. Le modèle est ensuite utilisé pour extraire les x-vectors des enregistrements de parole et ces vecteurs sont utilisés pour identifier les locuteurs ou vérifier leur identité. Nous avons vu dans le chapitre premier que plusieurs méthodes de deep learning pouvaient permettre d'entraîner des modèles de reconnaissance de locuteurs. Dans le cas de notre étude nous utiliserons les TDNN.

Le TDNN (Time-Delay Neural Networkest) est une méthode de traitement de la parole qui se concentre sur les caractéristiques temporelles de la parole plutôt que sur les caractéristiques fréquentielles. Le TDNN utilise une fenêtre temporelle glissante pour extraire des caractéristiques de la parole à partir de segments de durée fixe. Ces caractéristiques sont ensuite traitées par des couches de neurones, chaque couche traitant des caractéristiques de plus en plus abstraites. Le TDNN utilise également des couches de rétropropagation de l'erreur pour ajuster les poids des neurones en fonction de l'erreur de prédiction. Le TDNN est particulièrement utile pour la reconnaissance de lo-

cuteur car il est capable de traiter **les variations temporelles de la parole, telles que les pauses, les changements de ton et les variations de vitesse. Il peut également être utilisé pour détecter les caractéristiques distinctives de la voix d'un locuteur, telles que la hauteur, la durée et le débit de la parole.**

Le TDNN est également un type de réseau de neurones convolutif (CNN) qui prend en compte les caractéristiques temporelles dans les données d'entrée. Il utilise des convolutions sur des fenêtres de temps pour extraire des caractéristiques pertinentes de la parole. Le TDNN est ensuite suivi d'une couche d'agrégation globale pour résumer les caractéristiques de la parole dans un vecteur d'embedding de haut niveau (X-Vector).

En résumé, le TDNN est un outil puissant pour la reconnaissance de locuteur car il peut prendre en compte les informations temporelles de la parole et les variations de la parole dues aux différences de prononciation et de dialecte. Il peut être entraîné à reconnaître les locuteurs à partir d'enregistrements vocaux et peut être utilisé pour prédire le locuteur d'un enregistrement vocal inconnu.

2.3.2 Structure

Les TDNN sont constitués d'une couche d'entrée, de plusieurs couches cachées et d'une couche de sortie mais ils se différencient par l'organisation des liaisons inter-couches. Les TDNN introduisent des contraintes qui leur permettent d'avoir un certain degré d'invariance par décalage temporel et déformation. Celles-ci utilisent trois idées : poids partagés, fenêtre temporelle et délai. Les poids partagés permettent de réduire le nombre de paramètres du réseau neuronal et induisent ainsi une capacité de généralisation plus importante. Les poids sont partagés suivant la direction temporelle, c'est à dire que pour une caractéristique donnée, la fenêtre associée à celle-ci aura les mêmes poids selon la direction temporelle. De plus cette contrainte entraîne une capacité d'extraire les différences au fur et à mesure du balayage du signal. Le concept de fenêtre temporelle implique que chaque neurone de la couche $l+1$ n'est connecté qu'à un sous ensemble de la couche l (nous n'avons plus une connectivité totale). La taille de cette fenêtre est la même entre deux couches données. Cette fenêtre temporelle permet que chaque neurone n'est qu'une vision locale du signal, il peut être vu comme une unité de détection d'une caractéristique locale du signal. En plus des deux contraintes précédentes nous introduisons des délais entre deux fenêtre successive pour une couche donnée. De plus chaque couche a deux directions : une direction temporelle et une direction caractéristique [4].

2.3.3 Fonctionnement

Le but du TDNN est non pas d'apprendre basiquement le signal mais d'extraire les caractéristiques de celui-ci. La première couche acquiert le signal puis une ou plusieurs couches cachées transforment le signal en des vecteurs de caractéristiques(X-Vectors). Un neurone donné détecte une caractéristique locale de la variation de la courbe. Le champ de vision du neurone est restreint à une fenêtre temporelle limitée. Avec la contrainte des poids partagés, le même neurone est dupliqué dans la direction temps (la même matrice de poids dupliquée) pour détecter la présence ou l'absence de la même caractéristique à différentes places le long du signal. En utilisant plusieurs neurones à chaque position temporelle, le réseau de neurone effectue la détection de caractéristiques différentes : les sorties des différents neurones produisent un nouveau vecteur caractéristique pour la couche supérieure.

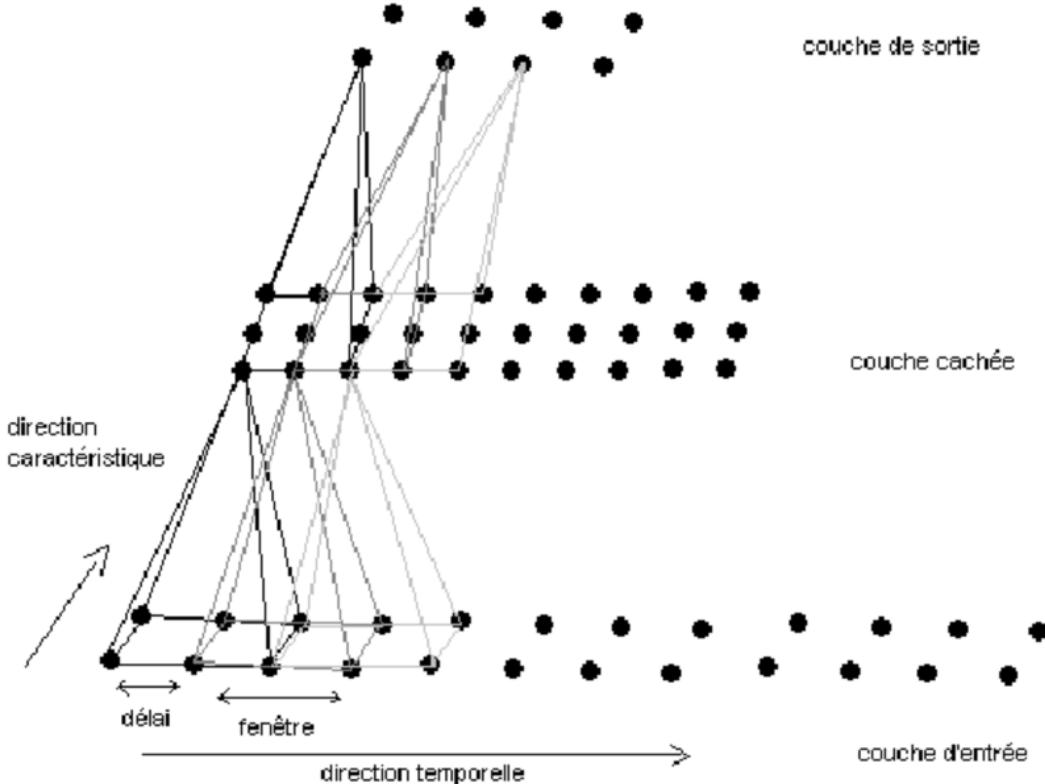


FIGURE 2.3 : Structure d'un RN de type TDNN

Source: The Marvin Project

La composante temporelle du signal d'origine est peu à peu éliminée au fur et à mesure de sa transformation en caractéristique par les couches supérieures, pour compenser cette perte d'information on augmente le nombre de neurones dans la direction caractéristique [11].

Conclusion

Dans ce chapitre, nous avons présenté l'architecture de notre système, les matériels ainsi que les différentes technologies utilisées dans le cadre de notre travail. En ce qui concerne le choix de la méthodologie nous avons opté pour le TDNN.

Il faut retenir que le TDNN (Time Delay Neural Network) est une technique d'apprentissage profond qui est largement utilisée pour la reconnaissance de locuteurs. Cette méthode a l'avantage de traiter les signaux acoustiques dans leur temporalité, ce qui permet de prendre en compte les variations temporelles des caractéristiques acoustiques des locuteurs. Le TDNN est également capable de traiter des données de haute dimensionnalité, telles que les signaux audios, en utilisant des réseaux de neurones profonds, ce qui en fait une méthode performante pour la reconnaissance de locuteurs. Dans notre implémentation, il nous permettra d'extraire les x-vectors à partir des caractéristiques acoustiques afin de faire l'identification ou l'authentification du locuteur.

MATERIELS ET METHODES

Introduction

Ce chapitre est consacré à la modélisation et à l'implémentation de notre model d'IA et du backend Django. Le model permettre d'offrir deux Rest API (un pour l'authentification et le second pour l'identification). Ainsi n'importe quel client web, mobile ou desktop pourra consommer les Apis. La modélisation fera l'analyse du contexte qui nous permet de sortir les différentes fonctionnalités que nous allons mettre en exergue à travers un diagramme de cas d'utilisation. Dans l'implémentation, nous abordons les détails relatifs à l'installation des outils de développement et de test, une esquisse du système, Enfin, nous présentons certaines parties du code source.

3.1 architecture du système

Le TDNN (Time Delay Neural Network) est un type de réseau de neurones récurrent utilisé pour la reconnaissance de la parole et la modélisation de la langue. Il est souvent utilisé dans les systèmes de reconnaissance automatique de la parole, tels que les systèmes de reconnaissance de locuteurs.

Les x-vectors sont des vecteurs de caractéristiques extraits de l'entrée audio, qui sont utilisés pour la reconnaissance de locuteurs. Les TDNN sont souvent utilisés pour extraire les x-vectors à partir de l'audio.

L'architecture d'un TDNN utilisé pour extraire les x-vectors est constituée de plusieurs couches de neurones, chacune avec une taille de fenêtre différente. Chaque couche est composée de plusieurs neurones qui sont connectés à des entrées audio avec un retard temporel.

Le retard temporel est utilisé pour permettre au réseau de traiter des informations audio qui se produisent à différents moments dans le temps. Par exemple, une couche peut être configurée pour traiter des informations audio qui ont été enregistrées 100 millisecondes avant l'instant actuel, tandis qu'une autre couche peut traiter des informations audio enregistrées 200 millisecondes avant l'instant actuel voir figure 2.3.

Les sorties de chaque couche sont ensuite combinées pour produire un vecteur de caractéristiques global, qui est utilisé comme x-vector pour la reconnaissance de locuteurs.

En résumé, l'architecture d'un TDNN utilisé pour extraire les x-vectors est généralement constituée de plusieurs couches de neurones, chacune avec une taille de fenêtre différente et connectée à des entrées audio avec un retard temporel. Les sorties de chaque couche sont combinées pour produire un vecteur de caractéristiques global, qui est utilisé comme x-vector pour la reconnaissance de locuteurs.

3.2 Le model

Tel que nous l'avons énoncé dans les parties précédentes, De nombreux modèles neuronaux peuvent être utilisés pour aborder ce type de tâche. Dans cette étude, nous nous concentrerons sur un classificateur TDNN (xvector).

Comme le montre la figure ci-dessous, l'architecture TDNN est basée sur la topologie à vecteur x (xvector), et introduit plusieurs améliorations pour les IA dans le domaine du traitement des signaux vocaux.

La couche de pooling utilise un mécanisme d'attention dépendant du canal et du contexte, ce qui permet au réseau d'attirer l'attention sur différents frames par canal. Des blocs SqueezeExcitation unidimensionnels (SE) redimensionnent les canaux des cartes de caractéristiques de niveau de frame intermédiaires pour insérer des informations de contexte global dans les blocs de convolution qui fonctionnent localement. Ensuite, l'intégration de blocs Res2 unidimensionnels améliore les performances tout en réduisant simultanément le nombre total de paramètres en utilisant des convolutions groupées de manière hiérarchique.

Enfin, l'agrégation de fonctionnalités multi-couches (MFA) fusionne des informations complémentaires avant le pooling de statistiques en concaténant la carte de caractéristiques de niveau de frame finale avec des cartes de caractéristiques intermédiaires de couches précédentes.

Le réseau est entraîné en optimisant la perte AAMsoftmax sur les identités de locuteurs dans le corpus d'entraînement. L'AAM-softmax est une amélioration puissante par rapport à la perte softmax régulière dans le contexte des problèmes de classification et de vérification fine. Il optimise directement la distance cosinus entre les embeddings de locuteur.

3.3 Les données d'entraînement

L'entraînement se fera avec un petit ensemble de données open source appelé minilibrispeech à laquelle nous avons ajouté des données vocales d'un certain nombre de personnalités politiques et médias. Tous les fichiers vocaux sont au format .flac, channel Mono et de fréquence 16 KHz.

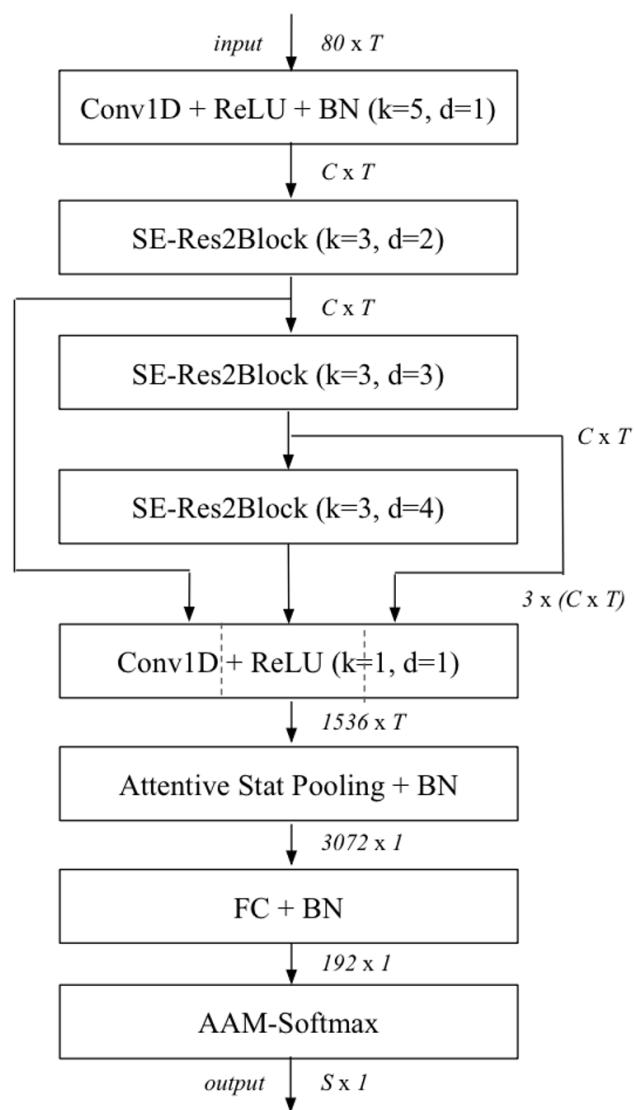


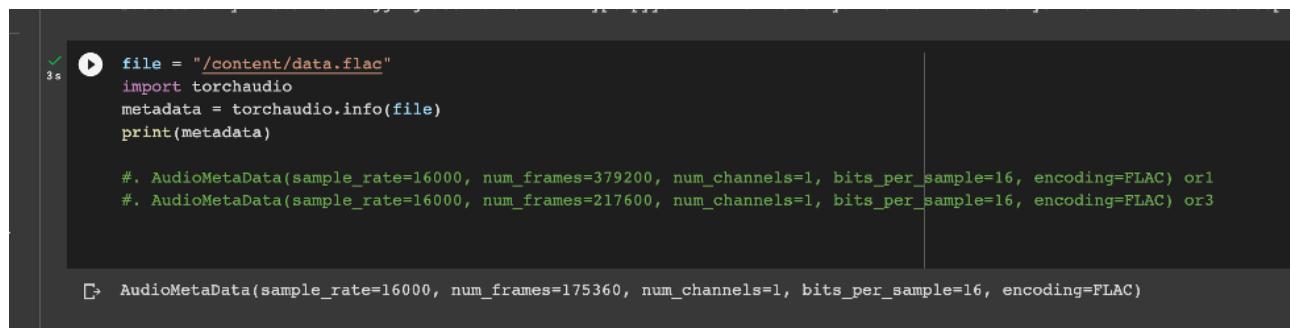
FIGURE 3.1 : Structure d'un RN de type TDNN

Nous avons ainsi, un dataset d'une trentaine d'empreinte vocale, qui ne contient que quelques heures de données de formation. Nous avons aussi mis en place un data pipeline qui à partir d'un API Rest peut servir à étendre le dataset.

3.4 Le code

Pour l'entraînement, on va cloner le répertoire github.com/kadersaka/minfo-speaker-identification. On y trouve un ensemble de fichiers et de dossiers, dont :

- **train.py** : le fichier de code principal, décrit l'ensemble du processus d'entraînement.
- **train.yaml** : le fichier d'hyperparamètres, définit tous les paramètres d'exécution.
- **custom-model.py** : un fichier contenant la définition d'un module PyTorch dont on aura besoin pour l'inférence.



```

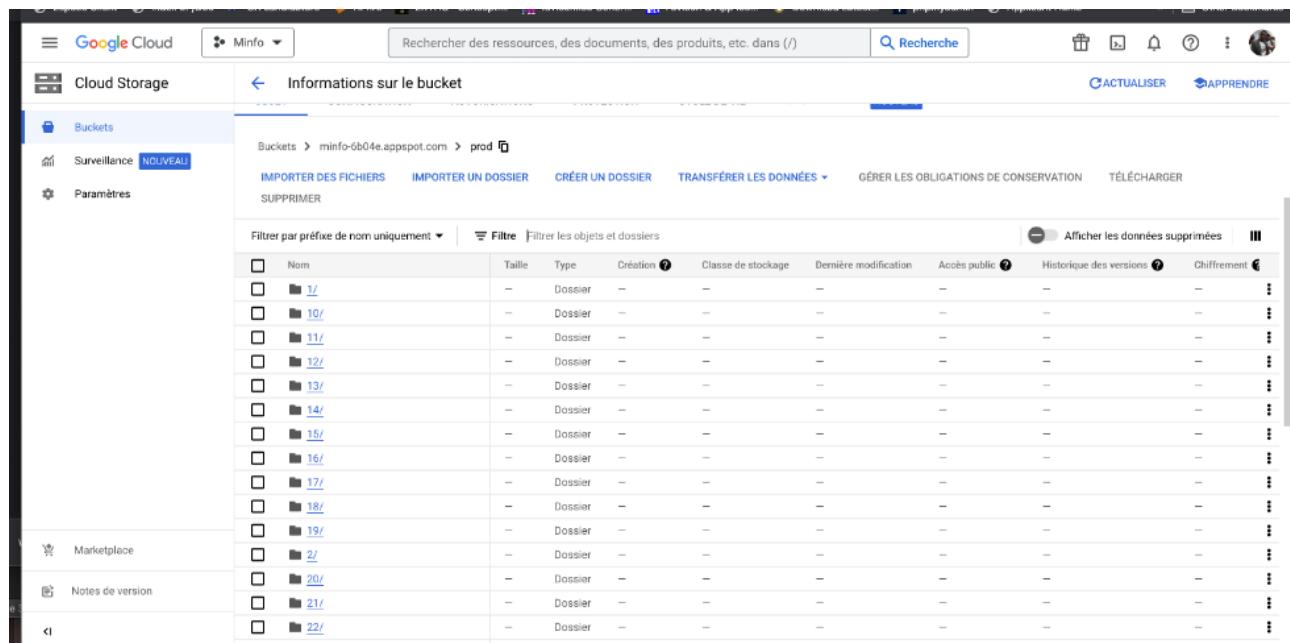
3s  file = "/content/data.flac"
    import torchaudio
    metadata = torchaudio.info(file)
    print(metadata)

    #. AudioMetaData(sample_rate=16000, num_frames=379200, num_channels=1, bits_per_sample=16, encoding=FLAC) or
    #. AudioMetaData(sample_rate=16000, num_frames=217600, num_channels=1, bits_per_sample=16, encoding=FLAC) or3

    ↵  AudioMetaData(sample_rate=16000, num_frames=175360, num_channels=1, bits_per_sample=16, encoding=FLAC)

```

FIGURE 3.2 : caractéristiques des enregistrements vocaux



The screenshot shows the Google Cloud Storage interface. On the left, there's a sidebar with 'Google Cloud' at the top, followed by 'Cloud Storage', 'Buckets', 'Surveillance', 'NOUVEAU', and 'Paramètres'. The main area is titled 'Informations sur le bucket' and shows a list of buckets under 'Buckets > minfo-6b04e.appspot.com > prod'. Below the buckets, there are several buttons: 'IMPORTER DES FICHIERS', 'IMPORTER UN DOSSIER', 'CRÉER UN DOSSIER', 'TRANSFÉRER LES DONNÉES', 'GÉRER LES OBLIGATIONS DE CONSERVATION', and 'TÉLÉCHARGER'. A large table lists objects and folders in the 'prod' bucket. The columns include 'Nom' (Name), 'Taille' (Size), 'Type' (Type), 'Création' (Creation), 'Classe de stockage' (Storage class), 'Dernière modification' (Last modified), 'Accès public' (Public access), 'Historique des versions' (Version history), and 'Chiffrement' (Encryption). Most entries are '1/' through '22/' which are folders. There are also entries for '2/' and '20/' which are files.

FIGURE 3.3 : Le dataset dans le bucket Cloud Storage

- **mini-librispeech-prepare.py** : si nécessaire, télécharge et prépare les manifestes de données si elles ne sont pas encore dans l'environnement d'entraînement.

Pour entraîner le modèle speaker-id, on va exécuter la commande :

3.5 Installations

La première chose à faire est de créer un notebook Colab. Ensuite, nous allons activer l'exécution en mode GPU car nous en auront besoin pour l'entraînement. Puis il faudra cloner le code du répertoire GitHub ci-dessus ; (voir figure 3.5). Ensuite il faudra installer la librairie SpeechBrain (voir figure 3.6).

3.6 Preparation des données

A cette étape la première chose à faire est de donner les autorisations au fichier Colab de pouvoir accéder au Bucket du Cloud Storage pour récupérer le dataset. Ensuite, nous utilisons GsUtils pour

Filtrer par préfixe de nom uniquement							Filtre Filtrer les objets et dossiers	Afficher les données supprimées	
	Nom	Taille	Type	Création	Classe de stockage	Dernière modification	Accès public	Copier l'URL	
<input type="checkbox"/>	kd-1.flac	144,4 Ko	audio/flac	11 déc. 2021, 15:20:49	Standard	11 déc. 2021, 15:20:49	▲ Publique sur Internet	Copier l'URL	
<input type="checkbox"/>	kd-10.flac	111 Ko	audio/flac	11 déc. 2021, 15:20:49	Standard	11 déc. 2021, 15:20:49	▲ Publique sur Internet	Copier l'URL	
<input type="checkbox"/>	kd-2.flac	135,6 Ko	audio/flac	11 déc. 2021, 15:20:49	Standard	11 déc. 2021, 15:20:49	▲ Publique sur Internet	Copier l'URL	
<input type="checkbox"/>	kd-3.flac	146,5 Ko	audio/flac	11 déc. 2021, 15:20:50	Standard	11 déc. 2021, 15:20:50	▲ Publique sur Internet	Copier l'URL	
<input type="checkbox"/>	kd-4.flac	113,9 Ko	audio/flac	11 déc. 2021, 15:20:50	Standard	11 déc. 2021, 15:20:50	▲ Publique sur Internet	Copier l'URL	
<input type="checkbox"/>	kd-5.flac	156,7 Ko	audio/flac	11 déc. 2021, 15:20:50	Standard	11 déc. 2021, 15:20:50	▲ Publique sur Internet	Copier l'URL	
<input type="checkbox"/>	kd-6.flac	159,6 Ko	audio/flac	11 déc. 2021, 15:20:50	Standard	11 déc. 2021, 15:20:50	▲ Publique sur Internet	Copier l'URL	
<input type="checkbox"/>	kd-7.flac	168,1 Ko	audio/flac	11 déc. 2021, 15:20:51	Standard	11 déc. 2021, 15:20:51	▲ Publique sur Internet	Copier l'URL	
<input type="checkbox"/>	kd-8.flac	153,4 Ko	audio/flac	11 déc. 2021, 15:20:51	Standard	11 déc. 2021, 15:20:51	▲ Publique sur Internet	Copier l'URL	
<input type="checkbox"/>	kd-9.flac	153,9 Ko	audio/flac	11 déc. 2021, 15:20:51	Standard	11 déc. 2021, 15:20:51	▲ Publique sur Internet	Copier l'URL	

FIGURE 3.4 : enregistrements vocaux d'un seul utilisateur

```
%%capture
# Let's import our project
!git clone https://github.com/kadersaka/minfo_speaker_identification
```

FIGURE 3.5 : enregistrements vocaux d'un seul utilisateur

la gestion des données du Bucket, notamment leur copie de cloud Storage vers l'environnement local du Colab (voir figure 3.7).

Après avoir récupéré le dataset, nous pouvons maintenant préparer les données. L'objectif de la préparation des données est de créer les fichiers manifestes de données. Ces fichiers indiquent à SpeechBrain où trouver les données audio et leur classification au niveau de l'énoncé correspondante. Ce sont des fichiers texte écrits dans les formats populaires CSV et JSON (voir figure 3.8).

Comme vous pouvez le voir, nous avons une structure hiérarchique spécifiant tous les champs nécessaires à la tâche adressée. Par exemple, nous rapportons le chemin de l'enregistrement vocal, sa durée en secondes (nécessaire si nous voulons trier les phrases avant de créer les mini-lots) et l'identité (ID) du locuteur dans l'enregistrement donné. Il faut préciser qu'à cette étape le dataset est divisé en trois partitions, chacun avec son fichier manifest :

- le premier pour les l'entrainement
- le second pour les test
- et le troisième pour la validation du model

3.7 Entrainement du modèle TDNN

Nous allons former le modèle basé sur TDNN utilisé pour les X-Vector. La mise en commun statistique est utilisée au-dessus des couches convolutives (CNN) pour convertir un vocal de longueur variable en incorporations de longueur fixe. Les paramètres sont sur le fichier train.yaml

Dans ce fichier (voir figure 3.9) :

```
pip install speechbrain
Collecting speechbrain
  Downloading speechbrain-0.5.10-py3-none-any.whl (393 kB)
    ██████████ | 393 kB 5.4 MB/s
Requirement already satisfied: torchaudio in /usr/local/lib/python3.7/dist-packages (from speechbrain) (0.10.0+cu111)
Collecting huggingface-hub
  Downloading huggingface_hub-0.1.2-py3-none-any.whl (59 kB)
    ██████████ | 59 kB 6.4 MB/s
Collecting hyperyaml
  Downloading HyperPyYAML-1.0.0-py3-none-any.whl (15 kB)
Requirement already satisfied: numpy in /usr/local/lib/python3.7/dist-packages (from speechbrain) (1.19.5)
Requirement already satisfied: torch in /usr/local/lib/python3.7/dist-packages (from speechbrain) (1.10.0+cu111)
Requirement already satisfied: packaging in /usr/local/lib/python3.7/dist-packages (from speechbrain) (21.3)
Requirement already satisfied: tqdm in /usr/local/lib/python3.7/dist-packages (from speechbrain) (4.62.3)
Requirement already satisfied: joblib in /usr/local/lib/python3.7/dist-packages (from speechbrain) (1.1.0)
Requirement already satisfied: scipy in /usr/local/lib/python3.7/dist-packages (from speechbrain) (1.4.1)
Collecting sentencepiece
  Downloading sentencepiece-0.1.96-cp37-cp37m-manylinux2_17_x86_64.manylinux2014_x86_64.whl (1.2 MB)
    ██████████ | 1.2 MB 39.9 MB/s
Requirement already satisfied: typing-extensions>=3.7.4.3 in /usr/local/lib/python3.7/dist-packages (from sentencepiece->speechbrain) (3.10.0.2)
Requirement already satisfied: importlib-metadata in /usr/local/lib/python3.7/dist-packages (from huggingface-hub->speechbrain) (4.8.2)
Requirement already satisfied: filelock in /usr/local/lib/python3.7/dist-packages (from huggingface-hub->speechbrain) (3.4.0)
Requirement already satisfied: pyyaml in /usr/local/lib/python3.7/dist-packages (from huggingface-hub->speechbrain) (3.13)
Requirement already satisfied: requests in /usr/local/lib/python3.7/dist-packages (from huggingface-hub->speechbrain) (2.23.0)
Requirement already satisfied: pyparsing>=3.0.5,>=2.0.2 in /usr/local/lib/python3.7/dist-packages (from packaging->speechbrain) (3.0.6)
Collecting pyyaml
  Downloading PyYAML-6.0-cp37-cp37m-manylinux_2_5_x86_64.manylinux1_x86_64.manylinux_2_12_x86_64.manylinux2010_x86_64.whl (596 kB)
    ██████████ | 596 kB 42.7 MB/s
Collecting ruamel.yaml>=0.15
  Downloading ruamel.yaml-0.17.17-py3-none-any.whl (109 kB)
    ██████████ | 109 kB 41.2 MB/s
```

FIGURE 3.6 : enregistrements vocaux d'un seul utilisateur

- Dans la première partie, nous spécifions quelques paramètres de base, tels que le seed, le chemin du dossier qui contiendra nos fichiers vocaux et le chemin du dossier de sortie (voir figure 3.10) :
- Nous spécifions ensuite le chemin des fichiers manifestes de données pour l'entraînement, la validation et le test voir figure 3.11) :
Ces fichiers seront créés automatiquement lors de l'appel du script de préparation des données (mini_librispeech_prepare.py) à partir du fichier d'entraînement (train.py).
- Ensuite, nous configurons le train_logger et déclarons les objets error_stats qui recueilleront des statistiques sur le taux d'erreur de classification (voir figure 3.12) :
- Nous pouvons maintenant spécifier certains hyperparamètres d'entraînement tels que le nombre d'époques, la taille du lot, le taux d'apprentissage, le nombre d'époques et la dimensionnalité d'intégration (voir figure 3.13).

La variable ckpt_interval_minutes peut être utilisée pour enregistrer des points de contrôle toutes les N minutes au cours d'une période d'entraînement. Dans certains cas, une époque peut prendre plusieurs heures, et enregistrer périodiquement le point de contrôle est une bonne pratique sûre. Nous pouvons maintenant définir les modules les plus importants qui sont nécessaires pour former notre modèle (voir figure 3.14)

La partie d'augmentation est basée à la fois sur env_corrupt (qui ajoute du bruit et de la réverbération) et augmentation (qui ajoute des pertes de temps/fréquence et un changement de vitesse). Nous terminons la spécification des hyperparamètres avec la déclaration de l'optimiseur, de l'ordonnanceur du taux d'apprentissage et du pointeur de contrôle :

Dans notre cas, nous utilisons Adam comme optimiseur et une décroissance linéaire du taux d'apprentissage sur les 15 époques. Adam (Adaptive Moment Estimation) est un optimiseur très populaire pour l'apprentissage profond, car il combine les avantages de deux autres optimiseurs, à savoir la descente de gradient stochastique (SGD) avec moment et la méthode d'estimation adaptative de la moyenne du carré des gradients (RMSProp). L'optimiseur Adam est très efficace pour résoudre des problèmes de gradient évanescents ou explosifs, car il utilise des moments d'ordre supérieur pour adapter le taux d'apprentissage pour chaque paramètre. En pratique, l'utilisation de la

```

project_id = 'minfo-6b04e'

[ ] from google.colab import auth
auth.authenticate_user()

[ ] from googleapiclient.discovery import build
gcs_service = build('storage', 'v1')

[ ] bucket_name = "minfo-6b04e.appspot.com"

- GSUtil connect to Google CLoous Storage

[ ] # project_id = '[your Cloud Platform project ID]'
!gcloud config set project {project_id}

Updated property [core/project]. 

[ ] !mkdir data

[ ] #!gsutil -m cp -r gs://{bucket_name}/prod/* gs://{bucket_name}/prod_denoised

[ ] #!gsutil -m cp -r gs://{bucket_name}/prod/* data/
!gsutil -m cp -r gs://{bucket_name}/prod_denoised/* data/

```

FIGURE 3.7 : récupération des données

décroissance linéaire du taux d'apprentissage avec l'optimiseur Adam peut améliorer la stabilité et la convergence du modèle, en particulier pour les problèmes de grande envergure avec des données complexes. Cependant, il est important de noter que le taux d'apprentissage doit être choisi avec soin, car une décroissance trop rapide peut conduire à une convergence lente ou à des minima locaux, tandis qu'une décroissance trop lente peut entraîner une instabilité ou une divergence. Voici donc comment les objets, les fonctions et les hyperparamètres déclarés dans le fichier yaml sont utilisés dans train.py pour implémenter le classifieur. Parcourrons maintenant le fichier principal train.py (voir figure 3.16) :

Nous effectuons ici quelques opérations préliminaires telles que l'analyse de la ligne de commande, l'initialisation des données distribuées en parallèle (nécessaire si plusieurs GPU sont utilisés), la création du dossier de sortie et la lecture du fichier yaml afin de récupérer les paramètres par défaut. Après lecture du fichier yaml avec load_hyperpyyaml, tous les objets déclarés dans les fichiers d'hyperparamètres sont initialisés et disponibles sous forme de dictionnaire (ainsi que les autres fonctions et paramètres rapportés dans le fichier yaml). Par exemple, nous aurons hparams['embedding_model'], hparams['classifier'], hparams['batch_size'], etc. Nous exécutons également le script de préparation des données prepare_mini_librispeech qui crée les fichiers manifestes de données (distribuées à 80% pour l'entraînement, 10% pour la validation, et 10% pour les tests). Il est encapsulé avec sb.utils.distributed.run_on_main car cette opération écrit les fichiers manifestes sur le disque.

Nous appelons ensuite la fonction qui crée les objets de l'ensemble de données pour la formation, la validation et le test.

La première partie est juste une déclaration de CategoricalEncoder qui sera utilisée pour convertir les étiquettes catégorielles en leurs index correspondants. Nous exposons ensuite les fonctions de traitement audio et label. La fonction audio_pipeline prend le chemin d'un signal audio (wav) et le lit. Il renvoie un objet tensor contenant la phrase vocale lue. L'entrée en paramètre de cette fonction

```

1  {
2      "10dc0961-3d9b-4b5f-962e-9fe871254deb": {
3          "wav": "{data_root}/28/10dc0961-3d9b-4b5f-962e-9fe871254deb.flac",
4          "length": 10.0,
5          "spk_id": "28"
6      },
7      "Roland1": {
8          "wav": "{data_root}/23/Roland1.flac",
9          "length": 10.96,
10         "spk_id": "23"
11     },
12     "537d918d-7672-4a9d-95b7-402a15f5e37e": {
13         "wav": "{data_root}/35/537d918d-7672-4a9d-95b7-402a15f5e37e.flac",
14         "length": 9.3686875,
15         "spk_id": "35"
16     },
17     "1088-134318-0001": {
18         "wav": "{data_root}/8/1088-134318-0001.flac",
19         "length": 12.235,
20         "spk_id": "8"
21     },
22     "25210228-6f59-4c50-8e9a-191ddb60b60d": {
23         "wav": "{data_root}/35/25210228-6f59-4c50-8e9a-191ddb60b60d.flac",
24         "length": 12.4520625,
25         "spk_id": "35"
26     },
27     "aea3c2eb-5bd7-4189-8e86-9eb2b2700f61": {
28         "wav": "{data_root}/28/aea3c2eb-5bd7-4189-8e86-9eb2b2700f61.flac",
29         "length": 10.0,
30         "spk_id": "28"
31     },

```

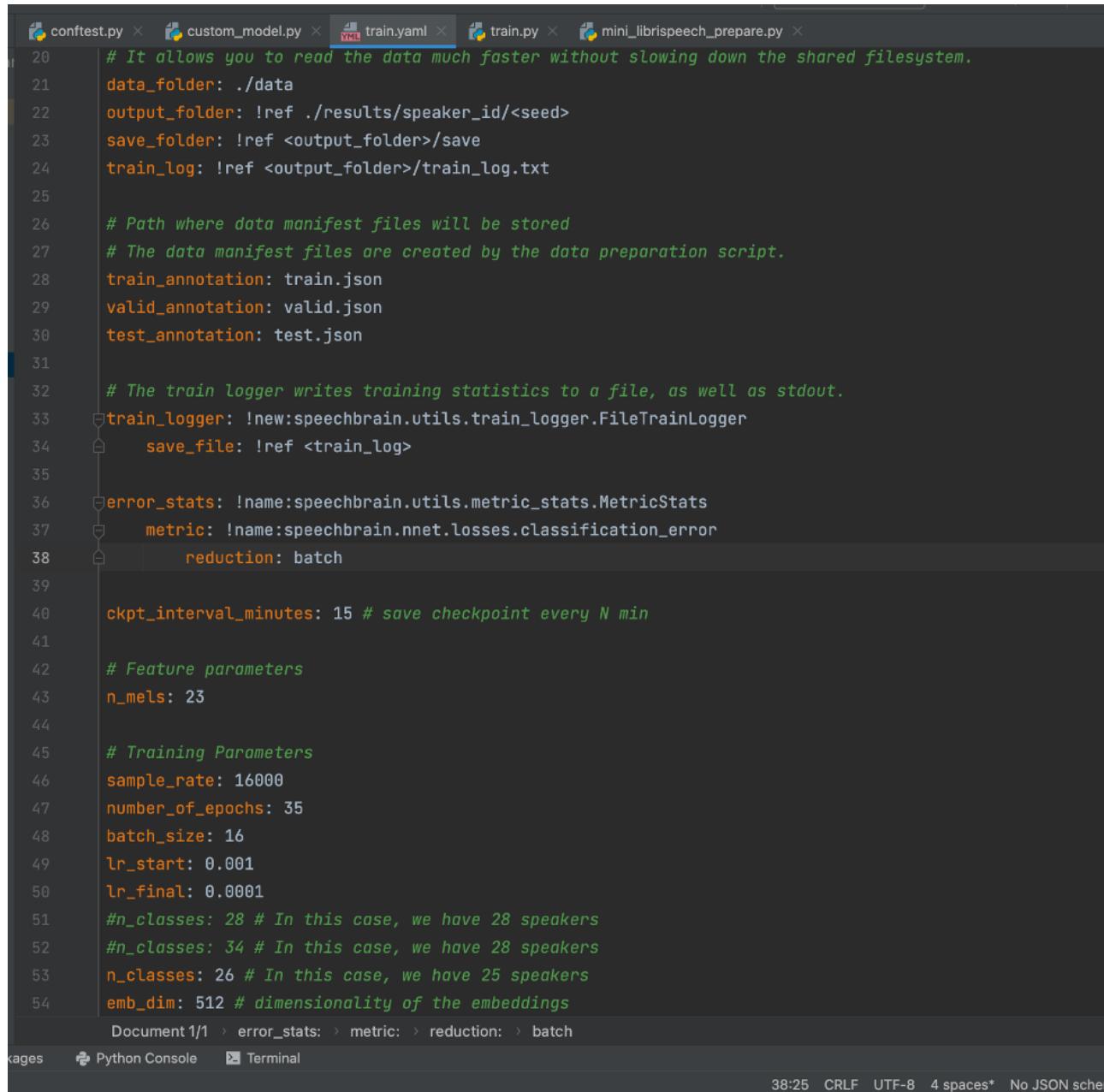
FIGURE 3.8 : Fichier manifestes d'entraînement

(c'est-à-dire le fichier wav) doit avoir le même nom que la clé correspondante dans le fichier manifeste de données généré : Ex :

De même, nous définissons une autre fonction appelée label_pipeline pour traiter les étiquettes au niveau de l'énoncé et les mettre dans un format utilisable par le modèle défini. La fonction lit la chaîne spk_id définie dans le fichier JSON et l'encode avec l'encodeur catégoriel. Nous créons ensuite un objet de type DynamicItemDataset et le connectons aux fonctions de traitement définies ci-dessus. Nous définissons les clés de sortie souhaitées à exposer. Ces clés seront disponibles dans la classe brain dans la variable batch comme :

- batch.id
- lot.sig
- lot.spk_id_encoded

La dernière partie de la fonction est dédiée à l'initialisation de l'encodeur d'étiquettes. L'encodeur d'étiquettes prend en entrée le jeu de données d'apprentissage et attribue un index différent à toutes les entrées spk_id trouvées. Ces index correspondront aux index de sortie du classifieur. Après la définition des jeux de données, la fonction main peut procéder à l'initialisation et à l'utilisation de la classe brain :



```

20  # It allows you to read the data much faster without slowing down the shared filesystem.
21  data_folder: ./data
22  output_folder: !ref ./results/speaker_id/<seed>
23  save_folder: !ref <output_folder>/save
24  train_log: !ref <output_folder>/train_log.txt
25
26  # Path where data manifest files will be stored
27  # The data manifest files are created by the data preparation script.
28  train_annotation: train.json
29  valid_annotation: valid.json
30  test_annotation: test.json
31
32  # The train logger writes training statistics to a file, as well as stdout.
33  train_logger: !new:speechbrain.utils.train_logger.FileTrainLogger
34    save_file: !ref <train_log>
35
36  error_stats: !name:speechbrain.utils.metric_stats.MetricStats
37    metric: !name:speechbrain.nnet.losses.classification_error
38      reduction: batch
39
40  ckpt_interval_minutes: 15 # save checkpoint every N min
41
42  # Feature parameters
43  n_mels: 23
44
45  # Training Parameters
46  sample_rate: 16000
47  number_of_epochs: 35
48  batch_size: 16
49  lr_start: 0.001
50  lr_final: 0.0001
51  #n_classes: 28 # In this case, we have 28 speakers
52  #n_classes: 34 # In this case, we have 28 speakers
53  n_classes: 26 # In this case, we have 25 speakers
54  emb_dim: 512 # dimensionality of the embeddings

```

Document 1/1 > error_stats: > metric: > reduction: > batch

Python Console Terminal

38:25 CRLF UTF-8 4 spaces* No JSON schema

FIGURE 3.9 : Fichier des paramètres d'entraînement

La méthode d'ajustement effectue l'entraînement, tandis que le test est effectué avec celui d'évaluation. Les chargeurs de données d'entraînement et de validation sont fournis en entrée de la méthode d'ajustement, tandis que l'ensemble de données de test est introduit dans la méthode d'évaluation. Examinons maintenant les méthodes les plus importantes définies dans la classe brain :

Les Calculateurs

Commençons par la fonction forward, qui définit tous les calculs nécessaires pour transformer l'audio d'entrée en prédictions de sortie.

Dans ce cas, la chaîne de calcul est très simple. Nous plaçons simplement le lot sur le bon appareil et calculons les caractéristiques acoustiques. Nous traitons ensuite les caractéristiques avec l'encodeur TDNN qui produit un objet tensor de taille fixe (notre X-Vector). Ce dernier alimente un classifieur qui sort les probabilités a posteriori sur les N classes (ici les 31 locuteurs). L'augmentation

```

# Seed needs to be set at top of yaml, before objects with parameters are made
seed: 1986
__set_seed: !!python/object/apply:torch.manual_seed [<seed>]

# If you plan to train a system on an HPC cluster with a big dataset,
# we strongly suggest doing the following:
# 1- Compress the dataset in a single tar or zip file.
# 2- Copy your dataset locally (i.e., the local disk of the computing node).
# 3- Uncompress the dataset in the local folder.
# 4- Set data_folder with the local path.
# Reading data from the local disk of the compute node (e.g. $SLURM_TMPDIR with SLURM-based clusters) is very important.
# It allows you to read the data much faster without slowing down the shared filesystem.
data_folder: ./data
output_folder: !ref ./results/speaker_id/<seed>
save_folder: !ref <output_folder>/save
train_log: !ref <output_folder>/train_log.txt

```

FIGURE 3.10 : paramètres de base d’entraînement

```

# Path where data manifest files will be stored
# The data manifest files are created by the data preparation script.
train_annotation: train.json
valid_annotation: valid.json
test_annotation: test.json

```

FIGURE 3.11 : paramètres des chemins des fichiers manifestes

des données est ajoutée dans la méthode `prepare_features` :

En particulier, lorsque la corruption (Fichiers parasites sonores) de l’environnement est déclarée dans le fichier yaml, nous concaténons dans le même lot la version propre et la version augmentée des signaux parasites. Cette approche double la taille du lot, mais elle implémente un régularisateur très puissant. Le fait d’avoir à la fois la version propre et la version bruyante du signal dans le même lot force le gradient à pointer dans une direction de l'espace des paramètres qui est robuste contre les distorsions du signal.

Objectifs de calcul

Intéressons-nous maintenant à la méthode `compute_objectives` qui prend en entrée les cibles, les prédictions, et estime une fonction de perte :

Les prédictions en entrée sont celles calculées dans la méthode directe. La fonction de coût est évaluée en comparant ces prédictions avec les étiquettes cibles.

Les autres méthodes

Au-delà de ces deux fonctions importantes, nous avons d’autres méthodes qui sont utilisées par l’IA. Le `on_state_starts` est appelé au début de chaque époque et il est utilisé pour configurer les trackers de statistiques. Le `on_stage_end` est appelé à la fin de chaque étape (par exemple, à la fin de chaque période d’entraînement) et s’occupe principalement de la gestion des statistiques, du recuit du taux d’apprentissage et des points de contrôle.

Pour démarrer l’entraînement nous allons exécuter le code suivant : `python train.py train.yaml –number_of_epochs=15`

nous pouvons y passer des paramètres qui vont surcharger les paramètres par défaut énoncés sur le fichier `train.yaml`. On observe sur la console que les pertes de validation et d’entraînement (`loss`

```
# The train logger writes training statistics to a file, as well as stdout.
train_logger: !new:speechbrain.utils.train_logger.FileTrainLogger
  save_file: !ref <train_log>

error_stats: !name:speechbrain.utils.metric_stats.MetricStats
  metric: !name:speechbrain.nnet.losses.classification_error
    reduction: batch
```

FIGURE 3.12 : paramètres des chemins des fichiers manifestes

```
ckpt_interval_minutes: 15 # save checkpoint every N min

# Feature parameters
n_mels: 23

# Training Parameters
sample_rate: 16000
number_of_epochs: 35
batch_size: 16
lr_start: 0.001
lr_final: 0.0001
n_classes: 28 # In this case, we have 28 speakers
emb_dim: 512 # dimensionality of the embeddings
dataloader_options:
  batch_size: !ref <batch_size>
```

FIGURE 3.13 : Autres paramètres

function) diminuent très rapidement dans les premières époques. Ensuite, nous voyons essentiellement quelques améliorations mineures et des oscillations de performances.

Voyons quels fichiers/dossiers sont générés dans le output_folder spécifié :

- **train_log.txt** : contient les statistiques (par exemple, train_loss, valid_loss) calculées à chaque époque.
- **log.txt** : est un enregistreur plus détaillé contenant les horodatages pour chaque opération de base.
- **env.log** : affiche toutes les dépendances utilisées avec leur version correspondante (utile pour la réplicabilité).
- **train.py, hyperparams.yaml** : sont une copie du fichier d'expérience avec les hyperparamètres correspondants
- **save** : est l'endroit où nous stockons le modèle généré.

Dans le dossier de sauvegarde, on a des sous-dossiers contenant les points de contrôle enregistrés lors de l'entraînement (au format CKPT+données+heure). Typiquement, on y trouve ici deux points de contrôle : le meilleur (c'est-à-dire le plus ancien) et dernier (c'est-à-dire le plus récent).

À l'intérieur de chaque point de contrôle, nous stockons toutes les informations nécessaires pour reprendre l'entraînement (par exemple, les modèles, les optimiseurs, les planificateurs, le compteur d'époques, etc.). Les paramètres des modèles d'intégration sont rapportés dans le fichier embedding_model.ckpt, tandis que ceux du classifieur sont dans classifier.ckpt. C'est juste un format binaire lisible avec torch.load. Le dossier de sauvegarde contient également l'encodeur d'étiquette (label_encoder.txt), qui mappe chaque entrée d'id de locuteur à leurs index correspondants. A la fin de l'entraînement, l'erreur de validation doit passer à zéro (ou très proche de zéro).

```

58     env_corrupt: !new:speechbrain.lobes.augment.EnvCorrupt
59         open_in_folder: !ref <data_folder>
60         babble_prob: 0.0
61         reverb_prob: 0.0
62         noise_prob: 1.0
63         noise_snr_low: 0
64         noise_snr_high: 15
65
66     augmentation: !new:speechbrain.lobes.augment.TimeDomainSpecAugment
67         sample_rate: !ref <sample_rate>
68         speeds: [95, 100, 105]
69
70     compute_features: !new:speechbrain.lobes.features.Fbank
71         n_mels: !ref <n_mels>
72
73     mean_var_norm: !new:speechbrain.processing.features.InputNormalization
74         norm_type: sentence
75         std_norm: False
76
77     embedding_model: !new:custom_model.Xvector
78         in_channels: !ref <n_mels>
79         activation: !name:torch.nn.LeakyReLU
80         tdmn_blocks: 5
81         tdmn_channels: [512, 512, 512, 512, 1500]
82         tdmn_kernel_sizes: [5, 3, 3, 1, 1]
83         tdmn_dilations: [1, 2, 3, 1, 1]
84         lin_neurons: !ref <emb_dim>
85
86     classifier: !new:custom_model.Classifier
87         input_shape: [null, null, !ref <emb_dim>]
88         activation: !name:torch.nn.LeakyReLU
89         lin_blocks: 1
90         lin_neurons: !ref <emb_dim>
91         out_neurons: !ref <n_classes>
92
93     epoch_counter: !new:speechbrain.utils.epoch_loop.EpochCounter
94         limit: !ref <number_of_epochs>
95
96     modules:
97         compute_features: !ref <compute_features>
98         env_corrupt: !ref <env_corrupt>
99         augmentation: !ref <augmentation>
100        embedding_model: !ref <embedding_model>
101        classifier: !ref <classifier>
102        mean_var_norm: !ref <mean_var_norm>
103

```

FIGURE 3.14 : Paramètres d’augmentation des données avec les bruits

3.8 Inférence

À ce stade, nous pouvons utiliser notre meilleur model pour effectuer des prédictions sur de nouvelles données. Speechbrain dispose de la classe EncoderClassifier qui facilite l’inférence. La classe peut également être utilisée pour extraire certains plongements à la sortie de l’encodeur. Voyons d’abord comment pouvons-nous l’utiliser pour charger notre meilleur modèle xvector pour effectuer une classification des locuteurs :

Mais comment cela fonctionne-t-il avec notre model que nous avons formé auparavant ?

Exemple d’inférence pour l’authentication. le résultat est contenue dans la variable prédiction (voir figure 3.28)(True = le locuteur est bien celui dont il s’agit, False=Le locuteur n’est pas celui qu’il prétend).

Ici nous avons effectué l’inférence pour l’identification de locuteur ce qui nous a retourné le locuteur avec l’ID 23 (voir figure 3.29).

3.9 Le développement des Apis

A cette étape, nous pouvons exporter notre model vers Cloud Storage afin de pouvoir l’utiliser sur notre backend Django. Le code se presnte comme suit :

- UpdateModel cette classe permet de recuperer le dernier modeldepui Cloud Storage pour l’utilisation a des fin de reconnaissance.

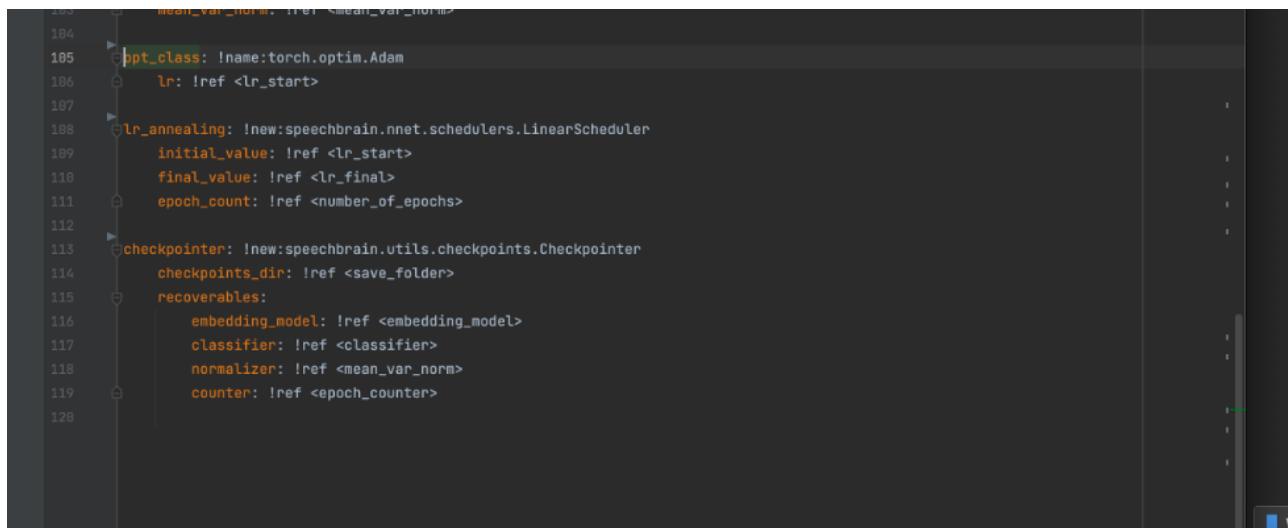


FIGURE 3.15 : Paramètres d'augmentation des données avec les bruits

Requête : GET**Route : https://api.dev.minfo.com/api/speaker/updatemodel**

```

1 class UpdateModel(APIView):
2     permission_classes = [AllowAny]
3
4     def get(self, request):
5         storage_client = storage.Client.from_service_account_json('minfo-6b04e-
6         ed4a5427f384.json')
7         buckets = list(storage_client.list_buckets())
8         print(buckets)
9         bc_name = "minfo-6b04e.appspot.com"
10        bucket = storage_client.get_bucket(bc_name)
11
12        blobs2 = bucket.list_blobs(prefix="models/")
13        print('-----blobs:')
14        print(blobs2)
15        for blob in blobs2:
16            print(blob.name)
17            # destination_uri = '{} / {}'.format("/best_model", blob.name)
18            blob.download_to_filename("speaker_identification/" + blob.name)
19
20        blobs3 = bucket.list_blobs(prefix="denoised_models/")
21        print('-----blobs3:')
22        print(blobs3)
23        for blob in blobs3:
24            print(blob.name)
25            # destination_uri = '{} / {}'.format("/best_model", blob.name)
26            blob.download_to_filename("speaker_identification/" + blob.name)
27
28        data = {
29            "success": True,
30            'message': "Best trained model successfully updated",
31        }
32        return Response(data, status=status.HTTP_200_OK)

```

```

# Recipe begins!
if __name__ == "__main__":
    # Reading command line arguments.
    hparams_file, run_opts, overrides = sb.parse_arguments(sys.argv[1:])

    # Initialize ddp (useful only for multi-GPU DDP training).
    sb.utils.distributed.ddp_init_group(run_opts)

    # Load hyperparameters file with command-line overrides.
    with open(hparams_file) as fin:
        hparams = load_hyppyyaml(fin, overrides)

    # Create experiment directory
    sb.create_experiment_directory(
        experiment_directory=hparams["output_folder"],
        hyperparams_to_save=hparams_file,
        overrides=overrides,
    )

    # Data preparation, to be run on only one process.
    sb.utils.distributed.run_on_main(
        prepare_mini_librispeech,
        kwargs={
            "data_folder": hparams["data_folder"],
            "save_json_train": hparams["train_annotation"],
            "save_json_valid": hparams["valid_annotation"],
            "save_json_test": hparams["test_annotation"],
            "split_ratio": [80, 10, 10],
        },
    )
)

```

FIGURE 3.16 : Code d'entraînement

```

# Create dataset objects "train", "valid", and "test".
datasets = dataio_prep(hparams)

```

FIGURE 3.17 : commande de séparation des données

- **IdentifySpeaker** : Cette classe permet d'identifier le locuteur.

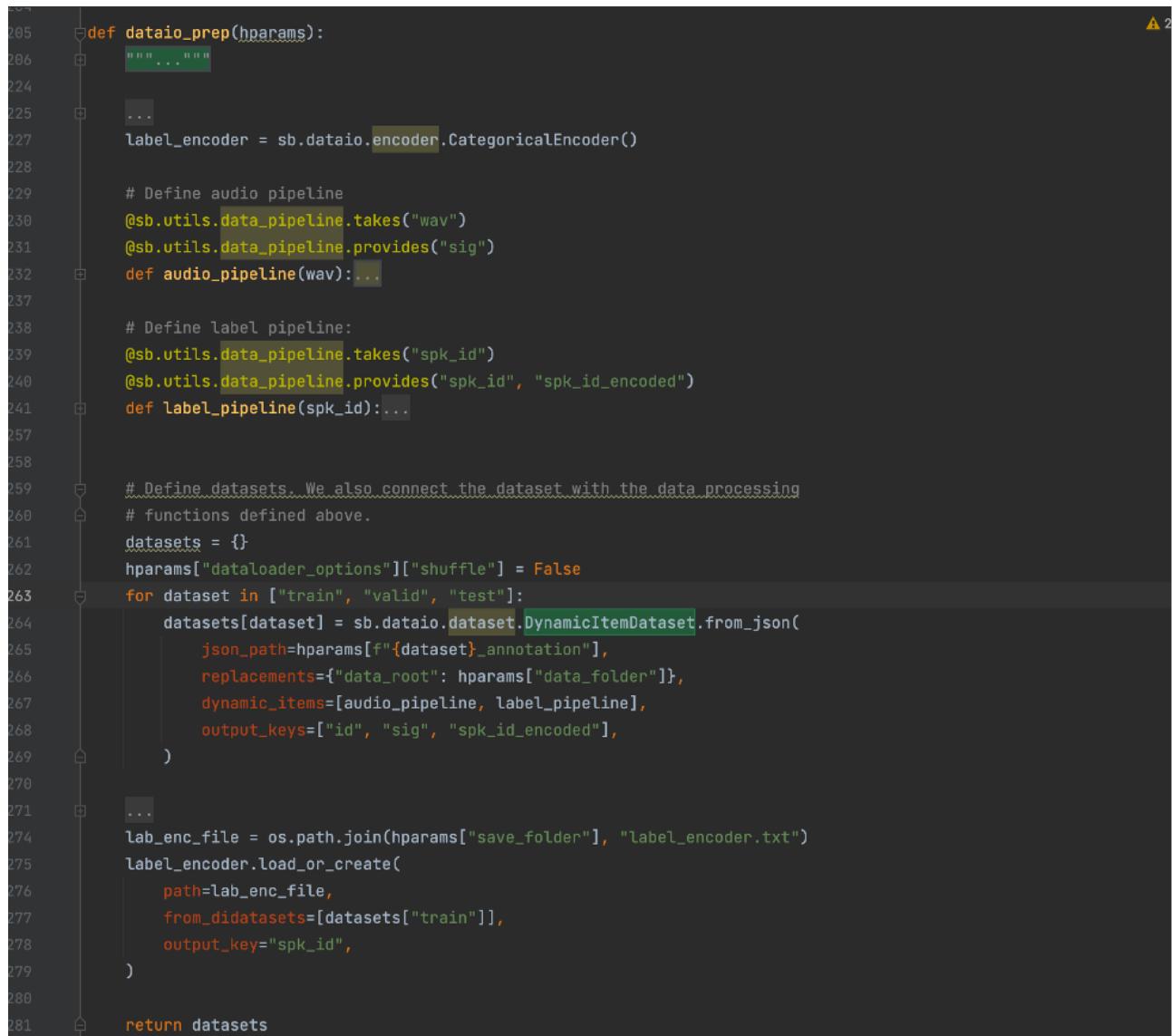
Requête : POST

Route : <https://api.dev.minfo.com/api/speaker/id>

```

1 class IdentifySpeaker(APIView):
2     permission_classes = [AllowAny]
3     def post(self, request):
4         voice = self.request.FILES['voice']
5         if voice:
6             id_data = identify_speaker(voice)
7             id_data_denoised = identify_speaker_denoised(id_data[5], id_data[4])
8             id = id_data[0]
9             accuracy = id_data[1]
10            id2 = id_data[2]
11            accuracy2 = id_data[3]
12            from django.core.files.base import ContentFile, File
13            denoised_voice = None
14
15            did = id_data_denoised[0]
16            daccuracy = id_data_denoised[1]
17            did2 = id_data_denoised[2]

```



```

205     def dataio_prep(hparams):
206         """
207         ...
208
209         label_encoder = sb.dataio.encoder.CategoricalEncoder()
210
211         # Define audio pipeline
212         @sb.utils.data_pipeline.takes("wav")
213         @sb.utils.data_pipeline.provides("sig")
214         def audio_pipeline(wav):...
215
216
217         # Define label pipeline:
218         @sb.utils.data_pipeline.takes("spk_id")
219         @sb.utils.data_pipeline.provides("spk_id", "spk_id_encoded")
220         def label_pipeline(spk_id):...
221
222
223
224
225         # Define datasets. We also connect the dataset with the data processing
226         # functions defined above.
227         datasets = {}
228         hparams["dataloader_options"]["shuffle"] = False
229         for dataset in ["train", "valid", "test"]:
230             datasets[dataset] = sb.dataio.dataset.DynamicItemDataset.from_json(
231                 json_path=hparams[f"{dataset}_annotation"],
232                 replacements={"data_root": hparams["data_folder"]},
233                 dynamic_items=[audio_pipeline, label_pipeline],
234                 output_keys=["id", "sig", "spk_id_encoded"],
235             )
236
237
238         ...
239         lab_enc_file = os.path.join(hparams["save_folder"], "label_encoder.txt")
240         label_encoder.load_or_create(
241             path=lab_enc_file,
242             from_didatasets=[datasets["train"]],
243             output_key="spk_id",
244         )
245
246
247         return datasets

```

FIGURE 3.18 : Code de séparation des données

```

18         daccuracy2 = id_data_denoised[3]
19
20         if id:
21             # id = int(id)+1 97073197 3600
22             pooi = Pooi.objects.filter(id=id).first()
23             pooi2 = Pooi.objects.filter(id=id2).first()
24
25             pooi_data = PooiSerializer(pooi).data
26             pooi2_data = PooiSerializer(pooi2).data
27
28             dpooi = Pooi.objects.filter(id=did).first()
29             dpooi2 = Pooi.objects.filter(id=did2).first()
30
31             dpooi_data = PooiSerializer(dpooi).data
32             dpooi2_data = PooiSerializer(dpooi2).data
33
34             # scanResult = ScanResult(voice=voice, score=float(str(accuracy)),
35             denoisedScore=float(str(accuracy2)), denoisedVoice=voice, pooi=pooi,
36             denoisedPooi=pooi2, created_by=request.user)

```

```

"10dc0961-3d9b-4b5f-962e-9fe871254deb": {
    "wav": "{data_root}/28/10dc0961-3d9b-4b5f-962e-9fe871254deb.flac",
    "length": 10.0,
    "spk_id": "28"
},

```

FIGURE 3.19 : ligne du fichier manifeste

```

320
321     spk_id_brain = SpkIdBrain(
322         modules=hparams["modules"],
323         opt_class=hparams["opt_class"],
324         hparams=hparams,
325         run_opts=run_opts,
326         checkpointer=hparams["checkpointer"],
327     )
328
329     spk_id_brain.fit(
330         epoch_counter=spk_id_brain.hparams.epoch_counter,
331         train_set=datasets["train"],
332         valid_set=datasets["valid"],
333         train_loader_kwargs=hparams["dataloader_options"],
334         valid_loader_kwargs=hparams["dataloader_options"],
335     )
336
337     test_stats = spk_id_brain.evaluate(
338         test_set=datasets["test"],
339         min_key="error",
340         test_loader_kwargs=hparams["dataloader_options"],
341     )
342

```

FIGURE 3.20 : Classe Brain

```

35             # scanResult = ScanResult(voice=voice, score=0.0, denoisedScore=0.0,
36             denoisedVoice=voice, pooli=pooli, denoisedPooli=pooli2, created_by=request.user)
37
38             with open(id_data[4], 'rb') as f:
39                 denoised_voice = File(f, name=os.path.basename(f.name))
40
41                 print('-----')
42                 print(type(denoised_voice))
43
44                 scanResult = ScanResult(voice=voice, denoisedVoice=
45                 denoised_voice, pooli=pooli, denoisedPooli=pooli2, created_by=request.user)
46                 dscanResult = ScanResult(voice=voice, denoisedVoice=
47                 denoised_voice, pooli=dpooli, denoisedPooli=dpooli2, created_by=request.user)
48
49                 print('-----brand logo type 0')
50                 if request.user is not None and request.user.is_authenticated:
51                     print('-----brand logo type 1')
52                     scanResult.created_by = request.user
53                     print('-----brand logo type 2')
54                     print('-----brand logo type 3')
55                     scanResult.save()
56                     print('-----brand logo type4')
57
# score
data = {
    'success': True,
    "score": accuracy,
    "id": id,
}

```



```

37     def compute_forward(self, batch, stage):
38         ...
39
40         # We first move the batch to the appropriate device.
41         batch = batch.to(self.device)
42
43
44         # Compute features, embeddings, and predictions
45         feats, lens = self.prepare_features(batch.sig, stage)
46         embeddings = self.modules.embedding_model(feats, lens)
47         predictions = self.modules.classifier(embeddings)
48
49
50     return predictions

```

FIGURE 3.21 : compute forward



```

63
64     def prepare_features(self, wavs, stage):
65         wavs, lens = wavs
66
67         ...
68
69         if stage == sb.Stage.TRAIN:
70             if hasattr(self.modules, "env_corrupt"):
71                 wavs_noise = self.modules.env_corrupt(wavs, lens)
72                 wavs = torch.cat([wavs, wavs_noise], dim=0)
73                 lens = torch.cat([lens, lens])
74
75
76             if hasattr(self.hparams, "augmentation"):
77                 wavs = self.hparams.augmentation(wavs, lens)
78
79
80         # Feature extraction and normalization
81         feats = self.modules.compute_features(wavs)
82         feats = self.modules.mean_var_norm(feats, lens)
83
84
85     return feats, lens

```

FIGURE 3.22 : augmentation des données

```

58             "pooi": pooi_data,
59             "denoised_score": accuracy2,
60             "denoised_pooi": pooi2_data,
61             "scanResult": ScanResultSerializer(scanResult).data,
62
63             "dscore": daccuracy,
64             "did": did,
65             "dpooi": dpooi_data,
66             "ddenoised_score": daccuracy2,
67             "ddenoised_pooi": dpooi2_data,
68
69         }
70
71         return Response(data, status=status.HTTP_200_OK)
72     else:
73         data = {
74             'message': 'POOI not found'
75         }
76
77         return Response(data, status=status.HTTP_400_BAD_REQUEST)
78     else:
79         data = {
80             'message': 'Voice file is required'
81         }
82
83         return Response(data, status=status.HTTP_400_BAD_REQUEST)

```

```

85
86     def compute_objectives(self, predictions, batch, stage):
87         _, lens = batch.sig
88         spkid, _ = batch.spk_id_encoded
89
90         if stage == sb.Stage.TRAIN and hasattr(self.modules, "env_corrupt"):
91             spkid = torch.cat([spkid, spkid], dim=0)
92             lens = torch.cat([lens, lens])
93
94         loss = sb.nnet.losses.nll_loss(predictions, spkid, lens)
95
96         self.loss_metric.append(
97             batch.id, predictions, spkid, lens, reduction="batch"
98         )
99
100        if stage != sb.Stage.TRAIN:
101            self.error_metrics.append(batch.id, predictions, spkid, lens)
102
103    return loss

```

FIGURE 3.23 : définition de la fonction de perte

- **MatchVoices** : Elle permet d'authentifier un locuteur.

Requête : POST

Route : <https://api.dev.minfo.com/api/speaker/match>

```

1 class MatchVoices(APIView):
2     permission_classes = [AllowAny]
3
4     def post(self, request):
5         from django.core.files.storage import FileSystemStorage
6         voice1 = self.request.FILES['voice1']
7         voice2 = self.request.FILES['voice2']
8
9         if voice2 and voice1:
10             fs = FileSystemStorage()
11
12             v1 = fs.save(voice1.name, voice1)
13             v1_path = fs.path(v1)
14
15             v2 = fs.save(voice2.name, voice2)
16             v2_path = fs.path(v2)
17
18             from speechbrain.pretrained import SpeakerRecognition
19             verification = SpeakerRecognition.from_hparams(source="speechbrain/
20             spkrec-ecapa-voxceleb",
21                                         savedir="",
22             pretrained_models="spkrec-ecapa-voxceleb")
23             score, prediction = verification.verify_files(v1_path, v2_path)
24             # score, prediction = verification.verify_files(file1, file2)
25             data = {
26                 # 'score': voice1.name,
27                 # 'prediction': voice2.name
28                 'score': score,
29                 'prediction': prediction
30             }
31             return Response(data, status=status.HTTP_200_OK)
32         else:
33             data = {

```

```

Chage Insérer Exécution Outils Aide Toutes les modifications ont été enregistrées
  + Code + Texte
  tensor([9])
  84% 38/45 [00:30<00:05,  1.39it/s, train_loss=0.299]-----speaker id is:21
  -----speaker id encoded 1:
  tensor([35])
  -----speaker id is:28
  -----speaker id encoded 1:
  tensor([0])
  -----speaker id is:28
  -----speaker id encoded 1:
  tensor([0])
  -----speaker id is:28
  -----speaker id encoded 1:
  tensor([0])
  -----speaker id is:43
  -----speaker id encoded 1:
  tensor([6])
  -----speaker id is:31
  -----speaker id encoded 1:
  tensor([10])
  -----speaker id is:33
  -----speaker id encoded 1:
  tensor([7])
  -----speaker id is:21
  -----speaker id encoded 1:
  tensor([35])
  -----speaker id is:31
  -----speaker id encoded 1:
  tensor([10])
  -----speaker id is:6
  -----speaker id encoded 1:
  tensor([29])
  -----speaker id is:6
  -----speaker id encoded 1:
  tensor([29])
  -----speaker id is:35
  -----speaker id encoded 1:
  tensor([2])
  -----speaker id is:41
  -----speaker id encoded 1:
  tensor([5])
  -----speaker id is:35
  -----speaker id encoded 1:
  tensor([2])
  -----speaker id is:33
  -----speaker id encoded 1:
  tensor([7])
  -----speaker id is:43
  -----speaker id encoded 1:
  tensor([6])
  87% 39/45 [00:31<00:04,  1.34it/s, train_loss=0.319]-----speaker id is:30
  -----speaker id encoded 1:

```

FIGURE 3.24 : Logs de l'entraînement variation de la fonction loss

```

32     'message': 'Hello Django REST API'
33 }
34 return Response(data, status=status.HTTP_400_BAD_REQUEST)

```

Conclusion

Dans cette partie, nous avons montré comment créer un modèle TDNN pour l'authentification et l'identification des locuteurs. Le système proposé contient tous les ingrédients de base pour développer un système de pointe (c'est-à-dire, l'augmentation de données, l'extraction de caractéristiques, l'encodage, la mise en pool statistique, le classificateur, etc.). Dans le prochain chapitre nous présenterons les résultats obtenus après notre implémentation.

```

+ Code + Texte
  ↗ 00:00:00.000000: speaker id is:28
  ↗ 00:00:00.000000: speaker id encoded 1:
  ↗ tensor([0])
  ↗ 83% 5/6 [00:00<00:00,  8.97it/s]-----speaker id is:31
  ↗ 00:00:00.000000: speaker id encoded 1:
  ↗ tensor([10])
  ↗ 00:00:00.000000: speaker id is:33
  ↗ 00:00:00.000000: speaker id encoded 1:
  ↗ tensor([7])
  ↗ 00:00:00.000000: speaker id is:35
  ↗ 00:00:00.000000: speaker id encoded 1:
  ↗ tensor([2])
  ↗ 00:00:00.000000: speaker id is:31
  ↗ 00:00:00.000000: speaker id encoded 1:
  ↗ tensor([10])
  ↗ 00:00:00.000000: speaker id is:33
  ↗ 00:00:00.000000: speaker id encoded 1:
  ↗ tensor([7])
  ↗ 00:00:00.000000: speaker id is:31
  ↗ 00:00:00.000000: speaker id encoded 1:
  ↗ tensor([10])
  ↗ 00:00:00.000000: speaker id is:30
  ↗ 00:00:00.000000: speaker id encoded 1:
  ↗ tensor([8])
  ↗ 00:00:00.000000: speaker id is:30
  ↗ 00:00:00.000000: speaker id encoded 1:
  ↗ tensor([8])
  ↗ 100% 6/6 [00:00<00:00,  9.54it/s]
speechbrain.net.schedulers - Changing lr from 0.00016 to 0.0001
speechbrain.utils.train_logger - Epoch: 14, lr: 1.64e-04 - train loss: 3.41e-01 - valid loss: 6.80e-01, valid error: 1.02e-01
speechbrain.utils.checkpoints - Saved an end-of-epoch checkpoint in results/speaker_id/1986/save/CKPT+2023-02-23+15-30-43+00
speechbrain.utils.checkpoints - Deleted checkpoint in results/speaker_id/1986/save/CKPT+2023-02-23+15-30-06+00
speechbrain.utils.epoch_loop - Going into epoch 15
  0% 0/45 [00:00:07, ?it/s]-----speaker id is:28
  ↗ 00:00:00.000000: speaker id encoded 1:
  ↗ tensor([0])
  ↗ 00:00:00.000000: speaker id is:23
  ↗ 00:00:00.000000: speaker id encoded 1:
  ↗ tensor([1])
  ↗ 00:00:00.000000: speaker id is:35
  ↗ 00:00:00.000000: speaker id encoded 1:
  ↗ tensor([2])
  ↗ 00:00:00.000000: speaker id is:8
  ↗ 00:00:00.000000: speaker id encoded 1:
  ↗ tensor([3])
  ↗ 00:00:00.000000: speaker id is:35
  ↗ 00:00:00.000000: speaker id encoded 1:
  ↗ tensor([2])
  ↗ 00:00:00.000000: speaker id is:28

```

FIGURE 3.25 : Logs de l'entraînement, passage de l'Epoch 14 à l'Epoch 15

minfo_speaker_identification

- data
- results
 - speaker_id
 - 1986
 - save
 - CKPT+2023-02-23+19-21-51+00
 - CKPT.yaml
 - brain.ckpt
 - classifier.ckpt
 - counter.ckpt
 - dataloader-TRAIN.ckpt
 - embedding_model.ckpt
 - normalizer.ckpt
 - optimizer.ckpt
 - CKPT+2023-02-23+19-23-45+00
 - CKPT.yaml
 - brain.ckpt
 - classifier.ckpt
 - counter.ckpt
 - dataloader-TRAIN.ckpt
 - embedding_model.ckpt
 - normalizer.ckpt
 - optimizer.ckpt
 - label_encoder.txt
 - env.log
 - hyperparams.yaml
 - log.txt
 - train.py
 - train_log.txt
 - samples

```

# Create folder for best model
# !mkdir /content/best_model/
!mkdir /content/model/
# Copy label encoder
# !cp results/speaker_id/1986/save/label_encoder.txt /content/best_model/
# Copy best model
# !cp "ls -td results/speaker_id/1986/*" /content/best_model/
!cp "ls -td results/speaker_id/1986/*" /content/best_model/
[ ] # !rm -r -f /content/best_model/hparams.yaml

we go use the trained classifier to perform prediction
EncoderClassifier one that can make inference easily

Let's see first how can we used it to load our best xvectors
embeddings and perform a speaker classification:

import torchaudio
from speechbrain.pretrained import EncoderClassifier
classifier = EncoderClassifier.from_hparams(
    filelist = '/content/data (3).flac',
    signal,
    fs = torchaudio.load(filelist)

# Compute speaker embeddings
embeddings = classifier.encode_batch(signals)

# Perform classification
output_probs, score, index, text_lab = classifier.classify(signals)

# Posterior log probabilities
print(output_probs)

# Score (i.e., max log posteriors)
print(score)

# Index of the predicted speaker
print(index)

```

FIGURE 3.26 : Les deux models

```

modifications ont été enregistrées
+ Code + Texte

import torchaudio
from speechbrain.pretrained import EncoderClassifier
# classifier = EncoderClassifier.from_hparams(source="speechbrain/spkrec-xvect-voxceleb")

classifier = EncoderClassifier.from_hparams(source="speechbrain/spkrec-xvect-voxceleb")
# signal, fs =torchaudio.load('/content/speechbrain/tests/samples/single-mic/example1.wav')

file1 = '/content/data (3).flac'
signal, fs =torchaudio.load(file1)

# Compute speaker embeddings
embeddings = classifier.encode_batch(signal)

# Perform classification
output_probs, score, index, text_lab = classifier.classify_batch(signal)

# Posterior log probabilities
print(output_probs)

# Score (i.e, max log posteriors)
print(score)

# Index of the predicted speaker
print(index)

# Text label of the predicted speaker
print(text_lab)

Downloading (...)ain/hyperparams.yaml: 100% [2.04k/2.04k [0:00<0:00, 60.3kB/s]
Downloading (...)bedding_model.ckpt: 100% [16.9M/16.9M [0:00<0:00, 87.8MB/s]
Downloading (...)an_var_norm_emb.ckpt: 100% [3.20k/3.20k [0:00<0:00, 124kB/s]
Downloading (...)classifier.ckpt: 100% [15.9W/15.9W [0:00<0:00, 84.5MB/s]
Downloading (...)inlabel_encoder.txt: 100% [129k/129k [0:00<0:00, 185kB/s]
tensor([-16.6096, -32.4139, -26.6173, ..., -21.0858, -17.6970, -14.3375])
tensor([-1.0456])
tensor([3551])
('id02289')

```

FIGURE 3.27 : Préparation du modèle

```

from speechbrain.pretrained import SpeakerRecognition
verification = SpeakerRecognition.from_hparams(source="speechbrain/spkrec-escapa-voxceleb", savedir="pretrained_models/spkrec-escapa-voxceleb")

file1 = '/content/data (3).flac'
file2 = '/content/data (1).flac'

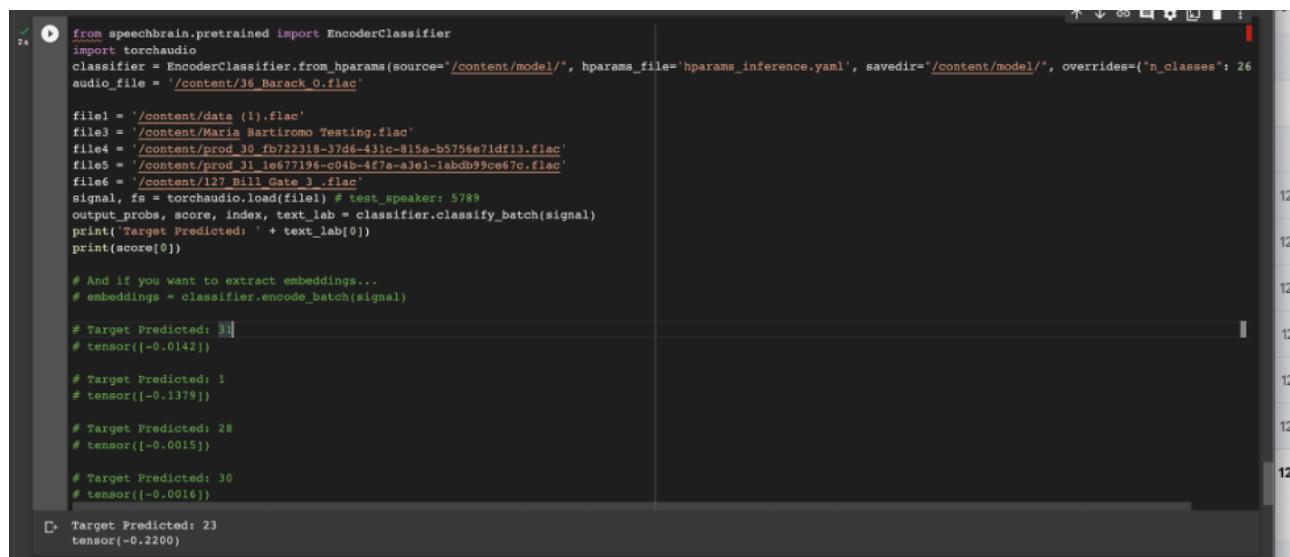
score, prediction = verification.verify_files(file1, file2)

print(score)
print(prediction) # True = same speaker, False=Different speakers

Downloading (...)ain/hyperparams.yaml: 100% [1.92k/1.92k [0:00<0:00, 81.1kB/s]
Downloading (...)bedding_model.ckpt: 100% [83.3M/83.3M [0:00<0:00, 355MB/s]
Downloading (...)an_var_norm_emb.ckpt: 100% [1.92k/1.92k [0:00<0:00, 51.5kB/s]
Downloading (...)classifier.ckpt: 100% [5.53M/5.53M [0:00<0:00, 81.8MB/s]
Downloading (...)inlabel_encoder.txt: 100% [129k/129k [0:00<0:00, 185kB/s]
tensor([10.8843])
tensor([True])

```

FIGURE 3.28 : Inférence sur le modèle pour l'authentification



```
from speechbrain.pretrained import EncoderClassifier
import torchaudio
classifier = EncoderClassifier.from_hparams(source="/content/model/", hparams_file='hparams_inference.yaml', savedir="/content/model/", overrides={"n_classes": 26}
audio_file = '/content/36 Barack_O.flac'

file1 = '/content/data/1.flac'
file3 = '/content/Maria_Bartromo_Testing.flac'
file4 = '/content/prod_30_fb722318-37d6-431c-815a-b5756e71df13.flac'
file5 = '/content/prod_31_1e677196-c04b-4f7a-a3e1-1abdb99ce67c.flac'
file6 = '/content/127_Bill_Gate_3.flac'
signal, fs = torchaudio.load(file1) # test_speaker: 5789
output_probs, score, index, text_lab = classifier.classify_batch(signal)
print('Target Predicted: ' + text_lab[0])
print(score[0])

# And if you want to extract embeddings...
# embeddings = classifier.encode_batch(signal)

# Target Predicted: 31
# tensor([-0.0142])

# Target Predicted: 1
# tensor([-0.1379])

# Target Predicted: 28
# tensor([-0.0015])

# Target Predicted: 30
# tensor([-0.0018])

D Target Predicted: 23
tensor(-0.2200)
```

FIGURE 3.29 : Inférence sur le modèle pour l'identification

Résultats et discussions

Introduction

Dans ce chapitre, nous présentons les différentes fonctionnalités de notre système. Les résultats obtenus sont assez satisfaisants. Ces résultats encourageants montrent des possibilités d'adoption de la reconnaissance de locuteur dans divers domaines, tels que la sécurité biométrique pour l'identification de personnes, la surveillance des appels téléphoniques, les centres d'appels et les interactions vocales dans les voitures connectées. Dans ce chapitre, nous repréciserons le cadre d'étude de notre système, ainsi que lesdits résultats, obtenus à l'issue de l'implémentation.

4.1 Présentation du cadre d'étude

Le présent système mis sur pieds a pour but la reconnaissance de locuteur. La reconnaissance de locuteur est une application de l'intelligence artificielle (IA) qui vise à identifier et à authentifier une personne en fonction de sa voix. Cette technologie peut être utilisée dans différents contextes tels que la sécurité, l'assistance vocale, la gestion de la relation client, etc. l'IA pourra être déployée et mise à la disposition des services de la criminologie pour servir lors des enquêtes judiciaires pour identifier des individus. Elle pourra également servir à sécuriser les systèmes en implémentant un MFA (Multi-Factor Authentication) afin d'accroître la sécurité des systèmes.

4.2 Résultats obtenus

Nous rappelons que la génération de nos modèles, à partir des données d'apprentissage a été possible grâce au service cloud de Google Colaboratory . Celui-ci est un outil intuitif donnant accès à des ressources informatiques permettant de travailler sur des projets en science des données. Par la suite ces modèles ont été sauvegardé sur Cloud Storage afin d'alimenter notre application Django qui offrira les APIs pour les opérations d'identification et d'authentification. L'utilisation de l'IA pour la reconnaissance de locuteur a connu une avancée significative au cours des dernières années grâce à l'amélioration des algorithmes d'apprentissage automatique et de traitement du signal vocal. Les

résultats obtenus sont encourageants et permettent d'identifier un locuteur avec une précision élevée. Les tests réalisés sur des enregistrements de voix ont montré que l'IA est capable de reconnaître un locuteur avec une précision pouvant varier en fonction de la qualité de l'enregistrement, du bruit ambiant et de la durée de la parole.

4.2.1 liste des locuteurs

La figure suivante présente la liste des locuteurs. Dans notre architecture, nous les avons appelé POI (Person Of Interest).

ID	FIRST NAME	LAST NAME	VOICE FILE VALIDATED	NAME	EMAIL	MOBILE PHONE	WEBSITE	WIKIPEDIA	AVIS	FICHIER	LINKEDIN	TWITTER	INSTAGRAM	YOUTUBE
1	Bill	Gate	0	Microsoft	bill@microsoft.com	-	https://microsoft.com	https://en.wikipedia.org	-	https://drive.google.com	https://microsoft.com	https://microsoft.com	https://microsoft.com	https://microsoft.com
2	Mark	Ortman	0	Spotify	mark@spotify.com	555-1234	https://spotify.com	https://en.wikipedia.org	-	https://drive.google.com	https://microsoft.com	https://microsoft.com	https://microsoft.com	https://microsoft.com
3	Mandy	Syler	0	Spotify	mandy@spotify.com	555-1235	https://spotify.com	https://en.wikipedia.org	-	https://drive.google.com	https://microsoft.com	https://microsoft.com	https://microsoft.com	https://microsoft.com
4	Winton	Syler	0	Spotify	winton@spotify.com	555-1236	https://spotify.com	https://en.wikipedia.org	-	https://drive.google.com	https://microsoft.com	https://microsoft.com	https://microsoft.com	https://microsoft.com
5	Pete	Fai	0	Spotify	pete@spotify.com	555-1237	https://spotify.com	https://en.wikipedia.org	-	https://drive.google.com	https://microsoft.com	https://microsoft.com	https://microsoft.com	https://microsoft.com
6	Pete	Fai	0	Spotify	pete@spotify.com	555-1238	https://spotify.com	https://en.wikipedia.org	-	https://drive.google.com	https://microsoft.com	https://microsoft.com	https://microsoft.com	https://microsoft.com
7	Pete	Fai	0	Spotify	pete@spotify.com	555-1239	https://spotify.com	https://en.wikipedia.org	-	https://drive.google.com	https://microsoft.com	https://microsoft.com	https://microsoft.com	https://microsoft.com
8	Pete	Fai	0	Spotify	pete@spotify.com	555-1240	https://spotify.com	https://en.wikipedia.org	-	https://drive.google.com	https://microsoft.com	https://microsoft.com	https://microsoft.com	https://microsoft.com
9	Pete	Fai	0	Spotify	pete@spotify.com	555-1241	https://spotify.com	https://en.wikipedia.org	-	https://drive.google.com	https://microsoft.com	https://microsoft.com	https://microsoft.com	https://microsoft.com
10	Pete	Fai	0	Spotify	pete@spotify.com	555-1242	https://spotify.com	https://en.wikipedia.org	-	https://drive.google.com	https://microsoft.com	https://microsoft.com	https://microsoft.com	https://microsoft.com
11	Pete	Fai	0	Spotify	pete@spotify.com	555-1243	https://spotify.com	https://en.wikipedia.org	-	https://drive.google.com	https://microsoft.com	https://microsoft.com	https://microsoft.com	https://microsoft.com
12	Pete	Fai	0	Spotify	pete@spotify.com	555-1244	https://spotify.com	https://en.wikipedia.org	-	https://drive.google.com	https://microsoft.com	https://microsoft.com	https://microsoft.com	https://microsoft.com
13	Pete	Fai	0	Spotify	pete@spotify.com	555-1245	https://spotify.com	https://en.wikipedia.org	-	https://drive.google.com	https://microsoft.com	https://microsoft.com	https://microsoft.com	https://microsoft.com
14	Pete	Fai	0	Spotify	pete@spotify.com	555-1246	https://spotify.com	https://en.wikipedia.org	-	https://drive.google.com	https://microsoft.com	https://microsoft.com	https://microsoft.com	https://microsoft.com
15	Pete	Fai	0	Spotify	pete@spotify.com	555-1247	https://spotify.com	https://en.wikipedia.org	-	https://drive.google.com	https://microsoft.com	https://microsoft.com	https://microsoft.com	https://microsoft.com
16	Pete	Fai	0	Spotify	pete@spotify.com	555-1248	https://spotify.com	https://en.wikipedia.org	-	https://drive.google.com	https://microsoft.com	https://microsoft.com	https://microsoft.com	https://microsoft.com
17	Pete	Fai	0	Spotify	pete@spotify.com	555-1249	https://spotify.com	https://en.wikipedia.org	-	https://drive.google.com	https://microsoft.com	https://microsoft.com	https://microsoft.com	https://microsoft.com
18	Pete	Fai	0	Spotify	pete@spotify.com	555-1250	https://spotify.com	https://en.wikipedia.org	-	https://drive.google.com	https://microsoft.com	https://microsoft.com	https://microsoft.com	https://microsoft.com
19	Pete	Fai	0	Spotify	pete@spotify.com	555-1251	https://spotify.com	https://en.wikipedia.org	-	https://drive.google.com	https://microsoft.com	https://microsoft.com	https://microsoft.com	https://microsoft.com
20	Pete	Fai	0	Spotify	pete@spotify.com	555-1252	https://spotify.com	https://en.wikipedia.org	-	https://drive.google.com	https://microsoft.com	https://microsoft.com	https://microsoft.com	https://microsoft.com
21	Pete	Fai	0	Spotify	pete@spotify.com	555-1253	https://spotify.com	https://en.wikipedia.org	-	https://drive.google.com	https://microsoft.com	https://microsoft.com	https://microsoft.com	https://microsoft.com
22	Pete	Fai	0	Spotify	pete@spotify.com	555-1254	https://spotify.com	https://en.wikipedia.org	-	https://drive.google.com	https://microsoft.com	https://microsoft.com	https://microsoft.com	https://microsoft.com
23	Pete	Fai	0	Spotify	pete@spotify.com	555-1255	https://spotify.com	https://en.wikipedia.org	-	https://drive.google.com	https://microsoft.com	https://microsoft.com	https://microsoft.com	https://microsoft.com
24	Pete	Fai	0	Spotify	pete@spotify.com	555-1256	https://spotify.com	https://en.wikipedia.org	-	https://drive.google.com	https://microsoft.com	https://microsoft.com	https://microsoft.com	https://microsoft.com
25	Pete	Fai	0	Spotify	pete@spotify.com	555-1257	https://spotify.com	https://en.wikipedia.org	-	https://drive.google.com	https://microsoft.com	https://microsoft.com	https://microsoft.com	https://microsoft.com
26	Pete	Fai	0	Spotify	pete@spotify.com	555-1258	https://spotify.com	https://en.wikipedia.org	-	https://drive.google.com	https://microsoft.com	https://microsoft.com	https://microsoft.com	https://microsoft.com
27	Pete	Fai	0	Spotify	pete@spotify.com	555-1259	https://spotify.com	https://en.wikipedia.org	-	https://drive.google.com	https://microsoft.com	https://microsoft.com	https://microsoft.com	https://microsoft.com
28	Pete	Fai	0	Spotify	pete@spotify.com	555-1260	https://spotify.com	https://en.wikipedia.org	-	https://drive.google.com	https://microsoft.com	https://microsoft.com	https://microsoft.com	https://microsoft.com
29	Pete	Fai	0	Spotify	pete@spotify.com	555-1261	https://spotify.com	https://en.wikipedia.org	-	https://drive.google.com	https://microsoft.com	https://microsoft.com	https://microsoft.com	https://microsoft.com
30	Pete	Fai	0	Spotify	pete@spotify.com	555-1262	https://spotify.com	https://en.wikipedia.org	-	https://drive.google.com	https://microsoft.com	https://microsoft.com	https://microsoft.com	https://microsoft.com
31	Pete	Fai	0	Spotify	pete@spotify.com	555-1263	https://spotify.com	https://en.wikipedia.org	-	https://drive.google.com	https://microsoft.com	https://microsoft.com	https://microsoft.com	https://microsoft.com
32	Pete	Fai	0	Spotify	pete@spotify.com	555-1264	https://spotify.com	https://en.wikipedia.org	-	https://drive.google.com	https://microsoft.com	https://microsoft.com	https://microsoft.com	https://microsoft.com
33	Pete	Fai	0	Spotify	pete@spotify.com	555-1265	https://spotify.com	https://en.wikipedia.org	-	https://drive.google.com	https://microsoft.com	https://microsoft.com	https://microsoft.com	https://microsoft.com
34	Pete	Fai	0	Spotify	pete@spotify.com	555-1266	https://spotify.com	https://en.wikipedia.org	-	https://drive.google.com	https://microsoft.com	https://microsoft.com	https://microsoft.com	https://microsoft.com
35	Pete	Fai	0	Spotify	pete@spotify.com	555-1267	https://spotify.com	https://en.wikipedia.org	-	https://drive.google.com	https://microsoft.com	https://microsoft.com	https://microsoft.com	https://microsoft.com
36	Pete	Fai	0	Spotify	pete@spotify.com	555-1268	https://spotify.com	https://en.wikipedia.org	-	https://drive.google.com	https://microsoft.com	https://microsoft.com	https://microsoft.com	https://microsoft.com
37	Pete	Fai	0	Spotify	pete@spotify.com	555-1269	https://spotify.com	https://en.wikipedia.org	-	https://drive.google.com	https://microsoft.com	https://microsoft.com	https://microsoft.com	https://microsoft.com
38	Pete	Fai	0	Spotify	pete@spotify.com	555-1270	https://spotify.com	https://en.wikipedia.org	-	https://drive.google.com	https://microsoft.com	https://microsoft.com	https://microsoft.com	https://microsoft.com

FIGURE 4.1 : Liste des POI

A partir de cette page nous pouvons facilement étendre la liste, ajouter des empreintes vocales pour les différents POIs.

4.2.2 mise à jour du model

Après l'entrainement sur Google Colab, le model est exporté vers Cloud Storage pour usage ultérieur.

Afin de récupérer ce model et l'utiliser en production dans notre projet Django, nous appelons la route `api/speaker/updatemodel`

4.2.3 Identification de locuteur

L'identification répond à la question : « Qui est cette personne ? ». Elle consiste à vérifier dans notre base de données à qui correspond une empreinte vocale.

Route : `/api/speaker/id`

Requête : POST

Paramètre : `voice(le fichier vocal)`

Le locuteur doit être dans la liste des POI.

The screenshot shows a list of objects in a Google Cloud Storage bucket named 'minfo-6b04e.appspot.com'. The objects are:

Nom	Taille	Type	Création	Classe de stockage	Dernière modification	Accès public
brain.ckpt	48 octets	application/octet-stream	8 août 2022, 11:26:59	Standard	8 août 2022, 11:26:59	Publique sur Internet
classifier.ckpt	1,1 Mo	application/octet-stream	8 août 2022, 11:26:59	Standard	8 août 2022, 11:26:59	Publique sur Internet
counter.ckpt	2 octets	application/octet-stream	8 août 2022, 11:26:58	Standard	8 août 2022, 11:26:58	Publique sur Internet
dataloader-TRAIN.ckpt	2 octets	application/octet-stream	8 août 2022, 11:26:59	Standard	8 août 2022, 11:26:59	Publique sur Internet
label_encoder.ckpt	416 octets	text/plain	8 août 2022, 11:26:59	Standard	8 août 2022, 11:26:59	Publique sur Internet
label_encoder.txt	416 octets	text/plain	8 août 2022, 11:26:59	Standard	8 août 2022, 11:26:59	Publique sur Internet
models.zip	47,2 Mo	application/zip	8 août 2022, 11:27:11	Standard	8 août 2022, 11:27:11	Publique sur Internet
normalizer.ckpt	1 Ko	application/octet-stream	8 août 2022, 11:26:59	Standard	8 août 2022, 11:26:59	Publique sur Internet
optimizer.ckpt	34,2 Mo	application/octet-stream	8 août 2022, 11:26:59	Standard	8 août 2022, 11:26:59	Publique sur Internet

FIGURE 4.2 : Meilleure modèle sauvegardé sur Cloud Storage

The screenshot shows a POST request in Postman to the endpoint `{{base_url}}/api/speaker/updatemodel`. The response status is 200 OK, and the response body is:

```

1  "success": true,
2  "message": "Best trained model successfully updated"
3
4

```

FIGURE 4.3 : récupération du meilleur modèle

4.2.4 authentification de locuteur

L'authentification permet de répondre à la question : « Le locuteur est-il vraiment celui qu'il prétend être ? ». Pour cet API on n'a pas besoin d'avoir les empreintes vocales du locuteur. Ainsi la plateforme qui intègre l'API peut conserver ses empreintes sur ses serveurs privés cela résoud en même temps les problématiques de priviléges d'accès aux données, de sécurités et de confidentialité.

route /api/speaker/match

Requête : POST

Paramètre : voice1 et voice2(le fichier vocal empreinte et le fichier vocale enregistré à l'authentification)

Le locuteur n'a pas besoin d'être enregistré dans la base des POI.

La prédition représente le résultat de l'authentication :

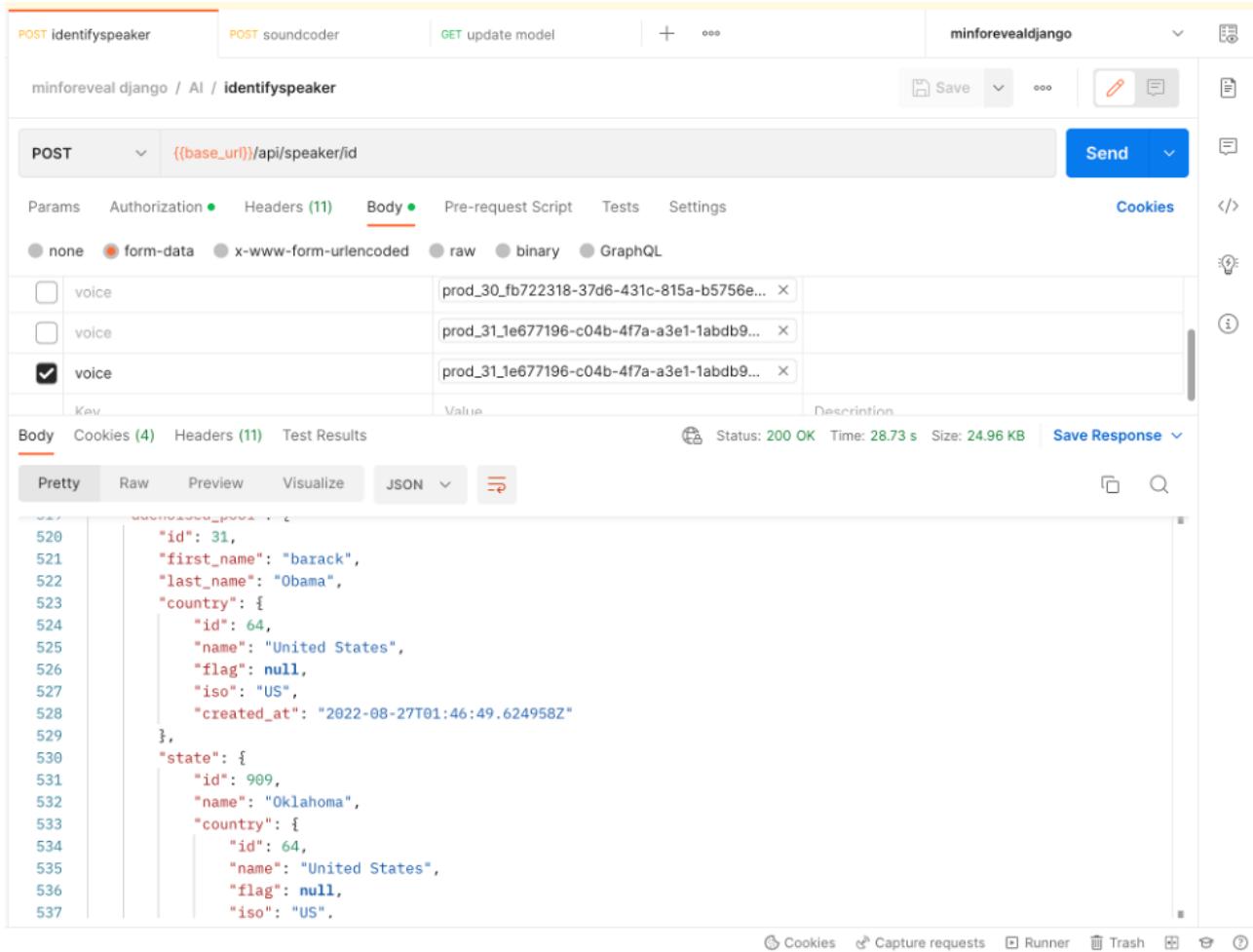


FIGURE 4.4 : Identification de locuteur

- True signifie que l'utilisateur est bien celui qu'il prétend être
- False signifie que l'utilisateur n'est pas le bon, c'est-à-dire que la voix ne correspond pas à l'empreinte vocale.

4.2.5 Analytic des connexions

Nous avons mise en place un service d'analytique pour le suivi et l'utilisation du système. Il enregistre la date, le fichier vocal et le résultat du scan.

4.3 avantage de notre architecture

Le TDNN (Time Delay Neural Network) est l'architecture de réseau de neurones artificiels qui a été utilisée dans. Le cas de notre étude pour la reconnaissance de locuteur. Concernant ses avantages, nous pouvons citer :

- Le TDNN est relativement rapide à entraîner et à utiliser une fois qu'il a été entraîné.
- Les TDNN peuvent prendre en compte les caractéristiques temporelles des signaux acoustiques. En effet, les TDNN sont capables de prendre en compte les données d'entrée passées et présentes, ce qui est important pour la reconnaissance de locuteur.

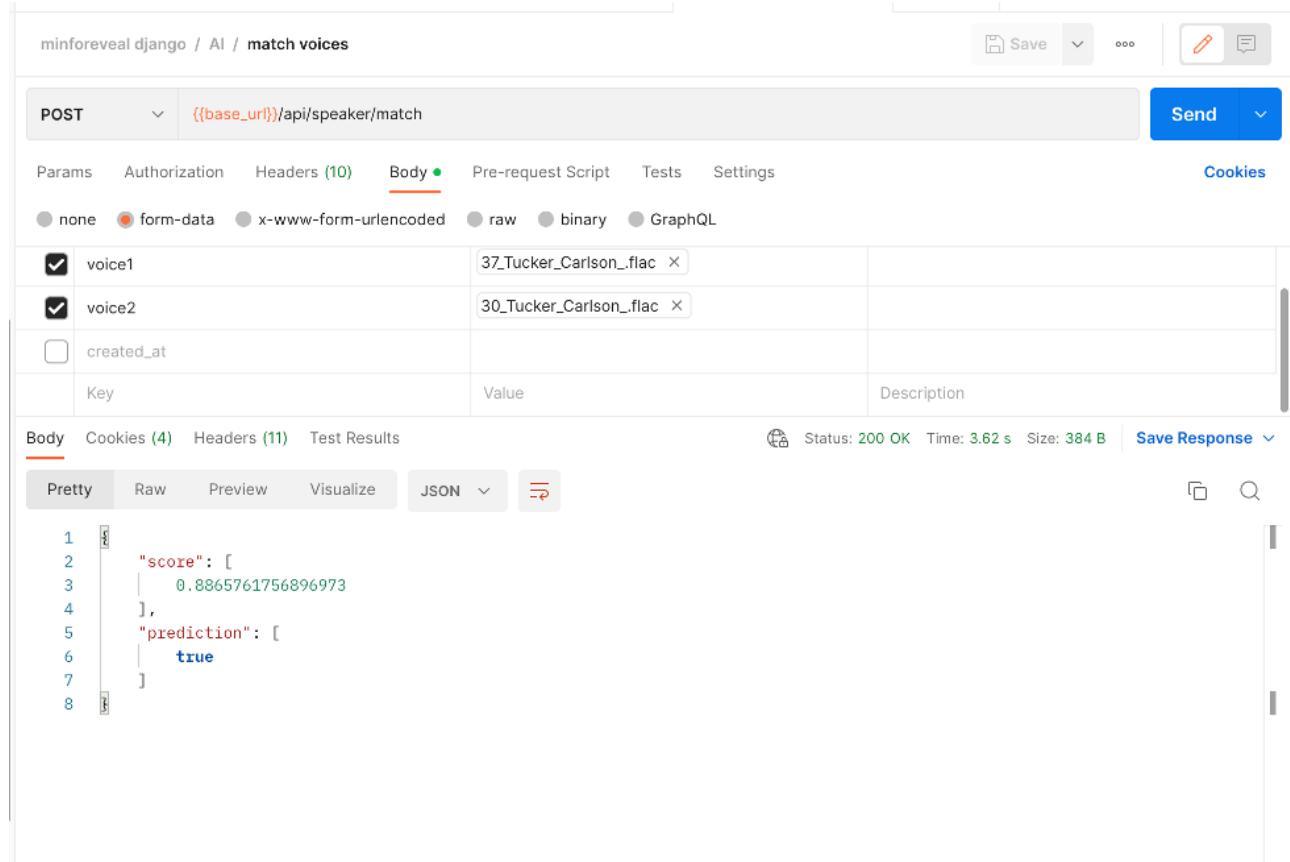


FIGURE 4.5 : Identification de locuteur

- Les TDNN peuvent être utilisés pour extraire des caractéristiques discriminantes des signaux acoustiques. Les TDNN peuvent apprendre des caractéristiques complexes des signaux acoustiques, ce qui est important pour la reconnaissance de locuteur.
- Les TDNN peuvent être entraînés sur de grandes quantités de données. Les TDNN sont capables d'apprendre à partir de grandes quantités de données, ce qui est important pour la reconnaissance de locuteur.

4.4 Les limites de notre architecture

Il convient de noter que les inconvénients du TDNN pour la reconnaissance de locuteurs peuvent varier en fonction du contexte d'utilisation et des caractéristiques des données de parole utilisées pour l'entraînement et le test. Les TDNN peuvent être sensibles aux variations de conditions d'enregistrement. Les TDNN peuvent être affectés par des facteurs tels que le bruit de fond, la qualité de l'enregistrement, la distance entre le locuteur et le microphone, etc. Les TDNN peuvent nécessiter une grande quantité de données d'entraînement. Pour obtenir des performances de reconnaissance de locuteur élevées, il peut être nécessaire de disposer d'une grande quantité de données d'entraînement. Les TDNN peuvent nécessiter une puissance de calcul importante. L'entraînement et l'utilisation de TDNN peut être coûteux en termes de puissance de calcul, surtout lorsqu'il s'agit de grandes quantités de données.

ID	VOICE	POOL	DENOISEDPPOOL
1	scanresult/data_5.flac	Fake Poi 20	Fake Poi 20
2	scanresult/prod_31_1e677196-c04b-4f7a-a3e1-1abdb99ce67c.flac	Fake Poi 20	Fake Poi 20
3	scanresult/data_3.flac	Fake Poi 20	Fake Poi 20
4	scanresult/prod_31_1e677196-c04b-4f7a-a3e1-1abdb99ce67c.flac	Fake Poi 20	Fake Poi 20
5	scanresult/132_Bill_Gate_3_.flac	Fake Poi 20	Fake Poi 20
6	scanresult/132_Bill_Gate_3_.flac	Fake Poi 20	Fake Poi 20
7	scanresult/data_4.flac	Donald Trump	Donald Trump
8	scanresult/prod_31_1e677196-c04b-4f7a-a3e1-1abdb99ce67c.flac	bill Gates	bill Gates

8 scan results

FIGURE 4.6 : Identification de locuteur

Conclusion

Dans ce chapitre, nous avons présenté les résultats obtenus à travers quelques figures présentant l'utilisation de nos différentes APIs. Dans le cas où l'IA est bien entraînée avec des données de haute qualité et utilisée dans un contexte approprié, elle peut fournir une reconnaissance de locuteur précise et fiable. Cela peut être utile dans diverses applications, telles que la sécurité de l'authentification vocale, l'analyse de la parole pour la recherche et la reconnaissance de locuteur dans les applications de communication. Cependant, l'utilisation d'une IA pour la reconnaissance de locuteur peut également soulever des préoccupations en matière de vie privée et de sécurité. Si les données vocales utilisées pour entraîner l'IA sont mal protégées, elles peuvent être utilisées pour identifier les utilisateurs sans leur consentement, ou même être utilisées à des fins malveillantes. Il est donc important de prendre des mesures pour protéger la vie privée des utilisateurs et de veiller à ce que l'utilisation de l'IA pour la reconnaissance de locuteur soit justifiée et éthique. C'est pour ça que notre implémentation de l'authentification nécessite une sauvegarde par notre système des empreintes vocales. En résumé, l'utilisation d'une IA pour la reconnaissance de locuteur peut être utile dans certains contextes, mais elle nécessite une attention particulière pour garantir la précision, la fiabilité et la protection de la vie privée des utilisateurs.

Conclusion Générale

La mise en place d'une IA pour la reconnaissance de locuteur est un domaine en constante évolution qui présente des opportunités importantes pour de nombreuses applications dans diverses industries. Les avancées dans les domaines de l'apprentissage automatique, du traitement du signal et de la reconnaissance de la parole ont permis de développer des systèmes de reconnaissance de locuteur de plus en plus précis et efficaces.

Dans le domaine de la sécurité, l'IA peut être utilisée pour identifier des individus spécifiques en fonction de leur voix, ce qui peut aider à identifier des criminels ou des suspects dans une enquête. Dans les centres d'appels, l'IA peut aider à identifier automatiquement les appelants enregistrés, ce qui peut aider à accélérer le processus d'identification et améliorer l'efficacité du centre d'appels.

Dans le domaine de l'assistance aux personnes âgées ou atteintes de troubles cognitifs, l'IA peut aider à identifier les membres de la famille ou les amis d'un patient par leur voix, ce qui peut aider à améliorer leur qualité de vie et leur sécurité.

Cependant, la mise en place de ces systèmes soulève également des questions éthiques et de confidentialité. Il est important de garantir que les données des utilisateurs soient traitées de manière éthique et sécurisée, conformément aux lois et réglementations applicables.

Pour l'avenir, il est possible que les systèmes de reconnaissance de locuteur soient intégrés dans une gamme encore plus large de dispositifs et de technologies, ce qui permettrait une interaction plus intuitive entre les humains et les machines. Par exemple, les voitures autonomes pourraient utiliser la reconnaissance vocale pour identifier les conducteurs et les passagers, ou les appareils médicaux pourraient être configurés pour reconnaître la voix d'un patient afin de personnaliser les soins.

En fin de compte, la mise en place d'une IA pour la reconnaissance de locuteur présente à la fois des opportunités et des défis, et il est important d'évaluer soigneusement les avantages et les risques potentiels avant de l'utiliser dans une application spécifique.

En perspectives, nous comptons ajouter d'autres fonctionnalités, poursuivre nos travaux de recherche afin de perfectionner au mieux le modèle pour aboutir à une solution beaucoup plus optimale. Il s'agira aussi de mettre en place avec d'autres méthodes et algorithme de Deep Learning pour comparer les performances.

Bibliographie

- [1] Aware. Biométrie vocale embarquée, authentification et identification.
- [2] T. Doctor. What is voice authentication ? (pros, cons, faqs);.
- [3] FRANÇOIS CHOLLET. *Deep Learning with Python*. -, 2020.
- [4] Frontiers. X-vectors : New quantitative biomarkers for early parkinson's disease detection from speech - <https://www.frontiersin.org/articles/10.3389/fninf.2021.578369/full>.
- [5] Idemia. Qu'est-ce que la biométrie ?
- [6] imageware. Biometric voice recognition – everything you should know.
- [7] D. S.-D. G.-R. G. S. D. P. S. Khudanpur. X-vectors : Robust dnn embeddings for speaker recognition. 2022.
- [8] Mathworks. Speaker recognition using x-vectors - <https://www.mathworks.com/help/audio/ug/speaker-recognition-using-x-vectors.html>.
- [9] Nice. Biométrie vocale embarquée, authentification et identification.
- [10] S. S. M. S. F. H. S. Ouni. Embeddings. 2022.
- [11] T. M. Project. Structures d'apprentissage & réseaux neuronaux reconnaissance de l'empreinte vocale - <http://themarvinproject.free.fr/final/node3.html>.
- [12] RCDev. Méthode d'authentification biométrique vocale améliorée.
- [13] Softjourn. Is voice authentication secure enough to be your new password ?
- [14] thalesgroup. La biométrie au service de l'identification et l'authentification.
- [15] Vivoka. Biométrie vocale embarquée, authentification et identification.
- [16] Wikipedia. Time delay neural network - https://en.wikipedia.org/wiki/time_delay_neural_network.

Table des matières

Dédicace	ii
Remerciements	iii
Résumé	iv
.....	iv
Abstract	v
.....	v
List of Figures	vi
Liste des Algorithmes	viii
Introduction	1
.....	1
.....	1
.....	1
.....	1
.....	1
.....	1
.....	2
.....	2
.....	2
.....	2

.....	2
.....	2
.....	2
.....	2
.....	2
.....	2
1 ETAT DE L'ART DE IDENTIFICATION PAR BIOMETRIE VOCALE	3
Introduction	3
1.1 Définitions	3
1.1.1 Biométrie	3
1.1.2 Biométrie vocale	4
1.1.3 Identification et authentification biométrique	4
.....	4
.....	5
1.1.4 Deep learning	5
1.2 historique de la biométrie vocale	5
.....	5
.....	5
.....	5
.....	5
.....	5
.....	5
.....	5
1.3 Types de reconnaissances vocales	6
1.4 Fonctionnement	6
1.5 cas d'utilisation de la biométrie vocale	7
.....	7
1.6 Les avantages de la biométrie vocale	8
1.7 Les inconvénients de la biométrie vocale	8
1.8 Apprentissage profond	9
1.8.1 Notion d'IA	9
.....	9
.....	10
.....	10
1.8.2 Notion d'apprentissage automatique	10
.....	10
.....	11
.....	11
.....	11
.....	11
.....	12
.....	12
.....	12
.....	12
.....	12

1.8.3	Notions de Deep learning	12
1.8.3.1	définition	12
	12
	13
	13
	13
	13
	13
	14
1.8.3.2	Fonctionnement	14
	14
	14
	14
	14
	15
	15
1.9	deep Learning pour la reconnaissance vocale	16
	16
	17
	17
1.9.1	que représentent les x-vector dans la reconnaissance vocale?	18
	18
	18
	18
1.9.2	que représentent les DNN	18
1.9.3	Les ECAPA-TDNN	18
	18
1.9.4	Les ECAPA-TDNN	19
	19
1.10	Les bonnes pratiques en matière de reconnaissance vocale	19
	19
1.11	À quel point la reconnaissance vocale est-elle sécurisée?	19
	19
	20
	20
	20
	20
Conclusion	21
	21
	21
	21
	21
	21

	21
2 MODELISATION ET LE DEVELOPPEMENT DE NOTRE MODEL D'IA		22
Introduction	22
2.1 Architecture du système	22
2.2 Matériels	23
2.3 Choix technologique	24
2.3.1 Définitions	24
.....	24
.....	24
.....	25
.....	25
2.3.2 Structure	25
2.3.3 Fonctionnement	25
.....	25
Conclusion	26
	26
	26
3 MATERIELS ET METHODES		27
Introduction	27
3.1 architecture du système	27
.....	27
.....	27
.....	27
.....	27
.....	28
.....	28
.....	28
.....	28
3.2 Le model	28
.....	28
.....	28
.....	28
.....	28
.....	28
3.3 Les données d'entraînement	28
.....	29
3.4 Le code	29
3.5 Installations	30
3.6 Préparation des données	30
3.7 Entraînement du modèle TDNN	31
3.8 Inférence	38

3.9 Le développement des Apis	38
Conclusion	45
4 Résultats et discussions	49
Introduction	49
4.1 Présentation du cadre d'étude	49
4.2 Résultats obtenus	49
4.2.1 liste des locuteurs	50
4.2.2 mise à jour du model	50
4.2.3 Identification de locuteur	50
4.2.4 authentification de locuteur	51
4.2.5 Analytic des connexions	52
4.3 avantage de notre architecture	52
4.4 Les limites de notre architecture	53
Conclusion	53
Conclusion	55
.	55
.	55
.	55
.	55
.	55
.	55
.	55
Bibliographie	56
Bibliographie	56
Table des matières	57

