



HAL
open science

Amélioration des systèmes de reconnaissance de la parole des personnes âgées

Juline Le Grand

► **To cite this version:**

Juline Le Grand. Amélioration des systèmes de reconnaissance de la parole des personnes âgées. Sciences de l'Homme et Société. 2012. dumas-00736504

HAL Id: dumas-00736504

<https://dumas.ccsd.cnrs.fr/dumas-00736504>

Submitted on 28 Sep 2012

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Amélioration des Systèmes de reconnaissance de la parole des personnes âgées

Filière Sciences du Langage

Master Industrie De la Langue

Parcours Traitement Automatique du Langage Naturel

Mémoire de Master 2 Recherche

Juline Le Grand

Sous la direction de Michel Vacher, Solange Rossato
michel.vacher@imag.fr,solange.rossato@imag.fr

Laboratoire LIG, Équipe : GETALP
BP 53

38041 Grenoble cedex 9

Année universitaire 2011-2012

Déclaration antiplagiat



Déclaration anti-plagiat
Document à scanner après signature
et à intégrer au mémoire électronique

DECLARATION

1. Ce travail est le fruit d'un travail personnel et constitue un document original.
2. Je sais que prétendre être l'auteur d'un travail écrit par une autre personne est une pratique sévèrement sanctionnée par la loi.
3. Personne d'autre que moi n'a le droit de faire valoir ce travail, en totalité ou en partie, comme le sien.
4. Les propos repris mot à mot à d'autres auteurs figurent entre guillemets (citations).
5. Les écrits sur lesquels je m'appuie dans ce mémoire sont systématiquement référencés selon un système de renvoi bibliographique clair et précis.

NOM : LE GRAND..... PRENOM : Juline.....

DATE : 20/09/2012.....

Résumé

La Reconnaissance Automatique de la Parole est une technologie en plein essor dont l'utilisation pour l'aide aux personnes fragiles apparaît comme un domaine novateur et porteur d'espoirs.

Ce mémoire a été réalisé dans ce contexte avec pour objectif d'évaluer l'état du système de reconnaissance de la parole destiné à des personnes âgées dans le cadre du projet CIRDO.

Nous avons étudié, à partir des travaux précédents sur le sujet, les différences entre la parole dite âgée et non âgée; les comportements du système, après adaptations, de manière approfondie sur le plan des phonèmes et classes de phonèmes ainsi que le décodage en mots pour de la parole âgée avec pour métrique le WER.

Ces recherches nous ont permis de mieux comprendre le fonctionnement du système Sphinx3 et de trouver un paramétrage qui permette d'obtenir des résultats intéressants sur un corpus de voix âgées de parole lue et de parole spontanée, thème sur lequel seuls de très rares travaux existent.

Mots-clés : Reconnaissance Automatique de la Parole, personnes âgées, adaptation MLLR, modèle acoustique, modèle de langage, Assistance à la Vie Autonome.

Abstract

Recognition Automatic Speech is a burgeoning technology whose use for assistance to frail appears as an innovative area and hope. This dissertation was made in this context with the aim of assessing the state of the system for speech recognition to elderly in the CIRDO project.

We studied from previous works on the subject, the differences between the non-elderly speaking and elderly speaking ; behavior of the system, after adjustments, thoroughly in terms of phonemes and classes of phonemes ; words decoding of speech for elderly with the WER metric.

This research helped us better understand how the system works and finding a Sphinx3 parameterization which allows to obtain interesting results on a corpus of read speech voice aged and spontaneous speech, topic on which only very few studies exist.

Keywords : Automatic Speech Recognition, Elderly, MLLR adaptation, Acoustic model, Language model, Ambient Assisted Living.

Remerciements

Cette étude a été financée par l'Agence Nationale de la Recherche dans le cadre du projet CIRDO-Recherche Industrielle (ANR-2010-TECS-012).

Je remercie tout d'abord toutes les personnes qui ont accepté de participer aux enregistrements, permettant ainsi le travail que nous avons pu réaliser dans le cadre de ce projet.

Je tiens à remercier Michel Vacher, Solange Rossato et François Portet pour m'avoir encadrée tout au long de ce mémoire, pour leur soutien et leur aide précieuse et très enrichissante.

Je remercie également l'équipe GETALP pour leur accueil très chaleureux et tous les doctorants et stagiaires du "2e étage du batB" qui ont rendu ces quelques mois inoubliables.

Je remercie tout particulièrement Jean-Louis Mas pour m'avoir secouru à maintes reprises face à mon ordinateur absolument pas coopératif.

Je remercie enfin Yuko Sasa, Rémus Dugheanu, Claude Aynaud pour leur importante contribution à ce projet en amont de mon travail et surtout Frédéric Aman sans qui rien de tout cela n'aurait été possible, merci infiniment pour ton aide.

” Dans la vieillesse de vos parents, souvenez-vous de votre enfance“

Ravignan

” Il ne faut pas reprocher aux gens leur vieillesse, puisque tous nous désirons y parvenir“

Bion de Boristhène

Table des matières

| | |
|--|-----------|
| Résumé | 5 |
| Abstract | 6 |
| Remerciements | 7 |
| Citations | 8 |
| Sigles et abréviations | 16 |
| Introduction | 17 |
| 1 Etat de l’art | 19 |
| Etat de l’art | 19 |
| 1.1 La reconnaissance de la parole (ou RAP) | 19 |
| 1.1.1 La parole | 19 |
| 1.1.2 Voix, parole & langage | 20 |
| 1.1.3 RAP & principes | 21 |
| 1.1.4 RAP & outils | 24 |
| 1.2 La voix âgée | 25 |
| 1.2.1 L’évolution de la voix | 26 |
| 1.2.2 Modifications morphologiques et physiologiques | 27 |
| 1.2.3 Modifications acoustiques | 27 |
| 1.2.4 Modifications cognitives et psychologiques | 29 |
| 1.2.5 Impact sur productions acoustiques | 30 |
| 1.2.6 Langage & compétences | 31 |
| 1.3 La RAP des voix âgées | 31 |
| 1.3.1 Dégradations des performances | 31 |
| 2 Contexte de l’étude | 34 |

| | |
|---|-----------|
| Contexte de l'étude | 34 |
| 2.1 L'équipe GETALP | 34 |
| 2.2 Le projet CIRDO | 35 |
| 2.3 E-lío | 36 |
| 3 Problématique | 38 |
| Problématique | 38 |
| 3.1 Répondre à un besoin | 38 |
| 3.2 Améliorer les performances | 38 |
| 3.2.1 Objectif de ce mémoire | 39 |
| 3.3 Attentes, hypothèses | 40 |
| 4 Méthodologie | 41 |
| Méthodologie | 41 |
| 4.1 Sphinx3 | 41 |
| 4.2 Données | 42 |
| 4.2.1 Corpus AD80 disponible au LIG | 42 |
| 4.2.2 Corpus enregistré ERES38 | 43 |
| 4.3 Tâches réalisées et leur protocole | 44 |
| 4.3.1 Alignement forcé phonémique | 44 |
| 4.3.2 Décodage en mots | 46 |
| 4.3.3 Tests sur de la parole spontanée | 47 |
| 4.4 Utilisation des corpus pour cette étude | 48 |
| 4.4.1 Adaptations | 48 |
| 4.4.2 Test | 48 |
| 4.5 Adaptations MLLR | 48 |
| 4.5.1 Technique utilisée pour l'adaptation des modèles acous- tiques | 48 |
| 4.5.2 Technique d'adaptation à un seul locuteur | 50 |
| 5 Evaluation | 52 |
| Evaluation | 52 |
| 5.1 Le Score d'alignement | 52 |
| 5.2 Le Taux d'Erreur Mot (TEM) | 53 |
| 6 Résultats | 55 |

| | |
|--|-----------|
| Résultats | 55 |
| 6.1 La parole lue | 55 |
| 6.1.1 Résultats des scores d'alignement | 55 |
| 6.1.2 Par classe phonémique | 59 |
| 6.1.3 Voix âgée vs non âgées | 62 |
| 6.2 Résultats du TEM | 63 |
| 6.2.1 Modèle de langage vs modèle acoustique | 63 |
| 6.2.2 La question du surapprentissage | 64 |
| 6.2.3 Voix âgée vs non âgées | 66 |
| 6.2.4 Différence homme/femme | 66 |
| 6.2.5 Conclusion | 68 |
| 6.3 La parole spontanée | 69 |
| 6.3.1 Résultats des scores d'alignement | 69 |
| 6.3.2 Par classe phonémique | 72 |
| 6.4 Résultats du TEM | 75 |
| 7 Conclusion & perspectives | 76 |
| Conclusion & perspectives | 76 |
| 7.1 Conclusion | 76 |
| 7.2 Perspectives | 77 |
| Bibliographie | 78 |
| Webographie | 82 |
| 8 Annexes | 83 |
| Annexes | 83 |
| 8.1 Tableau récapitulatif du corpus AD80 | 83 |
| 8.2 Tableau des correspondances API/SAMPA/SPHINX | 84 |
| 8.3 Document explicatif sur les triphones | 86 |
| 8.4 Tableau des occurrences des phonèmes pour le cas de la parole spontanée traité de ERES38 | 90 |
| 8.5 Graphiques des phonèmes montrés en section Résultats | 92 |
| 8.6 Diagrammes représentant pour chaque phonème et par locu- teur la courbe des scores d'alignements en fonction des modèles acoustiques Cas de la parole lue | 95 |
| 8.7 Diagrammes représentant pour chaque phonème et par locu- teur la courbe des scores d'alignements en fonction des modèles acoustiques Comparaison parole lue et spontanée | 107 |

| | | |
|------|---|-----|
| 8.8 | Diagrammes représentant pour chaque phonème et par locuteur la courbe des scores d'alignements en fonction des modèles acoustiques Cas de la parole spontanée | 122 |
| 8.9 | Diagrammes représentant les scores de l'alignement forcé pour la parole lue | 137 |
| 8.10 | Diagrammes représentant les scores de l'alignement forcé pour la parole spontané | 141 |
| 8.11 | Tableau complet ML vs MA | 145 |
| 8.12 | Histogrammes présentant les WER en fonction des modèles acoustiques et par locuteur pour la parole lue | 147 |
| 8.13 | Détail des résultats obtenus pour la parole spontanée sur une partie du corpus ERES38 | 149 |

Table des figures

| | | |
|-----|--|----|
| 1.1 | Schéma du fonctionnement de la communication | 20 |
| 1.2 | Schéma du fonctionnement d'un système de reconnaissance automatique de la parole | 22 |
| 1.3 | Caption for LOF | 28 |
| 1.4 | Audiogramme proposé sur le site bruit&societe.ca (http://www.bruit societe.ca/fr-ca/thematique_cat.aspx?catid=2&scatid=13) | 29 |
| 2.1 | Laboratoire du LIG, bat.B | 34 |
| 2.2 | Logo du projet CIRDO | 35 |
| 2.3 | Dispositif e-lio | 37 |
| 4.1 | Tableau récapitulatif des corpus | 44 |
| 4.2 | Tableau présentant les classes de phonèmes | 46 |
| 4.3 | Schéma des adaptations apportées en fonction des corpus et des tâches du traitement | 49 |
| 4.4 | Légende des figures 4.3 et 4.5 | 49 |
| 4.5 | Schéma des adaptations apportées pour l'adaptation au locuteur en fonction des corpus et des tâches du traitement | 51 |
| 5.1 | Schéma du fonctionnement de l'alignement phonémique et du système d'obtention du score d'alignement | 53 |
| 5.2 | expression du calcul du WER | 54 |
| 6.1 | Courbes présentant les scores d'alignement forcé en fonction des locuteurs et des modèles acoustiques pour la parole | 56 |
| 6.2 | Tableau présentant les phonèmes les moins bien reconnus | 58 |
| 6.3 | Diagrammes des scores de l'alignement forcé pour les consonnes par classe phonémique en fonction des modèles acoustiques pour la parole lue | 60 |

| | | |
|------|---|----|
| 6.4 | Diagrammes des scores de l’alignement forcé pour les voyelles par classe phonémique en fonction des modèles acoustiques pour la parole lue | 61 |
| 6.5 | Diagramme présentant les résultats de locuteurs âgés vs non âgés par classe de phonèmes | 63 |
| 6.6 | Tableau présentant les WER par modèle de langage en fonction du modèle acoustique | 64 |
| 6.7 | Tableau présentant les résultats du décodage en mots en fonction des modèles acoustiques pour une partie du corpus AD80 et comparaison avec les moyennes des scores du décodage en mots pour le corpus ERES38 | 65 |
| 6.8 | Histogramme présentant les résultats du test de généralisation | 66 |
| 6.9 | Tableau présentant les différences homme/femme pour le WER en parole lue et spontanée | 68 |
| 6.10 | Courbes présentant les scores d’alignement en fonction des locuteurs et des modèles acoustiques pour la parole lue vs spontané | 70 |
| 6.11 | Courbes présentant les scores d’alignement en fonction des locuteurs et des modèles acoustiques pour la parole spontanée . | 71 |
| 6.12 | Diagrammes des scores de l’alignement forcé pour les consonnes par classe phonémique en fonction des modèles acoustiques pour la parole spontanée | 73 |
| 6.13 | Diagrammes des scores de l’alignement forcé pour les voyelles par classe phonémique en fonction des modèles acoustiques pour la parole spontanée | 74 |
| 6.14 | tableau présentant les différences entre la parole lue et spontanée | 75 |
| 8.1 | Tableau récapitulatif des sessions composant le corpus AD80 . | 83 |

Sigles et abréviations

Par ordre alphabétique :

ANR Agence Nationale de la Recherche
CIRDO Compagnons Intelligent qui Réagit au Doigt et à l'Œil
ESTER Évaluation des Systèmes de Transcription enrichie d'Émissions Radiophonique
LIG Laboratoire d'Informatique de Grenoble
MA Modèle acoustique
ML Modèle de langage
MLLR Maximum Likelihood Linear Regression
MFCC Mel-scale Frequency Cepstral Coefficients
MMC Modèles de Markov Cachés
RAP Reconnaissance Automatique de la Parole
TAL Traitement Automatique des Langues
TAP Traitement Automatique de la Parole
TEM Taux d'Erreur Mot
TIC Technologies de l'Information et de la Communication
WER Word Error Rate (Taux d'Erreur Mot)

Introduction

Face à l'accroissement démographique important de la population âgée dans les années à venir, la restructuration de l'environnement adapté aux seniors est au cœur des projets. Le maintien à domicile grâce à des professionnels et des moyens adaptés semble être une alternative intéressante face au manque de disponibilités et au coût des infrastructures spécialisées.

La perte d'autonomie se présente suivant des degrés variés et évolutifs, il est donc possible d'intervenir jusqu'à un certain stade en assistant la personne dans son quotidien.

De nombreux projets ont émergé pour répondre à cette demande. En effet, d'une part, du point de vue social, on constate depuis les années 90 une forte hausse des services d'aide à la personne (âgée ou non). Ainsi, infirmières, aides ménagères ou encore livraison de repas sont des services qui se rendent au domicile des personnes pour des "tâches" qu'ils ne peuvent plus réaliser seuls ou avec difficultés. Par ailleurs, du point de vue technologique, on remarque que les TIC sont de plus en plus appliquées pour faciliter le quotidien des personnes fragiles. Détournées de leur but commercial premier, de nombreuses technologies sont aujourd'hui utilisées et adaptées à la population handicapée et âgée de plus en plus sollicitante à l'égard des technologies.

Très vite, la manipulation classique des différentes technologies (ordinateurs, téléphones, internet, systèmes domotiques) peut cependant s'avérer complexe pour des personnes non initiées ou handicapées. C'est alors que, bien que les performances ne soient pas tout à fait encore satisfaisantes, la RAP pour l'utilisation de ces outils s'est montrée très pertinente par son aspect naturel et intuitif.

On trouve ainsi tout un panel de technologies vocales pour le 'grand public' : bornes interactives, logiciels de dictée vocales, services vocaux interactifs, commande vocale pour appareils divers (ordinateur, téléphone, réglage du chauffage, volets, portes, éclairage, lit médicalisé ...)

De cette manière, un grand nombre de tâches sont alors réalisables seulement avec la voix : donner des commandes, entrer des données, créer des

documents textes, naviguer sur internet, aide à la rééducation (exercices de mémoire).

Dans le cadre de l'aide à l'autonomie de personnes ayant divers degrés de dépendance, les systèmes domotiques ont été développés et adaptés.

La grande majorité des personnes âgées ne souhaite pas quitter leur logement et faire appel à des services à domicile. Ainsi, permettant plus d'autonomie, les systèmes domotiques adaptés peuvent aider ces personnes à rester chez elles. En plus de l'interaction avec l'environnement par la voix, les systèmes domotiques offrent la possibilité d'être autonome tout en étant en sécurité grâce à un contact permanent avec l'entourage.

L'un des grands avantages des systèmes de RAP intégrés en domotique est qu'ils sont très personnalisables et adaptables sur le plan des habitudes de la personne dans son logement et sur le plan de son degré de dépendance tout en prenant en compte son évolution. L'objectif est d'améliorer la vie quotidienne de la personne âgée en la secondant dans l'accomplissement de tâches qui sont souvent difficiles pour elle, fatigantes ou dangereuses.

En plus d'un confort matériel, il ne faut pas exclure que ces systèmes pourraient apporter un réel bien-être psychologique. En effet il permet de diminuer l'état de stress lié à la peur d'accidents domestiques, chutes, oublis divers. La confiance que l'on porte aux technologies permettra aux personnes âgées, par l'apport d'un accompagnement continu et d'un renforcement des liens sociaux, d'être plus sereines face à leur isolement et leur vieillissement.

Notre travail sera présenté de la façon suivante : le premier chapitre retracera succinctement l'état de l'art correspondant au sujet ainsi qu'aux thèmes sous-jacents ; le chapitre 2 décrira le contexte de l'étude ; puis nous verrons dans le chapitre 3 la problématique à laquelle cette étude vise à apporter des réponses ; la partie 4 présentera quand à elle la méthodologie employée pour établir le protocole et obtenir des résultats ; le chapitre 5 exposera la phase d'évaluation et la partie 6 détaillera nos résultats ; enfin nous donnerons notre conclusion et nos perspectives sur ce travail.

Chapitre 1

Etat de l'art

1.1 La reconnaissance de la parole (ou RAP)

1.1.1 La parole

La parole est définie par Saussure dans son Cours de Linguistique Générale [de Saussure, 1975] [10] comme étant l'expression concrète de la langue. C'est un mode d'expression propre à l'Homme. La parole est réalisée en contexte d'énonciation, elle est donc plus ou moins dépendante de ce contexte d'énonciation. La parole est la réalisation individuelle du langage qui est lui un phénomène social. Saussure présente la parole comme étant une image auditive qui vient s'associer à un concept cognitif.

Sur le plan physique, la parole est le résultat d'une variation de la pression produite par l'émission d'un son par un locuteur. Il s'agit d'une onde sonore créée par le passage de l'air expulsé des poumons dans les appareils phonatoires et articulatoires du locuteur, ce qui provoque une modification de cette onde puis elle se propage dans l'air. Elle est caractérisée par quatre paramètres que sont la hauteur, la durée, l'intensité et le timbre.

La fréquence du son de la parole se mesure en hertz, l'intensité en décibels avec une moyenne pour la parole à 60dB par rapport au seuil de l'audition. La production de la parole est rapide : 150-300 mots/min. (Macley & Osgood, 1959), 3-5 syllabes/sec. (Deese, 1984), 10-15 phonèmes/sec.

En fonction des locuteurs, la parole est sujette à un très grand nombre de variations (langagières, sociétales, psychologiques, sociales...) qui appuient le fait que la parole soit individuelle et en lien intrinsèque avec la société et le mode de vie de l'individu et ses caractéristiques physiologiques et cognitives. La parole est basée sur l'utilisation des sons d'une langue, c'est-à-dire des phonèmes répertoriés et significatifs pour une langue donnée. Il faut savoir qu'il y a une inter-influence entre les phonèmes dans la parole, on parle de

coarticulation.

L'une des difficultés en traitement de la parole est l'absence de marqueurs entre les mots comme on a les espaces pour l'écrit ainsi que l'extrême variabilité des productions de parole possibles.

Le schéma suivant présente le fonctionnement de la communication [Shannon, 1948] entre deux locuteurs donnés.

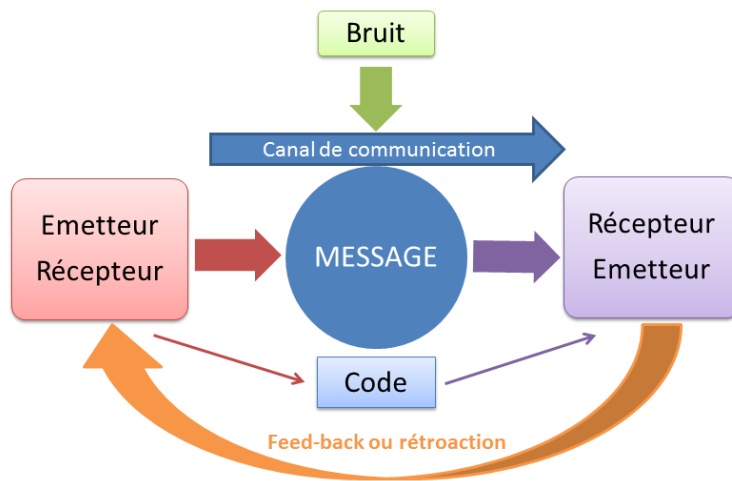


FIGURE 1.1 – Schéma du fonctionnement de la communication

1.1.2 Voix, parole & langage

Face à la pluridisciplinarité des domaines que touche ce travail, il est important de définir les termes utilisés et placer les dichotomies. Dans le domaine de la linguistique la différence entre ce que sont la voix, la parole et le langage se doit d'être précisé ici.

La *voix* correspond à l'ensemble des caractéristiques physiques du son de la parole perçu par un être humain. Elle met en œuvre entre autre le système respiratoire et les cordes vocales. La voix résulte de la morphologie, de la physiologie, de la santé, etc. du locuteur. Elle est ce que l'on entend mais ne porte aucune signification en elle-même, c'est le support acoustique de la parole.

La *parole* correspond à l'expression du langage par la voix. Nous avons déjà décrit la parole dans la section précédente.

Le *langage* est un processus cognitif qui a pour but de rendre possible la communication entre les Hommes. Le langage permet à un individu de traduire sa pensée en phénomène physique perceptible et interprétable par son entourage.

Ainsi on retrouve dans la littérature que le langage structure la pensée, il y a une inter-influence très importante à ce niveau.

Il existe une forte influence sociale et culturelle sur le langage, aussi le vécu de l'individu agit sur sa pensée donc son comportement et son langage.

La particularité du langage en laquelle réside toute la difficulté de son traitement par l'informatique est son caractère ambigu. En effet, tant à l'écrit qu'à l'oral, les productions comportent des éléments qui peuvent avoir plusieurs significations et dont le sens induit par l'émetteur, pourtant souvent bien interprété par l'homme, reste un problème pour le traitement automatique. De par son caractère référentiel le langage se rapporte souvent, et de manière implicite, à des éléments communs entre les locuteurs, faisant appel à des données en mémoire. Ce sont ces références qui permettent de restituer le contexte et le sens de l'énoncé dont ne dispose pas toujours la machine. Afin de pallier à cela, il faut trouver des méthodes qui vont permettre de contourner ce manque d'informations, par exemple des méthodes statistiques, probabiliste mais qui peuvent s'écarter des principes du langage dont il est pourtant question.

1.1.3 RAP & principes

Existantes depuis les années 70, les technologies vocales et le traitement automatique de la parole se sont révélés assez utilisées et développées dans de nombreux domaines ces dernières années. Cela est dû au fait que la parole soit un moyen naturel d'utilisation ou d'interaction avec la machine. Le problème qui subsiste est la gestion de l'immense variation qui existe entre des sujets qui parlent pourtant la même langue. Ces différences sont d'origines sociales, physiologiques, culturelles ou encore psychologiques. En cela et répondant au principe du langage, toute production est unique (locuteur, espace-temps, choix des mots, intonation, prononciation. . .) d'où la difficulté de traitement. Les méthodes actuelles les plus utilisées pour la reconnaissance de la parole sont les méthodes statistiques et le traitement du signal. Ce traitement est donc réalisé sur des critères récurrents à tous les locuteurs et peu en prenant compte des variations propres aux locuteurs qui sont pourtant fortement porteuses de sens. Il serait donc pertinent de réussir à analyser la parole dans son contexte complet. Mais toute la difficulté réside dans le fait de pouvoir attribuer une signification à des mesures physiques effectuées sur la parole dans son contexte d'énonciation.

La Reconnaissance Automatique de la Parole est un processus qui permet de passer d'un signal acoustique de parole à la transcription de ce signal en version écrite. Ce message peut ensuite être utilisé par divers traitements.

Dans le schémas 1.2 il est présenté l'architecture d'un système de RAP [Vacher et al., 2011]. On y voit quatre étapes. Tout d'abord l'interface audio où l'input (onde sonore) arrive et où sont calculés les paramètres acoustiques à partir du spectre. Ensuite viennent le module acoustique et l'extraction de phonèmes qui estiment les successions de phonèmes les plus probables grâce à un modèle acoustique où les paramètres cibles des phonèmes ont été modélisés. Puis arrive l'extraction de mots qui va décider les mots le plus probables en fonction des suites de phonèmes reconnus à partir d'un dictionnaire contenant les transcriptions de chaque mots en phonèmes. Enfin vient l'extraction de phrases qui détermine les suites de mots les plus probables grâce à un modèle de langage tri-gramme appris à partir des transcriptions spécifiques ou plus générales.

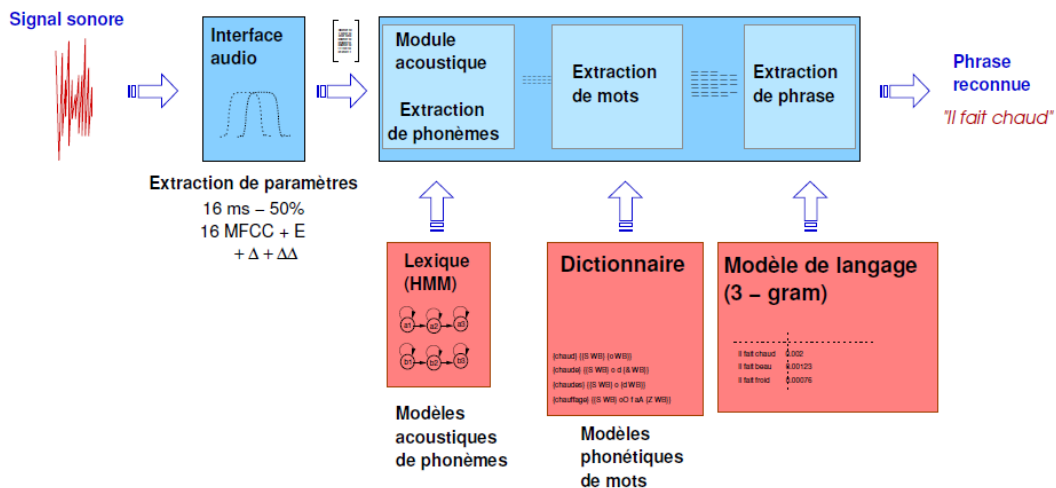


FIGURE 1.2 – Schéma du fonctionnement d'un système de reconnaissance automatique de la parole

On répertorie différentes activités en traitement automatique de la parole (TAP) [Haton et al., 1991] :

- codage et compression de parole (transmission et restitution de parole) ;
- synthèse de la parole (reproduction d'un signal vocal à partir d'un dictionnaire d'éléments phonétiques et/ou de règles de combinaisons) ;
- reconnaissance automatique de la parole (identification par une machine de mots ou phonèmes prononcés par un locuteur humain dans le cadre d'un contexte d'action) ;
- vérification et identification du locuteur (authentification d'une personne par sa voix et restitution d'éléments de parole à locuteur) ;

Les différentes activités vues ci-dessus peuvent être appliquées à divers

domaines [Allegre, 2003] :

- services vocaux interactifs (SVI), très utilisés par les entreprises pour les relations clients ou services clients sur internet ou par téléphone ;
- contrôle qualité et saisie de données, proposent aux entreprises un gain de temps et une liberté d'action dans leur travail de logistique ;
- aérospatial, naval, militaire, médecine, des domaines très demandeurs de ces technologies qui présentent un atout considérable pour des métiers qui demandent des actions supplémentaires quand les mains sont déjà occupées ;
- formation (à distance ou non), très demandé dans le cadre de l'apprentissage des langues étrangères pour la correction de la prononciation ;
- aide et assistance aux personnes fragiles ou dépendantes, de plus en plus utilisé pour l'aide à la vie autonomes des personnes handicapées et/ou âgées par un contrôle de l'environnement grâce à la voix ;
- dictée vocales, permet de se passer d'un clavier dans le cadre d'actions diverses ;

Depuis le début du TAP, plusieurs verrous technologiques ont sauté. De nombreuses tâches fonctionnent désormais très bien comme la reconnaissance de mots isolés (le plus souvent en monolocuteur) pour un vocabulaire de quelques centaines/dizaines de mots.

Depuis le premier système commercialisé en 1970, il y eu de gros progrès. La reconnaissance de mots isolé (petit vocabulaire) fonctionne bien et ce pour du multilocuteurs et dans des conditions estimées difficiles, c'est à dire en environnement bruité. En ce qui concerne la reconnaissance de grands vocabulaires (plusieurs milliers de mots isolés) en monolocuteurs et dans le cadre de domaines plus ou moins spécifiques des systèmes comme Dragon par IBM ou Speech System en 1991 sont assez performants, utilisant une modélisation stochastique de la parole. En 1991 la reconnaissance de la parole continue est utilisé pour des tâches très restreintes avec un vocabulaire et un des énoncés très limités (par exemple saisie de chiffres ou réservation de billets d'avion).

En 1991 comme 20 ans après, le dialogue homme-machine est encore trop difficile à mettre en place mais serait pourtant un enjeux important, en effet la perception et la cognition sont un enjeu fort de l'intelligence artificielle à laquelle s'intègre une partie du TAL et donc de la RAP. Il faut savoir que tous les grands groupes industriels et technologiques sont très axés sur le TAP qui leur propose, comme nous venons de le voir, des capacités supplémentaires qui viennent seconder ou assister l'être humain dans son travail ou ses tâches quotidiennes.

La grande difficulté de la RAP vient des spécificités de l'objet d'étude qu'est le signal vocal. La parole humaine est très complexe dans sa variation, en-

tre locuteurs, situations d'énonciation, signal bruité ou non, langue parlée, état émotionnel du locuteur. . . les paramètres sont extrêmement nombreux et sont susceptible de varier au cours d'un même énoncé. La segmentation du signal en unités du langage (phonèmes, mots, groupes de sens) par l'homme se fait naturellement grâce à un processus linguistique. Ce traitement par la machine apparaît de manière très complexe dans le sens où de nombreux paramètres liés à la structure du langage (acoustico-phonétiques, syntaxique, morphologique, pragmatique) doivent être pris en comptes et se présentent comme « imbriqués » les uns aux autres.

Les taux de reconnaissance dépendent fortement du locuteur, le caractère bruité du signal joue également un rôle puis des paramètres comme la tâche (taille du vocabulaire (domaine restreint, mots isolés, . . .), cependant le degrés selon lesquels ils influent l'un ou l'autre sont difficilement discernables.

1.1.4 RAP & outils

Les outils développés en RAP ont beaucoup évolués depuis le tout premier produit, un système de reconnaissance vocale de mots isolés, apparut dans les années 70. En effet les technologies sont de plus en plus abouties et de nombreux verrous technologiques ont été résolus permettant de nouvelles capacités et des objectifs prometteurs bien que ce domaine de la reconnaissance automatique de la parole ne soit pas le domaine le plus avancé en traitement automatique du langage.

Le système commercialisé le plus connu de reconnaissance automatique de la parole est sans doute Dragon développé par IBM. Il s'agit d'un système de reconnaissance vocale applicable sur micro-ordinateur ou téléphone mobiles comme Siri sur l'iphone.

Dans le domaine des SVI on trouve en France la société Vecsys depuis 1979 qui créa par exemple le serveur de la SNCF [1]. Un autre SVI très utilisé et assez performant en France est celui d'Orange.

On retrouve également la RAP utilisée pour le domaine de la formation. En effet, la RAP permet entre autre ici de vérifier la bonne prononciation des langues secondes. On voit sur le marché plusieurs logiciels d'apprentissage des langues proposer cette technologie : TELL ME MORE, Talk to me par Druide Informatique et beaucoup de plateformes de e-learning adoptent cet outil.

De nombreux systèmes de dictée vocale ont vu le jour, technologie assez bien maîtrisée, la plupart des grands noms de l'industrie de l'informatique fournissent ce service intégrés à d'autres systèmes ou en logiciels indépendants. Il existe par exemple ViaVoice par IBM, DragonSystem développé par Nuance et commercialisé chez IBM ou encore DictaLink de Mysoft. Des solutions

gratuites existent aussi sur internet comme TakIIITypeIt.

Pour le cas de l'aide aux personnes fragiles qui nous intéresse tout particulièrement, plusieurs outils existent depuis un certain nombre d'années. Sont créés des outils d'aide au contrôle de l'environnement en domotique, d'aide à la communication en grande majorité mais depuis quelques années on cherche également à créer des outils d'aide à la sécurité pour fournir plus d'autonomie aux personnes. On retrouve comme outils existant le système Nemo de Protéor [8] qui propose pour les personnes qui n'ont pas leur entière mobilité (handicapées ou âgées) un outil de contrôle de l'environnement avec un système de commande vocale avec auto apprentissage, Li1 par Ubiquiet [9] qui est destiné particulièrement aux séniors pour leur assurer « sécurité, lien social et télé-médecine ». Li1 dispose également d'un système de reconnaissance vocal cependant les informations collectées à ce sujet ne sont pas beaucoup plus précises.

1.2 La voix âgée

« La vieillesse est une période inévitable et naturelle de la vie humaine caractérisée par une baisse des fonctions physiques, la perte du rôle social joué comme adulte, des changements dans l'apparence physique et un acheminement graduel vers une diminution des capacités »
(B. R. Mishara, R. G. Riegel, *Le vieillissement*, Presses Universitaires de France, Paris, 1984).

Définit par le CHU de Jussieu, « le vieillissement correspond à l'ensemble des processus physiologiques et psychologiques qui modifient la structure et les fonctions de l'organisme à partir de l'âge mûr ». Ces processus sont liés à des facteurs génétiques et environnementaux agissants tout au long de la vie de l'individu. Les conséquences du vieillissement se doivent d'être discernées des manifestations de maladies qui s'additionnent bien souvent à l'avancée dans le grand âge [2]. Comme il est dit dans cet ouvrage, le vieillissement se traduit par une diminution des capacités fonctionnelles de l'organisme. Le vieillissement, indépendant de l'âge, présente une extrême variabilité entre les individus, pour un individu au cours du temps et en fonction des organes dont les capacités fonctionnelles sont modifiées.

Avec l'avancée dans l'âge, des troubles normaux du langage surviennent. Nous ne traiterons dans cette étude uniquement les cas non-pathologiques. Le vieillissement provoque une certaine dégénérescence des particularités de la voix. Cela s'explique par le fait que les conséquences de l'âge touchent chacune des parties du corps.

Le phénomène de vieillissement a plusieurs effets sur les caractéristiques de la voix. La voix devient plus grave et tremblotante, la prosodie se modifie, le rythme du débit de parole ralentit et l'articulation est plus lente, moins précise à cause de l'état de fatigue musculaire et de la dentition. En revanche, on ne note apparemment pas de changements sur le vocabulaire ou la syntaxe. S'ajoute à la parole la communication non verbale. Chez les personnes âgées, la raideur du corps, les problèmes de vue et des sens en général, ont un impact non négligeable sur la communication. Pour illustrer, l'audition et la vision permettent de déterminer les distances, or étant altérées, le sujet a des difficultés à appréhender son environnement et cela peut être source d'angoisse.

1.2.1 L'évolution de la voix

La référence majeure pour notre étude est celle menée par Vipperla [Vipperla et al., 2010] qui a analysé dans une étude longitudinale l'impacte du vieillissement sur la parole et ses caractéristiques.

De la même manière que les caractéristiques de la voix évoluent incontestablement durant l'enfance jusqu'à l'âge adulte, il en est de même tout au long de la vie. La période adulte de 20 à 65 ans est sans doute la période où la voix évolue de manière moins marquée, mais le phénomène du vieillissement provoque une certaine modifications de la voix.

En premier lieu, l'apparition de tremblements et la perte musculaire générale réduit la puissance et le contrôle de beaucoup de gestes. La compétence locutoire en fait partie, et le contrôle des cordes vocales, de la langue, des lèvres, du voisement ou autres fonctionnalités du langage en sont perturbés.

La parole est plus lente car la personne met plus de temps à structurer sa pensée et construire ses productions, il y a plus de pauses ou d'hésitations, et qui durent plus longtemps que chez un sujet non âgés.

D'autre part, le vieillissement provoque une certaine dégradation des fonctions cognitives [Gay, 2009a], [Gay, 2009b]. On observe des difficultés sur la prise de décision et les capacités de la mémoire de travail et à court terme, primordiales pour la structure fondamentale du langage, se trouvent diminués. Le discours en est souvent plus ambigu et moins cohérent. Cela se remarque par le fait que la personne âgée emploie des mots vagues qui sont représentatifs de problèmes d'accès au lexique.

De manière globale, on perçoit des difficultés à maintenir le volume de la voix et une intensité constante. Le discours est approximatif, au delà d'un certain seuil, il devient vite un effort physique et intellectuel qui provoque de la fatigue. Il faut noter cependant que la dégradation des fonctions (physiques

ou cognitives) n'est pas seulement une conséquence de l'âge mais plutôt du degré de dépendance de la personne.

1.2.2 Modifications morphologiques et physiologiques

Le corps est en évolution constante et avec l'âge le métabolisme devient plus faible, de nombreux facteurs internes et externes le rendent plus fragile jusqu'à atteindre une dépendance progressive lorsque le temps s'écoule année après année.

Dans notre étude sur la parole âgée, un point important était à faire afin de montrer en quoi la morphologie et la physiologie d'un individu avaient un rôle dans l'évolution de ses productions langagières.

La majorité des organes du corps ont un rôle dans le langage, ceux pour la respiration, ceux pour l'alimentation, ceux du mouvement, etc. Le langage n'émane pas uniquement de la bouche ou du visage, chaque partie du corps s'exprime dès que le locuteur veut signifier quelque chose - nous parlons ici tant du verbal que du non verbal.

Au cours du vieillissement, une perte musculaire et osseuse réduisent l'organisme. On observe une courbure de la colonne vertébrale et de 1,2 à 5 cm de réduction de sa taille [5], les poumons sont tassés, ou par exemple on peut moins lever la tête pour établir un contact visuel avec l'interlocuteur.

Aussi, les deux organes principaux dans la production de la parole, le larynx et les cordes vocales, s'atrophient [3], [Vipperla et al., 2010]. De par la perte musculaire, on observe un relâchement des muscles et des tissus induisant un contrôle moins précis des articulations et du contrôle des cordes vocales, ce qui a des répercussions inévitables sur la parole.

Les poumons ont un rôle important dans la production de la parole car c'est l'air expulsé des poumons qui va permettre, en passant par le larynx et les cordes vocales, de générer la voix. Le vieillissement entraîne une diminution des capacités respiratoires donc moins de puissance pour la production de la parole qui devient proportionnellement un effort [des Enseignants de Gériatrie, 2000].

1.2.3 Modifications acoustiques

De nombreux facteurs provoquent des changements dans la production des sons de la langue, même si le sujet ne subit pas de maladies particulières comme Parkinson ou Alzheimer. A cause des changements morpho-physiologiques, en particulier de l'augmentation de la sécrétion de salive, du port de prothèse dentaire ou de dents manquantes, les phonèmes ne sont plus

articulés et prononcés de la même manière [6].

La presbyacousie entraîne des problèmes avec les consonnes, en particulier f, s, *f*, t, p, en perception. On remarque également des difficultés avec les sons aigus ce qui entraîne des confusions de mots. Au début de la détérioration des sens, le cerveau compense les sons manquants grâce au sens de la phrase et au contexte. Mais il a de plus en plus de mal en vieillissant ou si le problème n'est pas traité rapidement. La lecture labiale joue un rôle important pour la communication avec les personnes âgées, cela serait à prendre en compte dans le cadre des habitats intelligent. En effet, de par leur diminution des performances d'un ou plusieurs des sens, ils pallient leurs handicaps par d'autres moyens. Par exemple, beaucoup de personnes âgées souffrant de troubles de l'audition n'entendent pas lorsqu'on leur parle si elles ne voient pas leur locuteur parler.

Ainsi certaines fréquences et certains phonèmes en particulier poseront d'avantage de problèmes que d'autres dans le sens où, sur le plan visuel il s'agit d'un phonème non labial (donc visuel) ou s'il s'agit d'un phonème dont les fréquences correspondent aux fréquences que la personne âgée ne peut plus très bien voire plus du tout percevoir.

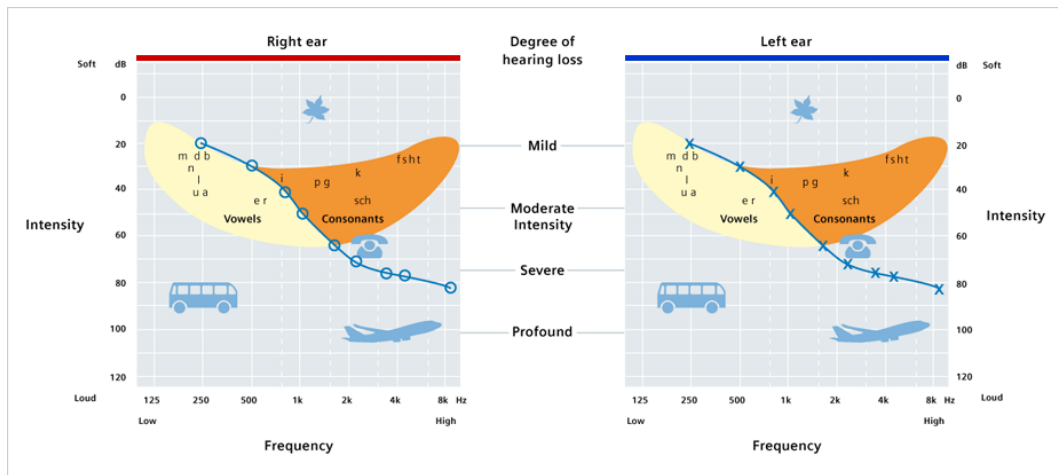


FIGURE 1.3 – Audiogramme proposé par Siemens Audios Solutions (<http://w1.hearing.siemens.com/fr/>)

Les phonèmes ont des spécificités acoustiques (fréquence, intensité) qui induit une perception différente. Ainsi certains phonèmes seront plus robustes face à la perte d'audition. Un indice qui peut montrer les différences inter-phonèmes est leur position sur un audiogramme.

Au regard de ces deux audiogrammes on peut voir que les phonèmes n'ont pas la même évolution face à la perte d'audition, en effet certains sont touchés

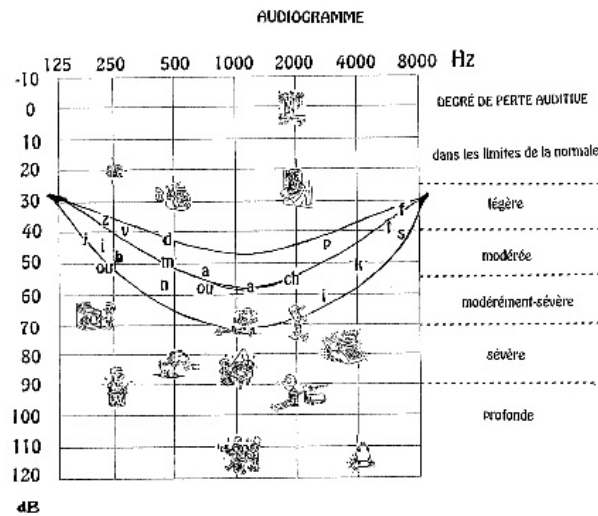


FIGURE 1.4 – Audiogramme proposé sur le site [bruit&societe.ca](http://www.bruit&societe.ca/fr-ca/thematique_cat.aspx?catid=2&scatid=13) (http://www.bruit&societe.ca/fr-ca/thematique_cat.aspx?catid=2&scatid=13)

plus tôt tandis que d'autre vont mieux résister au phénomène. On peut en déduire que les phonèmes, dans la boucle perception-production, ne soient pas tous équitable et que la perte des capacités à percevoir un phonème va dépendre de ses caractéristiques acoustiques.

Si l'on compare les deux audiogrammes, on voit que les théories ne sont pas encore sûres dans le sens où l'ordre des phonèmes est assez différent. Il est en effet très difficile de réaliser un tableau de référence sur les phonèmes affectés en premier par la surdité car il existe une extrême variabilité acoustiques pour un même phonème pourtant perçut comme étant identique par les locuteurs d'une langue donnée.

La majorité des paramètres de la parole sont modifiés par l'âge [Xue and Delisky, 2001]. En effet, à cause des modifications physiologiques et morphologiques, c'est toute la parole en son essence qui est bouleversée .

1.2.4 Modifications cognitives et psychologiques

Au cours du vieillissement, le cerveau perd 12% de sa masse [4], ce qui produit inévitablement un affaiblissement de la motricité, des fonctions d'apprentissage et de mémoire, de la structuration des idées, un ralentissement du fonctionnement cognitif, une altération des sensation (vue et ouïe en particulier) On relève une baisse de la rapidité mentale avec pour principale conséquence un ralentissement des réactions aux stimuli (intellectuels, physiques

ou psychiques), une perturbation de la disponibilité du lexique caractérisé par un manque de mot en particulier les noms de personnes. Les mots sont présents mais c'est leur accès qui est plus long.

Les personnes âgées ont une moindre capacité à organiser des stratégies, donc plus de difficulté à traiter l'information et à avoir de l'attention. Des troubles de la compréhension sont dus à la perte d'audition, à une diminution des capacités attentionnelles, et à une capacité limitée de la mémoire qui permet de suivre un raisonnement ou d'opérer des déductions.

Le facteur majeur qui va entraîner dans une grande majorité des cas la « dégénérescence » des capacités de l'individu est l'isolement. En effet, on remarque que dès lors qu'une personne n'est plus insérée dans la société (travail, famille ou tout autre rôle social) alors son corps vieillit plus rapidement et se laisse mourir [Hervy, 2003]. En outre des maladies liées à l'avancée dans l'âge, la psychologie c'est à dire l'état de stress ou l'anxiété, la valorisation qu'un individu a de lui-même, les rôles sociaux dans lesquels il se projette vont déterminer son maintien en activité.

« L'être humain est un peu comme un tabouret à trois pieds, bio-psycho-social. Tout au long de la vie du tabouret et tout au long de la vie de l'homme, selon l'usage, l'usure ou les accidents, tel ou tel pied aura besoin de réparation, de renfort, d'attention. Aucun pied ne peut être négligé : s'il manque un pied à un tabouret à trois pieds, il tombe. » [7]

1.2.5 Impact sur productions acoustiques

Comme nous l'avons vu dans la section 1.2.1 un très grand nombre de phénomènes peuvent entraîner la dégradation de la capacité à « bien parler ». L'impact de ces phénomènes sur la production des phonèmes s'expliquera en partie par le problème de la baisse de l'audition survenant avec l'âge. En effet, les locuteurs âgés subissant une perte d'audition ne s'entendent pas très bien parler donc ne peuvent pas s'autocorriger en cas de mauvaise prononciation de phonèmes. Ils considèrent donc leur articulation correcte alors qu'elle ne l'est pas. Dans le cadre de la RAP, les concepteurs doivent trouver un moyen de normaliser ces phénomènes et les intégrer aux systèmes en modélisant des règles propres à la parole âgée. C'est en adaptant les modèles acoustiques et les dictionnaires de prononciations des mots de la langue que les résultats obtenus par les systèmes de RAP pourront être plus prometteurs.

1.2.6 Langage & compétences

Du point de vue linguistique, il est indubitable que la capacité que nous avons à parler soit intrinsèquement liée à la pratique et l'usage que nous en faisons. En effet, c'est dans l'interaction avec d'autres sujets parlants que nous entretenons nos capacités locutoires. Dans le cas des personnes âgées, on remarquera que le fait d'être isolé provoque un renfermement du sujet. Aussi, la diminution de ses capacités cognitives dues au vieillissement entraînent une diminution des compétences langagières et de communication. Le fait de ne pas stimuler assez souvent les processus cognitifs a pour conséquences une régression de capacités. Comme nous avons pu le voir plus haut, l'accumulation des changements morphologiques, physiologiques, sociaux et cognitifs nuisent dans la majeure partie des situations aux capacités des individus. Ces phénomènes posent un réel verrou technologique pour la RAP de la voix âgée.

1.3 La RAP des voix âgées

1.3.1 Dégradations des performances

Si la RAP des voix non âgées fonctionne plutôt bien c'est parcequ'il y a beaucoup de corpus disponible de voix non âgée, elle est étudiée depuis de nombreuses années (40 ans dans le domaine de la RAP). C'est pourquoi beaucoup de problèmes auquel est confronté la parole âgée aujourd'hui ont déjà été surmontés pour la parole non âgée. Cependant il faut savoir que même pour la RAP de la parole non âgée, les performances ne sont pas parfaites et peuvent encore être améliorées.

Les conséquences du vieillissement sur les performances des systèmes de RAP sont inévitable. En effet l'état de l'art nous montre que chaque organe permettant la parole est modifié (de manière néfaste pour les systèmes) suivant divers degrés avec l'avancée dans l'âge et que cette évolution est extrêmement aléatoire d'un locuteur à l'autre et n'est en aucun cas corrélée à l'avancée dans l'âge. On relève dans la littérature une dégradation des performances des systèmes de l'ordre de 20% dès qu'il s'agit de parole âgée. Cette dégradation s'explique par le fait que les systèmes sont créés pour une RAP de voix générale soit testés -a priori- sur des voix-types d'utilisateurs lambda. Or comme peu de travaux montrent que la parole des personnes de plus de 70 ans est très différente de celle des locuteurs non âgés, le phénomène n'a donc pas été anticipé. La dégradation des performances viendrait simplement du fait que les concepteurs n'aient pas adapté les systèmes de RAP aux

spécificités acoustiques propres à la parole âgée. Cette chute de 20% des performances prouve qu'il est, pour le moment, nécessaire dans la conception d'un outil de RAP de le destiner soit à des locuteurs non âgés soit à des locuteurs âgés.

Un des phénomènes typiques de la parole âgée mais encore peu démontrée vis à vis de la RAP est la présence de bruits de bouches. En effet, comme il est dit dans l'état de l'art, les personnes âgées ont par exemple plus de salive dans la bouche ce qui entraîne des bruits de type aspiration, une dentition qui peut être défaillante voire remplacée par un dentier qui provoque parfois des claquements durant la parole, les muscles des cordes vocales sont relâchés ce qui engendre un voisement non contrôlé ou un sifflement durant la parole. Les systèmes de reconnaissance automatique de la parole ne sont pas adaptés à ces types de bruits qui se superposent aux phonèmes et gêne la reconnaissance des traits acoustiques.

Le problème qui subsiste est que le thème de la parole âgée est encore assez peu traité. En France très peu d'études ont été trouvées sur le sujet, mais c'est moins le cas pour les Etats-Unis où des études ont analysé l'évolution de la voix avec l'âge telle que l'étude du Vipperla mais sans pour autant appliquer cela au TAP. Si de telles études existaient, il serait plus facile de savoir où et comment adapter les systèmes pour qu'ils fonctionnent pour de la voix âgée.

La dégradation des performances des systèmes de RAP peut se combler uniquement grâce à des études très approfondies sur les aspects linguistiques, morphologiques, cognitifs et articulatoires liés à la parole des séniors. Malheureusement, les données nécessaires pour l'amélioration du système de RAP spécifique à la voix âgée, comme nous en aurions besoin dans le cadre du projet CIRDO, sont inexistantes. A cause de ces conditions, et comme nous ne pouvons être guidés par aucune étude, nous voyons bien que la tâche d'amélioration d'un système de RAP pour obtenir des résultats est exclusive.

Peu de corpus de parole existants sont spécifiques à la parole âgée, les travaux similaires que nous avons dans notre état de l'art sont seulement ceux de Vipperla [Vipperla et al., 2010]. En ce qui concerne d'autres travaux qui pourraient être assimilés nous retrouvons :

- le corpus « Parole d'Alzheimer » constitué et décrit par Melissa Barkat-Defradas et Hye Ran Lee de l'Université de Montpellier. Il s'agit d'un corpus de parole pathologique (Alzheimer) de personnes âgées ;
- l'Université de Lyon recueille, dans le cadre de l'une de leur formations en neuropsychologie, l'enregistrement de personnes âgées mais les pro-

- tocoles précis rendent les enregistrements inutilisables pour la RAP ;
- les pratiques et usages de la langue française parlée sont souvent étudiés dans de grandes études [Cappeau and Gadet, 2007] le plus souvent en sociolinguistique. On peut citer par exemple DGLFLF ou ESLO, mais ces corpus comportent de la parole de tout âges différents et ils ne sont pas destinés spécifiquement à l'étude de la parole âgée ;

Le fait que notre étude ait un caractère inédit ne veut pas dire qu'elle est impossible, au contraire, il s'agit de poser des hypothèses et tenter de répondre aux problèmes que pose la voix âgée en utilisant les données et outils adéquats. Pour pouvoir améliorer le système de RAP choisis dans le projet CIRDO nous allons donc devoir constituer un corpus adapté à nos conditions et réaliser chaque traitements (élaboration du protocole du corpus, enregistrements, transcriptions, segmentation audio, ...) ce qui sera très coûteux en temps et entièrement novateur vis à vis du domaine.

Par ailleurs, l'objectif final de notre projet de recherche est la détection de situation de détresse. Il n'a encore jamais été testé, avec notre système de RAP, une reconnaissance de situations de détresse réelle et encore moins lorsqu'il s'agit d'une personne âgée. Dans une situation de détresse le locuteur peut avoir de nombreux comportements différents et variables. Par exemple une personne dans une situation de ce type pourra être submergée par les émotions et cela provoquera sûrement une dégradation dans la qualité de la reconnaissance.

Chapitre 2

Contexte de l'étude

2.1 L'équipe GETALP

Le laboratoire d'accueil pour ce stage était le Laboratoire d'Informatique de Grenoble (LIG) et plus particulièrement l'équipe GETALP (Groupe d'Étude en Traduction / Traitement Automatisé des Langues et de la Parole).

Le laboratoire LIG effectue des recherches concernant l'ensemble des domaines de l'informatique. L'équipe au sein de laquelle j'ai effectué ce stage est donc l'équipe GETALP. Elle se concentre essentiellement sur six thèmes de recherche :



FIGURE 2.1 – Laboratoire du LIG, bat.B

- Thème 1 : Traduction Automatique (TA) et Automatisée (TAO),
- Thème 2 : Traitement Automatique des Langues (TALN) et plates-formes associées,
- Thème 3 : Collecte et construction de ressources linguistiques,
- Thème 4 : Multilinguisme dans les systèmes d'information,
- Thème 5 : Reconnaissance automatique de la parole, des locuteurs, des sons et des dialectes,

- Thème 6 : Analyse sonore et interaction dans les environnements perceptifs.

Les activités principales abordées à travers ces six thèmes de recherche sont :

- rendre l’informatique multilingue et « ubilingue »,
- informatiser les langues peu dotées et peu écrites en adaptant des ressources existantes,
- rendre la communication langagière multimodale (texte, parole, geste),
- trouver et implémenter des méthodes et outils d’évaluation liés à la tâche,
- utiliser l’interaction contributive pour collecter des ressources, améliorer des traductions et communiquer avec « sans garantie ».

2.2 Le projet CIRDO



FIGURE 2.2 – Logo du projet CIRDO

Dans le cadre de la politique actuelle en France, une restructuration des types d’hébergements spécialisés pour les personnes âgées est mise en œuvre. L’Insee prévoit qu’une personne sur trois aura plus de 60 ans en 2050. C’est donc un enjeu primordial que de trouver une alternative aux infrastructures existantes. Les maisons de retraites et foyers logements sont en nombre insuffisant et n’auront pas la capacité d’accueillir tous les demandeurs d’ici quelques années. Assurer les besoins, l’aide, les soins nécessaires et la sécurité de toutes les personnes âgées qui en ont besoin va s’avérer impossible d’autant plus que ces services sont de moins en moins accessibles financièrement. Il ne faut pas oublier non plus que beaucoup de personnes âgées ne souhaitent pas intégrer ces établissements mais rester, souvent seules, à leur propre domicile. L’évolution de la société met en danger ces personnes isolées de tout. Dans ce contexte plusieurs projets ont été mis en œuvres pour rendre concrètes de nouvelles solutions. L’Assistance à la Vie Autonome (AVA) appuie l’idée de favoriser le maintien à domicile grâce au développement de la

télé-assistance ou des habitats intelligents pour la santé, apportant confort, sécurité, relations avec l'extérieur.

Par ailleurs, en ce qui concerne l'interaction entre le résidant et l'habitat intelligent, l'interface vocale est un moyen plus naturel et surtout plus facilement utilisable par opposition à des interfaces tactiles lorsque la personne à par exemple des difficultés à se mouvoir[Vacher et al., 2006].

C'est dans ce contexte que la DGCIS (Direction Générale de la Compétitivité, de l'Industrie et des Services) et l'ANSP (Agence Nationale des Services à la Personne) ont fait un appel à projet intitulé « Services à la personne : Innover pour développer l'offre de services » en 2008. Le projet CIRDO-Formation (Compagnon Intelligent Réagissant au Doigt et à l'Œil) proposé par un consortium mené par la société Technosens a été retenue lors de cet appel. Ce projet qui s'est déroulé entre 2008 et 2010 a permis à des professionnels du secteur de la télésanté et des services à la personne, à des industriels et à des laboratoires de recherche d'étudier et de proposer un cahier des charges pour une solution adaptée, autour du système télélien social (E-lio), mettant en œuvre des techniques d'analyse sonore et de traitements de l'image[Debeaux and Vacher, 2010].

Suite aux résultats de ce premier projet, le consortium a proposé une réponse à l'appel à projet TECSAN en 2010. Le projet CIRDO - Recherche Industrielle à été retenu et à commencé officiellement le 1^{er} décembre 2010. Le but de ce projet est de développer un produit de télélien social augmenté et automatisé par l'intégration de services innovants (reconnaissance automatique de la parole, analyse de situations (scènes) dans un environnement complexe non contrôlé) visant à favoriser l'autonomie et la prise en charge par les aidants, des patients atteints de maladies chroniques ou de la maladie Alzheimer ou apparentées. Les travaux de l'équipe GETALP ont surtout pour l'instant concerné l'étude de la reconnaissance automatique de la parole des personnes âgées [Vacher et al., 2012][Ama, 2012b][Ama, 2012a][LeG, 2012].

2.3 E-lio

E-lio¹ est l'outil de lien social qui va être le point central du projet CIRDO. Il s'agit d'un système de communication en visiophonie adapté à la perte d'autonomie des personnes âgées. Il peut être utilisé en établissement ou à domicile. Ce dispositif permet d'établir un lien entre les établissements, les familles et les résidents.

Grâce à une utilisation simplifiée et adapté aux personnes âgées, E-lio se révèle plutôt pertinent de la part des concepteurs vis à vis de l'usage et

1. www.technosens.fr/index.php

des attentes dans le domaine du maintiens du liens social. A l'aide d'une caméra fixée sur la télévision ou l'ordinateur de la personne et une télécommande/téléphone composée uniquement de trois boutons pour naviguer et sélectionner sur la plateforme e-lio.

De leur côté, les familles et les professionnels peuvent envoyer des données ou entrer en contact avec les personnes âgées et ainsi avoir un moyen de communication peut-être plus adapté à leur relation et plus directe et simplifiée.



FIGURE 2.3 – Dispositif e-lio

Chapitre 3

Problématique

3.1 Répondre à un besoin

Les objectifs pour ce projet sont nombreux et il en est de même pour le travail de recherche que nous menons au LIG. Lorsque l'on s'affaire à un tel projet, il est primordial de se demander pourquoi ce projet et en quoi il nécessite un travail de recherche afin de cibler précisément les objectifs à atteindre. Dans notre cas, nous cherchons à répondre à la nécessité de pouvoir proposer une alternative aux personnes vieillissantes qui leur permettent :

- d'avoir une option abordable financièrement entre l'isolement à leur domicile et l'entrée en infrastructures spécialisées ;
- de proposer une solution qui les respecte dans leur intégrité physique et morale liée à leur âge ;
- de trouver un moyen justifié sur le plan ergonomique qui permette de consolider ou même tout simplement créer les liens sociaux primordiaux pour leur santé et leur sécurité.

La crainte de l'accident domestique est une source de stress très importante qu'il faudra prendre en compte pour répondre aux attentes. C'est cette crainte qui motive les utilisateurs à rechercher une alternative qui leur convienne.

3.2 Améliorer les performances

Le problème qui se pose aujourd'hui est que les systèmes de reconnaissance de parole (continue ou mots isolés) fonctionnent de manière acceptable mais que les taux d'erreurs augmentent de manière significative dès lors qu'il s'agit de sujets de plus de 65 ans [Dugheanu, 2011]. Hormis l'étude de

Vipperla [Vipperla et al., 2010], peu de travaux ont été menées sur la parole des personnes âgées dans le cadre du traitement automatique de la parole, or il semble que le facteur de l'âge ait des conséquences non négligeables à prendre en compte dans les traitements. Des mesures physiques ont été effectuées mais aucune étude (à grande échelle) à l'heure actuelle ne permet de caractériser précisément les différentes données de la parole des personnes âgées ni d'expliquer à quoi est due précisément cette évolution provoquant un fort taux d'erreur de la part des systèmes. Pour pouvoir expliquer et tenter de corriger ces aspects, il est nécessaire d'étudier de manière plus précise les caractéristiques et des productions des locuteurs et apporter des modifications nécessaires pour adapter les outils.

3.2.1 Objectif de ce mémoire

Nous allons tenter de cerner les éléments qui posent problèmes pour la qualité de la reconnaissance de la parole des personnes âgées et ainsi essayer de permettre l'optimisation des résultats. Ce travail se fera par la mise en place d'une série de tests qui a pour objet de mettre en évidence le fonctionnement du système de reconnaissance automatique de la parole et les paramètres sur lesquels il faut jouer pour obtenir des résultats satisfaisants. L'état de l'art a montré que les modifications sur la parole des personnes âgées se situent majoritairement au niveau phonétique (hauteur, timbre) et que l'impact du vieillissement sur les niveaux morphologiques, lexicaux et syntaxiques est moindre. Cela a pour conséquences de dégrader grandement l'efficacité des systèmes de RAP.

Notre objectif est d'analyser les performances d'un système de RAP sur la voix âgée et l'impact des adaptations successives des modèles acoustiques et des modèles de langage sur ces performances.

Nous donnerons priorité aux aspects relevant des modèles acoustiques étant donné que ce niveau apparaît comme prépondérant dans l'évolution de la parole des séniors. Pour cela nous nous appuierons sur différents corpus dont l'un comporte de la parole spontanée exploitable et comparable à une situation de maintiens à domicile.

Cela ouvre l'opportunité de comparer le comportement d'un système de reconnaissance automatique de la parole, et principalement celui des modèles acoustiques, entre la parole lue et la parole spontanée, plus sujette aux phénomènes de réduction, d'assimilation, d'interruption, soit plus imprévisible et irrégulière.

3.3 Attentes, hypothèses

Nous avons plusieurs attentes et hypothèses induites des travaux précédents, certaines sont parfois venue au cours du traitement.

Tout d'abord, en ce qui concerne les scores d'alignement, les travaux effectués dans [Dugheanu, 2011] [Lefol, 2010] ont montré que certains phonèmes posent problème dans la parole âgée. Le corpus ERES38 ayant été constitué dans cet objectif dans [Sasa and Grand, 2011]. Les tests d'alignement forcé devraient confirmer que les plosives et fricatives (p-t-k-b-d-g-f-s-ʒ) ont un score d'alignement moins bons que les autres phonèmes, consonnes et voyelles confondues, pour la parole âgée et montrer que l'adaptation du modèle acoustique permet d'améliorer ces scores.

Selon cette étude, les résultats du score d'alignement étaient utiles (entre autre) pour comparer la voix âgée à la voix non âgée.

Ensuite, comme les tests vont être réalisé à plus grande échelle (plus de locuteur et plus d'occurrences) on va pouvoir généraliser les conclusions des études phonémiques qui ont précédées ce travail et évaluer le bien-fondé d'une adaptation à partir d'un corpus de développement de voix âgées ainsi que le rôle d'une adaptation au locuteur classique. L'impact en terme de WER de ces adaptations devrait ressortir des tests effectués et montrer qu'une adaptation du modèle de langage le plus spécifique présente les meilleurs résultats.

Chapitre 4

Méthodologie

4.1 Sphinx3

L'outil de RAP utilisé dans le cadre du projet CIRDO est le système Syphinx3 développé par The Carnegie Mellon University de Pittsburg¹. Créé en 2001, Sphinx3 est très connu dans le domaine de la RAP.

Cet outil a été sélectionné parmi les différents systèmes possibles pour les raisons suivantes :

- il était déjà utilisé dans l'équipe donc une certaine maîtrise, non négligeable, du fonctionnement de l'outil s'est avérée être un atout pour ce projet ;
- il présente de bonnes performances en comparaison aux autres systèmes existants ;
- la paramétrisation du système permet de tester des aspects précis dans les corpus ;
- il existe une « communauté » ce qui permet de pouvoir entrer en contact facilement en cas de difficultés particulières en relation avec l'outil, ce qui n'existe pas dans tous les cas ;
- c'est un système ouvert ce qui, en plus des trois avantages précédents, a fait pencher la balance.

En plus de cela, Sphinx est de manière générale robuste et facile d'utilisation et propose plusieurs scripts déjà prêts pour effectuer divers traitements [Ameur, 2011].

1. cmusphinx.sourceforge.net

4.2 Données

4.2.1 Corpus AD80 disponible au LIG

La majorité des corpus ont été enregistré à partir de 2009 spécialement pour ce projet. D'autres corpus étaient disponibles au LIG mais aucun ne correspondait suffisamment à la tâche, seul un corpus de parole non âgée enregistré en 2004, décrit dans la section suivante a été réutilisé pour servir de référence de parole non âgée pour ce projet.

En traitement automatique de la parole, il est primordiale, du point de vue linguistique mais également pour une optimisation totale des performances de l'outil, que les corpus de travail (apprentissage comme test) soient adaptés à la tâche.

Des corpus disponibles et adaptés au projet CIRDO pour le développement d'un outil de reconnaissance vocale adapté à la parole âgée étaient inexistant, c'est pourquoi le travail des chercheurs au début du projet à vite débouché sur une constitution de corpus approprié, donc recueillir de la parole de personne de plus de 65 ans. Les corpus de parole âgée sont très rares à l'heure actuelle, le domaine d'application de la RAP aux personnes de plus de 65 ans est très peu développé, ce qui explique une certaine difficulté à trouver des ressources.

Le corpus AD80 a été enregistré sur plusieurs étapes afin d'évaluer dans un premier temps le comportement d'un système de RAP « classique » (soit non adapté à de la parole âgée). A ce stade, le corpus enregistré en 2010 reprenait une partie d'un corpus de voix non âgées enregistré en 2004 à Grenoble par le LIG pour une étude portant sur la reconnaissance d'appels de détresse, cela permettait d'avoir une référence de parole non âgées qui donnerait un élément de comparaison justifié du point de vue de la tâche.

Anodin-Détresse

Le corpus Anodin Détresse a été enregistré en deux étapes, la première constituant un corpus de locuteurs non âgés et la seconde un corpus de locuteurs âgés.

Le corpus AD non âgé a été enregistré en 2004 au laboratoire CLIPS de Grenoble par l'équipe GEOD [Vacher et al., 2006]. Ce corpus est constitué de 2646 phrases lues (soit 126 par locuteur) issues de la vie quotidienne pour 66 phrases sur 126 et de situations de détresse pour les autres. Pour donner des exemples « J'ai bu ma tisane », « Où est le sel », « Bonjour », « La porte est ouverte! » ou « Au secours », « Un médecin vite », « A l'aide! » « Ap-

pelez le SAMU! » pour ce qui est des phrases correspondant aux situations de détresse. 21 locuteurs ont été enregistrés, 10 femmes et 11 hommes âgés de 20 à 65 ans. La totalité de ce corpus fait 38min d'enregistrements.

Le corpus AD âgé a été enregistré en 2010 [Ama, 2012a] sur le même protocole que le corpus non âgé. 36 locuteurs ont été enregistrés, 25 femmes et 11 hommes âgés de 62 à 94 ans.

Voice-Age

La partie enregistrée en 2009 (Voice Age) comportait 7 locuteurs âgés enregistrés au CHU de Grenoble (2 pers) et à leur domicile (5 pers) en Normandie, âgés de 70 à 89 ans. Les conditions d'enregistrement et le protocole se sont avérés inadaptés pour pouvoir poursuivre les enregistrements, ce qui explique le nombre réduit de locuteurs [Aynaud, 2009] [Lefol, 2010]. Il s'agit du même type de phrases lues que pour Anodin-Détresse soit 5441 phrases. Ce corpus fait au total 4h8min d'enregistrements.

Corpus en cours d'enregistrement

Le corpus AD80 est toujours en cours d'enregistrement, l'équipe vise par ce corpus à obtenir des enregistrements de 40 locuteurs âgés et 40 locuteurs non âgés. Une partie complémentaire à AD80 a été apporté au projet en juin 2012 comportant 29 locuteurs enregistrés dans une maison de retraite dans le Gard. Cette partie comporte 8 hommes et 21 femmes âgées de 62 à 94 ans.

4.2.2 Corpus enregistré ERES38

Le corpus ERES38 a été enregistré en 2011 par le LIG [Sasa and Grand, 2011]. Il est composé de 24 locuteurs 16 femmes et 8 hommes âgés de 68 à 98 ans. Ce corpus est constitué de deux parties (l'une étant imbriquée dans l'autre), une partie de récit de vie au sein de laquelle une tâche de lecture a été demandée au locuteur, seulement deux sujets ont refusé de la réaliser à cause de problèmes de vue. Au total le corpus fait 17h43min.

Il faut noter que ce corpus a été constitué (plus précisément la tâche de lecture) à partir de l'état d'avancement du projet CIRDO afin d'apporter des données qui permettrait de tester ou répondre aux hypothèses restées jusqu'alors sans réponses.

Dans [Dugheanu, 2011] on trouve que certains phonèmes et/ou classes de phonèmes pourraient poser plus de problèmes que d'autres c'est à dire que leur score d'alignement phonémique serait particulièrement plus mauvais que

les autres phonèmes. Il a donc été proposé un texte à faire lire aux personnes enregistrées dans lequel ces phonèmes très précisément seraient inclus une ou plusieurs fois et dans des contextes différents (début, milieu, fin de mot et en contexte vocalique extrême [aiu]).

| | AD80 | | | | ERES38 | |
|----------------------|--------------------------------------|--------------------------------------|--------------------------------------|---|-----------------------|-----------------------|
| | AD non âgé | AD âgé 1 | AD âgé 2 | VA | lecture | spontané |
| nb. loc | 11H 10F | 11H 25F | 21F 8H | 4F 3H | 14H 8H | 5F 2H |
| âges min/max | 20 65 | 62 94 | 62 94 | 70 98 | 62 98 | 67 98 |
| année enreg. | 2004 | 2010 | 2012 | 2009 | 2011 | 2011 |
| lieu enreg. | labo. LIG | Isère à domicile | Gard en résidence | CHU Grenoble, domicile Normandie | Isère en résidence | Isère en résidence |
| durée totale | 38min | 1h 25min | env. 1h 40min | env. 2h | 48min | env. 15min |
| nb. phrases | 2646 | 4536 | 3654 | 5441 | 14 | estimation 105 |
| vocab. en nb.mots | 299 | 299 | 299 | env. 300 | 179 | env. 300 |
| utili- sation | test, comparaisons âgé/non âgé | test, comparaisons âgé/non âgé | test, comparaisons âgé/non âgé | test, comparaisons âgé/non âgé | adaptation, test | adaptation test |

FIGURE 4.1 – Tableau récapitulatif des corpus

4.3 Tâches réalisées et leur protocole

4.3.1 Alignement forcé phonémique

A partir des études de Rémus Dugheanu et Frédéric Aman dont ce mémoire est la suite, nous avons choisi de réutiliser comme dans [Dugheanu, 2011] l'alignement forcé qui permet d'effectuer une étude phonémique précise du corpus.

Dans le logiciel Sphinx3, un script permet d'effectuer cette tâche. Ce script prend en entrées un fichier audio (ou plusieurs), le fichier texte des transcriptions de références correspondantes ainsi qu'un modèle acoustique. Ce script impose un inconvénient, en effet, il est impossible d'accéder aux paramètres modifiés au cours du traitement, il agit comme une « boîte noire » donc nous n'avons pas pu mesurer les modifications des coefficients lors des différentes adaptations. Le processus est décrit en section 5.1.

Le processus d'alignement phonémique (ou alignement forcé) est différent du décodage acoustico-phonétique. Cette tâche consiste à aligner la suite de phonèmes cible sur celle de référence avec pour chaque phonème des moyennes de paramètres acquises par apprentissage qui constituent le modèle acoustique. Le décodage acoustico-phonétique est une étape importante d'un système de RAP. Il s'agit de segmenter le flux acoustique en unités phonémiques grâce à un modèle acoustique appris à partir d'un corpus de la langue en question. Contrairement à l'alignement forcé, le décodage acoustico-phonétique ne propose pas de suite de phonèmes de référence mais cherche simplement à reconnaître la suite de phonème cible grâce au modèle acoustique de référence créé suivant le même modèle que pour celui de l'alignement forcé. Les performances d'un système de RAP sont directement liées à la phase de décodage acoustico-phonétique [Fohr et al., 1994], en effet, c'est dans la mesure où les unités phonémiques de la langue sont bien reconnues que les éléments supérieurs tels que les mots puis phrases vont pouvoir l'être.

Nous avons lancé l'alignement forcé sur les 48 locuteurs des corpus Anodin-Détresse, ERES38 et Voice Age. En plus de cela nous avons testé l'alignement forcé sur ERES38 avec trois modèles acoustiques différents :

- un modèle acoustique « générique » pour de la RAP non adapté à de la parole âgée construit à partir de Bref120 [Gauvain et al., 1990] ;
- un modèle adapté à la voix âgée, un modèle acoustique issu de Bref120 auquel est appliquée une adaptation grâce au corpus ERES38 ;
- un modèle adapté au locuteur, réalisé locuteur par locuteur générant un modèle acoustique par locuteur à partir du modèle générique ;

Pour la tâche d'alignement forcé nous avons utilisé seulement un modèle de langage, le modèle restreint aux mots du texte de la lecture.

Grâce à divers scripts disponibles, nous avons calculé les scores d'alignements par phonème par locuteur ce qui permettrait de comparer la reconnaissance d'un phonème par rapport aux autres ou d'étudier l'évolution d'un phonème suivant des paramètres tels que les modèles acoustiques. Ensuite nous avons calculé des moyennes par phonèmes et regroupés en classes phonémiques. Cette étude effectuée par classe de phonèmes permettra de

faire ressortir des comportements spécifiques en fonction de traits acoustico-phonétiques et établir des éventuels liens avec le vieillissement. Les résultats sont présentés dans le chapitre 6.

| Consonnes | | | | | | | |
|----------------------|------------------|--------------------|------------------|--------------------|-----------------|-------------------|---------------|
| classes phonémiques | plosives sourdes | fricatives sourdes | plosives sonores | fricatives sonores | nasales nasales | liquides liquides | glides glides |
| phonèmes | p | f | b | v | m | l | ɥ |
| par | t | s | d | z | n | ʁ | j |
| classes ² | k | ʃ | g | ʒ | | | w |

| Voyelles | | | | |
|----------------------|---------|----------|---------|-----------|
| classes phonémiques | fermées | ouvertes | nasales | centrales |
| | fermées | ouvertes | nasales | moyennes |
| phonèmes | i | a | ã | ø-œ-øœ |
| par | u | ɑ | ẽ-õ | e-eɛ-ɛ-ə |
| classes ³ | y | aɑ | õ | o-ɔ-oɔ |

FIGURE 4.2 – Tableau présentant les classes de phonèmes

4.3.2 Décodage en mots

La tâche de décodage en mot est une tâche qui relève de l'aspect lexico-syntaxique sur le plan linguistique. Grâce à un corpus d'apprentissage un modèle de langage sera créé et modélisera ainsi un modèle statistique des suites de mots du corpus. En fonction de la probabilité des successions des mots dans le corpus d'apprentissage, le décodage va sélectionner sur le modèle des HMMs la suite de trois mots qui a le plus de probabilité d'apparaître dans cet ordre. Le résultat de ce traitement est calculé en pourcentage de mot bien reconnus.

Le WER est la mesure la plus connue et la plus utilisée pour évaluer un système de RAP. Pour pouvoir comparer à pied d'égalité nos résultats avec les études en RAP et surtout avec les études du projet CIRDO et sur la reconnaissance de la voix âgée, nous avons utilisé le WER.

Le logiciel Sphinx3 propose des scripts qui permettent de réaliser ce traitement.

Nous avons dans un premier temps adapté le modèle de langage à notre corpus et nous avons utilisé les mêmes modèles acoustiques que pour l'alignement forcé.

A coté de cela nous avons créé trois modèles de langages distincts :

- Lectures « pures », qui modélise uniquement les mots présents dans le texte de la lecture.
- Lectures « réelles », qui modélise tous les mots prononcés par les locuteurs durant la lecture, par exemple des mots tels que ben, donc, alors heu, ah, ...
- Générique, qui modélise un vocabulaire très large adapté à partir des corpus ESTER et PFC en plus des mots modélisés dans les lectures réelles.

4.3.3 Tests sur de la parole spontanée

Une fois l'ensemble des tests ci-dessus réalisés, nous avons tenté d'effectuer les mêmes tâches mais sur de la parole spontanée. A notre connaissance, aucun test de reconnaissance de parole spontanée n'aurait été fait pour de la parole âgée dans le cadre de l'assistance à la vie autonome. L'un des axes majeur du projet étant la détection d'appels de détresse, la parole lue (sur laquelle la majorité des outils de RAP sont développés) ne semble en aucun cas représentative de situation réelle où la voix sera modifiée par des émotions et non contrôlée et prévisible comme c'est le cas de la parole lue. Nous avons donc voulu observer le comportement de l'outil avec un corpus de parole spontanée (donc différent des corpus utilisé pour les tâches ci-dessus) mais avec des modèles acoustiques identiques (puisqu'il s'agit des mêmes locuteurs).

Les locuteurs que nous avons choisis pour tester la parole spontanée (7 des 24 locuteurs de ERES38) ont été choisis en fonction des transcriptions déjà réalisées au cours de notre travail en 2011 [[Sasa and Grand, 2011](#)]. Nous avons transcrit 50% des récits de vie. La tâche de transcription est faite manuellement et est un travail très coûteux en temps, une minute d'audio à transcrire correspond à minimum 10 minutes de travail. En 2011 nous avons transcrit 60% du corpus, soit le récit de vie de 18 locuteurs. Par manque de temps nous avons sélectionné environ 2 minutes d'enregistrement de récit de vie (parole spontanée), cela correspond environ au temps d'une lecture par locuteur soit un temps de parole total de 11,07 min. Puis nous avons effectué une normalisation des transcriptions car certains caractères étaient incompatibles et certaines notations non homogènes.

4.4 Utilisation des corpus pour cette étude

Chacune des parties de ce corpus a permis de réaliser une tâche plus ou moins différente. En Traitement Automatique de la Parole on distingue une partie du corpus pour la phase d'apprentissage et une autre, plus importante quantitativement, pour les phases de test.

4.4.1 Adaptations

Le corpus ERES38 a été enregistré en 2011 afin d'apporter au projet CIRDO un corpus de travail d'adaptation du système de RAP. Ce corpus a permis de fournir des données de référence pour l'adaptation des modèles acoustiques. Comme on le voit dans la figure 4.3 le traitement de l'adaptation se fait en boucle : on évalue dans un premier temps le comportement du système de RAP sans adaptations sur le corpus AD80, on réalise les adaptations nécessaires et possible à partir d'un corpus distinct (pour cette étude ERES38) puis on reteste le comportement du système afin d'évaluer l'impact des adaptations apportées, etc.

Les différentes adaptations portent sur les modèles acoustiques, les modèles de langage ou encore sur les dictionnaires de prononciations. Ainsi grâce aux enregistrements et leurs transcriptions on peut générer les divers modèles et dictionnaires nécessaires.

4.4.2 Test

Dans le cadre de ce mémoire, le corpus d'adaptation ERES38 a dans un premier temps servi de corpus de tests. Cependant, afin d'estimer l'effet éventuel de surapprentissage entraînant des résultats non représentatifs de la réalité si l'adaptation et les tests portent sur exactement le même corpus, nous avons utilisé le corpus 4.2.1 enregistré au cours de cette étude. Le corpus principal de test pour le projet CIRDO est le corpus AD80. Les figures 4.3 et 4.5 montrent parfaitement ce qui est testé au regard de l'adaptation.

4.5 Adaptations MLLR

4.5.1 Technique utilisée pour l'adaptation des modèles acoustiques

Plusieurs techniques sont connues pour effectuer une adaptation au locuteur. Dans notre cas cette adaptation est justifiée par le fait que la voix des

personnes âgées, à qui est destinée cet outil de RAP, est très différente de celle pour lesquelles les systèmes fonctionnent habituellement. En effet, on observe une dégradation de 20% des résultats avec des systèmes classiques. Une adaptation est donc nécessaire.

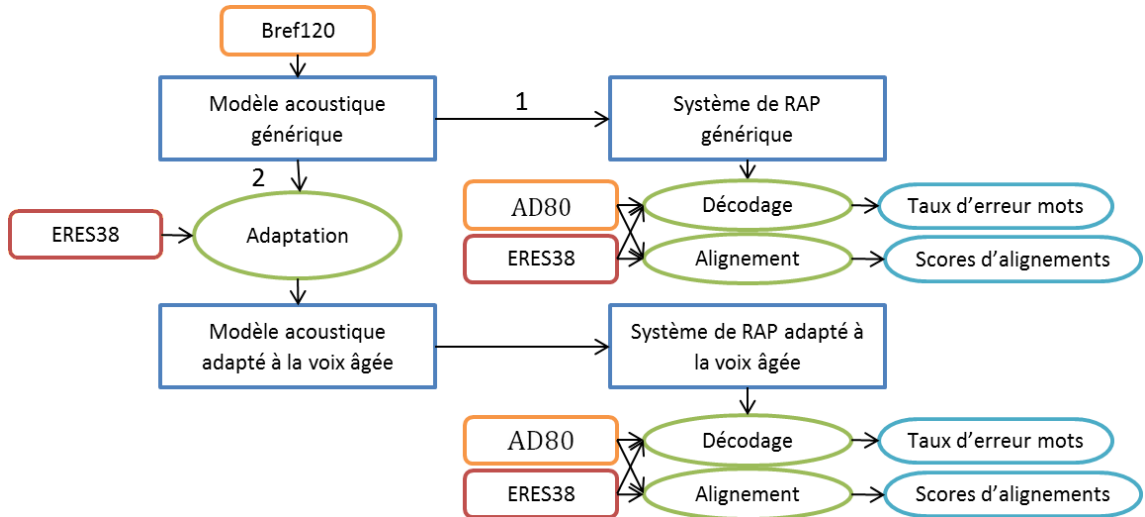


FIGURE 4.3 – Schéma des adaptations apportées en fonction des corpus et des tâches du traitement

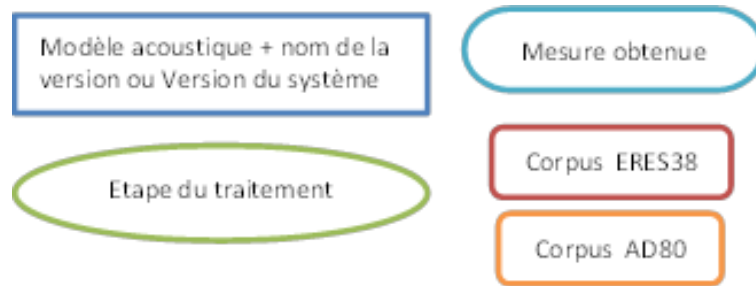


FIGURE 4.4 – Légende des figures 4.3 et 4.5

Technique MLLR

Une technique MLLR (Maximum Likelihood Linear Regression) est une technique qui permet d'adapter les paramètres d'un modèle acoustique d'un système de RAP. Une adaptation par MLLR permet de mieux prendre en compte les caractéristiques acoustiques liées aux spécificités de la voix âgée.

Grâce à cette technique on obtient de meilleurs résultats. En effet, les résultats obtenus sont généralement indépendants du locuteur mais souvent fonction de la qualité du signal et du corpus de développement utilisé qui se rapproche (dans notre cas) d'avantage des conditions de test.

La technique MLLR est un script qui va transformer et adapter de manière automatique les paramètres du modèle acoustique choisi, réduisant ainsi les différences entre les données d'apprentissage et les conditions de test [Tissier, 2011].

Le principe général consiste à « estimer la transformation W qui modélise les différences supposées linéaires entre les conditions d'apprentissage et les conditions de test. Le système adapté est obtenu en appliquant cette transformation linéaire W à un ensemble de paramètres des modèles issus de l'apprentissage. »

4.5.2 Technique d'adaptation à un seul locuteur

Dans le domaine de la Reconnaissance Automatique de la Parole, la restriction (du point de vue du vocabulaire) envers le sujet, thème ou contexte de l'énonciation permet d'optimiser la qualité du traitement. De la même manière que pour une adaptation à un groupe de voix âgées, nous avons réalisé une adaptation du système à l'unique voix du locuteur.

En modélisant les paramètres acoustique pour un seul locuteur on peut générer un modèle acoustique qui lui est propre. Ainsi le système de RAP est directement adapté à la production phonémique de ce locuteur.

L'adaptation du système peut se faire, non pas sur un corpus de développement regroupant plusieurs locuteurs, mais sur des données de paroles d'un seul locuteur. L'adaptation se fait ainsi à l'unique voix du locuteur.

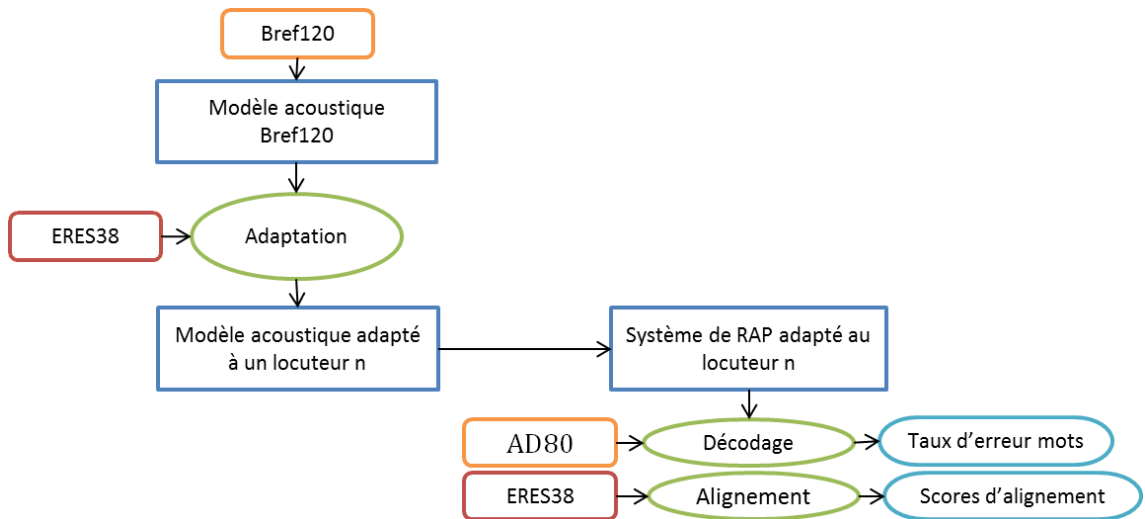


FIGURE 4.5 – Schéma des adaptations apportées pour l’adaptation au locuteur en fonction des corpus et des tâches du traitement

Chapitre 5

Evaluation

Pour évaluer notre système, estimer précisément son comportement dans le cadre de la voix âgée et afin d’avoir des éléments de comparaison pour les différentes tâches réalisées, nous avons choisi deux métriques. Ces deux métriques sont une référence pour l’évaluation des systèmes de RAP, il s’agit du « Score d’alignement » et du « Word Error Rate ».

L’intérêt de ces deux métriques est qu’elles permettent chacune d’évaluer un aspect linguistique différent. En effet, le score d’alignement est axée sur la dimension acoustico-phonétique tandis que le taux d’erreur mot (ou word error rate) traite de la dimension lexicale. Ces deux métriques et leur fonctionnement sont décrites dans les deux sections suivantes.

5.1 Le Score d’alignement

La première métrique utilisée pour l’évaluation de l’outil sur lequel nous avons choisi de travailler est le *Score d’alignement*. Il s’agit d’appliquer au corpus la technique de l’alignement forcé sur les phonèmes et d’en extraire les scores d’alignement.

Un premier travail avait permis grâce à cette technique d’évaluer les phonèmes les plus problématiques de la parole âgée. Le système d’alignement forcé utilisé de Sphinx prend en entrée les fichiers audio ainsi que les transcriptions textuelles de référence, il convertit ces transcriptions en suites de phonèmes et propose la décomposition phonémique du segment audio la plus proche de celle de référence. A partir de ce traitement, on calcule le score acoustique. Autrement dit, on associe pour chaque couple de fragment (source, cible) un score qui reflète la probabilité que la source soit produite comme la cible [Lardilleux et al., 2011].

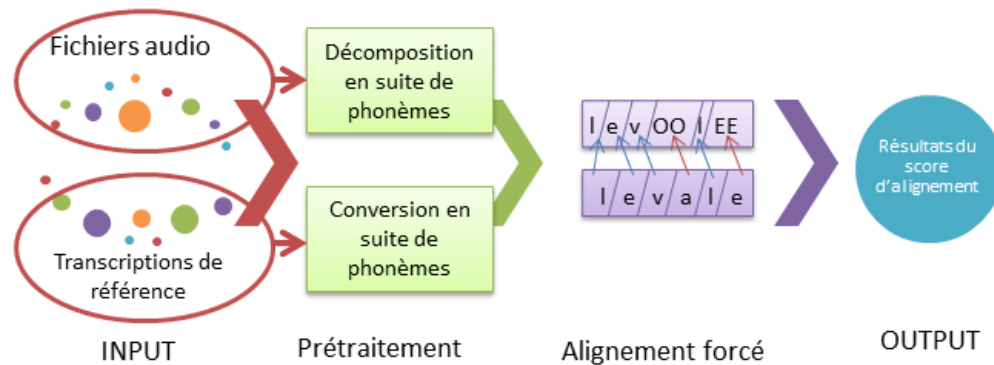


FIGURE 5.1 – Schéma du fonctionnement de l’alignement phonémique et du système d’obtention du score d’alignement

Le score d’alignement ainsi obtenu est normalisé sur le nombre de trames pour prendre en compte la durée du phonème. On obtient par locuteur et pour chaque phonème le nombre d’occurrences ainsi que sa moyenne logarithmique des scores d’alignement obtenus. Malheureusement, dans l’état de l’art, nous n’avons trouvé aucun travail utilisant cette métrique de façon similaire, les seules comparaisons possibles ont été réalisées au sein de ce projet à travers les différentes études effectuées.

5.2 Le Taux d’Erreur Mot (TEM)

La seconde métrique que nous avons utilisé est le Taux d’Erreur Mot (Word Error Rate ou WER). Il s’agit de la mesure de référence dans le domaine de la Reconnaissance Automatique de la Parole.

Le WER permet de mesurer le taux d’erreur des mots, autrement dit il mesure le nombre de mots qui diffèrent entre l’hypothèse et la référence [Evermann, 1999]. Puisque la « longueur » de l’hypothèse peut être différente de celle de référence, le WER est calculé par rapport à ces deux séquences de mots.

On distingue quatre situations possibles :

- Correcte : l’hypothèse et la référence sont identiques (ou équivalents selon un ensemble de règles) ;
- Substitution : le mot de référence est aligné sur un mot hypothèse différent ;
- Insertion : des mots supplémentaires par rapport à la référence sont insérés et l’alignement ne peut être effectué sur la référence ;

- Suppression : un mot de la référence ne figure pas dans l'hypothèse ;

La figure suivante présente le calcul réalisé :

$$WER = \frac{S + D + I}{N}$$

FIGURE 5.2 – expression du calcul du WER

où S = nombre de substitutions, D = nombre de suppressions, I = nombre d'insertions et N = nombre de mots de référence.

L'alignement optimal est calculé grâce à une procédure de programmation dynamique qui va minimiser la distance de Levenshtein (somme pondérée des erreurs d'insertion, de suppression et de substitution) des deux chaînes de mots.

Le taux d'erreurs mots est ainsi donné, en pourcentage, comme le nombre d'erreurs sur le nombre de mots de référence. En raison du nombre d'insertion possible, le WER peut être supérieur à 100%.

Pour se donner une idée globale des WER moyens par domaine nous avons lu que le WER pour des textes lus de type dictée vocale, système monolocuteur était de 5%, 10% pour de la RAP sur des journaux radio et TV et jusqu'à 40% pour des conversations téléphoniques informelles¹.

1. source. Wikipedia

Chapitre 6

Résultats

6.1 La parole lue

6.1.1 Résultats des scores d'alignement

Deux traitements étaient intéressants du point de vue du score d'alignement, le traitement par phonème qui permet de voir le comportement des phonèmes suivant les modifications que l'on effectue sur le système et d'un autre côté de tester le comportement des phonèmes par classes faisant ressortir des tendances liées aux spécificités des classes de phonèmes.

Par phonème seul

Grâce aux résultats des scores d'alignement, nous allons pouvoir analyser quel modèle acoustique est le plus adapté au type de tâche. Suivant notre protocole nous avons obtenu un score de -16472,6 (tous phonèmes confondus) pour le modèle acoustique générique, de -9246,44 pour le modèle adapté à la voix âgée et de -7815,02 pour le modèle adapté au locuteur. Autrement dit ces adaptations nous apportent une amélioration de près de 43% lorsque l'on adapte le modèle générique à la voix âgée et jusqu'à 52% lorsqu'on l'adapte à un seul locuteur. Le traitement par phonème seul montre parfaitement comment il évolue en fonction des modèles acoustiques. Le graphique 6.1¹ montre les résultats obtenus pour le phonème [b]. Vous trouverez l'ensemble des phonèmes en annexe 8.6.

Le graphique 6.1 montre les scores d'alignement pour les différents locuteurs en fonction des modèles acoustiques utilisés, comparés à la « référence » pour de la parole non âgée. On obtient grâce à ces graphiques une représentation visuelle très explicite de l'état de la reconnaissance des phonèmes de la parole

1. En annexe 8.5 une version plus grande est disponible.

âgée.

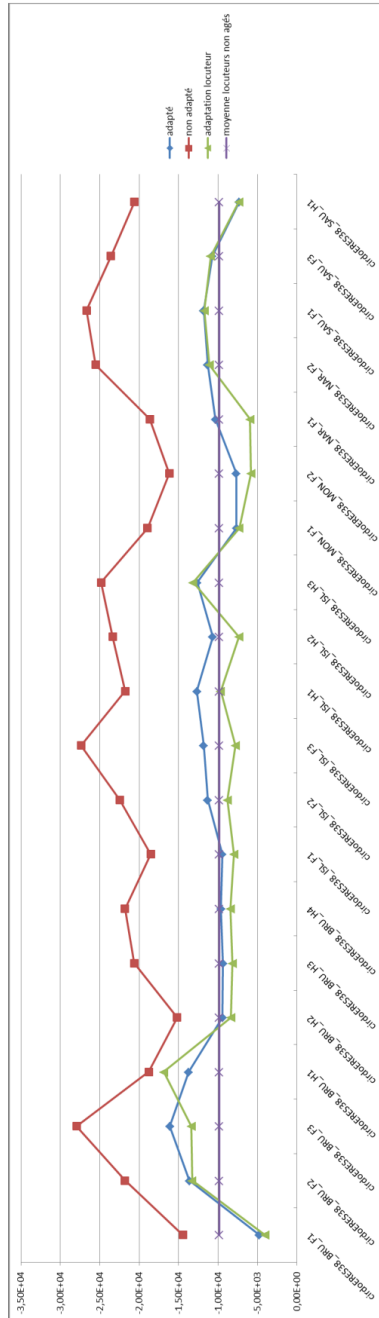


FIGURE 6.1 – Courbes présentant les scores d’alignement forcé en fonction des locuteurs et des modèles acoustiques pour la parole

L'une des remarques importante à faire est l'extrême variabilité d'un locuteur à l'autre. De plus, un même locuteur ne se situera pas toujours de la même manière par rapport aux autres : pour un phonème, ses scores seront très bons tandis que pour un autre phonème le score de ce même locuteur sera très mauvais. Par exemple, en regardant l'annexe 8.6, les phonèmes [2] (qui correspond au symbole [ø] en API, cf 8.2) et [a] pour le locuteur ISL_H1 pour le modèle non adapté : le score du phonème [2] est de -8535,87 ce qui est l'un des meilleurs scores pour le phonème en revanche, pour le phonème [a] le score, toujours pour le même modèle acoustique est de -31099,98 ce qui est vraiment très élevé.

Les études précédentes sur le projet annonçaient que certains phonèmes étaient moins bien reconnus que d'autres en particulier [p-t-k-b-d-g-f-s-ʒ] le tableau 6.1 présente le classement des phonème par moyenne de score d'alignement par phonèmes et par modèles acoustiques générique et adapté à la voix âgée. Les phonèmes sont classés par moins bon taux de reconnaissance et ceux grisés sont les phonèmes mis en avant comme étant les moins biens reconnus.

Ce tableau montre que pour le cas du corpus ERES38, les phonèmes annoncés (grisés ou roses dans la figure 6.2) présentent des résultats similaires en ce qui concerne les plosives mais que le cas des fricatives soit moins vérifié. Nous n'avons pas testé d'autres corpus, le raisonnement est donc à poursuivre car il existe une différence probable entre les corpus.

| Phonème | moyenne adaptée voix âgée | | Phonème | moyenne non adaptée |
|---------|------------------------------|--|---------|------------------------|
| p | -14741,34501 | | p | -27344,75892 |
| k | -12339,21832 | | k | -26236,86558 |
| un | -12089,43268 | | t | -23601,00657 |
| g | -11803,91763 | | oOO | -23345,72638 |
| h | -11712,6102 | | b | -21453,49883 |
| 2 | -10865,33221 | | a | -19630 |
| t | -10761,01985 | | un | -19355,4374 |
| b | -10615,77434 | | d | -19028,79718 |
| d | -10372,23969 | | g | -18772,99613 |
| z | -10356,9555 | | OO | -18749,19914 |
| RR | -9562,786122 | | an | -18570,7518 |
| aeA | -9502,221854 | | aAA | -18156,13373 |
| eEE | -9362,155649 | | EE | -18018,74068 |
| u | -9270,903587 | | 9 | -17140,3293 |
| aAA | -9229,991539 | | RR | -16003,18117 |
| OO | -9111,813863 | | aeA | -15472,76374 |
| l | -9048,568965 | | z | -15360,51335 |
| NJ | -8888,501352 | | o | -15316,88162 |
| i | -8809,924788 | | eEE | -15248,78656 |
| n | -8637,72662 | | u | -15096,34578 |
| oOO | -8527,49875 | | w | -14425,0942 |
| a | -8477 | | on | -14258,92794 |
| y | -8455,851464 | | y | -13920,85268 |
| f | -8233,597912 | | m | -13655,32125 |
| EE | -8232,910845 | | h | -13426,33893 |
| v | -8209,545887 | | ZZ | -13422,14076 |
| in | -8079,225481 | | HH | -13412,10953 |
| an | -8078,365589 | | in | -13399,27929 |
| w | -8024,22974 | | l | -13358,60949 |
| e | -7853,560907 | | 2 | -13204,05186 |
| ZZ | -7815,925418 | | v | -13151,7277 |
| j | -7771,849234 | | SS | -13034,94102 |
| s | -7735,338578 | | e | -13030,69887 |
| on | -7727,621252 | | s | -12986,30184 |
| HH | -7320,793472 | | i | -12863,21668 |
| m | -7277,776199 | | j | -12798,64265 |
| o | -7135,791331 | | n | -12465,47252 |
| SS | -7118,832561 | | f | -11129,75917 |
| 9 | -5912,722836 | | NJ | -10636,93395 |
| AA | -1106,996667 | | AA | -1140,643628 |
| 29 | 0 | | 29 | -1050,842857 |

FIGURE 6.2 – Tableau présentant les phonèmes les moins bien reconnus

La critique que nous pourrions faire ici est que ces résultats excluent totalement le fait que les phonèmes apparaissent en contexte et sont différents en fonction de ce dernier. Ici on effectue une moyenne sur le phonème seul, cela est donc partiellement représentatif de ce qui se passe réellement dans la langue. On entend par contexte le phonème plus celui qui précède et celui qui suit. Il faut savoir qu'il existe 318 028 contextes possibles dans la langue française soit 43 phonèmes donc 79 507 triphones que Sphinx modélise 4 états qui sont début de mot, milieu de mot, seul et fin de mot. Étudier chaque contexte se révèle donc impossible.

6.1.2 Par classe phonémique

Les résultats par locuteur et par phonèmes permettent d'avoir une étude détaillée mais complexe à analyser. Le choix d'une classification par mode d'articulation des phonèmes s'est imposée pour rendre plus facile la lecture des résultats. Nous avons donc regroupé les phonèmes suivant le tableau 4.3.1.

Afin d'étudier le comportement par classes de phonèmes et pouvoir observer si des tendances pouvaient être relevées comme il était dit dans des études précédentes [[Ama, 2012b](#)] [[Ama, 2012a](#)] [[Dugheanu, 2011](#)] nous avons mis ces résultats sous forme de diagrammes. Les diagrammes 6.1.2 et 6.1.2 pour de la parole lue (vous trouverez le détail des résultats en annexe 8.9) :

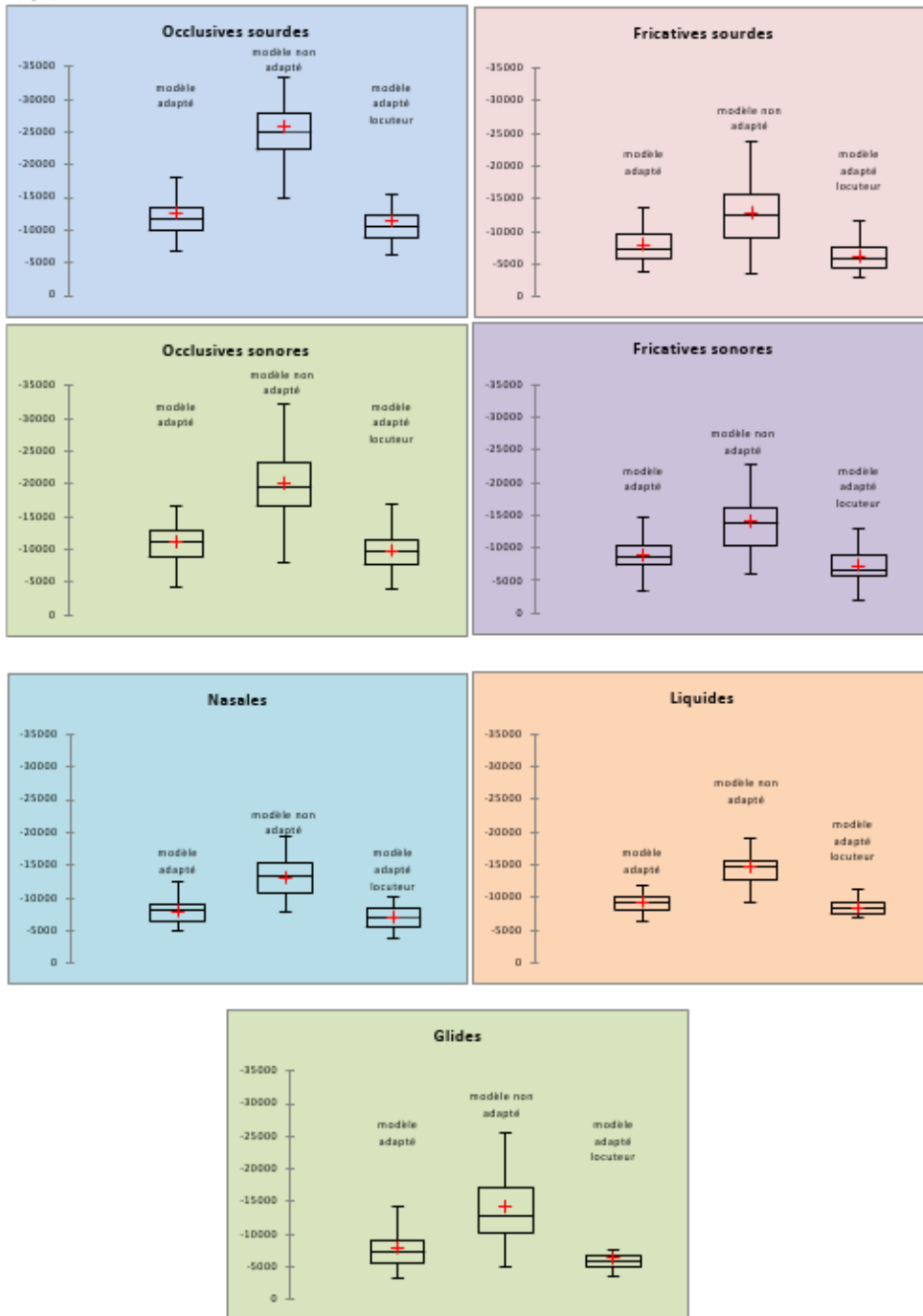


FIGURE 6.3 – Diagrammes des scores de l’alignement forcé pour les consonnes par classe phonémique en fonction des modèles acoustiques pour la parole lue

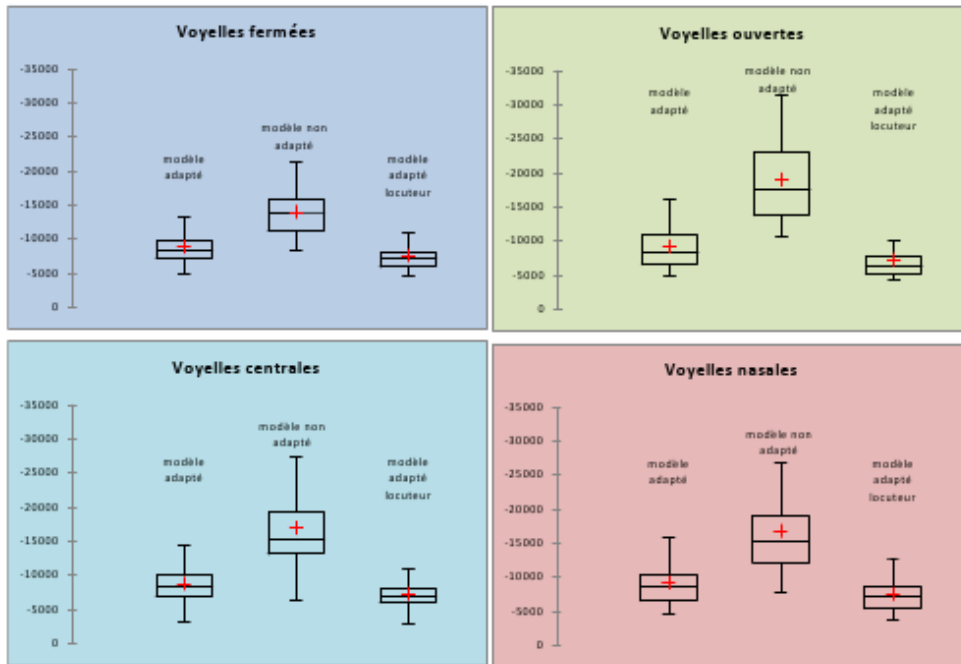


FIGURE 6.4 – Diagrammes des scores de l’alignement forcé pour les voyelles par classe phonémique en fonction des modèles acoustiques pour la parole lue

Les plosives sourdes (score de -16585,2) et sonores (-13662), ainsi que les voyelles nasales (-1120,6), moyennes (-10961,9) et ouvertes (-11793,5) obtiennent les scores d'alignement les plus élevés c'est à dire que ces classes sont, pour le cas de notre corpus, les classes les moins bien reconnues. Ces résultats correspondent parfaitement avec ceux obtenus dans [Dugheanu, 2011] page 35 pour le corpus Voice Age. Nous pouvons donc dire que ces classes de phonèmes seraient moins bien reconnues dès lors qu'il s'agit de voix âgée indépendamment du modèles acoustique.

L'adaptation grâce à un corpus de développement constitué de l'ensembles des voix de locuteurs âgés permet une nette amélioration des scores pour l'ensemble des classes étudiés. L'adaptation au seul locuteur permet à chaque fois d'obtenir les scores les meilleurs.

Ces diagrammes montrent par ailleurs le comportement des différentes classes de phonèmes et les répercussions des différentes adaptations des modèles acoustiques :

- La représentation correspondant au modèle acoustique non adapté à la voix âgée reste le plus irrégulier avec un minimum et un maximum très étendus et très élevés indiquant combien les scores sont mauvais pour ce modèle acoustique ;
- On peut observer que le voisement ne semble pas agir mais que le mode d'articulation reformerait des classes plus larges et plus justifiées. Ainsi on aurait comme classes occlusives, fricatives, nasales-liquides en une seule classe, glides pour les consonnes. Nous ne donnerons pas d'avis sur les voyelles dont les graphiques ne montrent qu'une vague distinction fermée autres voyelles ;
- Bien que la différence soit légère par endroit, le modèle acoustique adapté au locuteur présente toujours les meilleurs résultats ;

6.1.3 Voix âgée vs non âgées

Nous avons comparé les scores d'alignement forcé obtenus sur un corpus de voix âgée et ceux issus d'un corpus de voix non âgée. Cette comparaison nous permettra de situer les résultats obtenus sur le niveau phonémique à ceux obtenus sur un corpus de voix non âgée que l'on avait estimé comme étant satisfaisants.

Le diagramme 6.5 présente les résultats des scores d'alignement par classes phonémiques mettant en apposition les scores obtenus pour la parole âgée de ERES38 et ceux de la parole non âgée de AD80. On voit que de manière générale aucun des deux types de parole n'est mieux reconnu que l'autre. En revanche on peut y voir que l'écart se creuse pour certaines classes

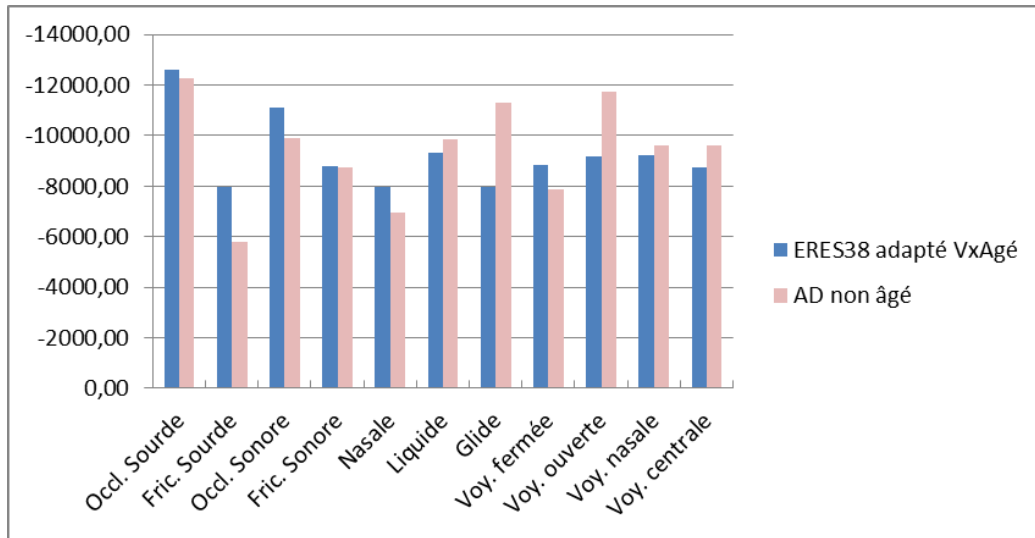


FIGURE 6.5 – Diagramme présentant les résultats de locuteurs âgés vs non âgés par classe de phonèmes

phonémiques :

- Sont à 'égalité' les plosives sonores avec un score d'environ -12400, fricatives sonores avec -8770 et les voyelles nasales avec -9400 ; ce qui indique la prise en compte des spécificités acoustiques de la parole âgée.
- La voix âgée présente de moins bons résultats que la voix non âgée pour les fricatives sourdes avec un score supérieur de 27%, les occlusives sonores +11% de dégradations, les nasales +13% et les voyelles fermées +11% ; ce qui montre que l'adaptation a été insuffisante pour ces phonèmes.
- A l'inverse, les autres classes ont un score d'alignement moins bon du côté de la parole non âgée, cela indique que les différentes adaptations des systèmes ont permis d'obtenir un meilleur taux que pour de la parole non âgée donc que l'adaptation a bien fonctionné.

6.2 Résultats du TEM

6.2.1 Modèle de langage vs modèle acoustique

Nous comparons les WER en fonction des locuteurs, des modèles acoustiques utilisés et de modèles de langages générique et comportant uniquement

| | Modèle de langage 1 | Modèle de langage 2 | Modèle de langage 3 |
|-------------------------|---------------------|---------------------|---------------------|
| Modèle bref20 | 60,40454545 | 50,71363636 | 65,58181818 |
| Modèle adapté voix âgée | 26,76363636 | 18,50909091 | 44,96818182 |
| Modèle adapté locuteur | 21,59090909 | 14,71818182 | 41,42727273 |

FIGURE 6.6 – Tableau présentant les WER par modèle de langage en fonction du modèle acoustique

les transcriptions des séquences de récits de vie utilisées². Le tableau 6.2.1 montre les résultats :

Ce tableau présente les résultats moyens sur l'ensemble des locuteurs. Les modèles de langage ont un impact énorme sur les taux d'erreur-mot. L'adaptation au contenu thématique est primordiale, permettant de passer d'un WER moyen de 50,66% pour un modèle générique à un WER moyen de 32,18% pour un modèle adapté. Cependant, cette unique modification se révèle insuffisante si elle n'est pas accompagnée d'une adaptation du modèle acoustique. En effet, pour un modèle acoustique générique, le taux d'erreur-mot peut passer de 65% à 60% pour une adaptation à la lecture et jusqu'à 50% lorsque les « petits mots » d'hésitation et d'interruption propres à l'oral sont pris en compte dans le modèle de langage.

6.2.2 La question du surapprentissage

Durant notre recherche nous avons réalisé certaines adaptations et tests sur le même corpus. Nous nous sommes alors posé la question d'un éventuel effet de surapprentissage, c'est à dire que les résultats obtenus seraient meilleurs car l'apprentissage permet dans ce cas d'anticiper chaque élément du corpus testé, ce qui n'est pas représentatif de conditions réelles d'utilisation. Pour pouvoir justifier notre démarche et vérifier le degré de surapprentissage sur nos résultats, nous avons testé le décodage en mots avec un corpus de pa-

². Les modèles de langages et acoustiques sont décrits dans la section « Tâches réalisées et leur protocole » 4.3

role âgée de 17 locuteurs issu de AD80 en utilisant les modèles acoustiques adaptés sur ERES38. Les modèles de langages seront les mêmes que pour ERES38 car le corpus est constitué à partir de la même lecture.

Dans le tableau qui suit, on voit que les résultats obtenus sur le corpus AD80 sont bien meilleurs, 29,6% en moyenne contre 42,6 pour ERES38. On peut donc dire que le phénomène de surapprentissage est assez peu influent donc il semble possible de faire une adaptation des système de RAP à partir d'un vaste corpus de parole âgée, lue et spontanée.

| locuteurs | modèle de langage 1 (lectures pures) | | modèle de langage 2 (product ^o des lect) | | modèle de langage 3 (Ester, lectures, PFC) | |
|---------------------|--------------------------------------|--------------|---|--------------|--|--------------|
| | MA bref120 | MA adaptVxAg | MA bref120 | MA adaptVxAg | MA bref120 | MA adaptVxAg |
| LD1 | 32,2 | 11,1 | 48,8 | 10 | 67,8 | 39,1 |
| LD2 | 46,9 | 26,2 | 41,3 | 29,9 | 56,4 | 41,6 |
| LD3 | 28,9 | 11,7 | 21,4 | 10,7 | 51,3 | 34,4 |
| LD4 | 24,4 | 10,2 | 22,4 | 4,6 | 47,9 | 32,7 |
| LD5 | 34,8 | 26,4 | 42,6 | 12,9 | 60,7 | 48,2 |
| LD6 | 34 | 9,3 | 23,8 | 8,5 | 47,7 | 36,6 |
| LD7 | 6,5 | 2,3 | 4,6 | 2 | 27,1 | 26,8 |
| LD8 | 21,3 | 6,6 | 19,3 | 9,6 | 39 | 30,2 |
| LD9 | 36,5 | 16,3 | 32,9 | 9,2 | 69,3 | 37,2 |
| L10 | 32,9 | 6,8 | 27,7 | 5,5 | 49,4 | 32,3 |
| L11 | 64,9 | 30,8 | 38,6 | 23,2 | 73,2 | 47,4 |
| L12 | 25,3 | 7,1 | 18,8 | 2,9 | 45,5 | 28,2 |
| L13 | 13,8 | 3,9 | 9,2 | 2 | 39,7 | 27,9 |
| L14 | 38,8 | 17,4 | 44,3 | 9,8 | 65,6 | 36,7 |
| L16 | 22,4 | 12,8 | 17,6 | 8,9 | 42,5 | 31,6 |
| L17 | 48,9 | 12,2 | 43,7 | 6,4 | 56,9 | 33,1 |
| TOTAL AD80 | 36,72 | 13,21 | 31,19 | 9,01 | 52,50 | 34,94 |
| TOTAL ERES38 | 58,95 | 23,15 | 49,52 | 15,41 | 65,42 | 43,35 |

FIGURE 6.7 – Tableau présentant les résultats du décodage en mots en fonction des modèles acoustiques pour une partie du corpus AD80 et comparaison avec les moyennes des scores du décodage en mots pour le corpus ERES38

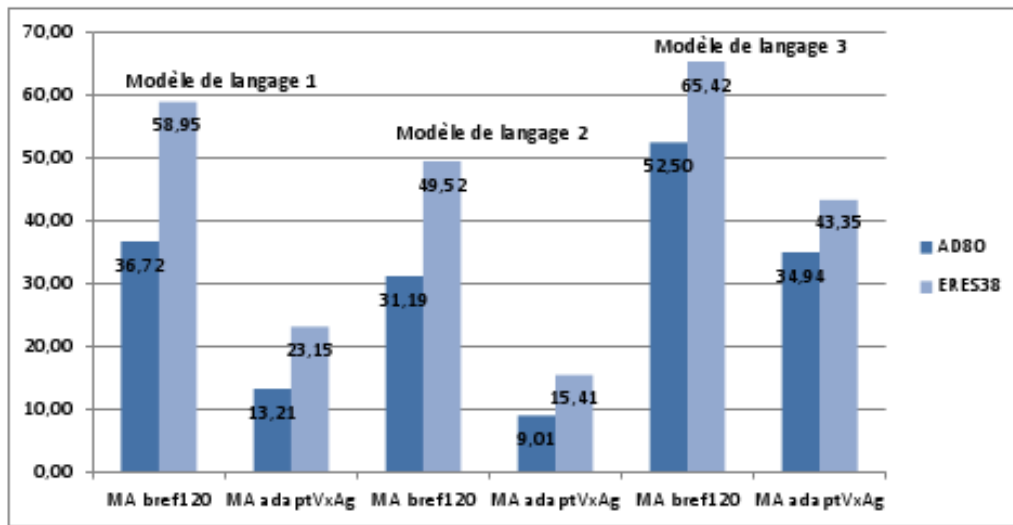


FIGURE 6.8 – Histogramme présentant les résultats du test de généralisation

6.2.3 Voix âgée vs non âgées

A l'heure actuelle, les systèmes de RAP pour de la parole « classique » donc non âgée présentent des taux de WER assez acceptables. Dans notre étude nous avons tenté d'adapter le système afin que les résultats obtenus pour de la voix âgée se rapprochent le plus possible voire soient meilleurs que la référence donnée dans l'état de l'art. Au début de l'étude, dans [Dugheanu, 2011] l'état du décodage mot montrait que le système sphinx3 (avec adaptation du modèle acoustique et du modèle de langage sur le corpus Voice Age) fonctionnait plutôt bien pour de la voix non âgée et présentait, pour la parole âgée, un taux de WER de +20,23% par rapport au taux obtenu en voix non âgée soit un WER de 27,56%. Le taux de 7,33% enregistré en 2011 constitue notre référence est le WER obtenu pour de la parole non âgée issue de AD80. Les tableaux de la section précédente le montrent, dans des conditions identiques (lecture et environnement légèrement bruité) nous avons obtenu pour de la parole âgée - avec un modèle acoustique adapté à celle-ci et un modèle de langage comportant les productions de tous les locuteurs (vocabulaire étendu) - un WER de 15,41% et 13,1% si l'on adapte au locuteur. On passe ainsi de +20,23 à seulement +8,08 points en adaptant le système au type de parole.

6.2.4 Différence homme/femme

Dans la littérature [Vigouroux et al., 2005], nous relevons qu'il existerait une différence homme/femme, liée aux spécificités de la voix âgée, au niveau

du WER soit 10% de plus chez les hommes. Nous avons testé cette hypothèse avec nos résultats obtenus sur la parole lue et spontanée. Le tableau 6.2.4 montre ces résultats.

Nous voyons bien ici qu'il existe une réelle différence entre les hommes et les femmes, cependant nous tenons à notifier que nos corpus, en particulier celui de la parole spontanée ne comporte que peu de locuteurs, trop peu pour pouvoir généraliser. La tendance reste tout de même formelle avec 13,5 (moyenne femmes) vs 20,5 (moyenne homme) soit +7 points relatifs en parole lue sur 22 locuteurs dont huit hommes ; et 17 vs 38,5 en parole spontanée soit +21,5 points pour la moyenne homme.

| WER par locuteur par lieu d'enregistrement | | | | |
|--|------------|---------|------------------|-------|
| locuteur | parole lue | | parole spontanée | |
| BRU_F1 | 11,96 | | — | |
| BRU_F2 | | | 16,5 | |
| BRU_F3 | | | — | |
| BRU_F4 | | | — | |
| BRU_F5 | | | — | |
| BRU_H1 | 19,95 | | — | |
| BRU_H2 | | | 48,9 | |
| BRU_H3 | | | — | |
| BRU_H4 | | | — | |
| ISL_F1 | 12,8 | | 18,8 | |
| ISL_F2 | | | — | |
| ISL_F3 | | | — | |
| ISL_H1 | 15,6 | | — | |
| ISL_H2 | | | 27,7 | |
| ISL_H3 | | | — | |
| MON_F1 | 12,7 | | 11,4 | |
| MON_F2 | | | — | |
| NAR_F1 | 17,3 | | — | |
| NAR_F2 | | | 21,1 | |
| SAU_F1 | 12,7 | | — | |
| SAU_F3 | | | — | |
| SAU_H1 | 26,2 | | — | |
| TOTAL | femme | homme | femme | homme |
| | 13,492 | 20,5833 | 16,95 | 38,3 |

FIGURE 6.9 – Tableau présentant les différences homme/femme pour le WER en parole lue et spontanée

6.2.5 Conclusion

La notion de WER est bien connue dans le domaine de la RAP cependant on remarque que les résultats obtenus dépendent énormément des concepts de modèles de langage et de modèle acoustique. Ainsi, la balance entre les deux va dépendre de la tâche dans laquelle l'outil se projette. Les adaptations des différents modèles vont permettre de trouver un paramétrage optimal avec un WER le meilleur possible à condition d'avoir un corpus répondant à la tâche et qu'il soit de bonne qualité du point de vue acoustique.

6.3 La parole spontanée

6.3.1 Résultats des scores d'alignement

Deux analyses étaient intéressantes du point de vue du score d'alignement, le traitement par phonème qui permet de voir le comportement des phonèmes suivant les modifications que l'on effectue sur le système et d'un autre côté de tester le comportement des phonèmes par classes faisant ressortir des tendances liées aux spécificités des classes de phonèmes. Ces analyses ont été menées pour chaque locuteur indépendamment afin de ne pas effacer l'effet de la variabilité entre locuteurs.

Par phonème seul

De la même manière que pour la parole lue, le modèle acoustique a un impact fort sur les résultats des scores d'alignement forcé pour les phonèmes. De façon générale nous avons obtenu pour la parole spontanée un score phonémique moyen de -16585,26 pour le modèle acoustique générique, -9206,51 pour le modèle adapté à la voix âgée et -7744,31 lorsque l'on adapte le modèle générique à un seul locuteur. Une adaptation portant sur le modèle acoustique a donc permis, dans le cadre de cette étude, d'obtenir une amélioration de la reconnaissance des phonèmes de 44,5% lorsque l'on adapte le système générique à de la parole âgée et jusqu'à +53,3% d'amélioration des résultats lorsque l'on adapte le système générique au locuteur.

Les deux graphiques³ suivants montrent dans le premier cas (graphique 2) la première courbe obtenue pour de la parole spontanée juxtaposée aux courbes de la parole lue ainsi qu'à la courbe moyenne pour le phonème en question en parole lue non âgée.

Et dans le graphique 6.11 on peut observer les scores d'alignement pour la parole spontanée. On y voit très bien que, dans le cas de l'adaptation à la voix âgée les scores se situent aux alentours de la moyenne non âgée réalisée alors qu'il s'agit de parole lue. Cela nous mène à penser que le traitement de la parole spontanée pour de la voix âgée est fortement envisageable pour un domaine précis et avec un corpus adapté.

3. En annexe 8.5 voir le graphique plus gros

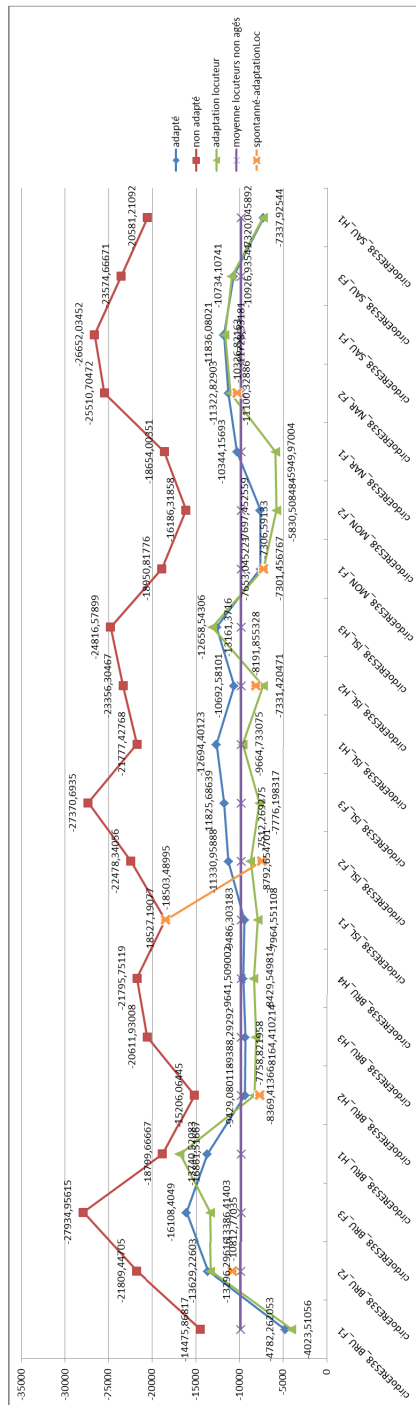


FIGURE 6.10 – Courbes présentant les scores d’alignement en fonction des locuteurs et des modèles acoustiques pour la parole lue vs spontané

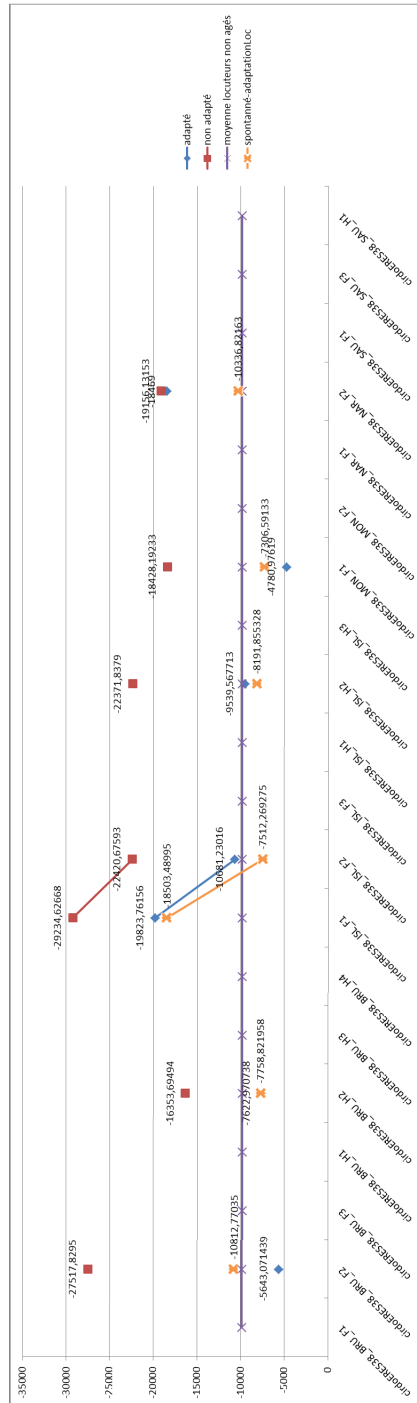


FIGURE 6.11 – Courbes présentant les scores d’alignement en fonction des locuteurs et des modèles acoustiques pour la parole spontanée

6.3.2 Par classe phonémique

Toujours dans l'idée de comparer notre travail avec ceux effectués auparavant, nous avons regroupé les scores phonémiques en classes de phonèmes. Les diagrammes 6.3.2 et 6.3.2montrent les mêmes calculs pour de la parole spontanée (vous trouverez le détail des résultats en annexe 8.10 :

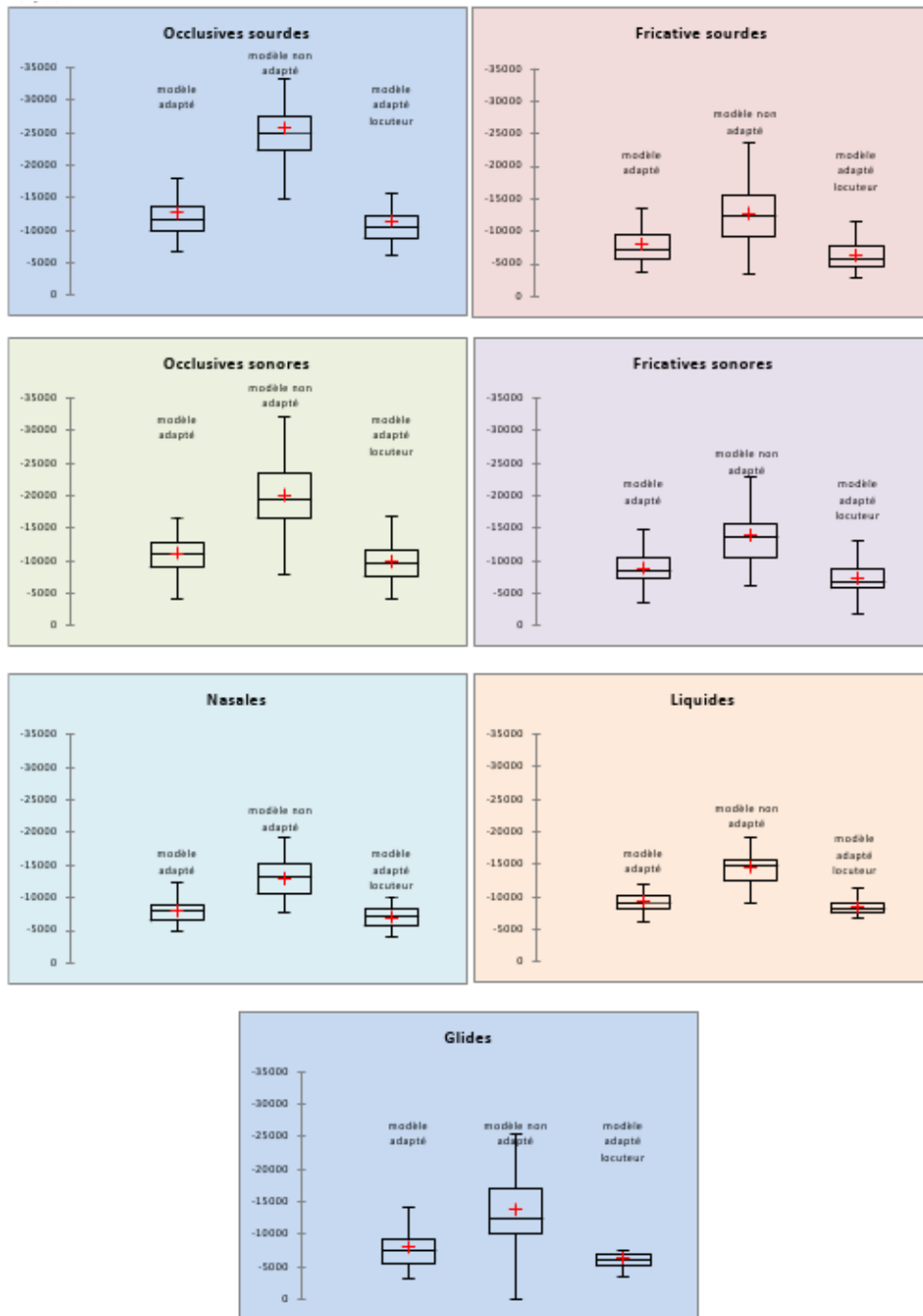


FIGURE 6.12 – Diagrammes des scores de l’alignement forcé pour les consonnes par classe phonémique en fonction des modèles acoustiques pour la parole spontanée

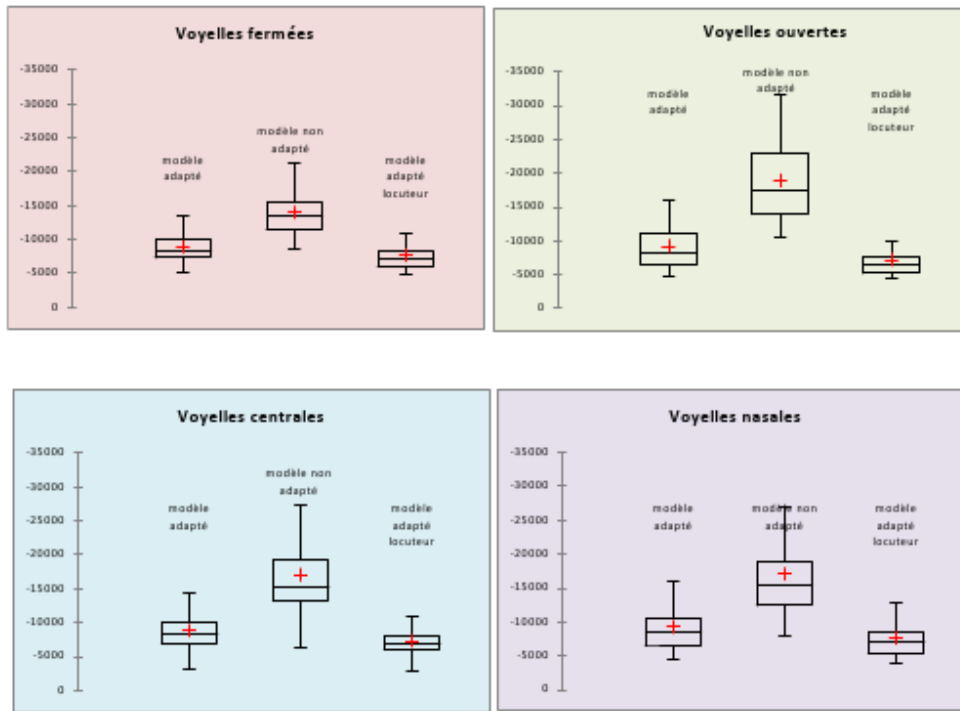


FIGURE 6.13 – Diagrammes des scores de l’alignement forcé pour les voyelles par classe phonémique en fonction des modèles acoustiques pour la parole spontanée

6.4 Résultats du TEM

Afin de mieux mettre en évidence notre propos, nous allons présenter ici les résultats du décodage-mot réalisé sur la parole spontanée en comparaison avec les résultats obtenus pour de la parole lue vue ci-dessus. Nous obtenons, pour un modèle acoustique générique et avec un modèle de langage très restreint, un WER de 55% et pour un modèle acoustique adapté nous obtenons 22% de WER, ce qui est encourageant ; et pour un modèle de langage générique et avec un modèle acoustique générique nous obtenons un WER de 68% tandis que avec un modèle adapté nous obtenons 48% ce qui reste trop élevé.

Le tableau 6.4 présente tous les résultats en WER montrant les différences entre la parole spontanée et la parole lue du corpus ERES38. Les modèles acoustiques sont les mêmes mais les modèles de langages sont propre à chaque tâche.

| locuteurs | modèle de langage 1 (productions pures) | | | modèle de langage 2 (productions+ESTER+PCF) | | |
|---|---|--------------|-------------|---|--------------|-------------|
| | MA bref120 | MA adaptVxAg | MA adaptLoc | MA bref120 | MA adaptVxAg | MA adaptLoc |
| BRU_F2 | 48,2 | 16,5 | 13,2 | 64 | 35,6 | 41,6 |
| BRU_H2 | 73,5 | 48,9 | 33,2 | 81 | 69,5 | 57,8 |
| ISL_F1 | 40,4 | 15,8 | 16,7 | 63,2 | 34,2 | 37,7 |
| ISL_F2 | 59,7 | 21,8 | 28,7 | 72 | 42,2 | 47,2 |
| ISL_H2 | 61,9 | 27,7 | 14,5 | 70,7 | 65,2 | 54,8 |
| MON_F1 | 36,1 | 11,4 | 16,7 | 51,9 | 43,5 | 40,1 |
| NAR_F2 | 68 | 21,1 | 24,5 | 79,6 | 53,1 | 55,8 |
| TOTAL ERES38_spontané | 55,40 | 23,31 | 21,07 | 68,91 | 49,04 | 47,86 |
| TOTAL ERES38_lecture | 58,95 | 23,15 | 18,93 | 65,42 | 43,35 | 40,74 |
| différence en % entre spontané et lecture | -6,40364449 | 0,69444444 | 10,1694915 | 5,07186291 | 11,6127779 | 14,8756219 |

FIGURE 6.14 – tableau présentant les différences entre la parole lue et spontanée

On relève une augmentation de 6 points en moyenne pour la parole spontanée par rapport à la parole lue. Le meilleur taux obtenu pour la parole spontanée est 21,07%, comparé aux 40% ou 28,5% annoncés dans la littérature pour les conversations téléphoniques informelles, le système DEMON [Bousquet-Vernhettes, 2002] ou dans [Dufour, 2010] (références se rapprochant le plus de notre étude).

Chapitre 7

Conclusion & perspectives

7.1 Conclusion

Ce travail ouvre de nombreuses réflexions de par l'aspect novateur du domaine. Très peu d'études existent sur l'étude de la reconnaissance automatique de la parole spécifique à la voix âgée, c'est pourquoi les tests effectués durant ce mémoire suscitent beaucoup d'interrogation et/ou de remise en cause d'hypothèses.

Grâce à ce travail on voit que l'étude de la voix âgée et de sa reconnaissance automatique sont possibles et entièrement justifiées comme domaine à part entière étant donné les spécificités qu'elle présente. L'étude que nous avons menée ici a permis de montrer une certaine relation entre le vieillissement du métabolisme (tant sur le plan physiologique que cognitif) et l'évolution de la performance langagière. Nous avons également pointé le fait que dans la démarche de l'amélioration d'un système de RAP, bien que la plupart des chercheurs jouent sur l'adaptation du modèle de langage, dans notre cas, l'adaptation du modèle acoustique s'est montrée plus efficace.

Par ailleurs, extrêmement peu de chercheurs ont étudié l'aspect parole lue vs parole spontanée. Tout d'abord nous avons vu qu'il est possible d'adapter un système avec un corpus de parole lue et d'appliquer les modèles acoustiques ainsi adaptés pour de la RAP de parole spontanée, les résultats sont satisfaisants bien que nous n'ayons pas eu le temps de comparer ces résultats avec des modèles acoustiques adaptés sur de la parole spontanée.

En ce qui concerne l'alignement phonémique, en reconnaissance par classe de phonème nous avons obtenus des résultats moyens par classes de phonèmes équivalentes entre la parole lue et la parole spontanée, nous en avons conclu que la reconnaissance des phonèmes n'était pas problématique qu'il s'agisse de parole lue ou de parole spontanée et qu'une fois que le modèle acoustique

donnait un score satisfaisant, alors on devait se pencher sur le modèle de langage. Cette étude a permis de montrer une partie des spécificités de la voix âgée en comparaison à la voix classique non âgée, démontrant qu'elle devait être traitée différemment de la parole non âgée. Aussi, nous avons présenté le comportement du logiciel Sphinx3 avant et après adaptations à la voix âgée, prouvant que dès lors que la tâche de RAP veut être vouée à des personnes de plus de 70 ans, alors il est indispensable d'adapter ce système.

7.2 Perspectives

Ce qui manque dans cette étude est une typologie des erreurs typiques provoquant cette hausse du taux d'erreur tant pour les phonèmes que pour le décodage en mots. En effet, nous avons remarqué qu'il y avait un certain nombre d'inversion sur le plan phonémique ce qui nuit au décodage en mots, une matrice de confusions réalisée sur un très grand nombre de locuteurs permettrait de générer une règle qui viendrait compléter le système de RAP et lui permettrait ainsi d'anticiper les erreurs et y pallier.

L'analyse phonémique a montré que le vieillissement avait effectivement des répercussions sur la reconnaissance des phonèmes et que certains phonèmes étaient moins bien reconnus. Cependant, les deux études ayant cherché à classer les phonèmes de la sorte ne sont pas arrivés aux mêmes conclusions sur tous les points. Il faudrait donc étudier sur plus de locuteurs quels sont, statistiquement les phonèmes les plus problématiques.

Au cours de cette étude nous avons constaté que le phonème [h] non compris dans les phonèmes du français, apparaissait dans l'alignement forcé comme étant prononcé dans le corpus de parole âgée et jamais dans le corpus de parole non âgée. Il serait intéressant de vérifier cette réflexion en enregistrant plus de locuteurs non âgés suivant le même protocole que pour Anodin Détresse par exemple.

Notre étude montre qu'il est possible d'adapter un système (ici à de la parole âgée) destiné à de la reconnaissance de parole spontanée en effectuant les adaptations sur un corpus correspondant de parole lue. Il se peut que l'adaptation à partir de parole lue pour des test sur de la parole spontanée ne soit pas efficace dans toutes les circonstances.

Très peu d'études ont été réalisées sur de la parole âgée, il serait donc nécessaire d'étudier, d'un point de vue linguistique, de manière longitudinale l'évolution de la parole subissant les effets du vieillissement pour tenter d'observer des tendances communes, types de confusions, taux d'erreur plus élevé pour telle ou telle classe de phonèmes, etc. Ce qui permettrait probablement d'obtenir de meilleurs résultats en traitement automatique.

Bibliographie

- [Gay, 2009a] (2009a). *Le vieillissement normal et pathologique du langage : étude comparative des discours oraux*, Lorient, France. 6emes journées internationales de Linguistique de Corpus.
- [Gay, 2009b] (2009b). *Pauses et hésitations dans le discours de patients Alzheimer et chez la personne saine*, Aix-en-Provence, France. 3emes Journées de Phonétique Clinique.
- [Ama, 2012a] (2012a). *Contribution à l'étude de la variabilité de la voix des personnes âgées en reconnaissance automatique de la parole*, volume 1 : JEP, Grenoble, France.
- [LeG, 2012] (2012). *Utilisation de la Reconnaissance Automatique de la Parole pour l'aide à l'autonomie des personnes âgées*, Castres-Mazamet, France.
- [Ama, 2012b] (2012b). *Étude de la performance des modèles acoustiques pour des voix de personnes âgées en vue de l'adaptation des systèmes de RAP*, Grenoble, France.
- [Allegre, 2003] Allegre, J. (2003). Approche de la reconnaissance automatique de la parole. Master's thesis, Conservatoire National des Arts et Métiers.
- [Ameur, 2011] Ameur, M. (2011). Caractérisation des accents étrangers en utilisant des techniques de reconnaissance automatique de la parole. Master's thesis, Université Grenoble3.
- [Aynaud, 2009] Aynaud, C. (2009). Reconnaissance de la parole chez les personnes âgées, adaptation d'un système de reconnaissance. Master's thesis, INP Grenoble.
- [Bousquet-Vernhettes, 2002] Bousquet-Vernhettes, C. (2002). *Compréhension robuste de la parole spontanée dans le dialogue oral homme-machine-Décodage conceptuel stochastique*. PhD thesis, Université Toulouse III - Paul Sabatier.

- [Cappeau and Gadet, 2007] Cappeau, P. and Gadet, F. (2007). 'Document 3' Où en sont les corpus sur les français parlés ?, volume 2007/1, pages 129–133. CAIRN.INFO.
- [de Saussure, 1975] de Saussure, F. (1975). *Cours de Linguistique Générale (1906-1911)*. Editions Payot.
- [Debeaux and Vacher, 2010] Debeaux, C. and Vacher, M. (2010). Evaluation du décodeur pocket-sphinx sur la plateforme e-lio. Technical report.
- [des Enseignants de Gériatrie, 2000] des Enseignants de Gériatrie, C. N. (2000). *Corpus de Gériatrie - Accès aux textes*, volume Tome 1. Editions 2000.
- [Dufour, 2010] Dufour, R. (2010). *Transcription automatique de la parole spontanée*. PhD thesis, Université du Maine.
- [Dugheanu, 2011] Dugheanu, R. (2011). Evaluation des outils pour la reconnaissance automatique de la parole adaptée aux personnes âgées. Master's thesis, Université Grenoble3.
- [Evermann, 1999] Evermann, G. (1999). *Minimum Word Error Rate Decoding*. PhD thesis, University of Cambridge.
- [Fohr et al., 1994] Fohr, D., Haton, J. P., and Laprie, Y. (1994). Knowledge-based techniques in acoustic-phonetic decoding of speech : Interest and limitations. *IJPRAI*, 8(1) :133–153.
- [Gauvain et al., 1990] Gauvain, J.-L., Lamel, L., and Eskenazi, M. (1990). Design considerations and text selection for bref, a large french read-speech corpus. In *ICSLP*.
- [Haton et al., 1991] Haton, J.-P., Caelen, J., Pierrel, J.-M., Perennou, G., and Gauvain, J.-L. (1991). *Reconnaissance Automatique de la Parole*. Dunod.
- [Hervy, 2003] Hervy, B. (2003). Animateur en gérontologie : des tisseurs de liens. *La Santé de l'Homme*, (363) :23–26.
- [Lardilleux et al., 2011] Lardilleux, A., Yvon, F., and Lepage, Y. (2011). Généralisation de l'alignement sous-phrastique par échantillonnage.
- [Lefol, 2010] Lefol, Q. (2010). Acquisition de corpus pour la génération de modèles acoustiques adaptés à la voix des personnes âgées. Master's thesis, INP Grenoble.
- [Sasa and Grand, 2011] Sasa, Y. and Grand, J. L. (2011). Corpus parole de personnes Âgées grenoble 2011 - explications, légendes et observations des transcriptions et dispositifs utilisés pour les entretiens.
- [Shannon, 1948] Shannon, C. E. (1948). A mathematical theory of communication. *The Bell system technical journal*, 27 :379–423.

- [Tissier, 2011] Tissier, F. (2011). Adapter le vocabulaire d'un système de transcription automatique de la parole aux thèmes abordés. Master's thesis, Ecole Supérieure d'Ingénieur de Rennes.
- [Vacher et al., 2011] Vacher, M., Portet, F., Fleury, A., and Noury, N. (2011). Development of audio sensing technology for ambient assisted living : Applications and challenges. *International Journal of E-Health and Medical Communications*, 2(1) :35 – 54.
- [Vacher et al., 2012] Vacher, M., Portet, F., Rossato, S., Aman, F., Golanski, C., and Dugheanu, R. (2012). Speech-based interaction in an aal context. *Gerontechnology*, 11 :310.
- [Vacher et al., 2006] Vacher, M., Serignat, J., Chaillol, S., Istrate, D., and Popescu, V. (2006). Speech and sound use in a remote monitoring system for health care. 4188/2006 :711–718.
- [Vigouroux et al., 2005] Vigouroux, N., Privat, R., and Truillet, P. (2005). Etude de l'effet du vieillissement sur les productions langagières et sur les performances en reconnaissance automatique de la parole. *Revue Parole*, 2004-3&-32 :281–318.
- [Vipperla et al., 2010] Vippera, R., Renals, S., and Frankel, J. (2010). Ageing voices : The effect of changes in voice parameters on asr performance. *EURASIP J. Audio, Speech and Music Processing*, 2010.
- [Xue and Delisky, 2001] Xue, S. A. and Delisky, D. (2001). *Effects of aging on selected acoustic voice parameters : preliminary normative data and educational implications*, volume 27, Issue 2.

Webographie

- [1] <http://www.predim.org/IMG/pdf/FicheSATIM.pdf>, consulté le 24/08/2012
- [2] http://www.chups.jussieu.fr/polys/geriatrie/tome1/01_vieillissement.pdf consulté le 25/08/2012
- [3] www.voicemedecine.com/aging.htm, consulté le 26/08/2012
- [4] www.lecorpshumain.fr/lecorpshumain/5-vieillissement.html, consulté le 25/08/2012
- [5] papidoc.chic-cm.fr/580villiphysio.html, consulté le 25/08/2012
- [6] www.e-sante.fr/essentiel-orthophoniste-sujet-age/actualite/20717, consulté le 26/08/2012
- [7] www.inpes.sante.fr/slh/articles/363/02.htm# consulté le 25/08/2012
- [8] aides.electroniques.proteor.fr/scripts/produit.php?id_pdt=327 consulté le 02/06/2012
- [9] www.ubiquiet.com/fr/page/pr%C3%A9sentation-ubiquiet-full.htm consulté le 02/06/2012
- [10] www.ac-grenoble.fr/PhiloSophie/logphil/textes/textesm/saussu3m.htm consulté le 22/05/2012

Chapitre 8

Annexes

8.1 Tableau récapitulatif du corpus AD80

| Corpus | Nb. locuteurs | Age min/max | Durée Durée | Nb. phrases | Lieux d'enregistrement |
|-----------|---------------|-------------|-------------|-------------|------------------------|
| AD | 21 | 22-64 | 38min | 2646 | studio LIG |
| VA | 7 | 62-94 | 1h25 | 5451 | CHU Grenoble, domicile |
| ERES38 | 24 | 68-98 | 17h44 | env. 17438 | Mais. retraite, foyers |
| AD80part2 | 36 | 62-94 | 1h25 | 4202 | Mais. retraite |

FIGURE 8.1 – Tableau récapitulatif des sessions composant le corpus AD80

8.2 Tableau des correspondances API/SAM-PA/SPHINX

| Index | Sphinx | Sampa | TextGrid(API) | API |
|-------|--------|-------|---------------|------|
| 1 | p | p | p | p |
| 2 | b | b | b | b |
| 3 | t | t | t | t |
| 4 | d | d | d | d |
| 5 | k | k | k | k |
| 6 | g | g | \gs | g |
| 7 | f | f | f | f |
| 8 | v | v | v | v |
| 9 | s | s | s | s |
| 10 | z | z | z | z |
| 11 | SS | S | \sh | ʃ |
| 12 | ZZ | Z | \zh | ʒ |
| 13 | j | j | i | j |
| 14 | m | m | m | m |
| 15 | n | n | n | n |
| 16 | NJ | J | \ng | ŋ |
| 17 | NG | N | \n. | ŋ |
| 18 | l | l | l | l |
| 19 | RR | R | \ri | ɻ |
| 20 | w | w | w | w |
| 21 | HH | H | \ht | ɥ |
| 22 | j | j | j | j |
| 23 | i | i | i | i |
| 24 | e | e | e | e |
| 25 | EE | E | \ef | ɛ |
| 26 | a | a | a | a |
| 27 | AA | A | \as | ɑ |
| 28 | OO | O | \ct | ɔ |
| 29 | o | o | o | o |
| 30 | u | u | u | u |
| 31 | y | y | y | y |
| 32 | 2 | 2 | \o/ | ø |
| 33 | 9 | 9 | \oe | œ |
| 34 | aeA | @ | \sw | ə |
| 35 | in | e~ | e ~^ | e~ |
| 36 | an | a~ | \as ~^ | a~ |
| 37 | on | o~ | o ~^ | o~ |
| 38 | un | 9~ | \oe ~^ | œ~ |
| 39 | eEE | E/ | e ef | eɛ |
| 40 | aAA | A/ | a as | aɑ |
| 41 | 29 | &/ | \o/\oe | øœ |
| 42 | oOO | O/ | o ct | oɔ |
| 43 | in un | U~/ | e ~^ \oe ~^ | e~œ~ |
| 44 | SIL | / | / | / |

8.3 Document explicatif sur les triphones

Traitement des triphones sous Sphinx3

Juline Le Grand

Contenu

| | |
|---|---|
| Introduction..... | 2 |
| 1 Les triphones | 2 |
| 1.1 Définition..... | 2 |
| 1.2 Les triphones dans la langue..... | 2 |
| 2 Pertinence d'un traitement par triphones..... | 3 |
| 3 Fonctionnement du traitement (tri)phonémique dans Sphinx3 | 3 |
| 3.1 Sélection des triphones dans Sphinx3 | 3 |
| 4 Les états de transitions | 3 |
| 5 Références | 3 |

Introduction

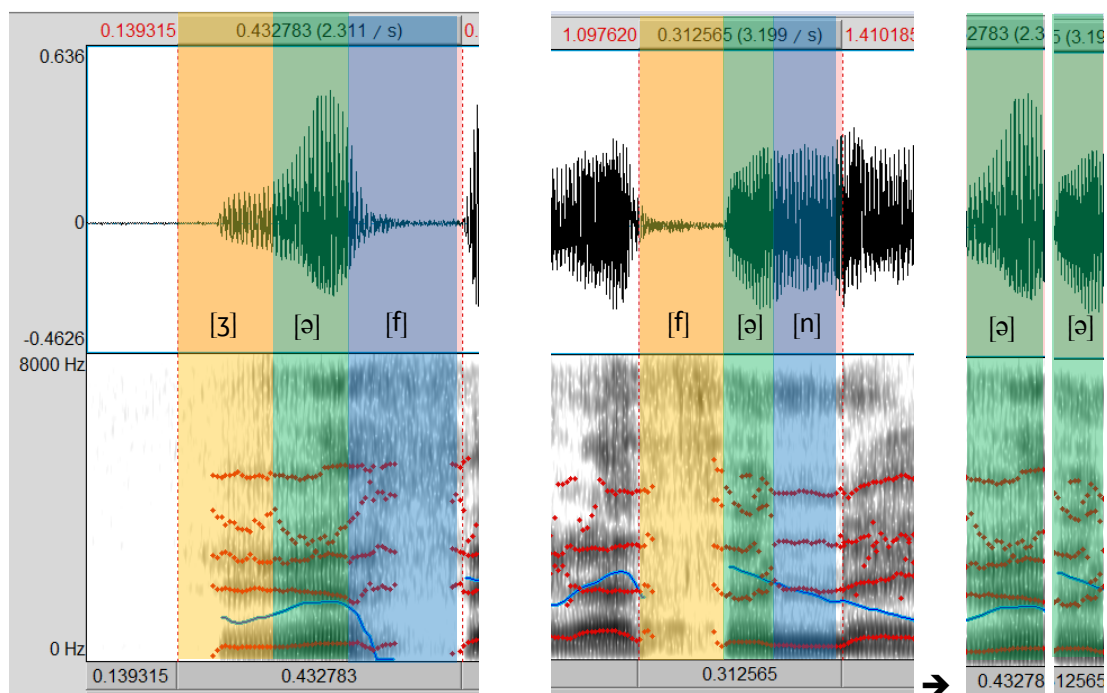
Le traitement de reconnaissance de la parole continue effectuée par le logiciel Sphinx3 est un traitement réalisé par triphone. Plusieurs aspects de cette technique sont éclaircis par le présent document.

1 Les triphones

1.1 Définition

Un triphone est un ensemble de 3 phonèmes¹. En effet il s'agit d'un phonème pour lequel on prend en compte le contexte gauche et le contexte droit. Un phonème n'est en aucun cas un élément indépendant, la raison de cela est que du point de vue acoustique, les contextes ont une influence significative sur la réalisation et sur les caractéristiques du phonème central.

Par exemple, pour les triphones suivants [ʒəf] et [fən] on voit que les [ə] centraux sont différents, influencés par leurs contextes respectifs.



1.2 Les triphones dans la langue

Selon les dictionnaires, il existe 39 phonèmes pour le français et dénombre près de 6000 triphones.

[p b t d k g f v s z ʃ ʒ j m n ŋ ŋ l ʁ w ʁ j i e ε a o u y ø œ ə ɛ̃ ä ö œ̃] sont les phonèmes répertoriés du français auxquels il faut ajouter le « silence ».

A partir de ces 40 éléments, on dénombre alors 64 000 combinaisons possibles de phonèmes en triphones. Or ce nombre est restreint à environ 6000 triphones par le fait certains cas soient exclus pour des raisons linguistiques. En effet, du point de vue de la langue, toutes les possibilités ne sont pas réalisables ni acceptables ou du moins n'apparaissent pas dans la langue.

¹ Phonème : plus petite unité sonore significative dans la langue.

2 Pertinence d'un traitement par triphones

En prenant en compte chaque phonème et son contexte grâce au triphone, on obtient des modèles pour toutes les réalisations possibles. En traitant par triphone, on optimise la reconnaissance du phonème et on obtient également des indices pour le phonème qui précède ou celui qui suit. En effet, si le phonème précédent à mal été reconnu, le triphone va permettre d'obtenir des informations supplémentaires en feed back ou prédiction.

3 Fonctionnement du traitement (tri)phonémique dans Sphinx3

3.1 Sélection des triphones dans Sphinx3

La CMU² dénombre pour le logiciel 43 phonèmes, ainsi 79 507 triphones sont possibles. Cependant ils ont intégré la notion de contexte :

- b = commence
- s = seule
- i = milieu
- e = fin

Ce critère supplémentaire se justifie par le fait qu'en fonction de la position du triphone dans la séquence sonore. De même que tous les triphones n'existent pas, tous les triphones n'apparaissent pas dans chacun des contextes.

On obtient donc une liste de 125 644 triphones qui répertorient tous les cas possibles en français selon la CMU.

Le modèle acoustique pour Sphinx place le triphone comme unité acoustique de base, donnant des informations précise, adapté pour l'apprentissage et que l'on peut généraliser. Cette unité est adéquate pour le travail sur un grand vocabulaire en parole continue.

Ils ont également traité par triphones car cela permet de gérer les mots inconnus.

4 Les états de transitions

Le traitement par triphone de Sphinx3 inclue également la notion d'« états de transitions ». Il s'agit des transitions entre les phonèmes dans les triphones. On dénombre dans le logiciel 4122 états de transitions.

Dans une phase d'adaptation il est donc important d'améliorer les phonèmes, les triphones et les états de transitions

5 Références

<http://www.speech.cs.cmu.edu/cgi-bin/cmudict>

http://www.google.fr/url?sa=t&rct=j&q=&esrc=s&source=web&cd=1&ved=0CCkQFjAA&url=http%3A%2F%2Fwww.liacs.nl%2F~erwin%2FSR2003%2FStudents%2F10_SR_Triphones.ppt&ei=a_ZqT9uEMi1hAfc08CPBw&usq=AFQjCNFaLluDhamCmxCjZVRQJngRYlkXvw&sig2=pvOcMfwru2nBn70qrYi3sw

² Carnegie Mellon University

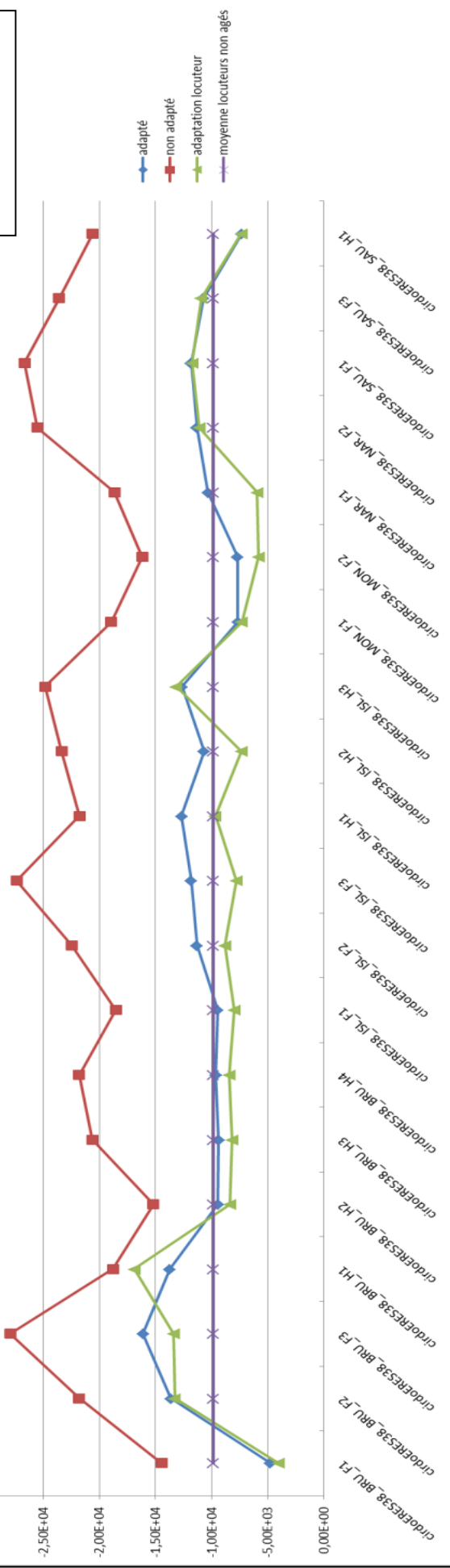
8.4 Tableau des occurrences des phonèmes pour le cas de la parole spontanée traité de ERES38

Apperçut des occurrences par phonème. Exemple pour la parole spontanée. (env. 2 min par locuteur)

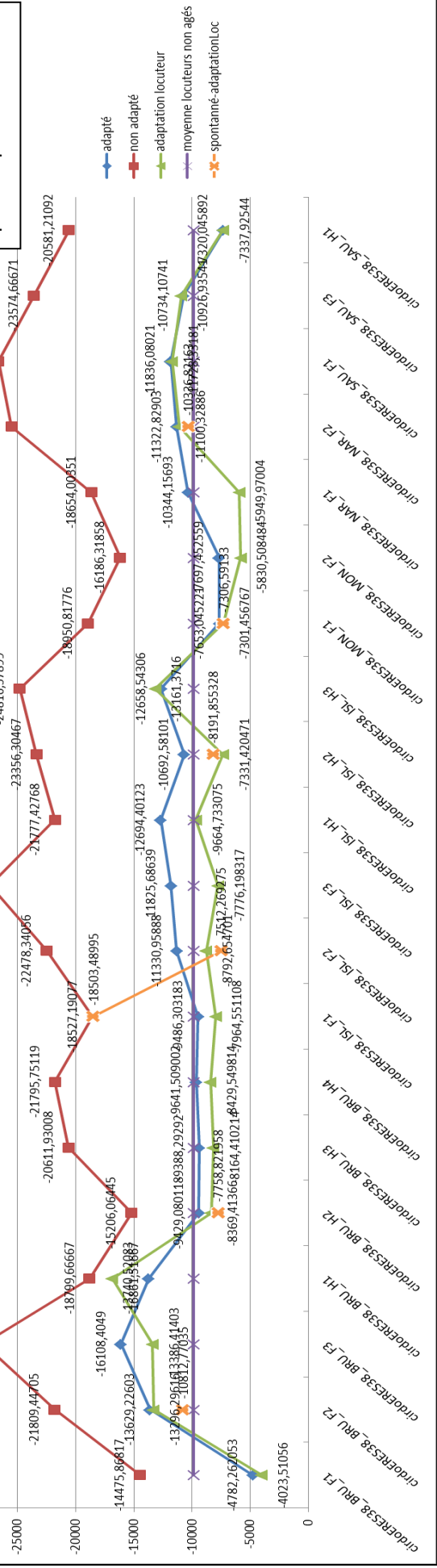
| phonème | BRU_F2 | BRU_H2 | ISL_F1 | ISL_F2 | ISL_H2 | MON_F1 | NAR_F2 | total occurrence | moyenne |
|---------|--------|--------|--------|--------|--------|--------|--------|------------------|---------------|
| z | 13 | 3 | 5 | 15 | 10 | 8 | 5 | 59 | 8,428 57 143 |
| 9 | 6 | 5 | 4 | 6 | 6 | 4 | 2 | 33 | 4,7 1428 57 1 |
| z9 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 0,28 57 1429 |
| a | 61 | 77 | 24 | 34 | 63 | 43 | 49 | 351 | 50,1428 57 1 |
| AA | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 2 | 0,28 57 1429 |
| aAA | 6 | 18 | 7 | 10 | 19 | 16 | 20 | 96 | 13,7 1428 57 |
| aaA | 22 | 31 | 7 | 25 | 29 | 23 | 17 | 154 | 22 |
| an | 32 | 61 | 13 | 18 | 38 | 28 | 12 | 202 | 28,8 57 1429 |
| b | 6 | 12 | 9 | 3 | 16 | 12 | 13 | 71 | 10,1428 57 1 |
| d | 23 | 42 | 16 | 21 | 42 | 33 | 14 | 191 | 27,28 57 143 |
| e | 43 | 55 | 18 | 28 | 65 | 69 | 32 | 310 | 44,28 57 143 |
| EE | 21 | 48 | 8 | 32 | 46 | 58 | 28 | 236 | 33,7 1428 57 |
| eEE | 4 | 8 | 2 | 3 | 7 | 4 | 9 | 37 | 5,28 57 1429 |
| f | 13 | 14 | 4 | 21 | 11 | 10 | 5 | 78 | 11,1428 57 1 |
| g | 8 | 6 | 5 | 6 | 13 | 4 | 0 | 42 | 6 |
| h | 5 | 0 | 1 | 0 | 2 | 2 | 4 | 14 | 2 |
| HH | 4 | 2 | 6 | 2 | 7 | 3 | 9 | 33 | 4,7 1428 57 1 |
| i | 38 | 41 | 17 | 38 | 64 | 51 | 38 | 287 | 41 |
| in | 14 | 9 | 7 | 11 | 17 | 7 | 9 | 74 | 10,57 14286 |
| j | 15 | 8 | 6 | 14 | 17 | 9 | 10 | 79 | 11,28 57 143 |
| k | 18 | 32 | 7 | 24 | 40 | 23 | 27 | 171 | 24,428 57 14 |
| l | 34 | 52 | 13 | 32 | 65 | 44 | 46 | 288 | 41,1428 57 1 |
| m | 24 | 43 | 13 | 18 | 36 | 29 | 29 | 192 | 27,428 57 14 |
| n | 22 | 25 | 10 | 25 | 35 | 31 | 14 | 162 | 23,1428 57 1 |
| NG | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| NU | 0 | 0 | 0 | 0 | 2 | 0 | 1 | 3 | 0,428 57 143 |
| o | 4 | 7 | 4 | 7 | 15 | 16 | 8 | 61 | 8,7 1428 57 1 |
| on | 8 | 11 | 5 | 14 | 15 | 16 | 17 | 86 | 12,28 57 143 |
| OO | 7 | 14 | 7 | 6 | 19 | 9 | 15 | 77 | 11 |
| oOO | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 4 | 0,57 1428 57 |
| p | 37 | 28 | 11 | 24 | 26 | 22 | 26 | 174 | 24,8 57 1429 |
| RR | 49 | 62 | 16 | 29 | 66 | 45 | 32 | 299 | 42,7 1428 57 |
| s | 26 | 43 | 13 | 16 | 41 | 34 | 27 | 200 | 28,57 14286 |
| SS | 5 | 20 | 7 | 6 | 16 | 10 | 13 | 77 | 11 |
| t | 37 | 45 | 15 | 39 | 64 | 48 | 19 | 268 | 38,28 57 143 |
| u | 4 | 13 | 5 | 13 | 22 | 19 | 7 | 83 | 11,8 57 1429 |
| un | 8 | 1 | 1 | 4 | 1 | 1 | 3 | 19 | 2,7 1428 57 1 |
| v | 10 | 15 | 9 | 9 | 20 | 23 | 13 | 99 | 14,1428 57 1 |
| w | 9 | 5 | 4 | 5 | 10 | 6 | 7 | 46 | 6,57 1428 57 |
| y | 12 | 16 | 4 | 21 | 14 | 26 | 16 | 109 | 15,57 14286 |
| z | 12 | 10 | 5 | 14 | 17 | 34 | 9 | 101 | 14,428 57 14 |
| ZZ | 13 | 27 | 10 | 18 | 30 | 18 | 9 | 125 | 17,8 57 1429 |

8.5 Graphiques des phonèmes montrés en section Résultats

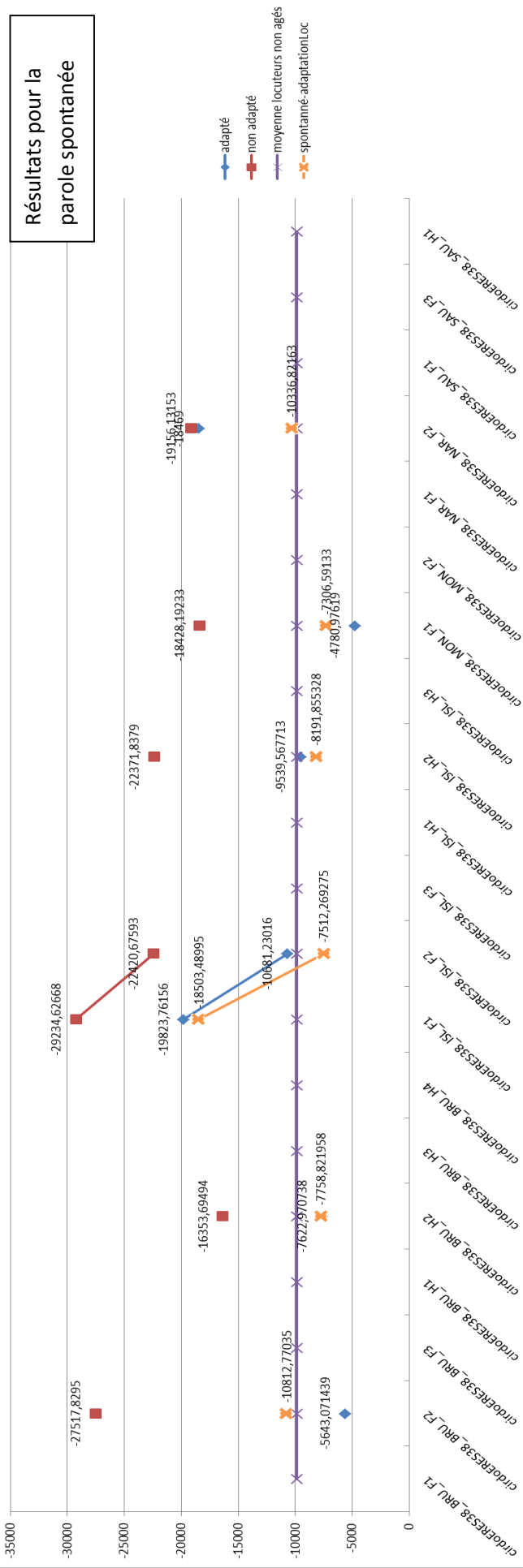
Résultats pour la tâche de lecture



Comparaison résultats de la lecture et de la parole spontanée



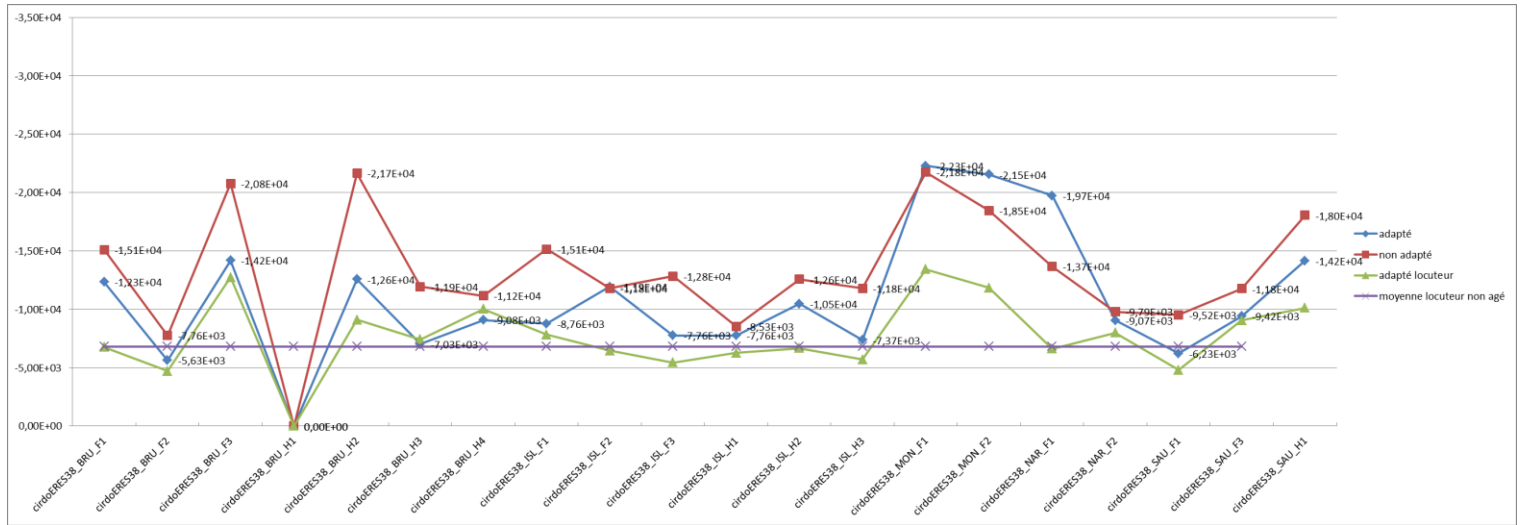
Résultats pour la parole spontanée



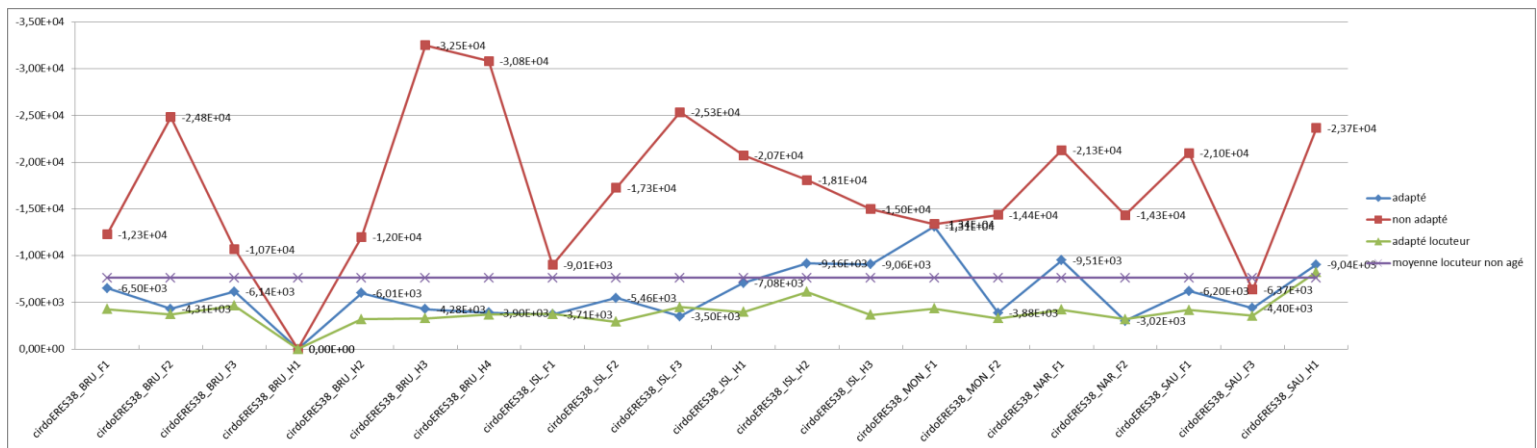
**8.6 Diagrammes représentant pour chaque phonème
et par locuteur la courbe des scores d'alignements
en fonction des modèles acoustiques
Cas de la parole lue**

Graphiques représentant pour chaque phonème les moyennes des scores d'alignements pour chacun des 3 modèles acoustiques ainsi que la moyenne pour des locuteurs non agés.

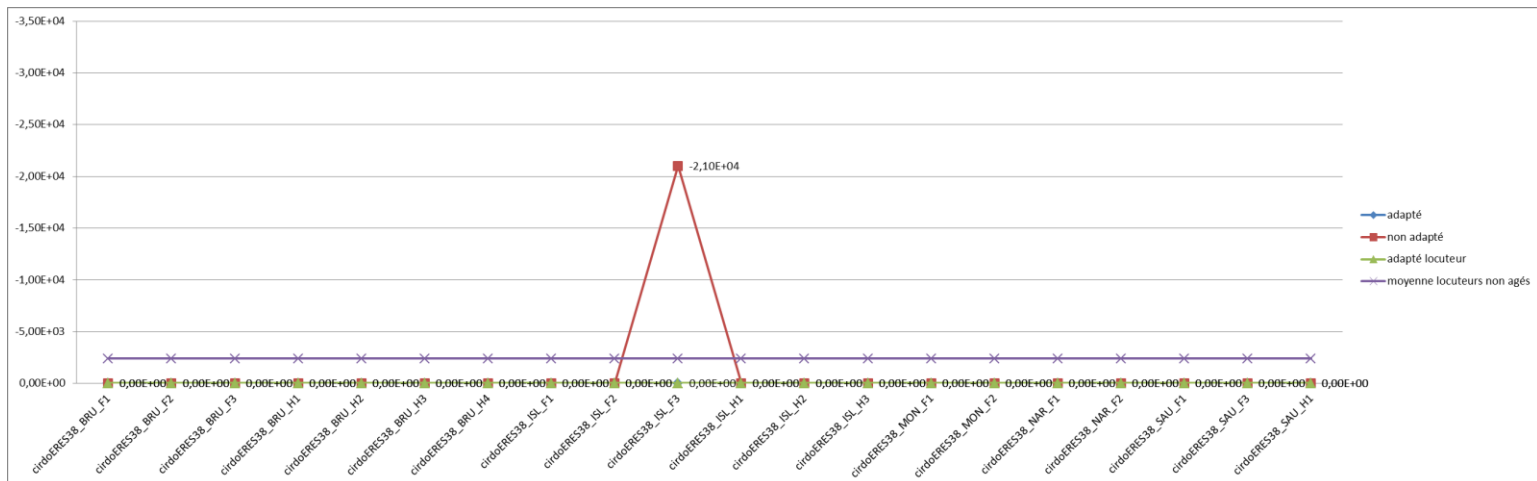
2



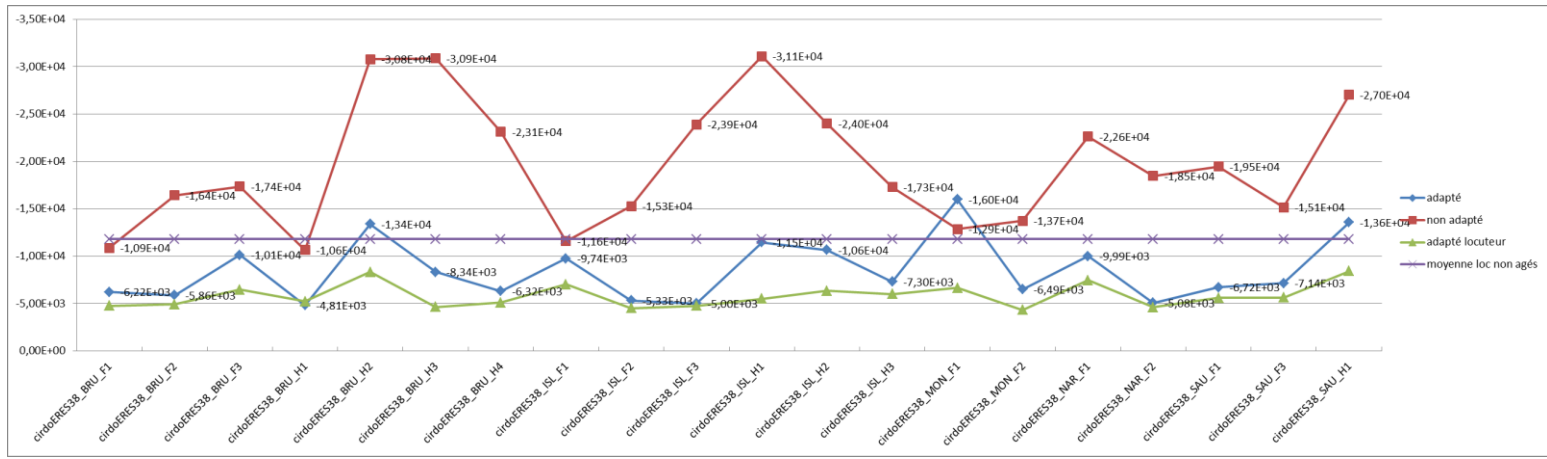
9



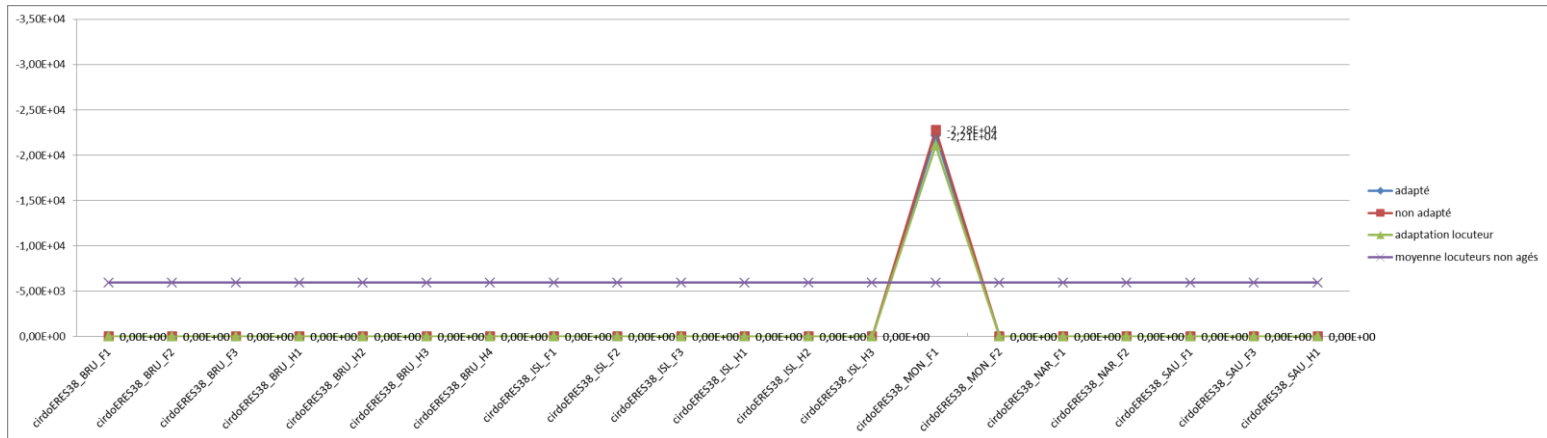
29



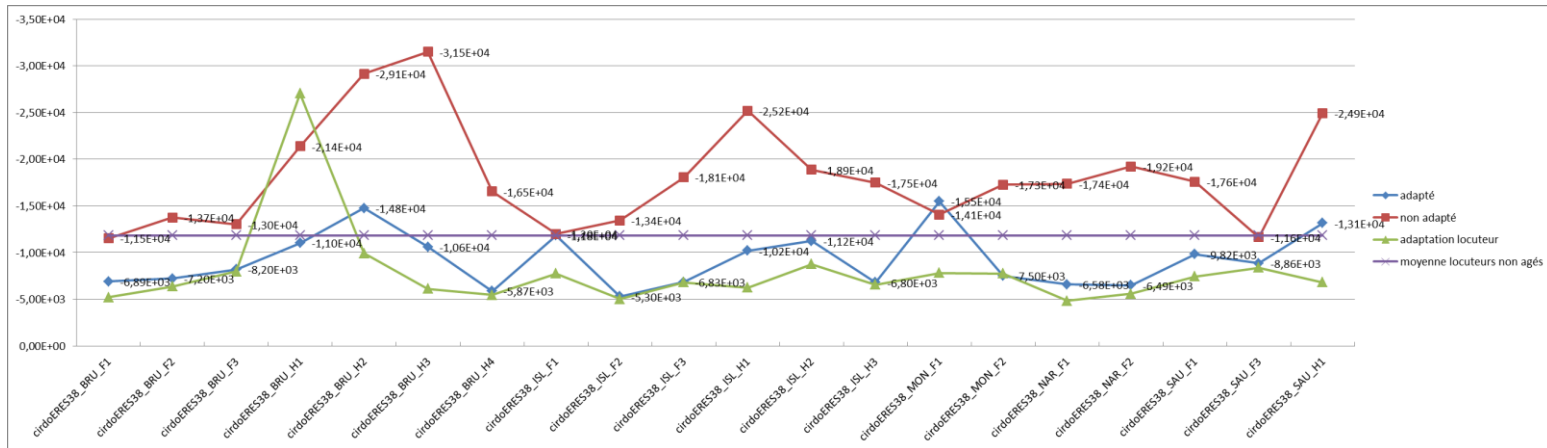
a



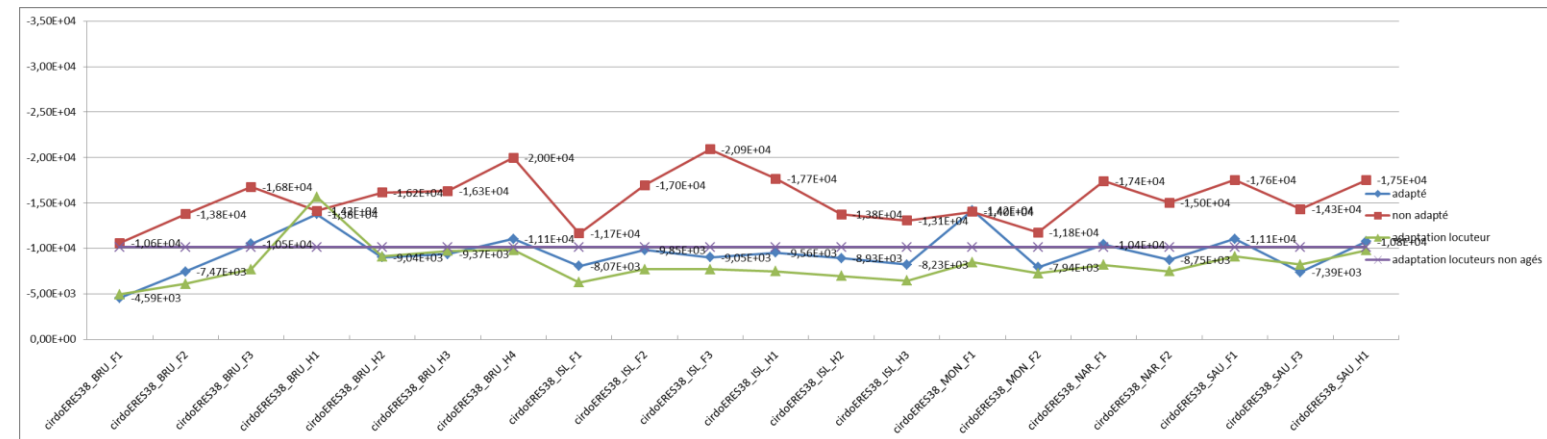
AA



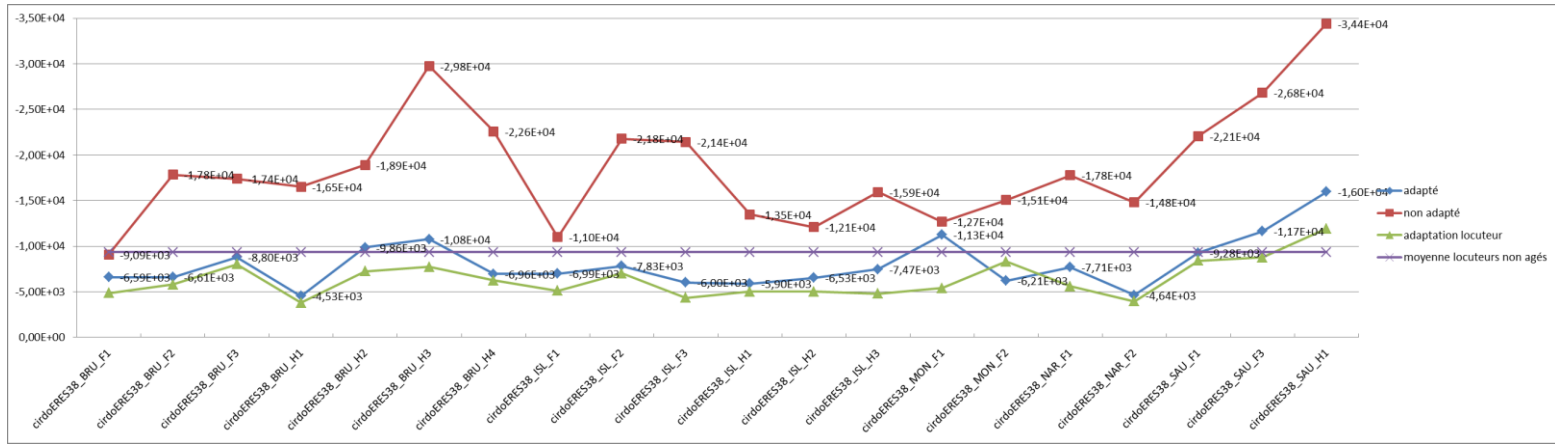
aaA



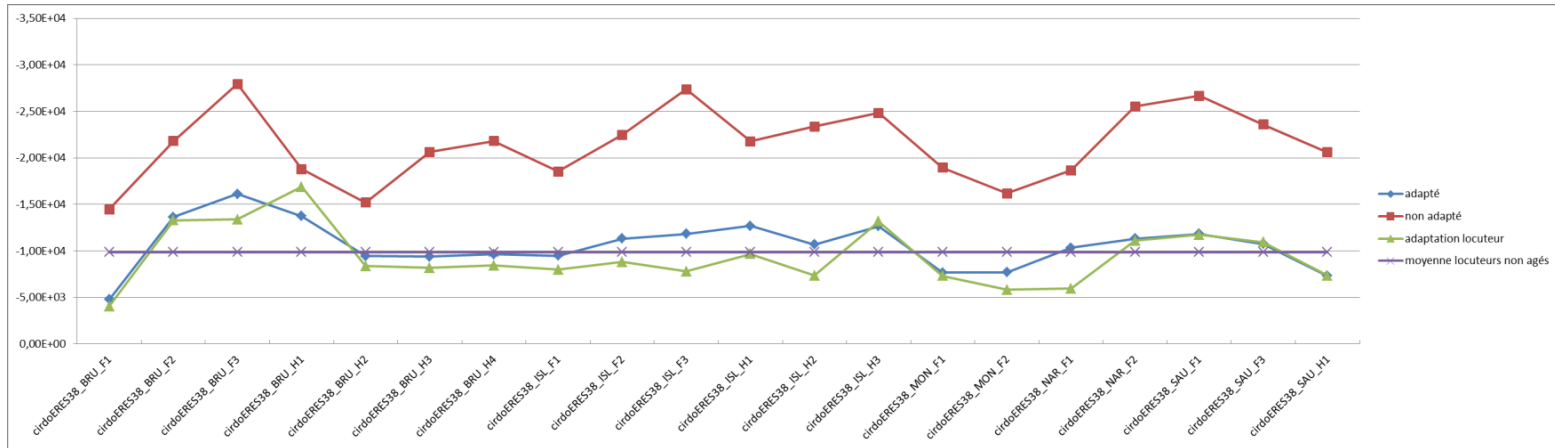
aeA



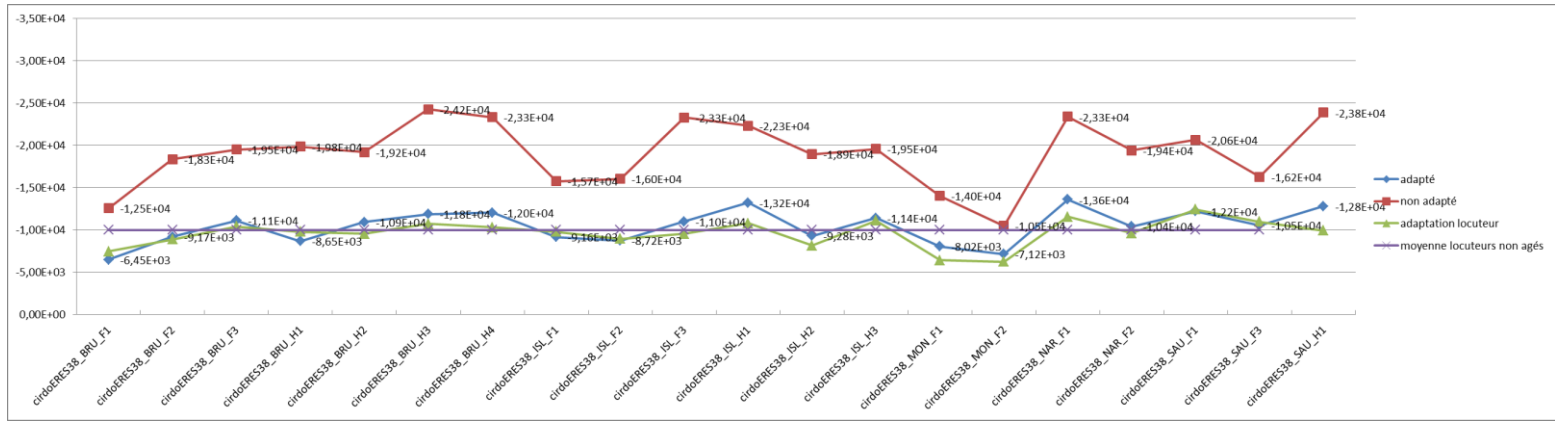
an



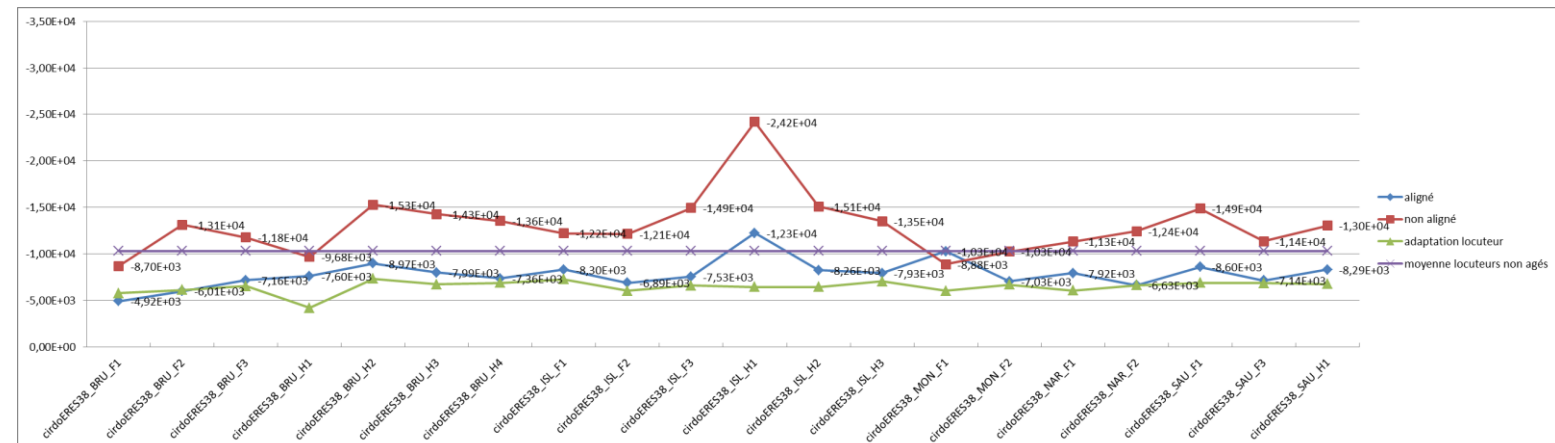
b



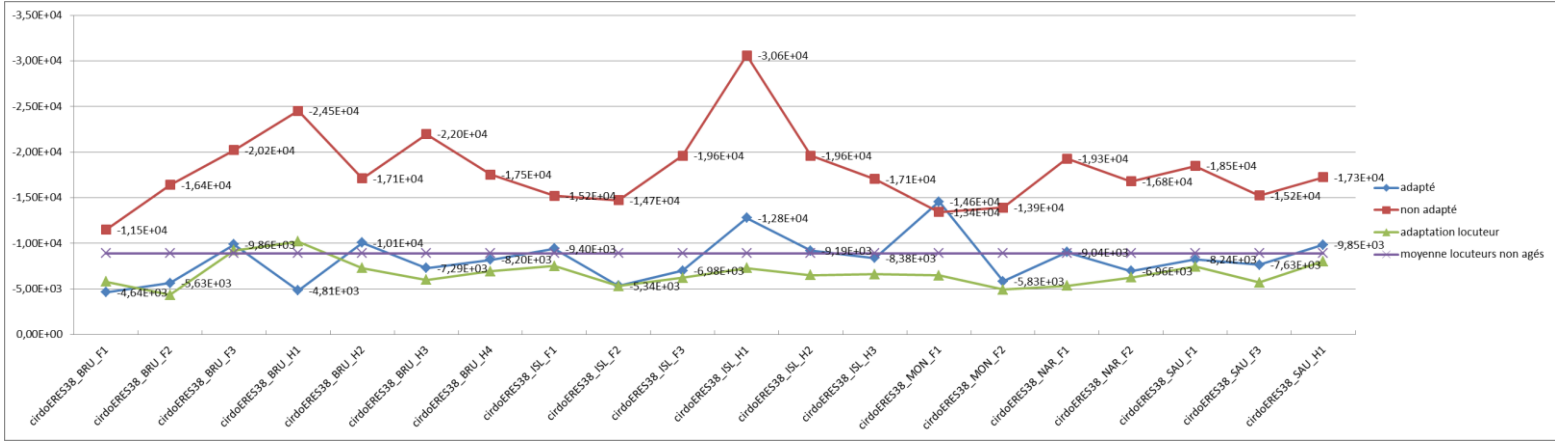
d



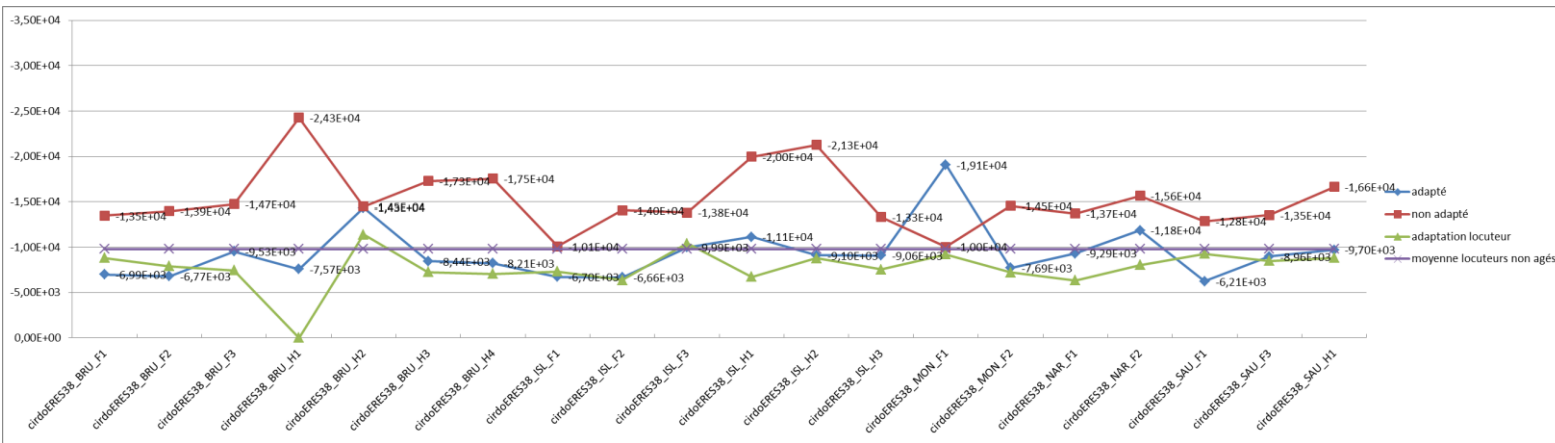
e



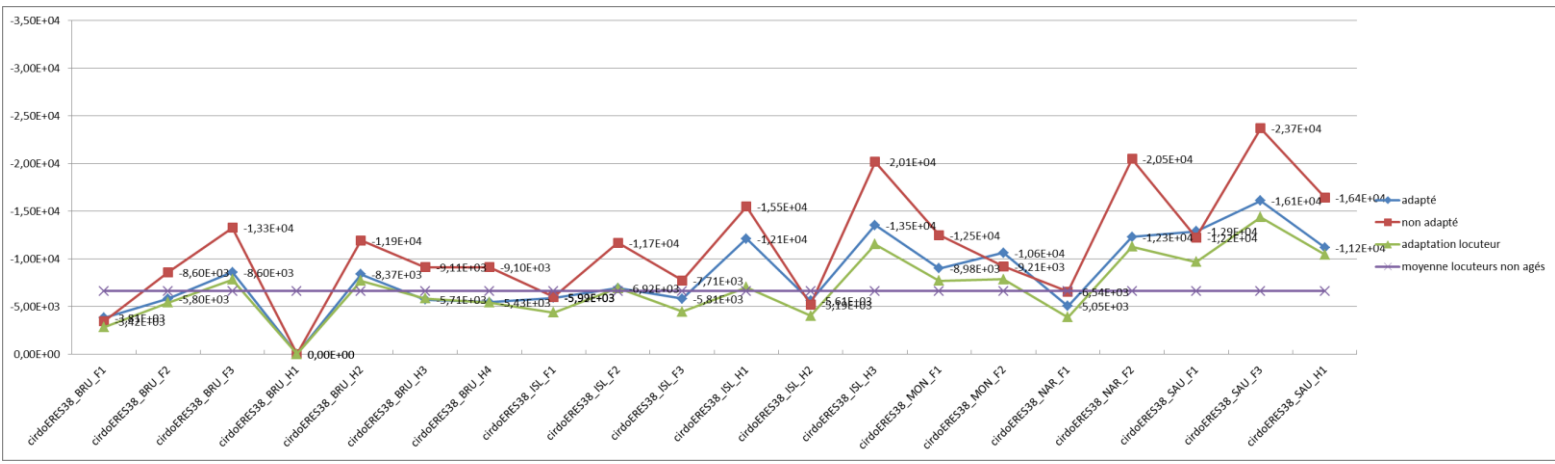
EE



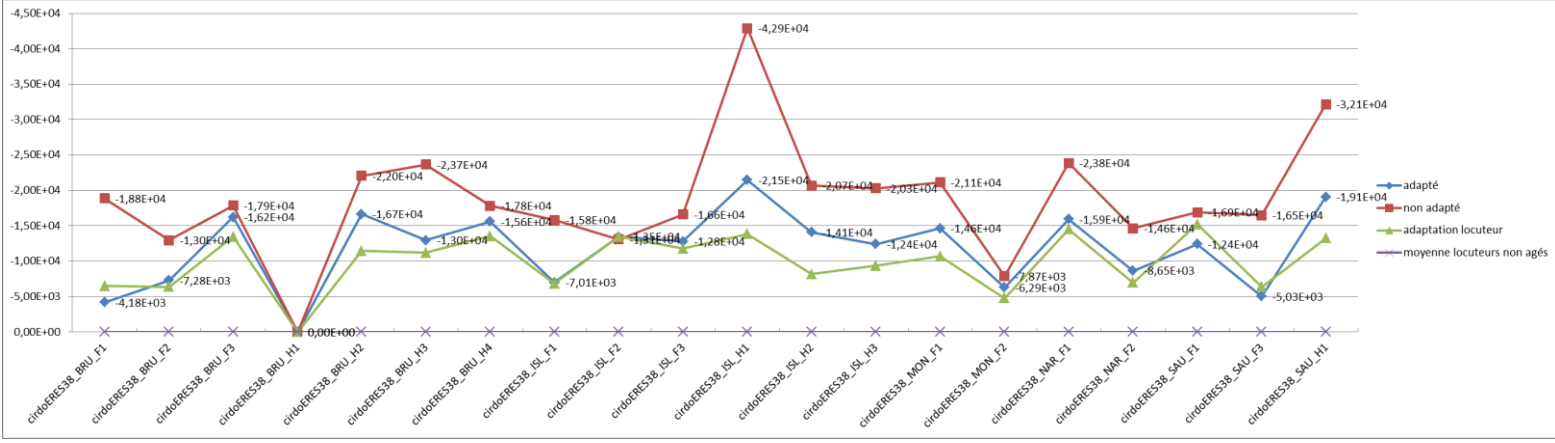
eEE



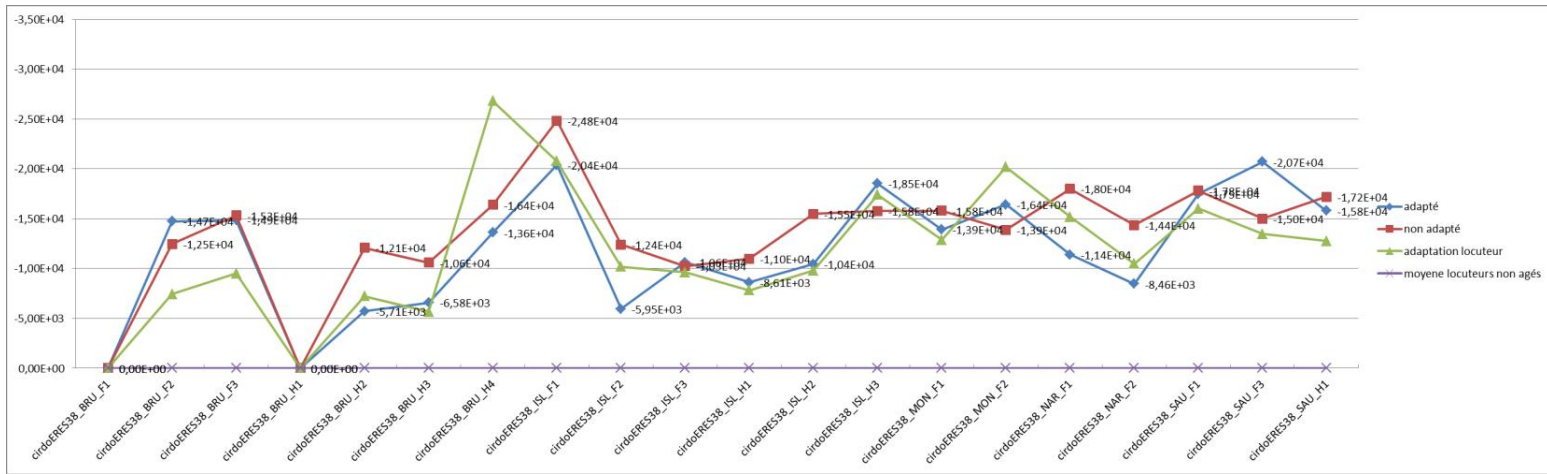
f



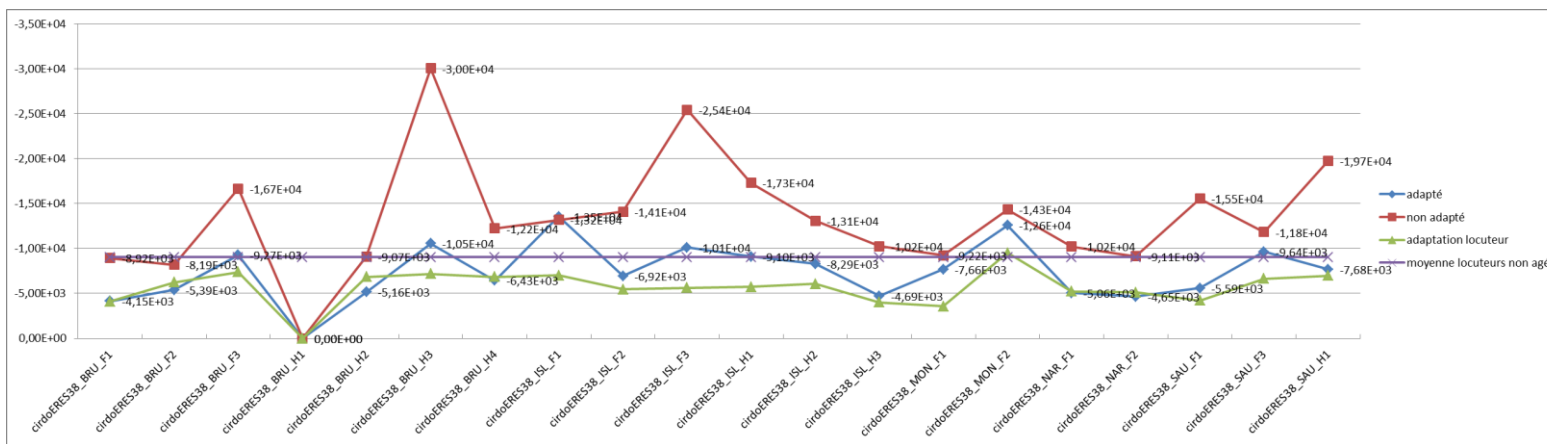
g



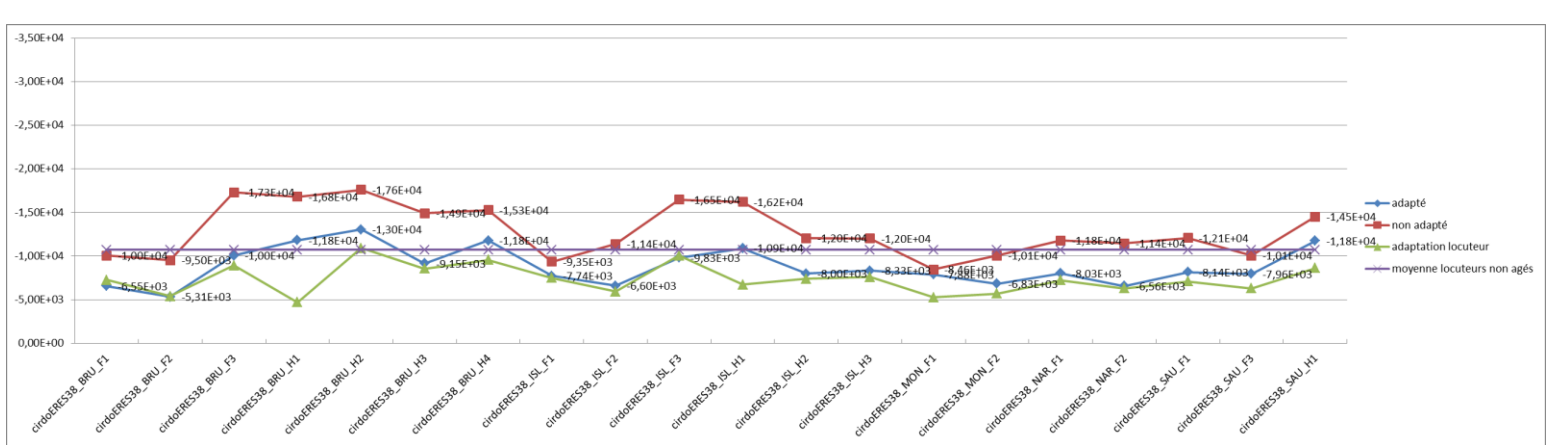
h



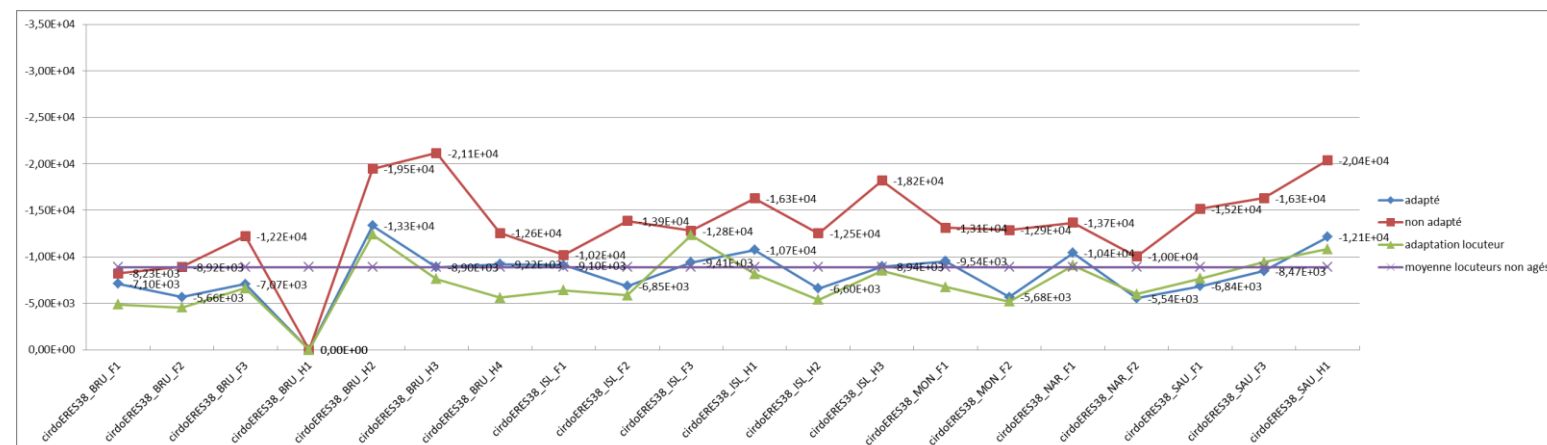
HH



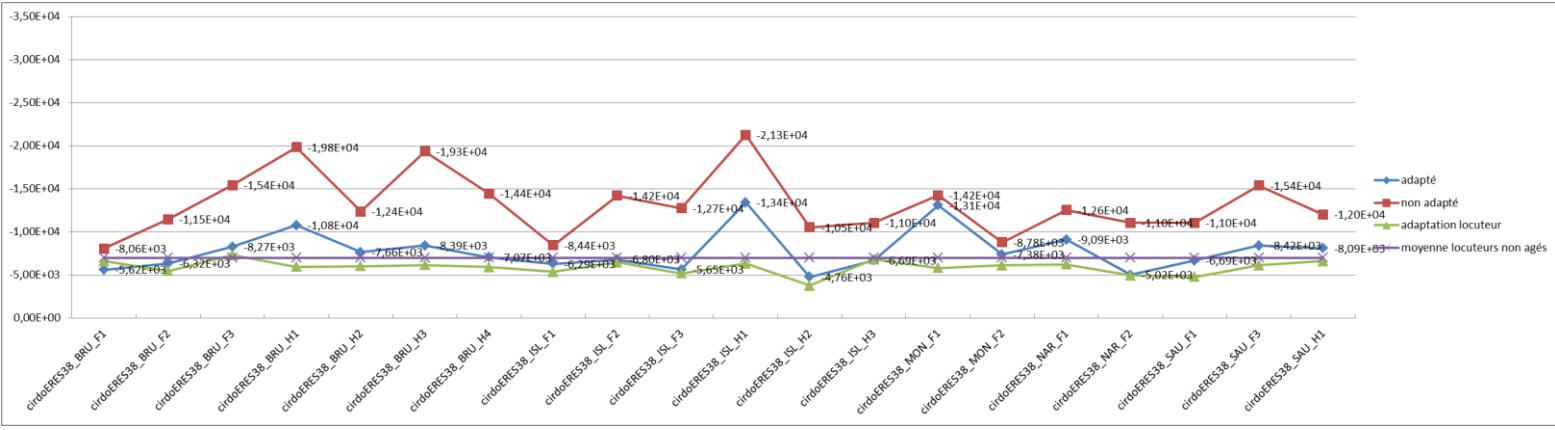
i



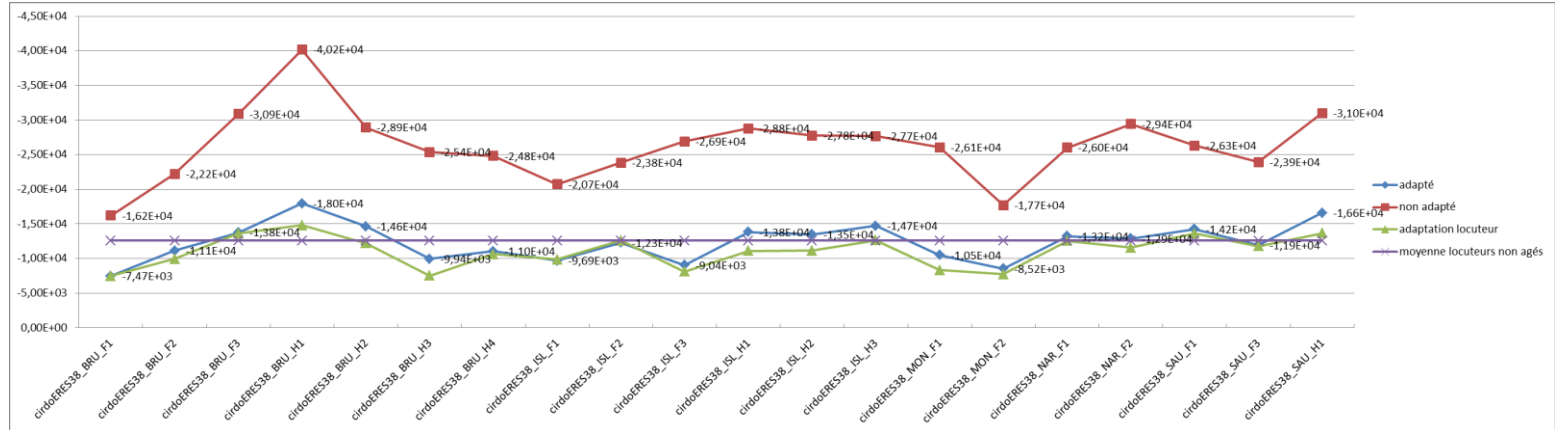
in



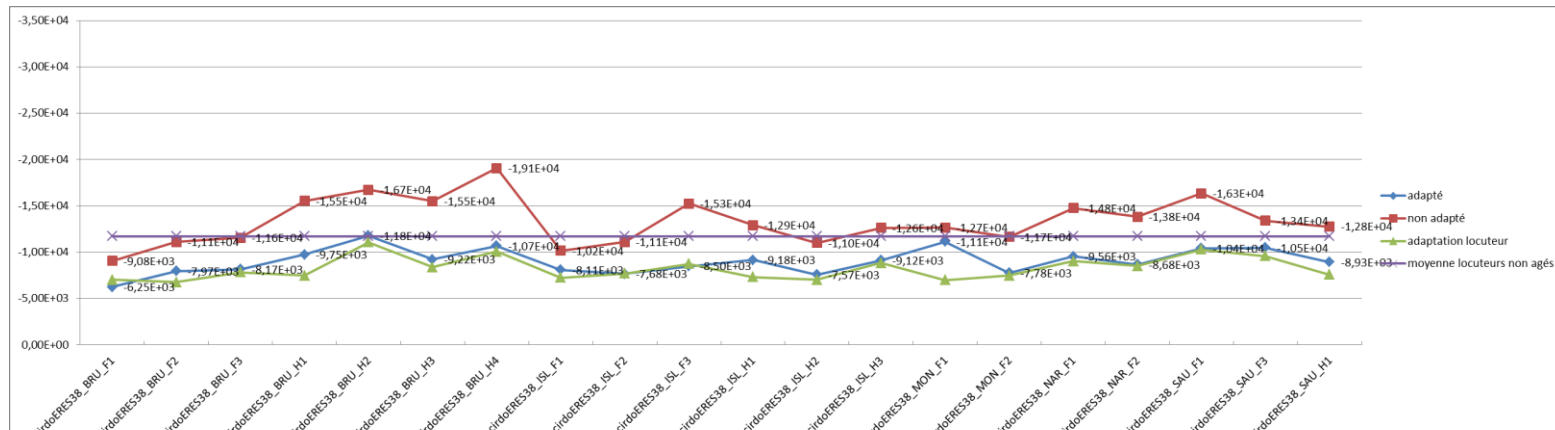
j



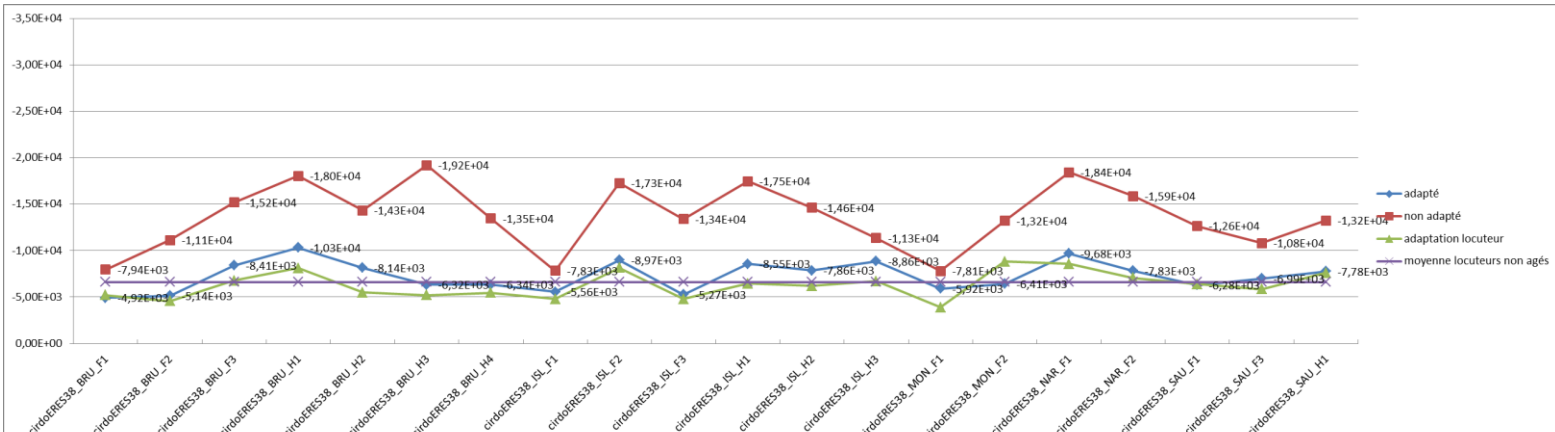
k



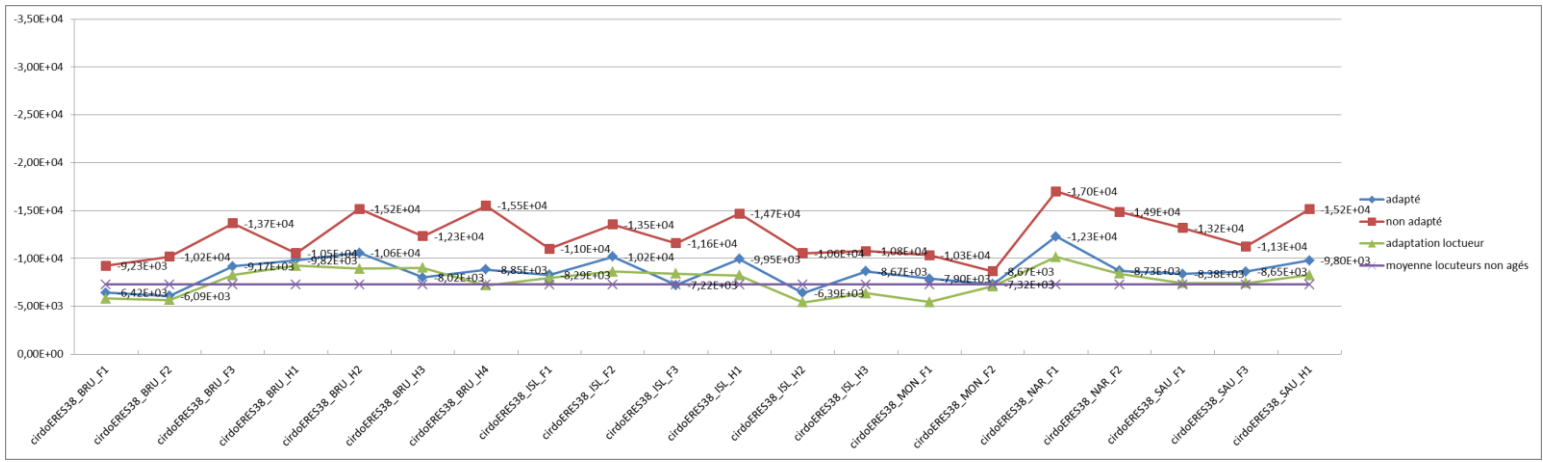
l



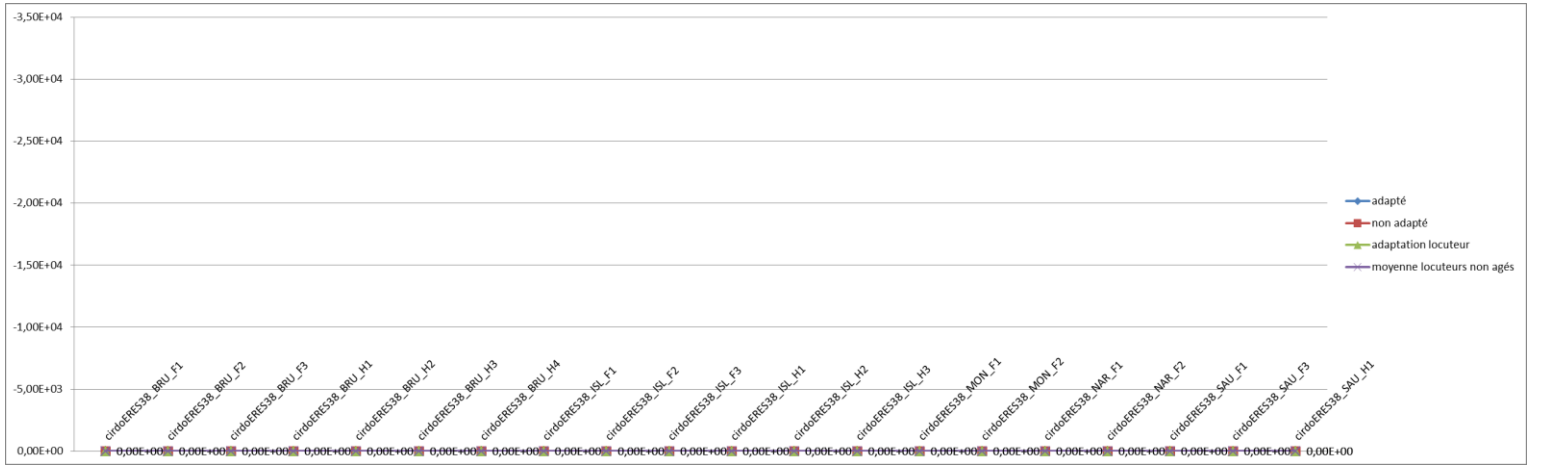
m



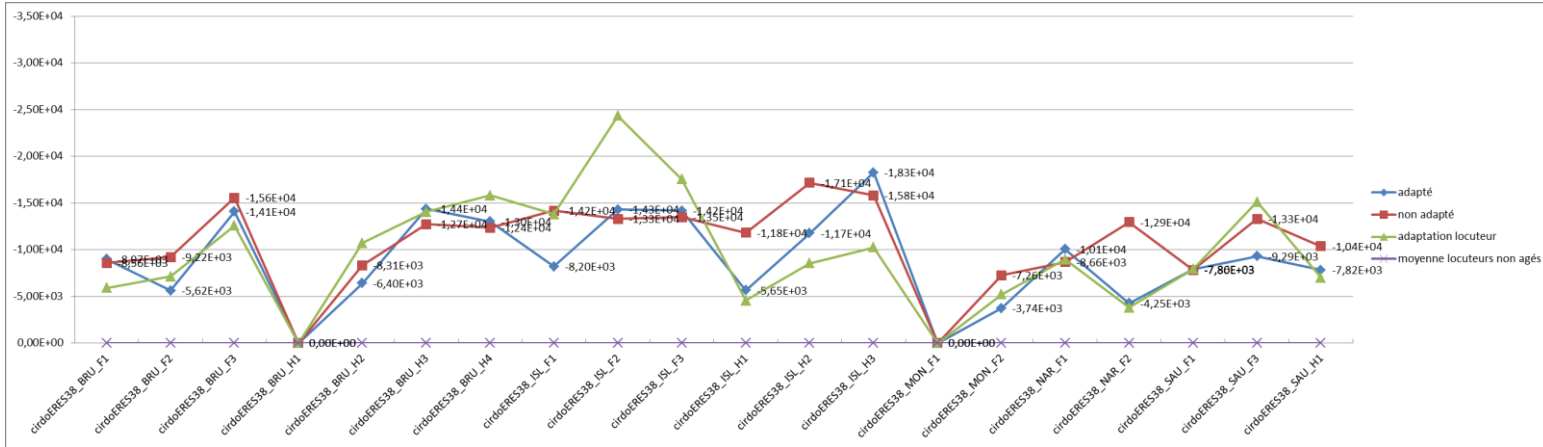
n



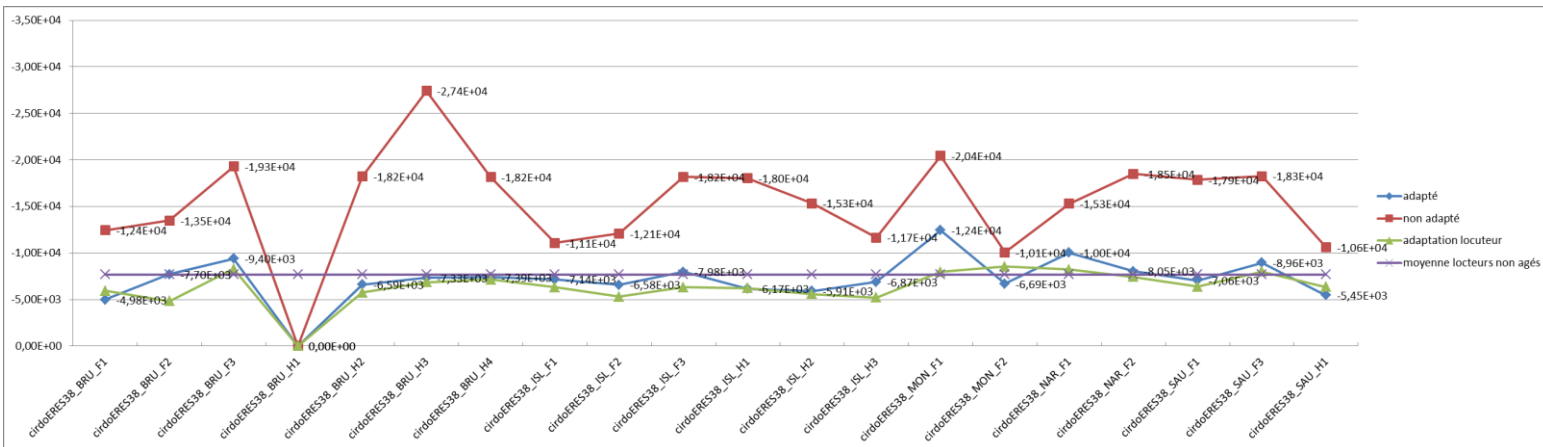
NG



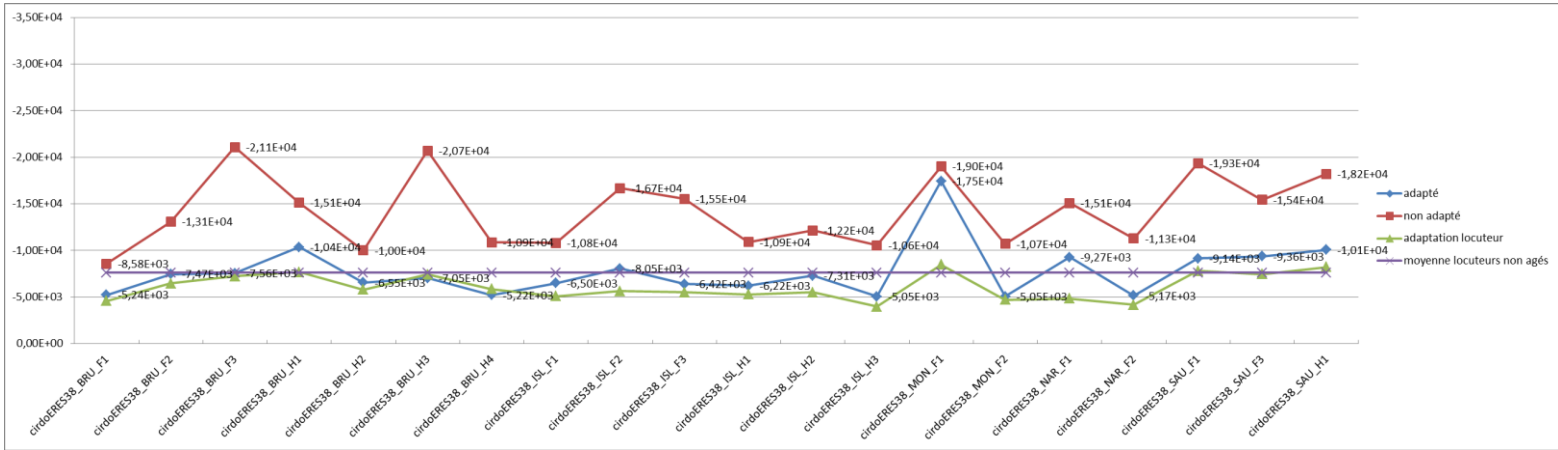
NJ



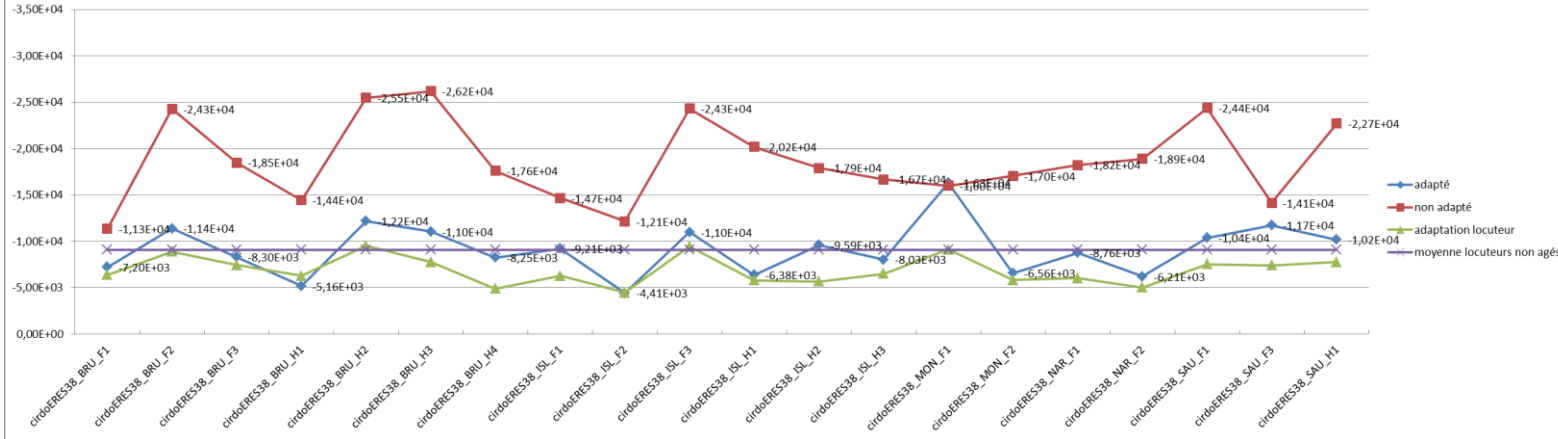
o



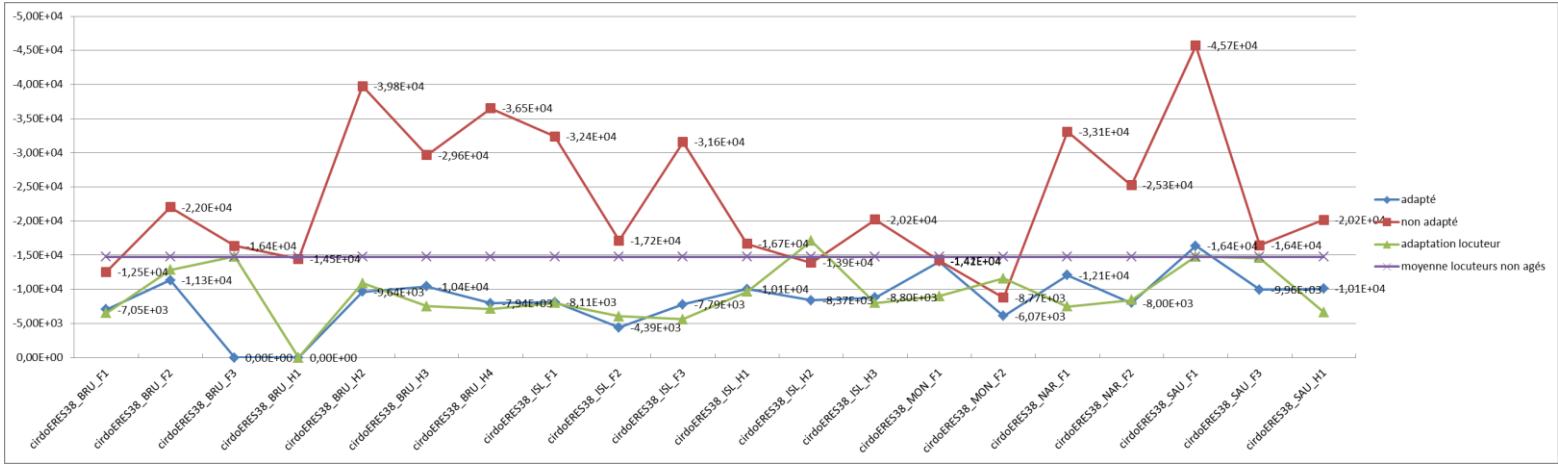
on



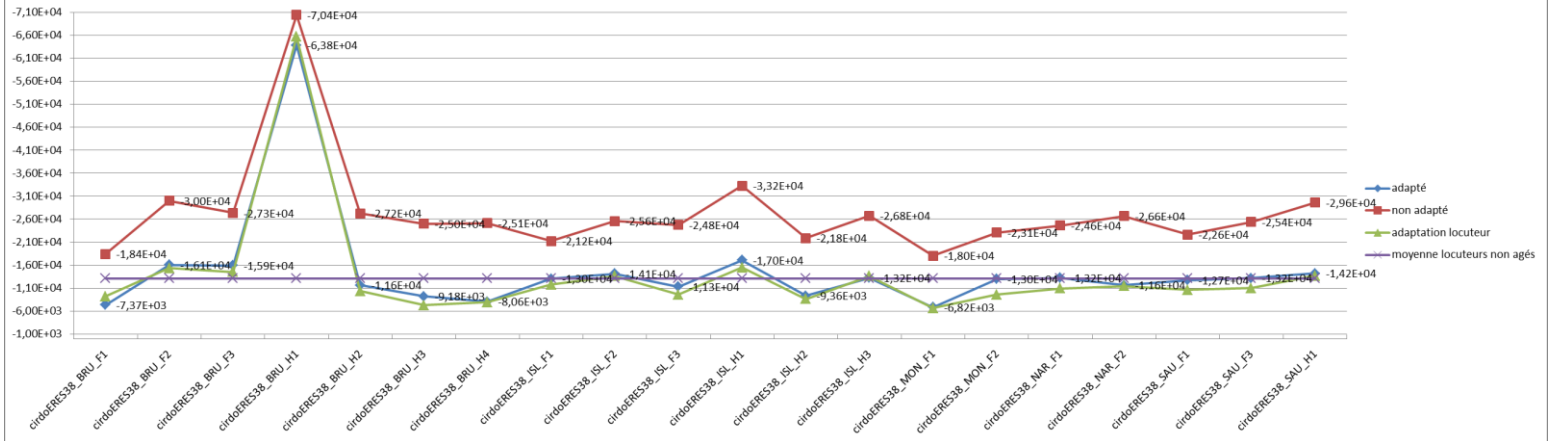
oo



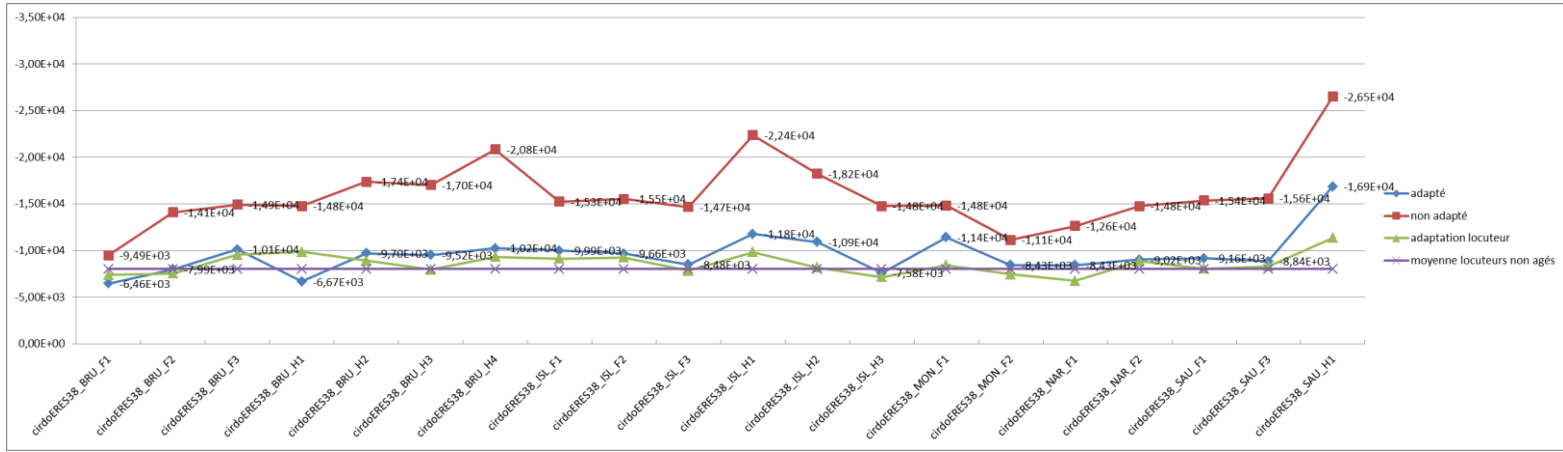
ooo



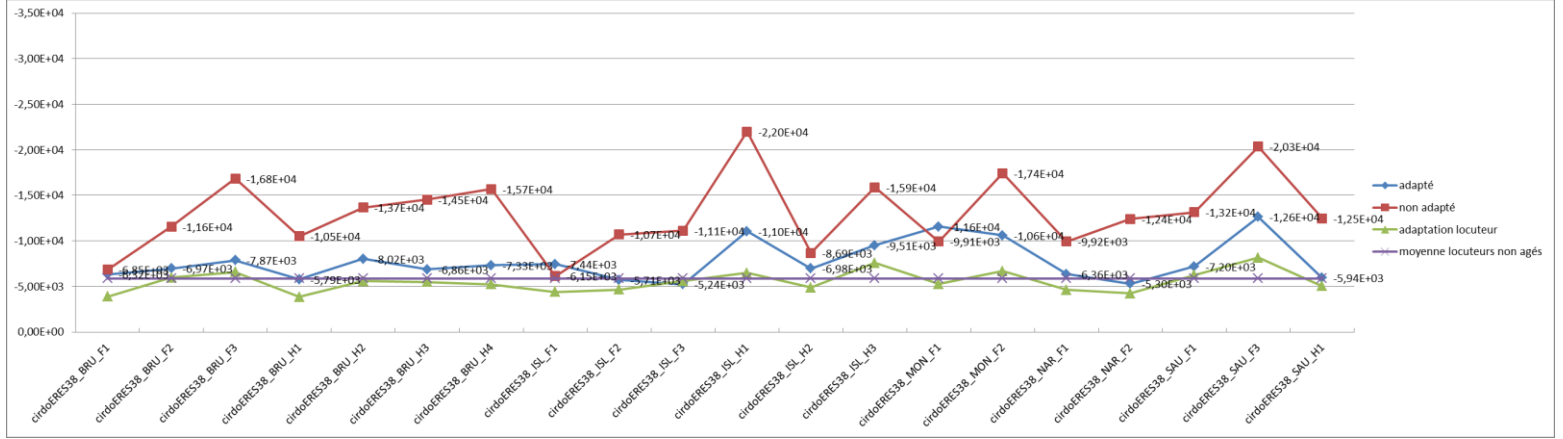
p



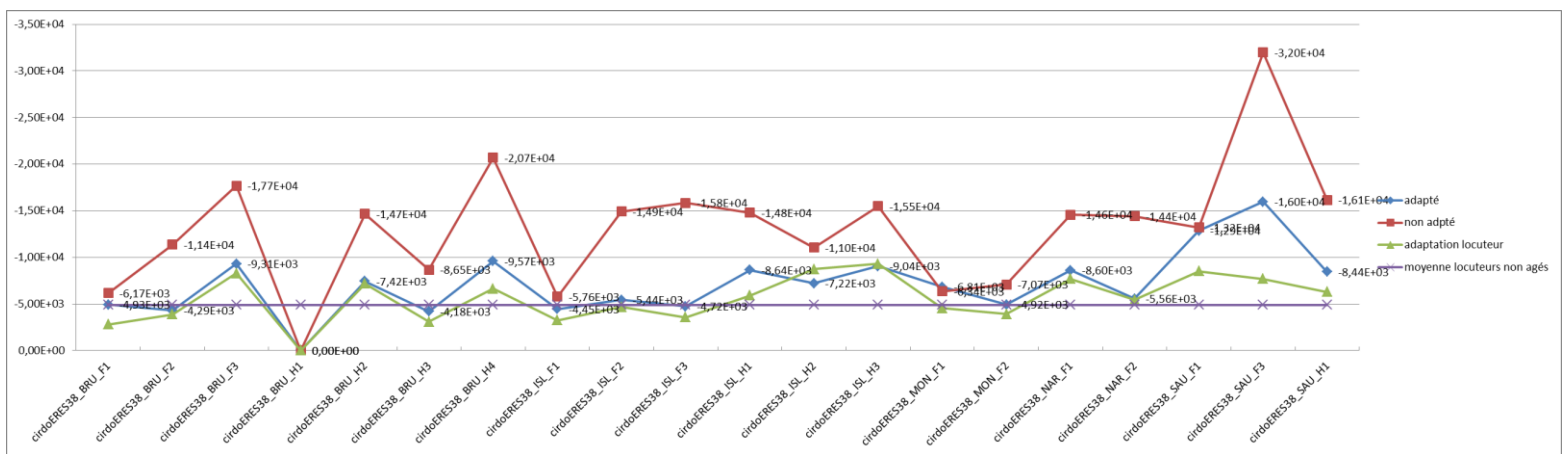
RR



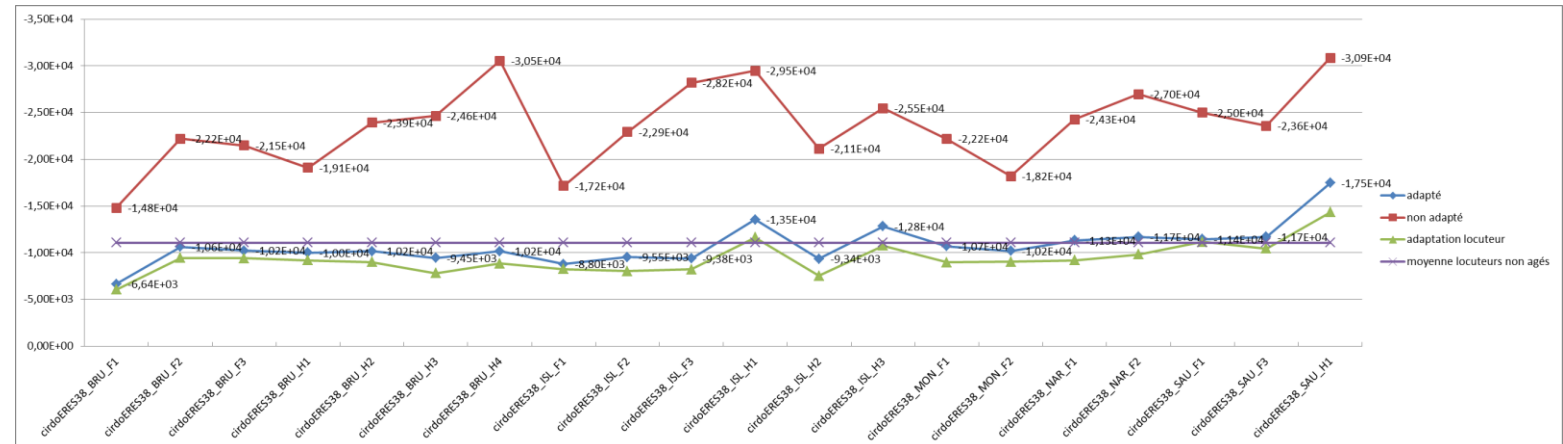
S



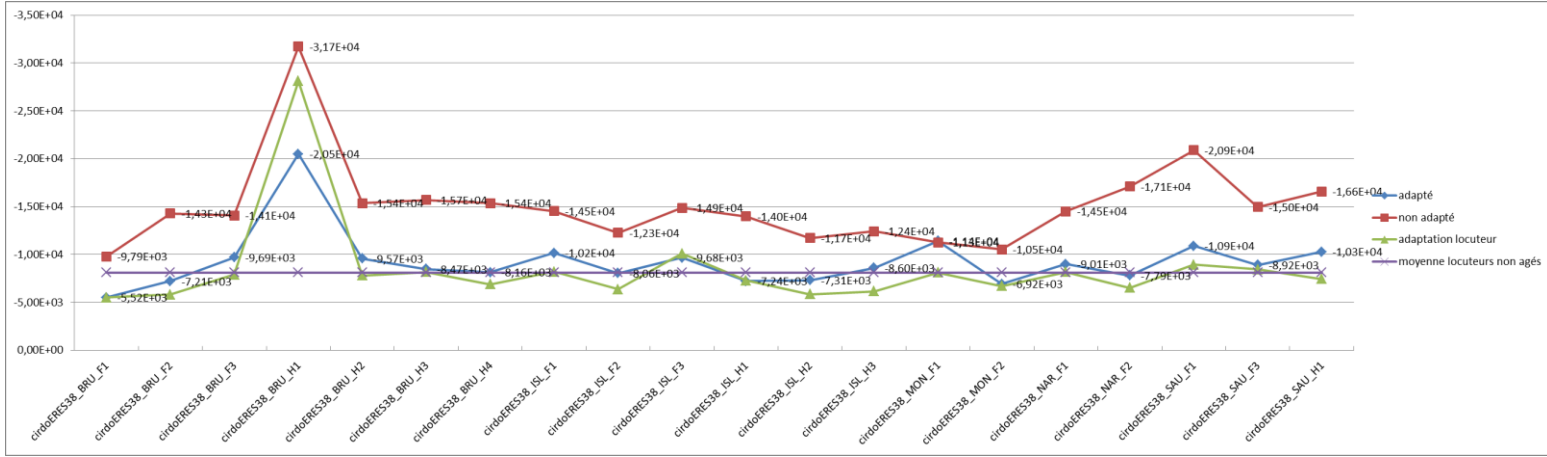
SS



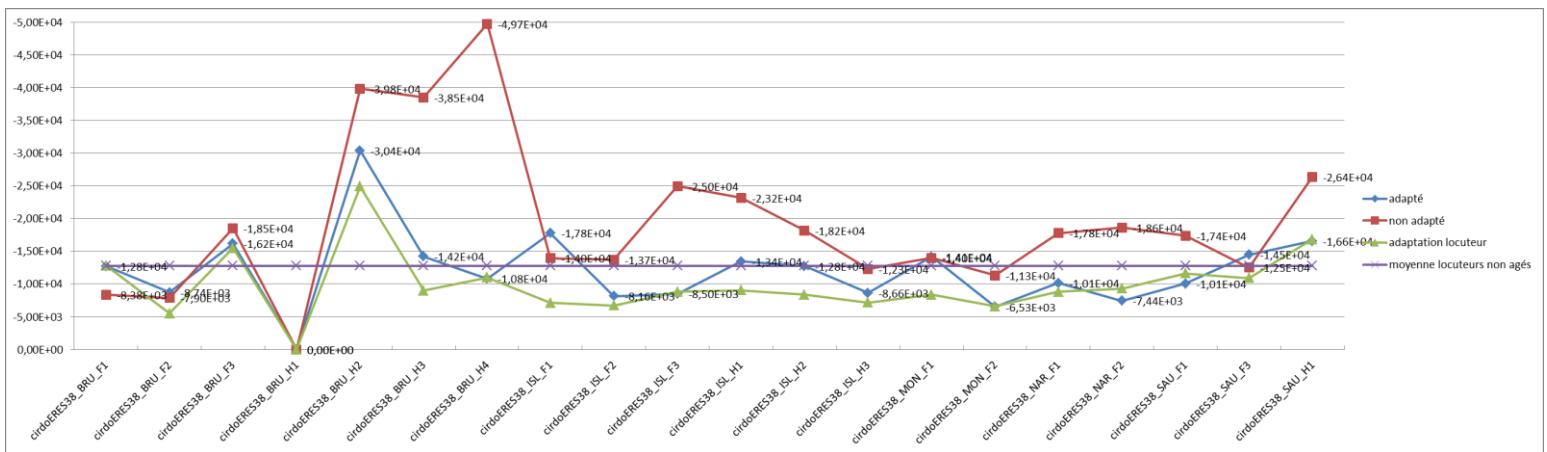
t



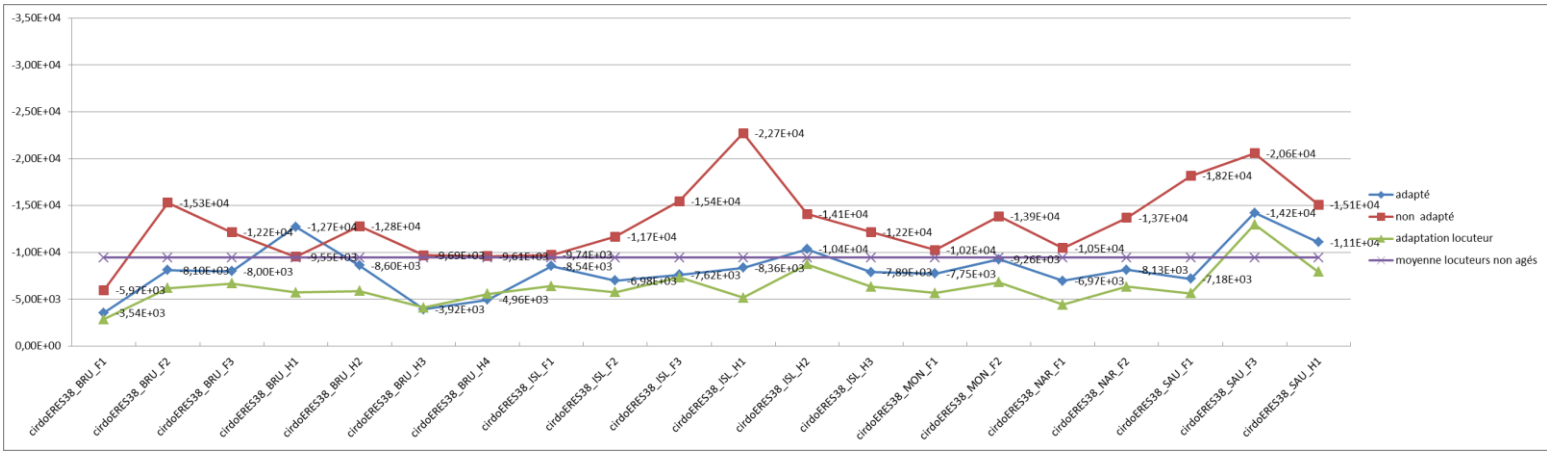
u



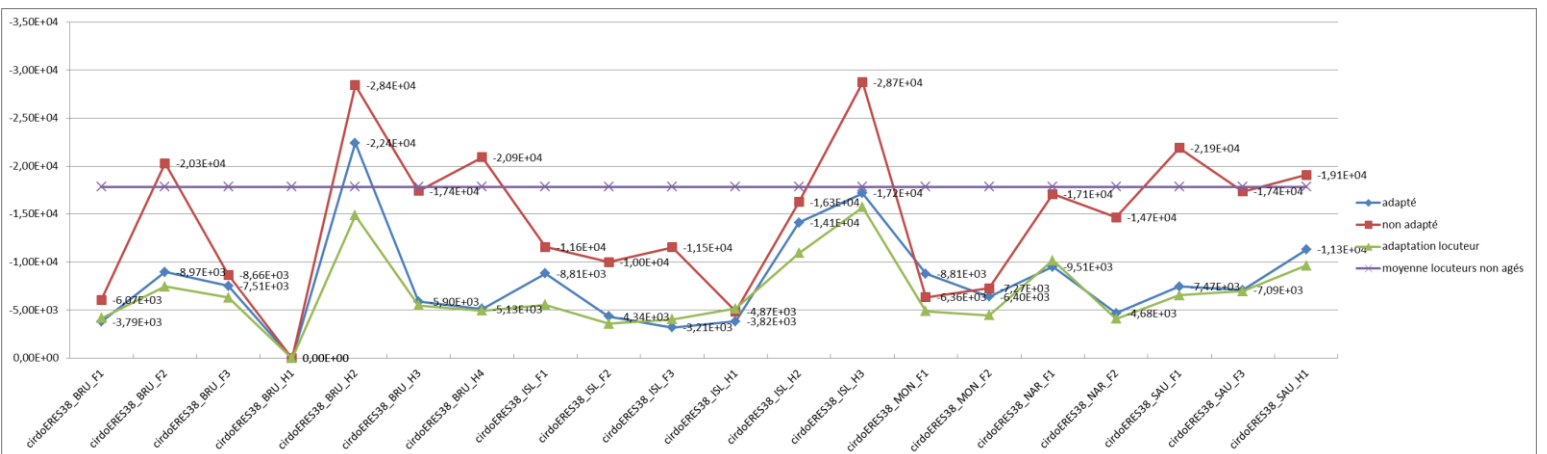
un



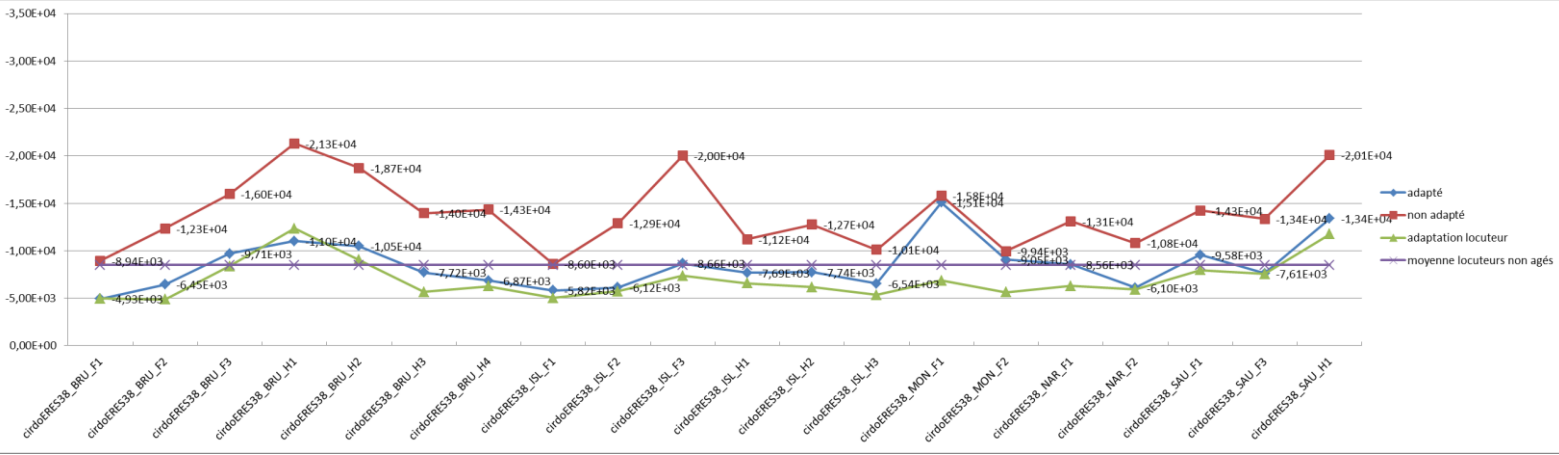
v



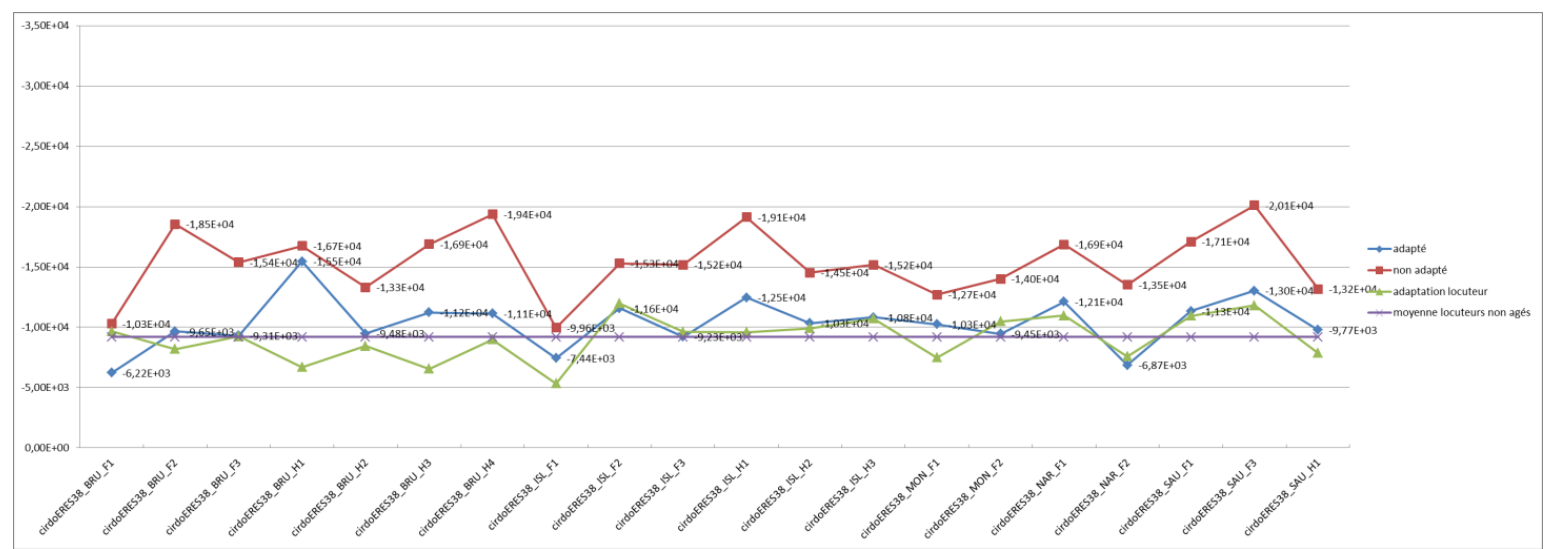
w



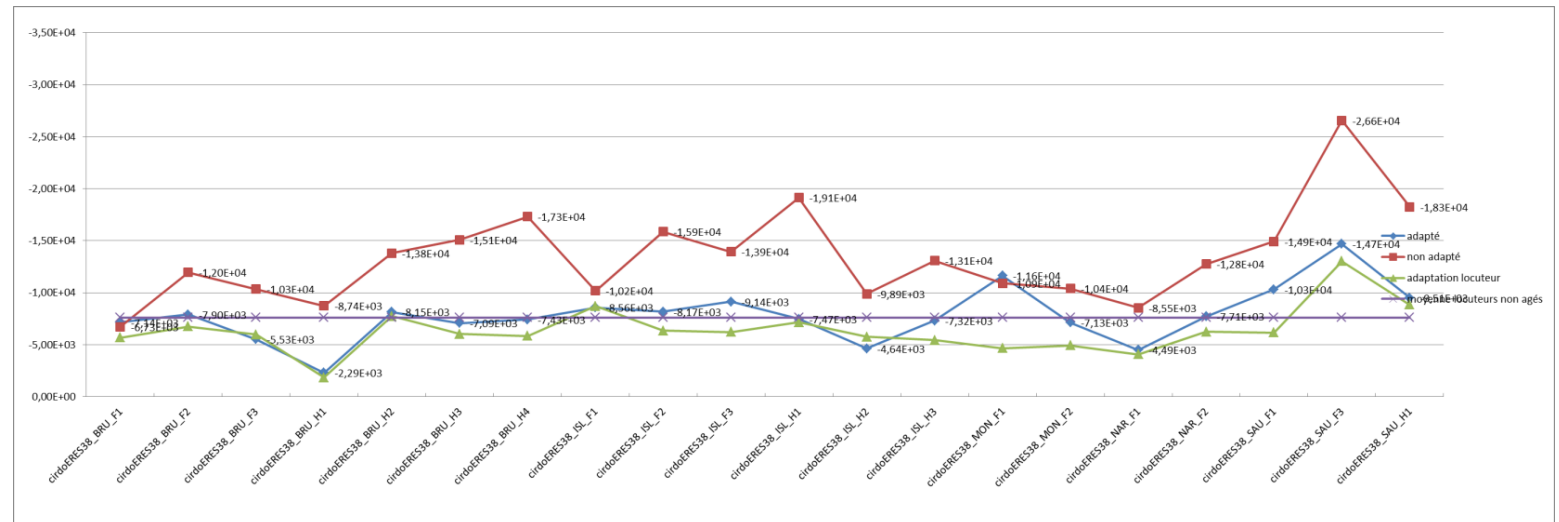
y



z



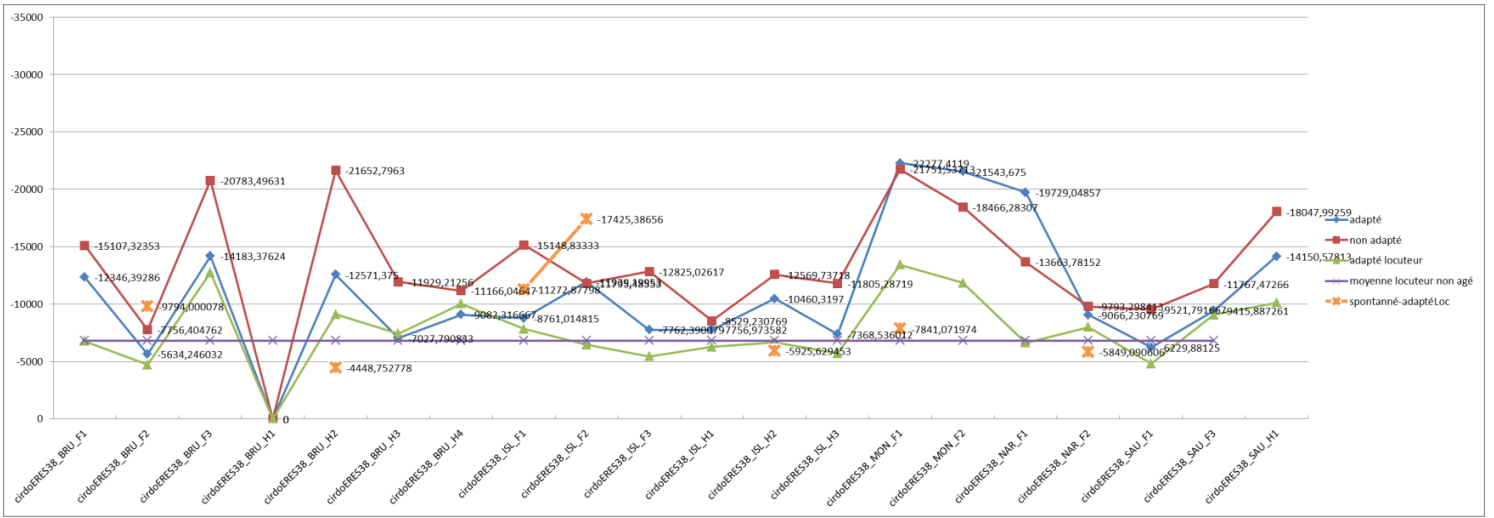
zz



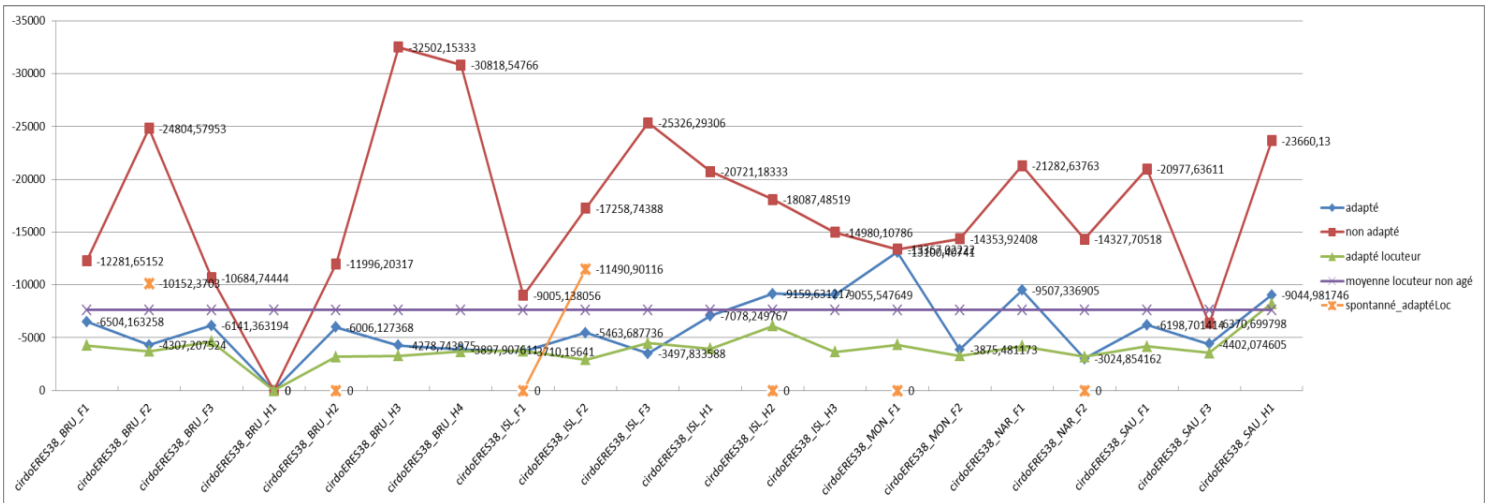
**8.7 Diagrammes représentant pour chaque phonème
et par locuteur la courbe des scores d'alignements
en fonction des modèles acoustiques
Comparaison parole lue et spontanée**

Comparaison parole lue vs parole spontanée. ATTENTION, lorsque la valeur = 0 le score n'est pas parfait mais nul (absence de données)

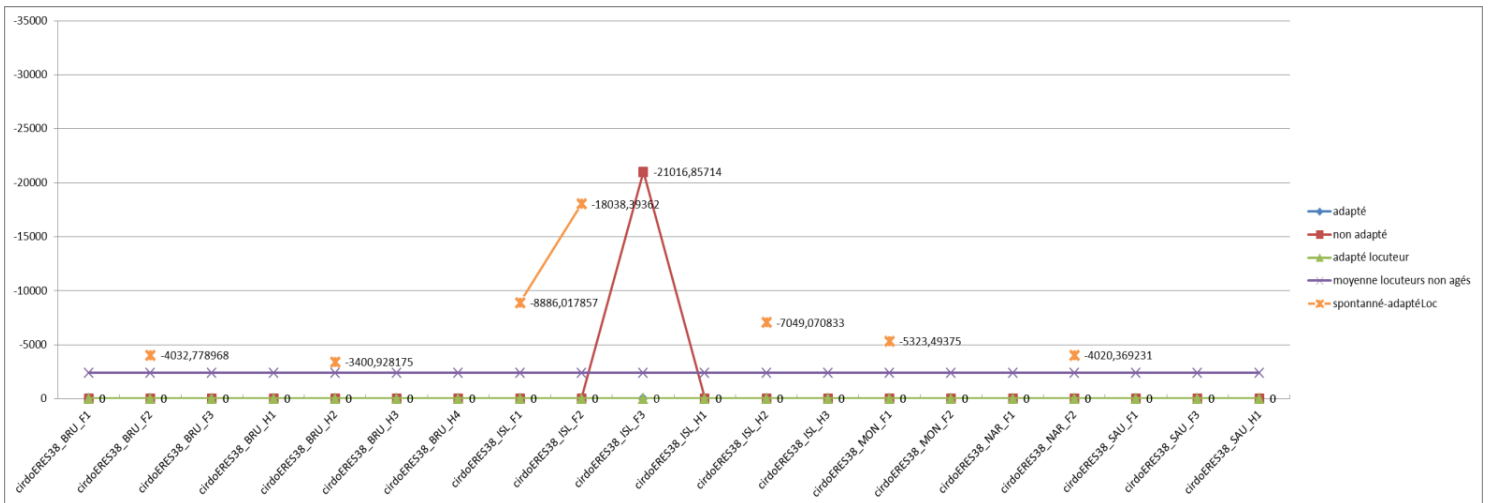
2



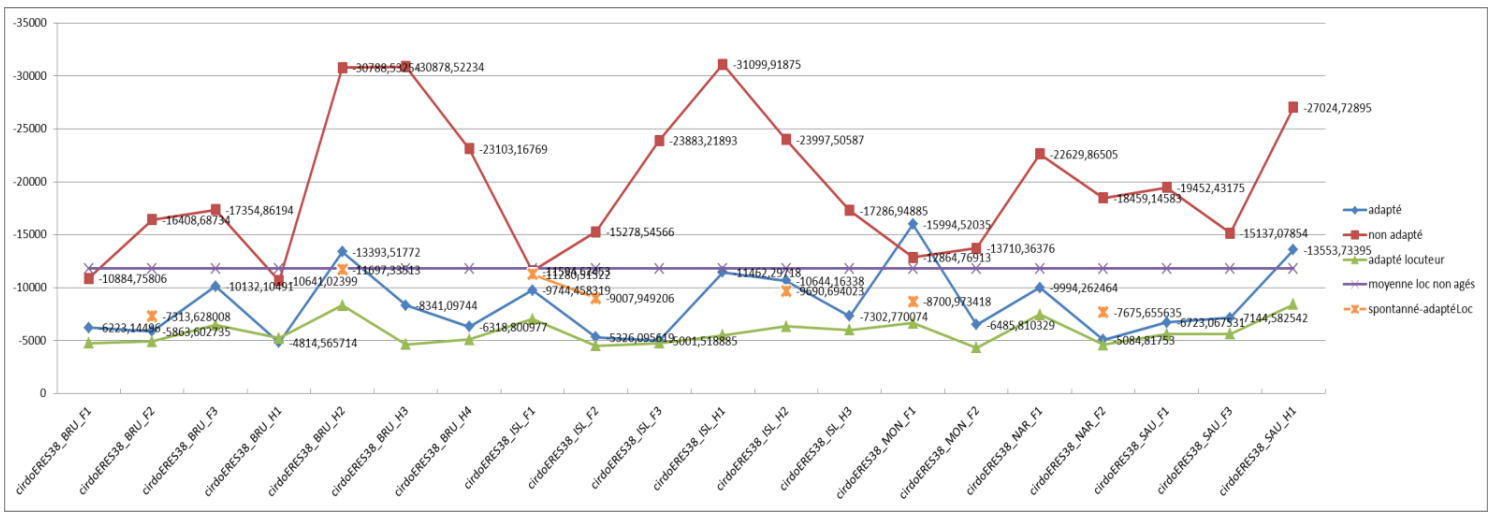
9



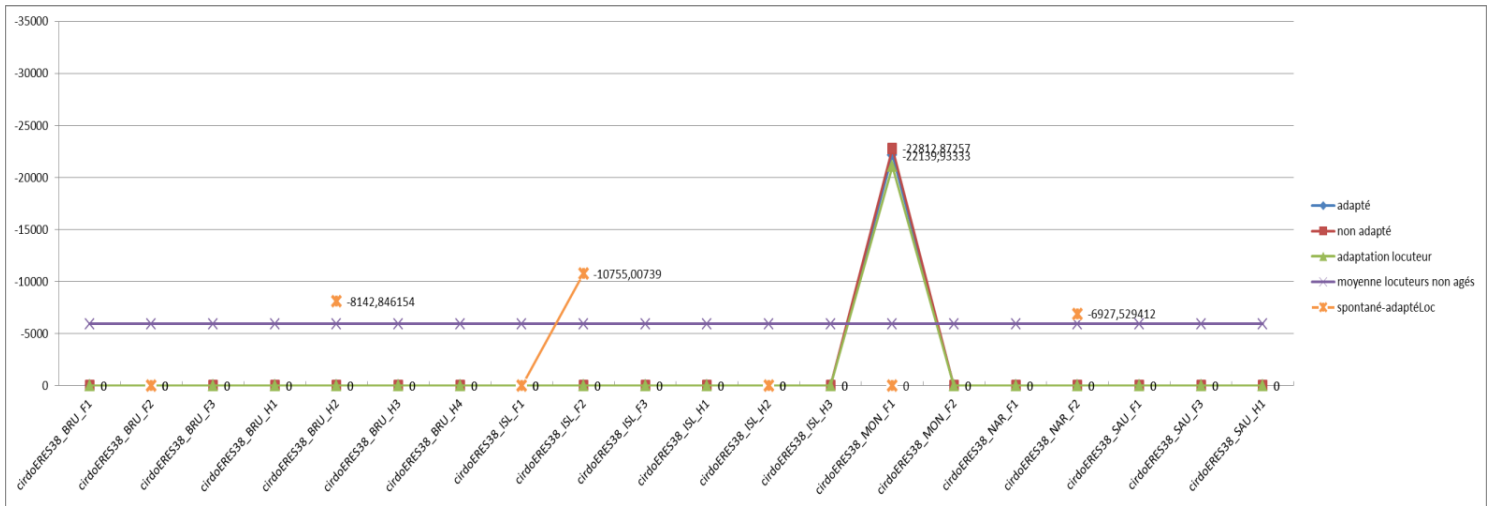
29



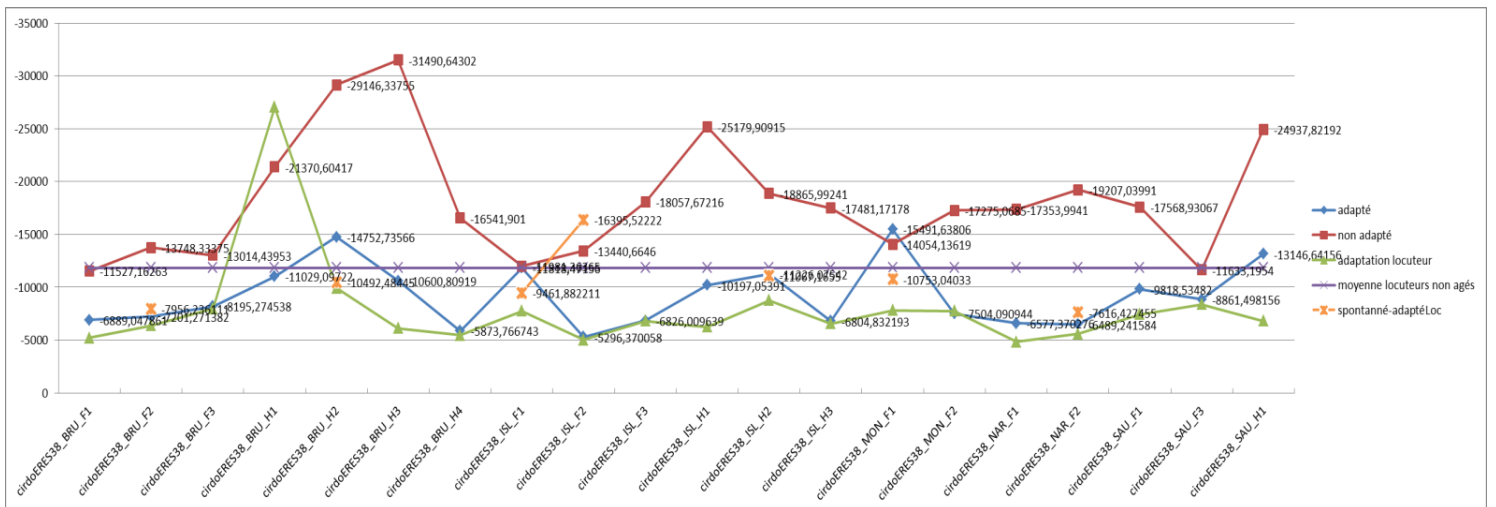
a



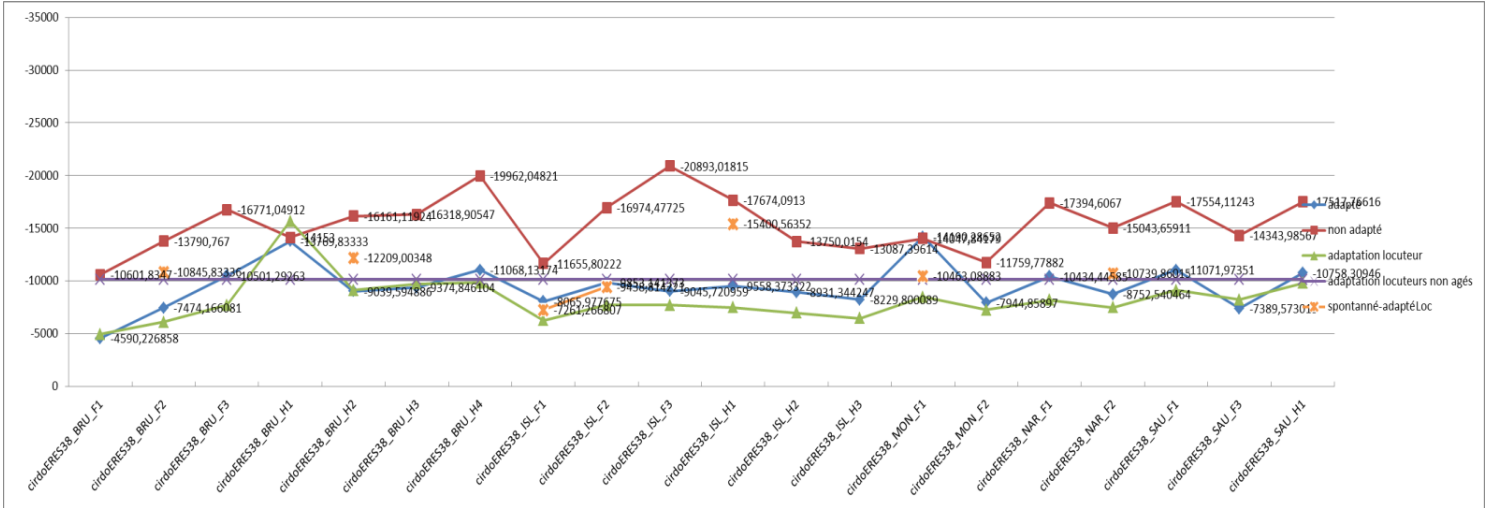
AA



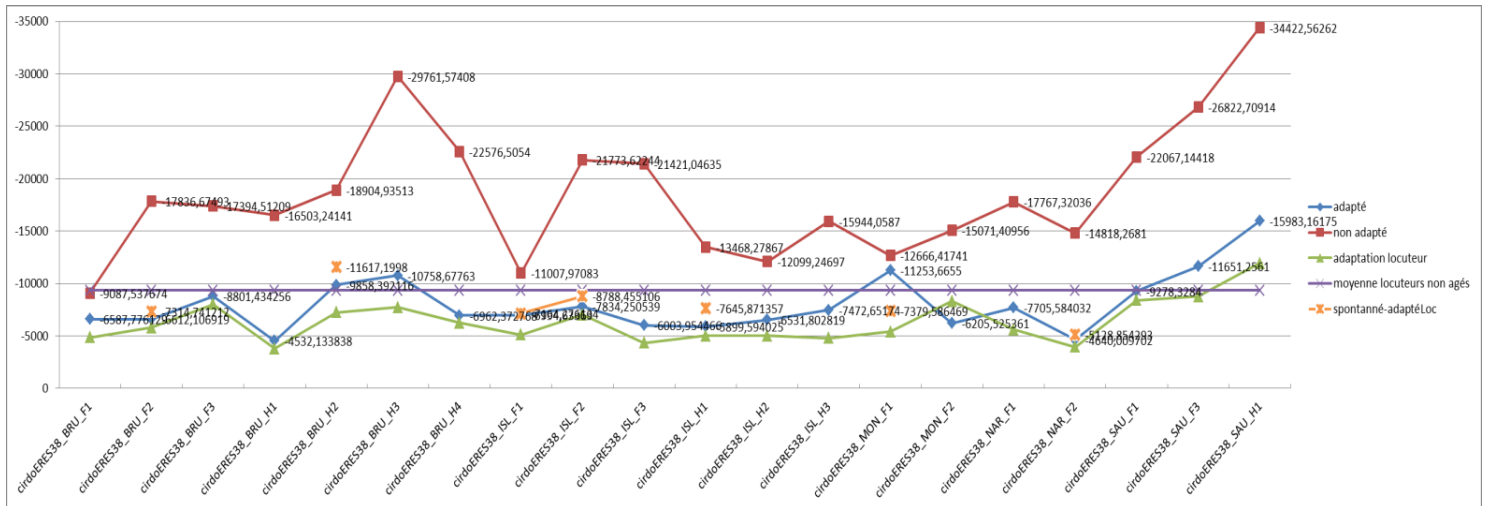
aAA



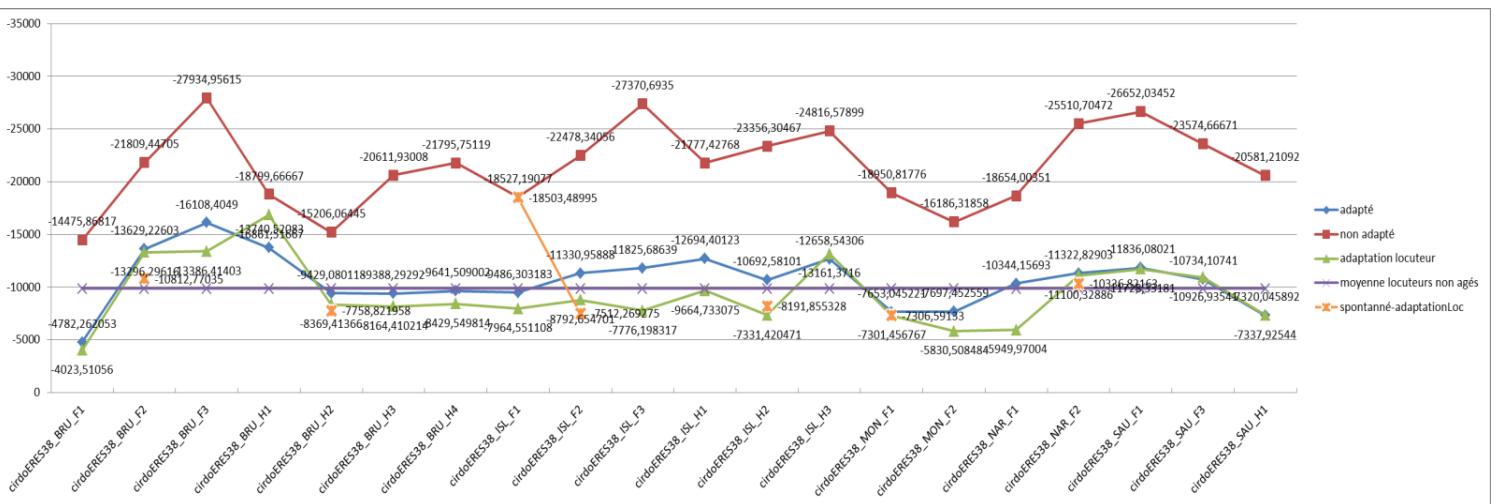
aeA



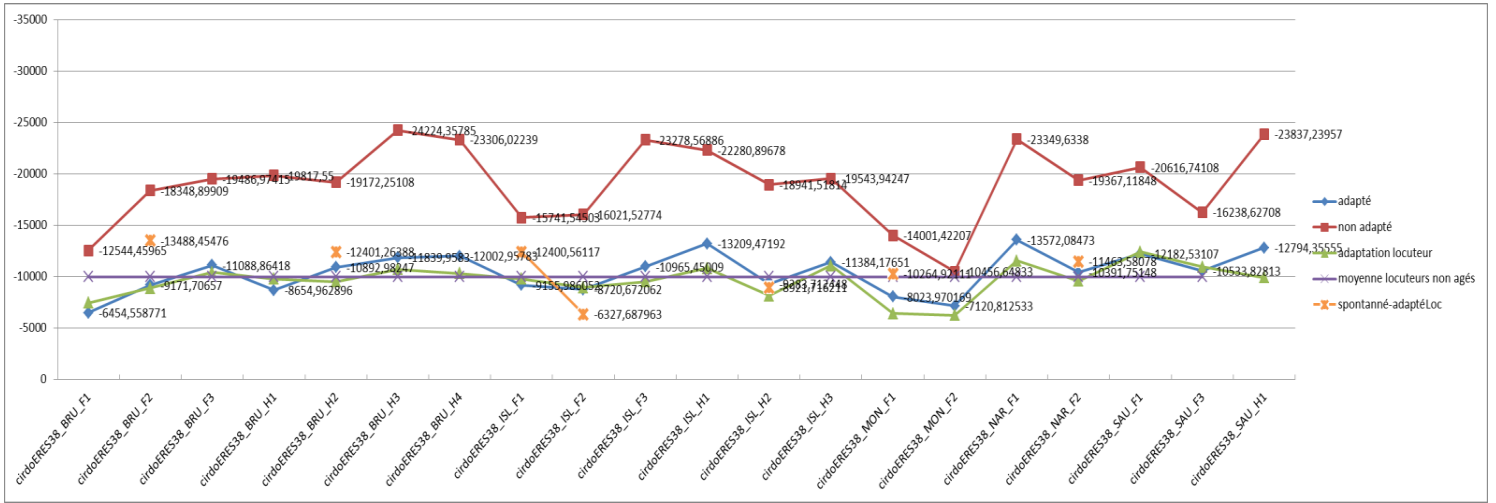
an



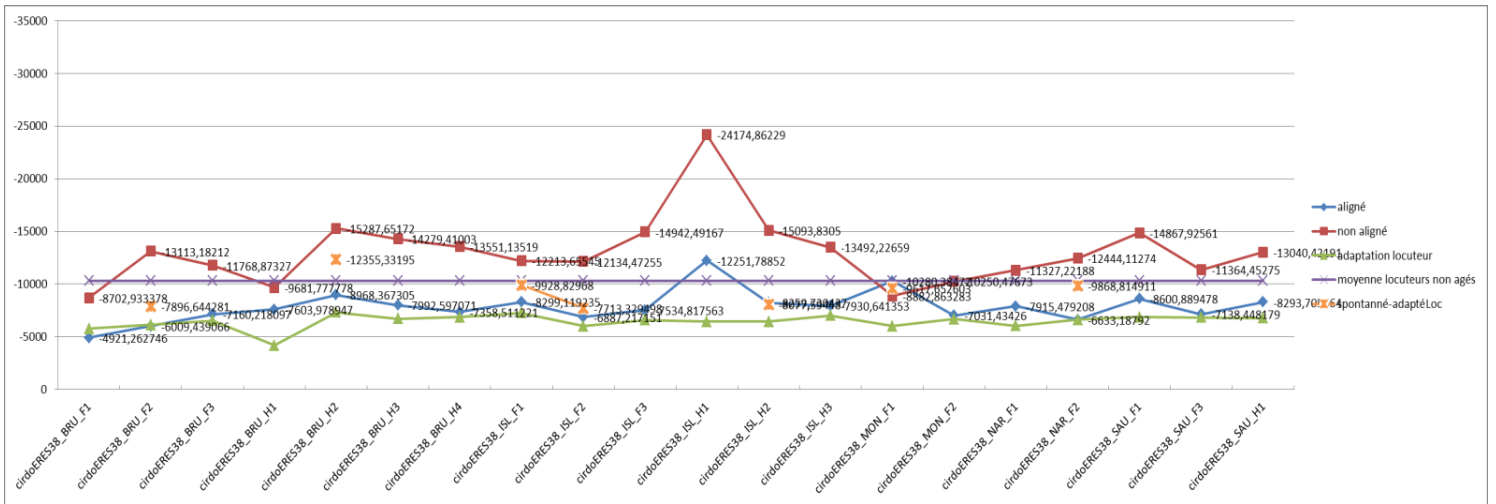
b



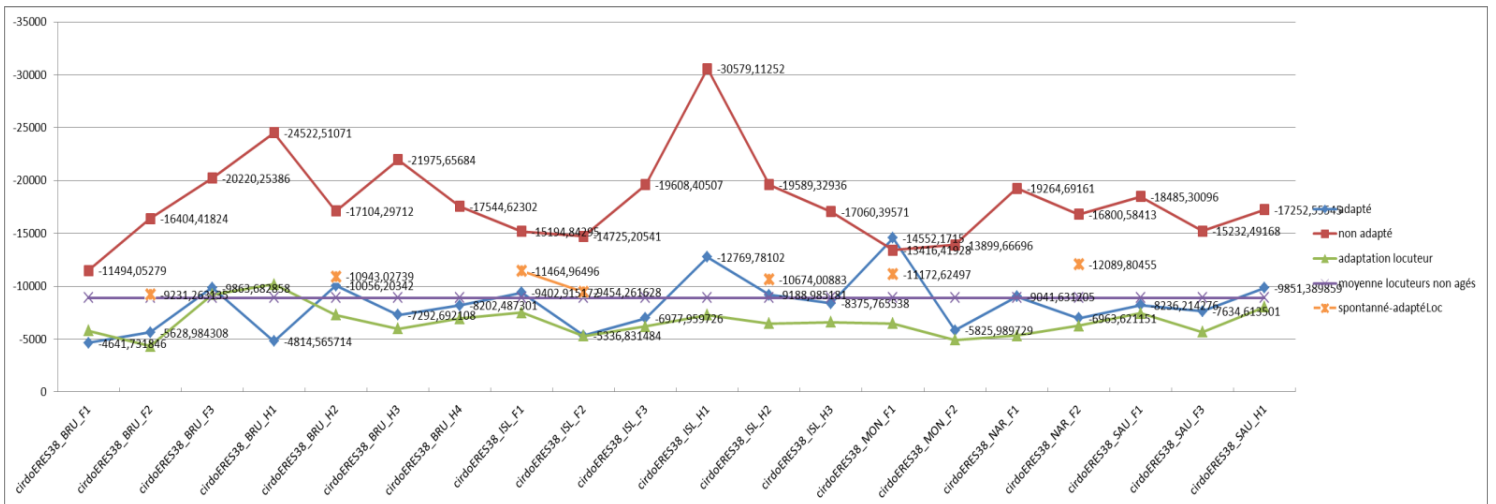
d



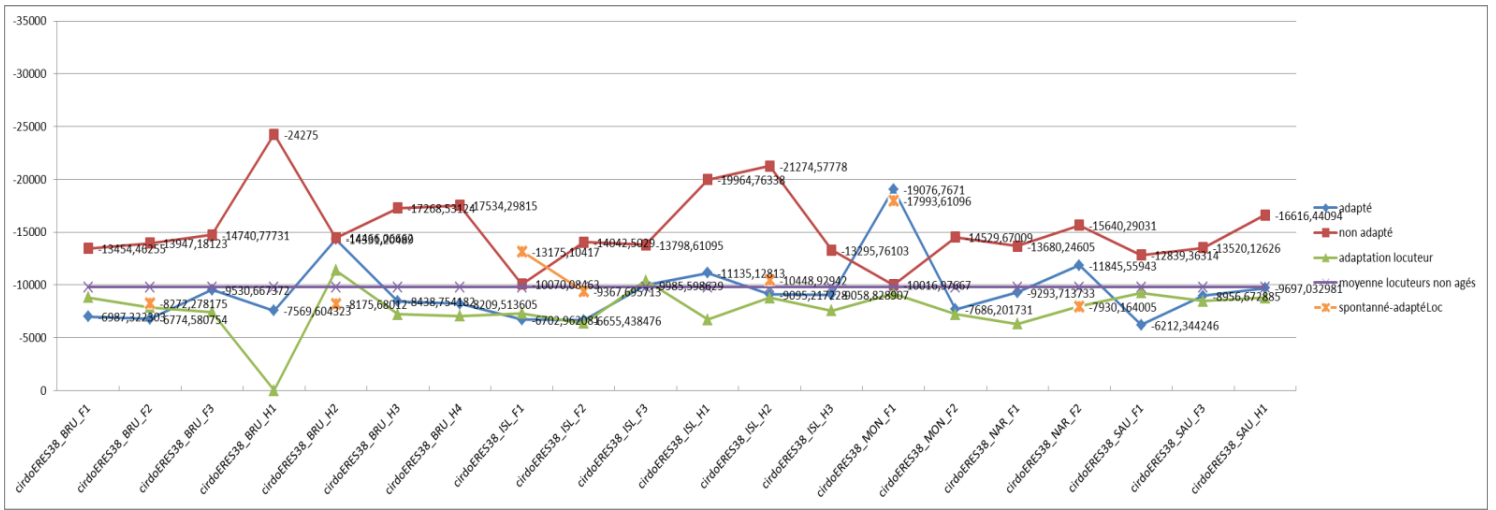
e



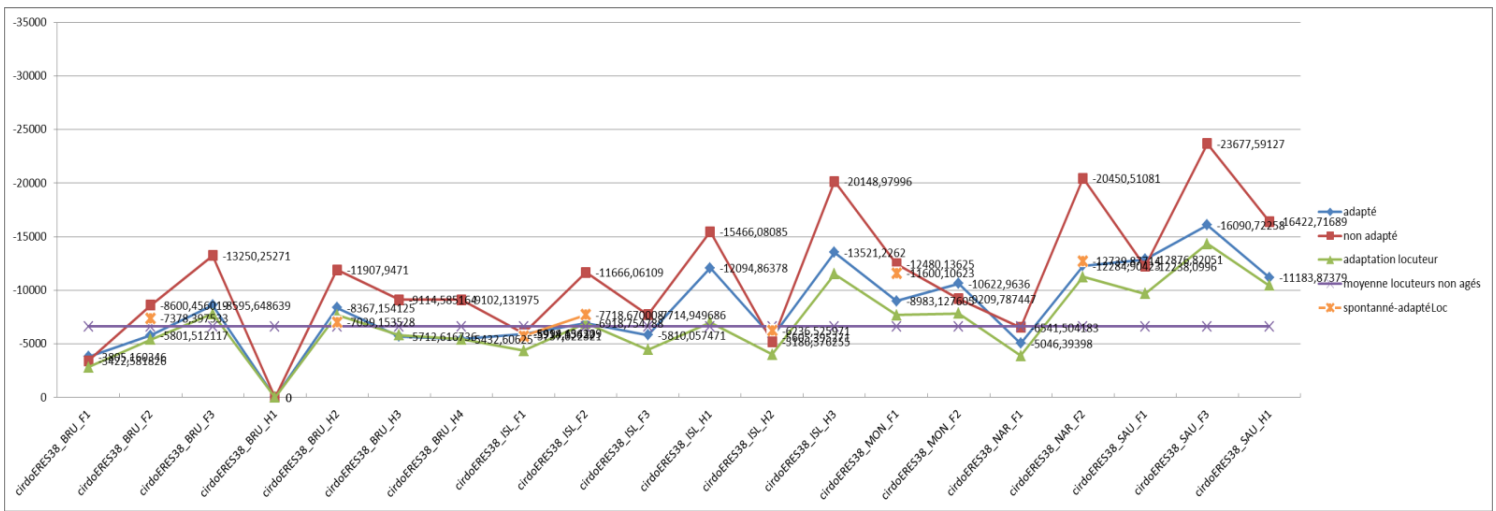
EE



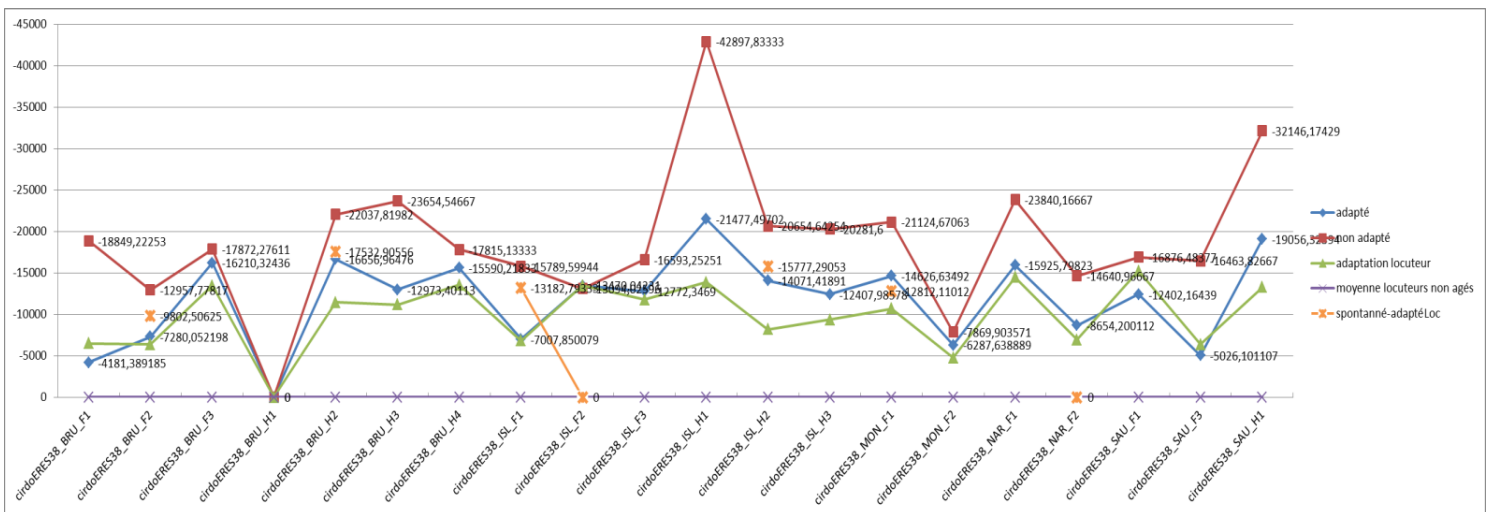
eEE



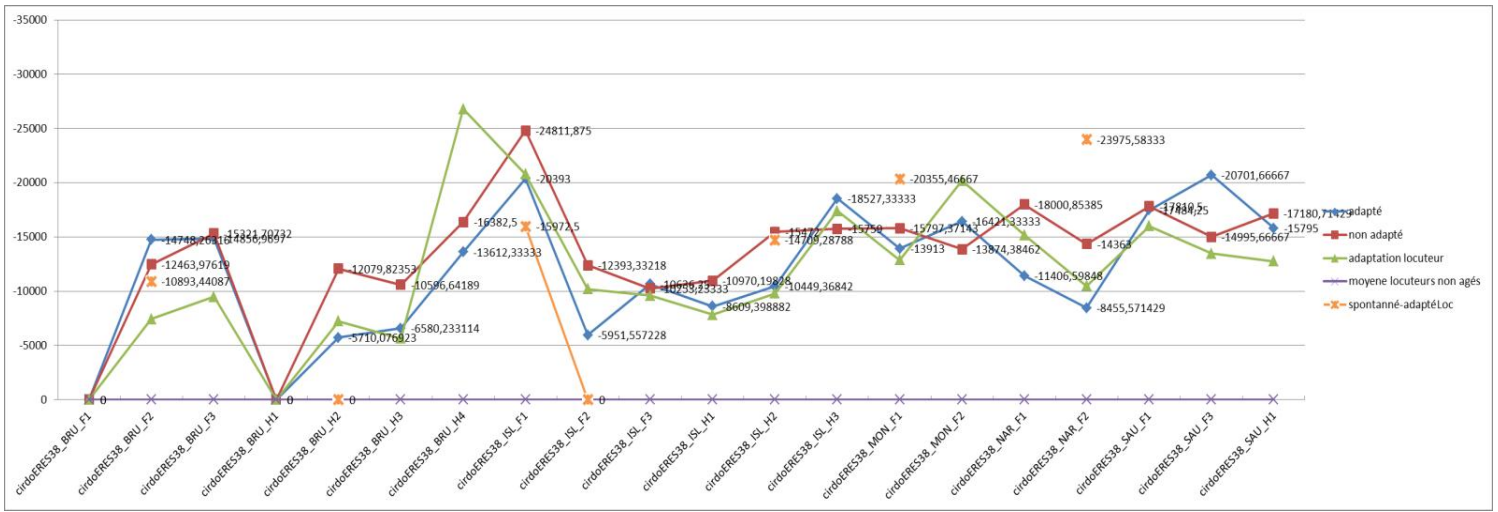
f



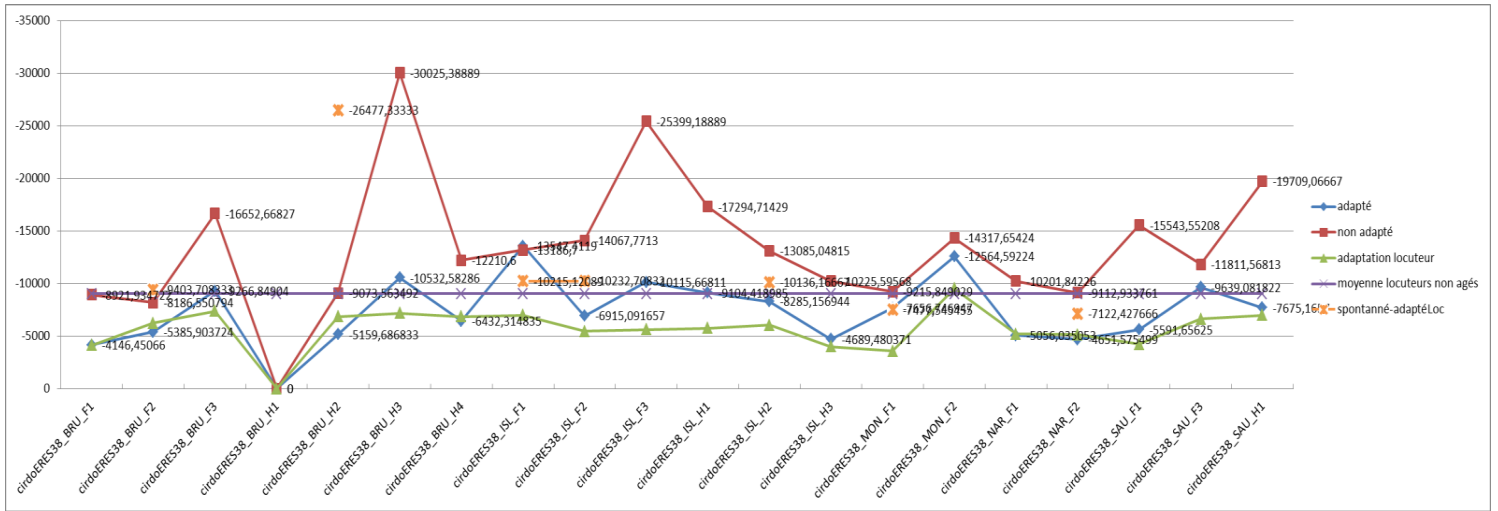
g



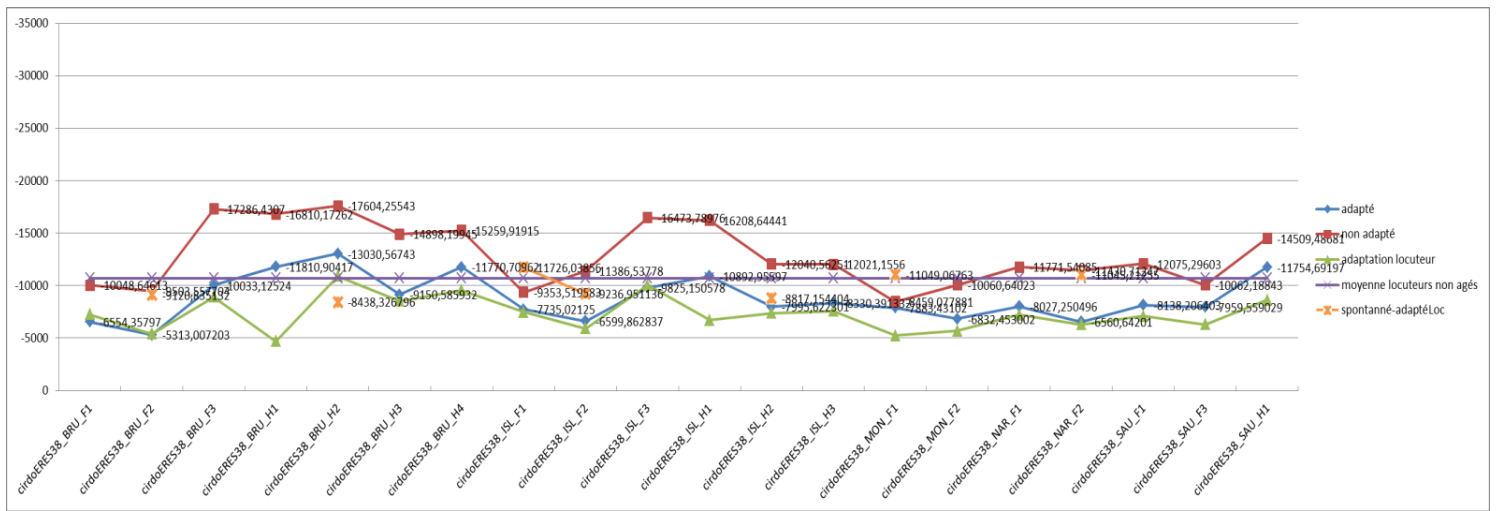
h



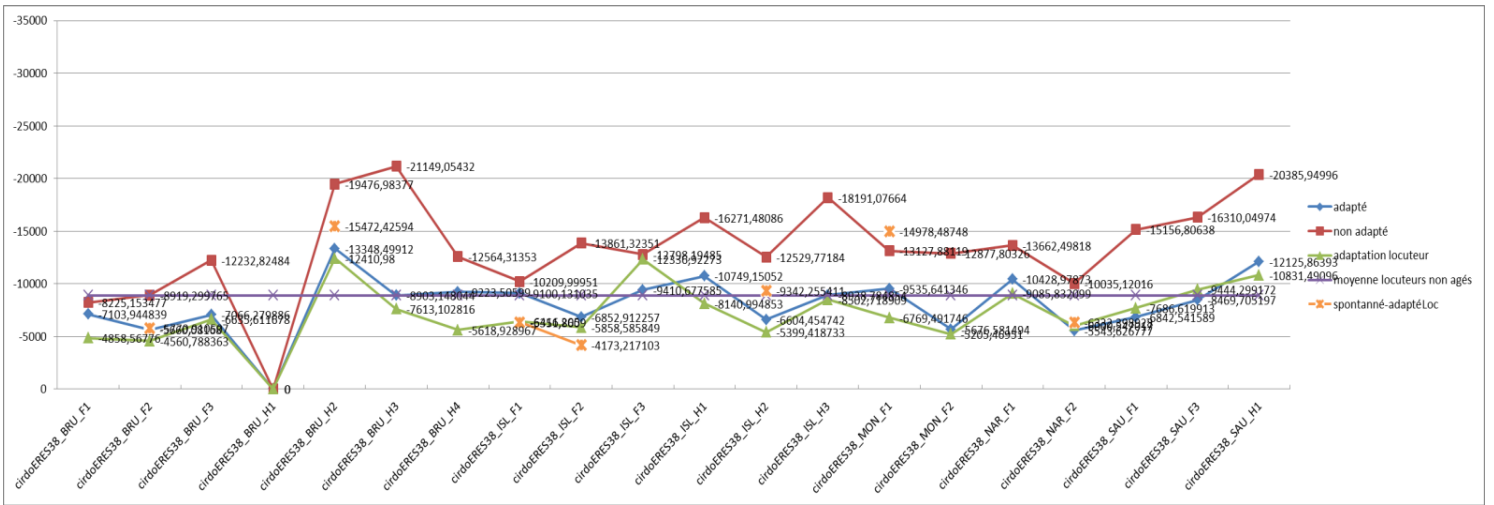
HH



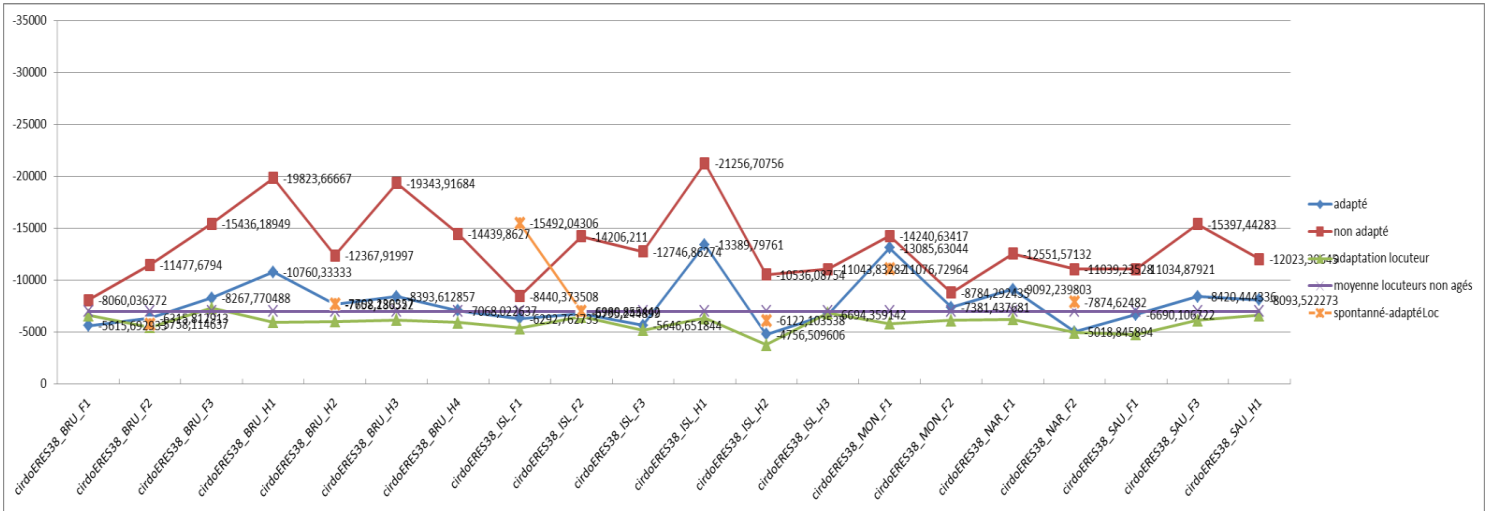
i



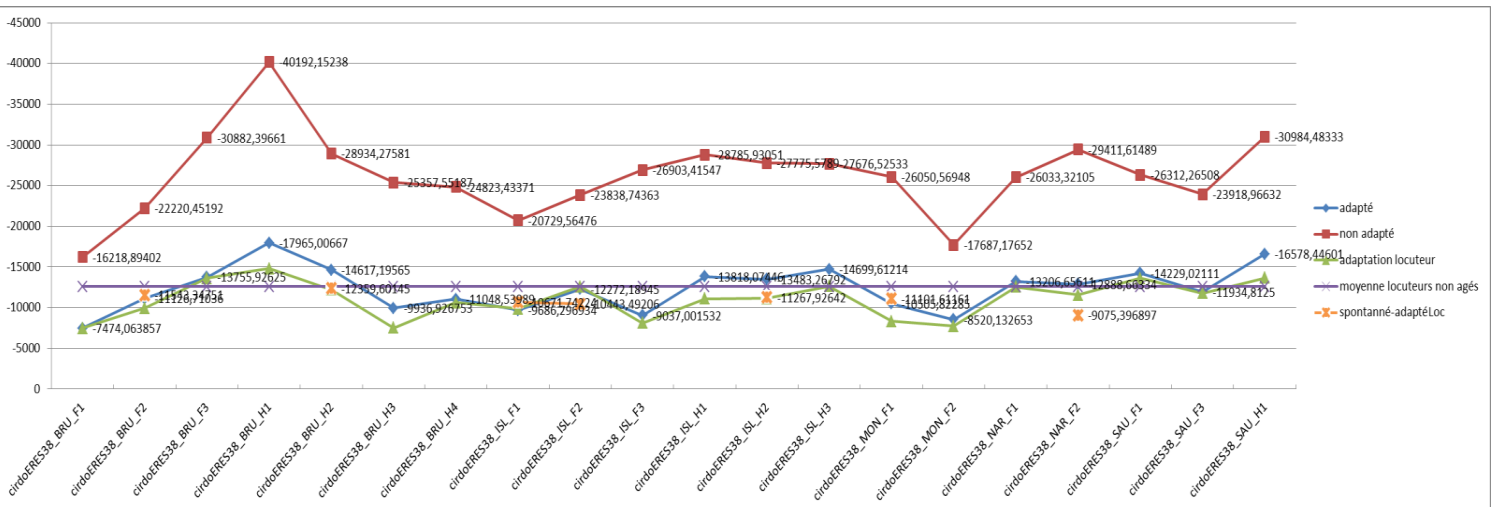
in



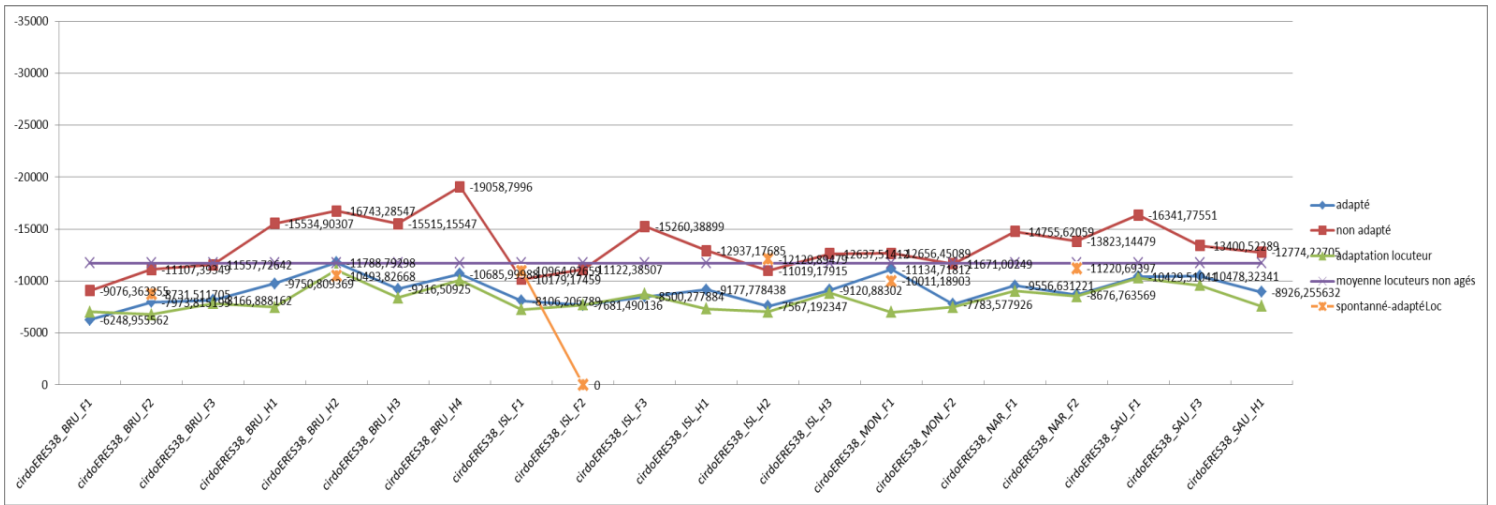
j



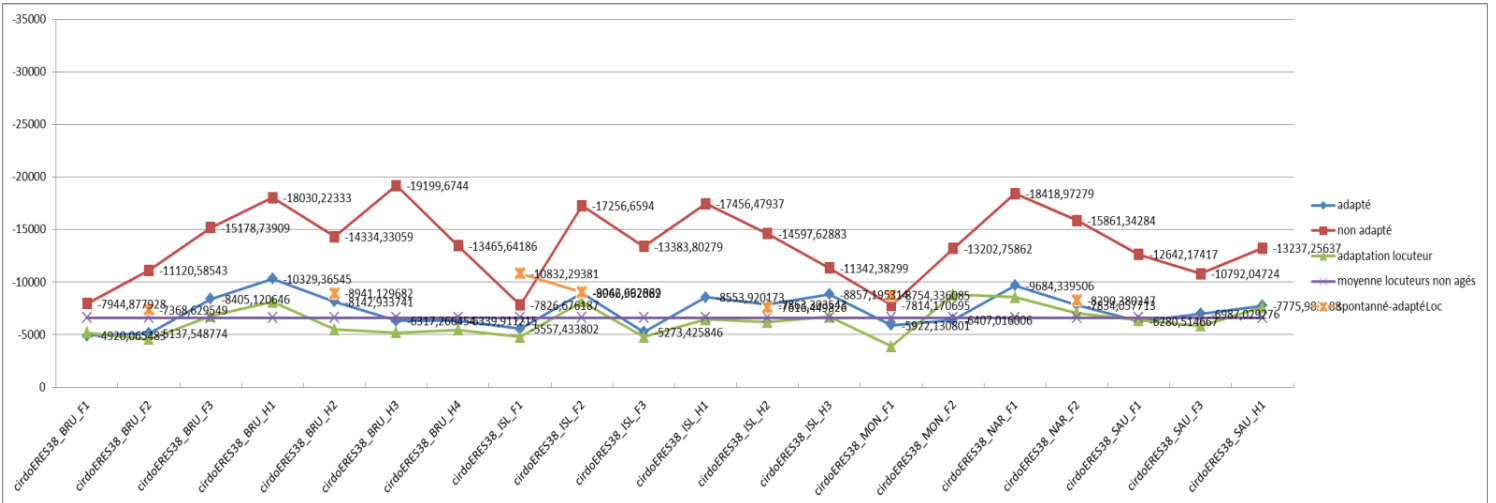
k



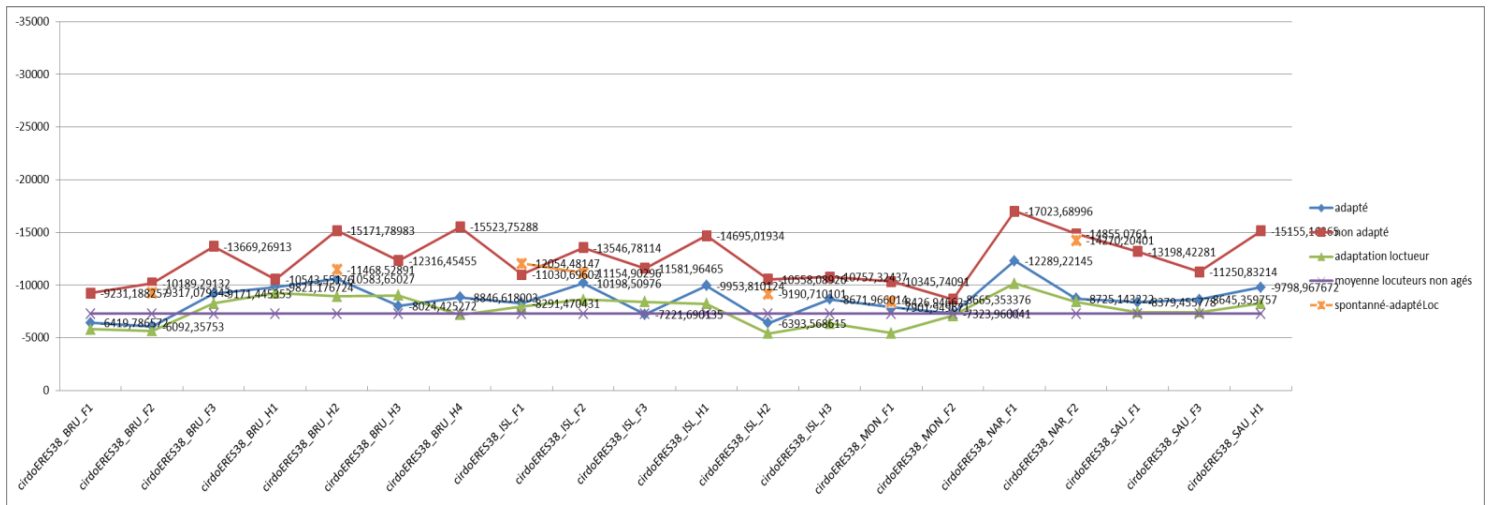
l



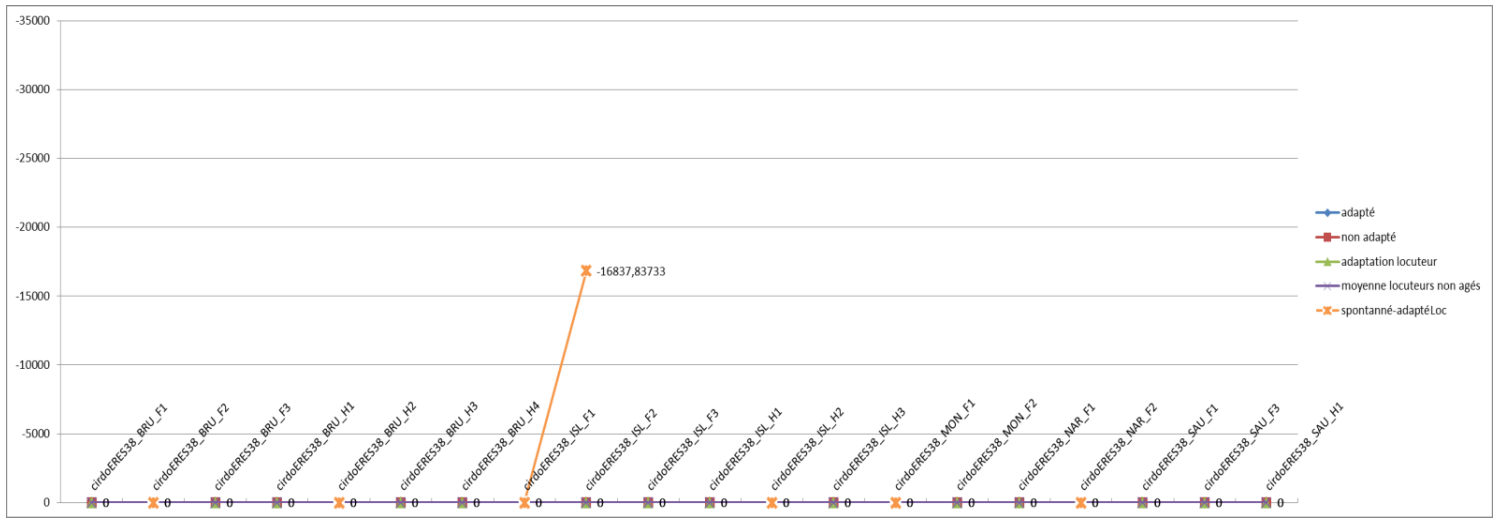
m



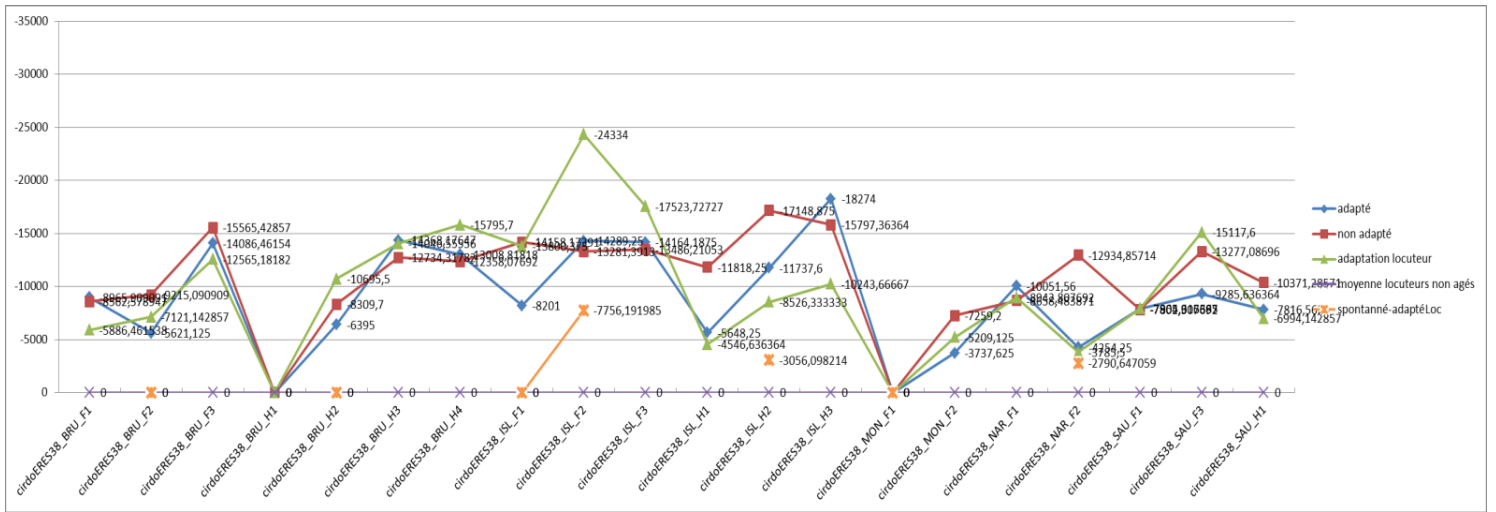
n



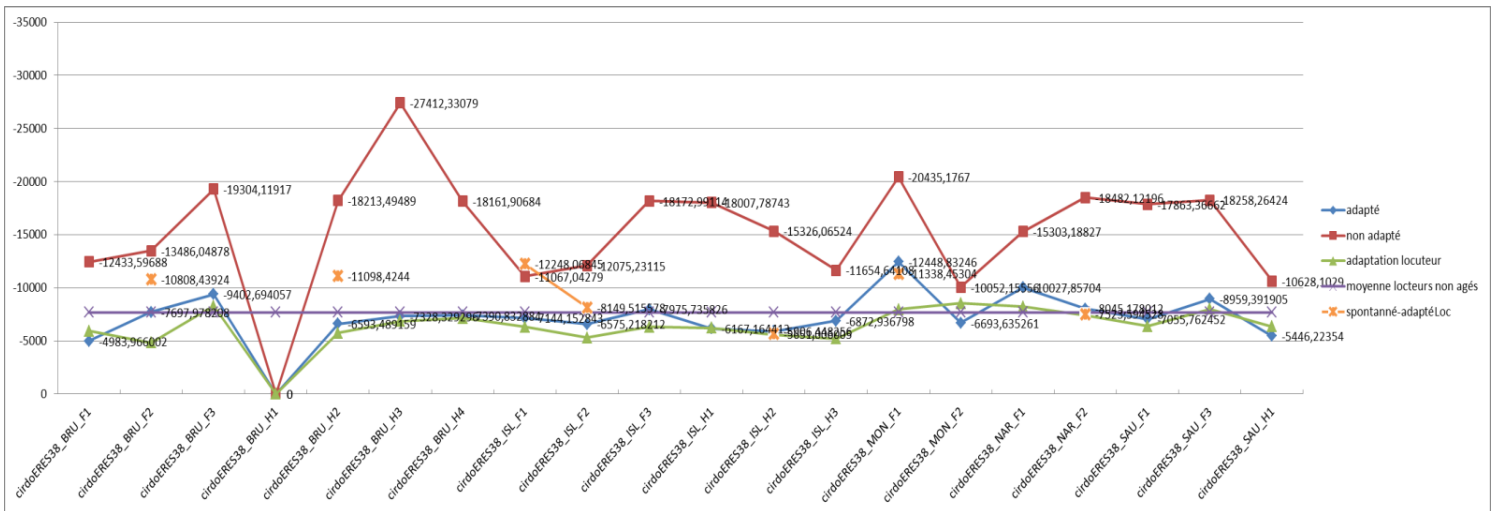
NG



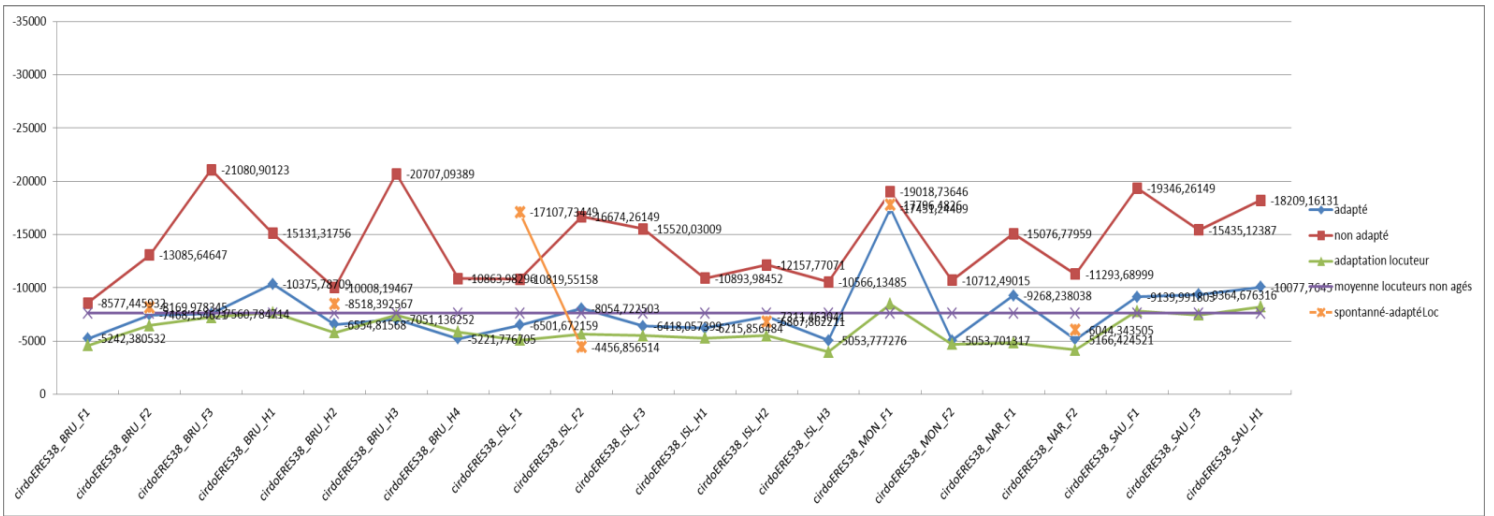
NJ



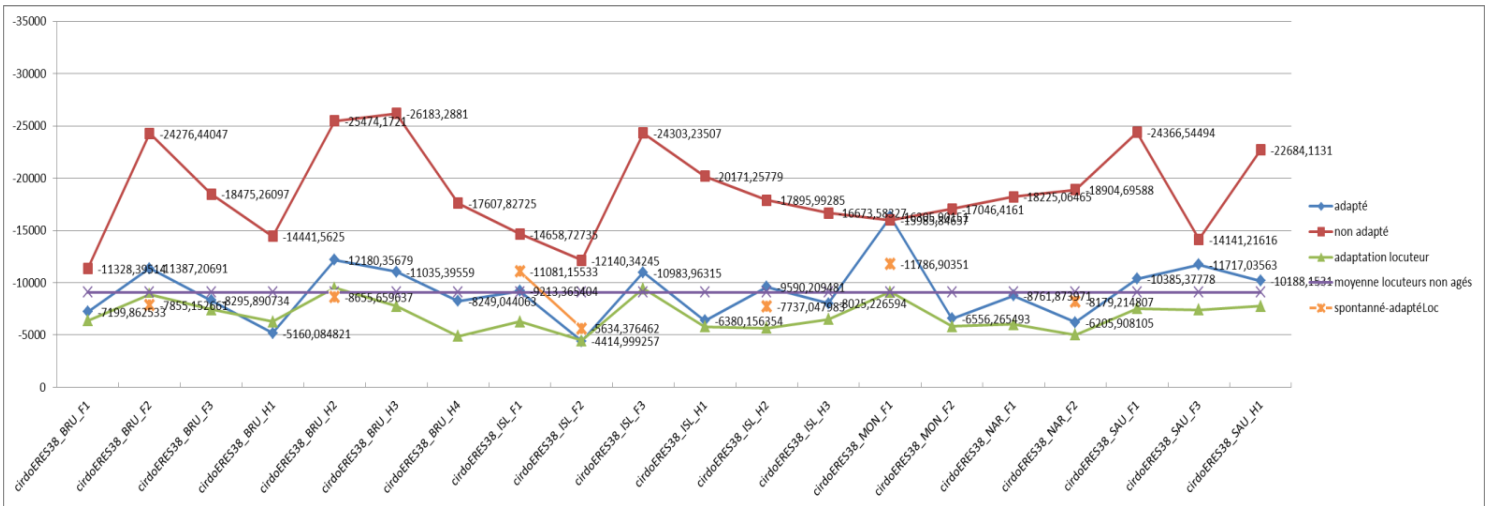
O



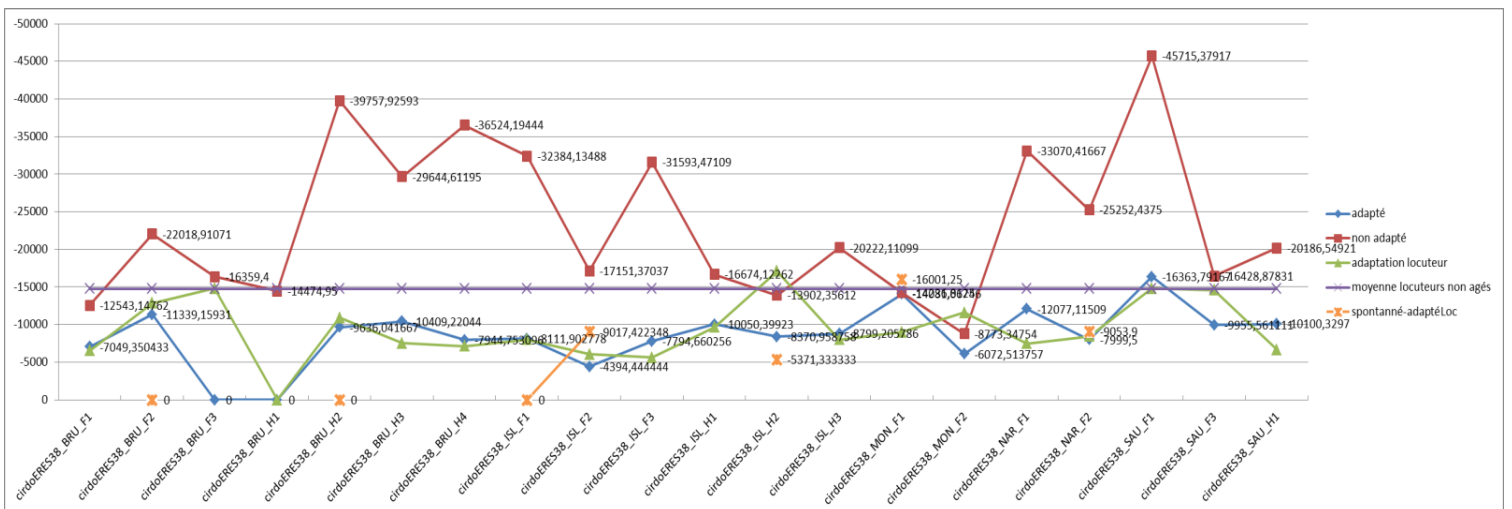
on



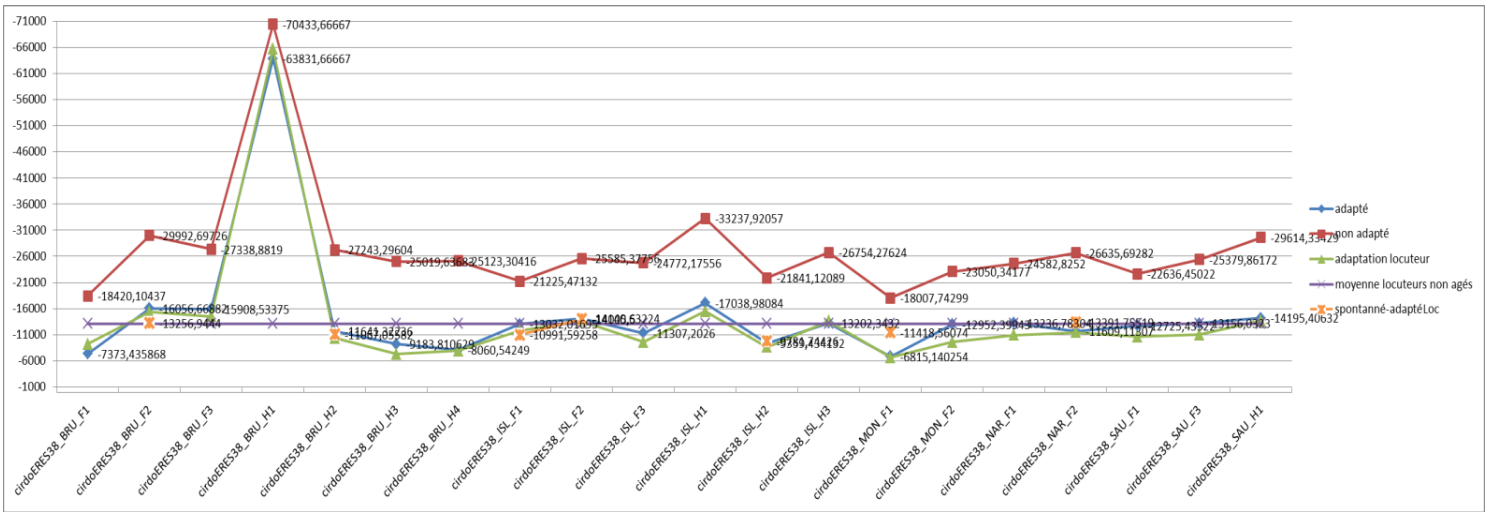
oo



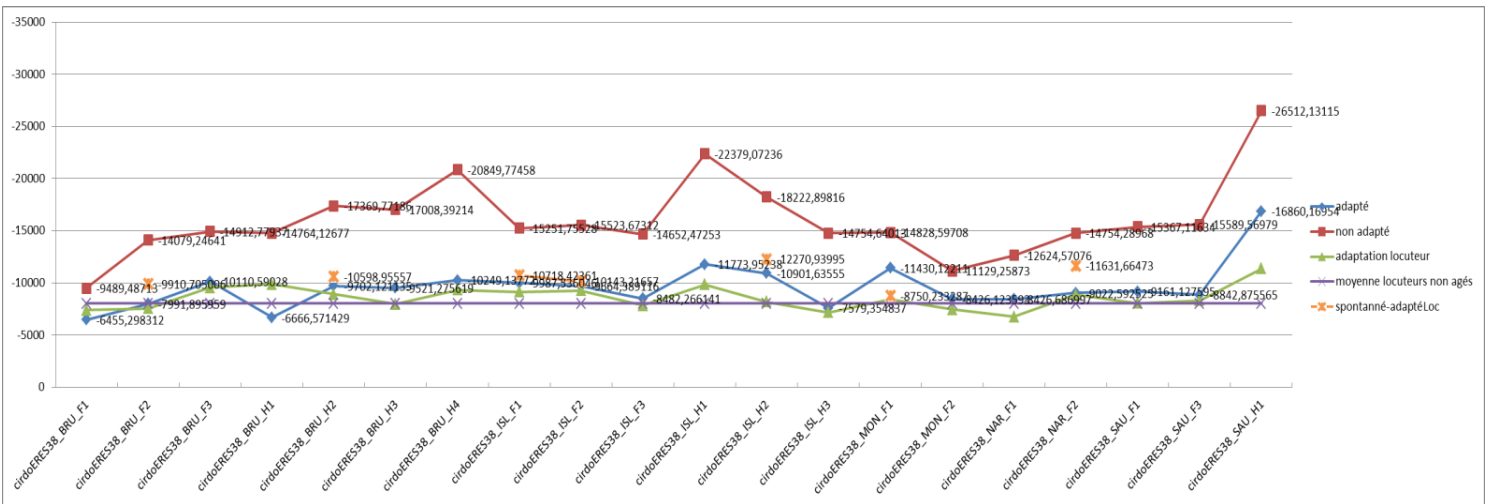
ooo



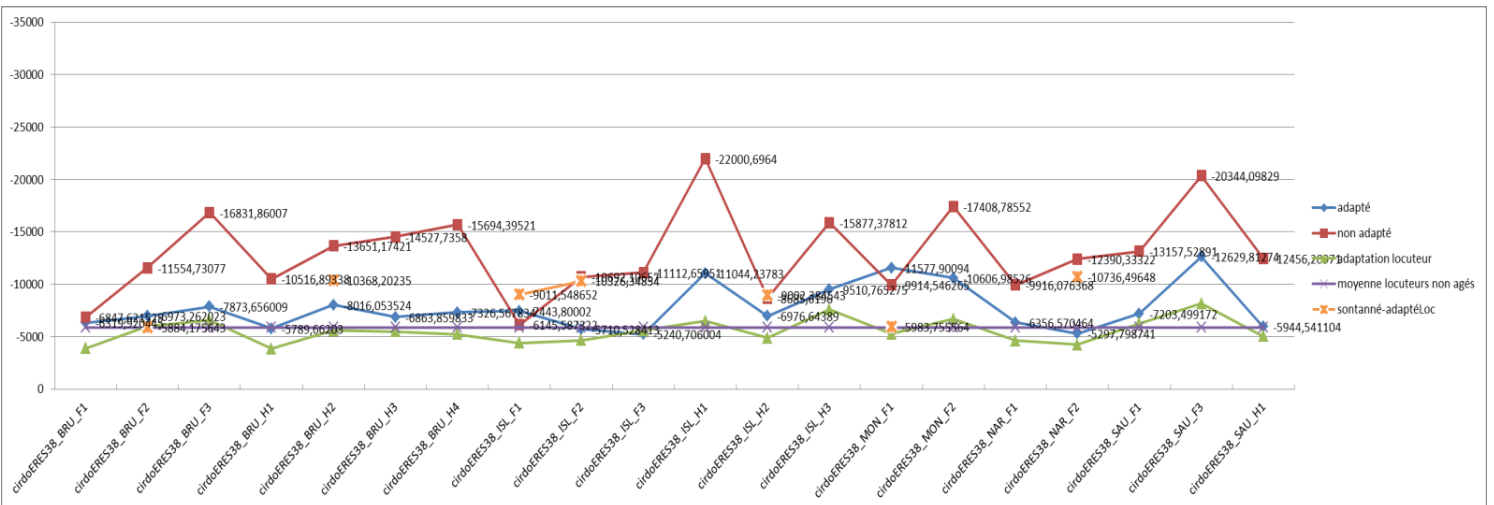
p



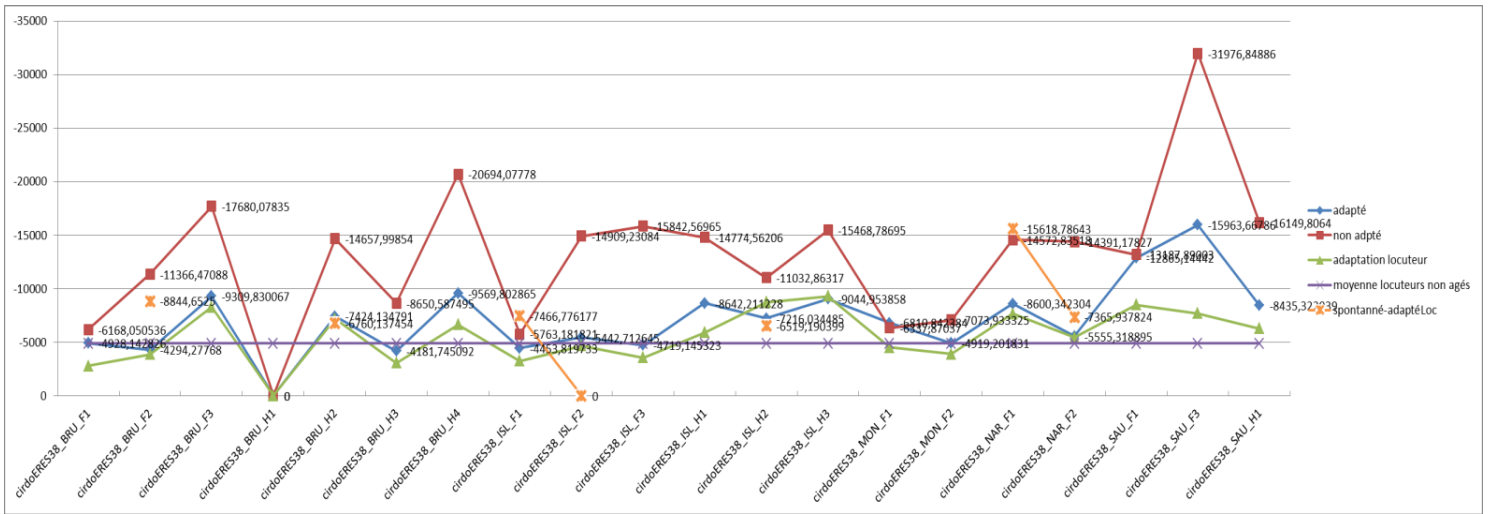
RR



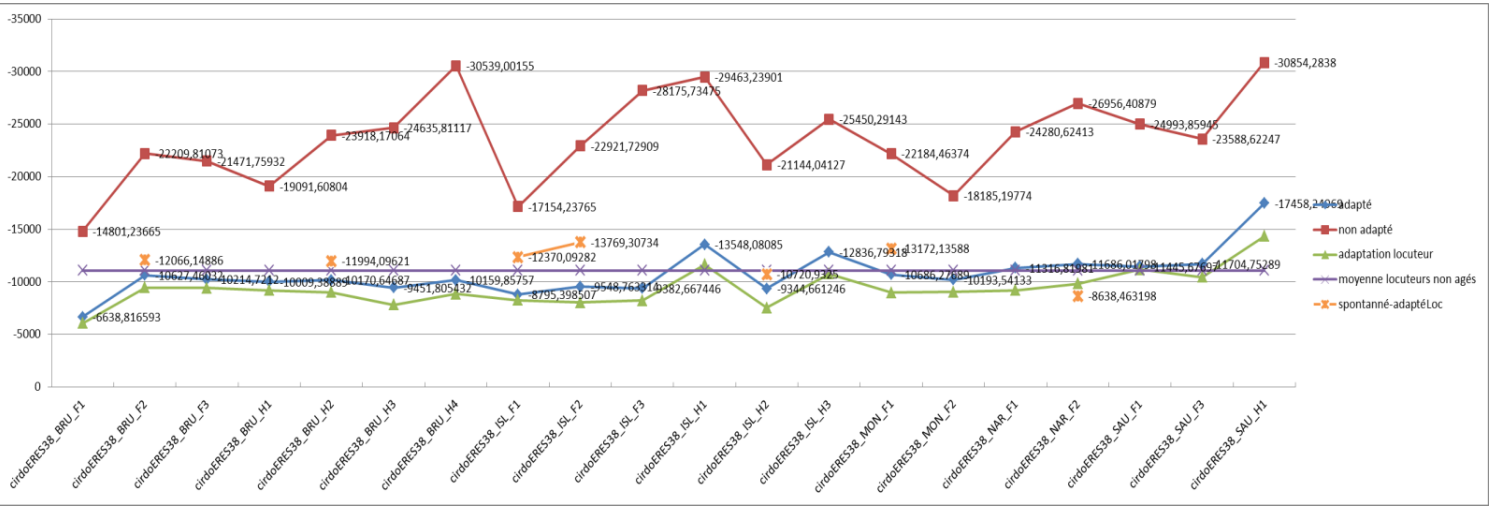
S



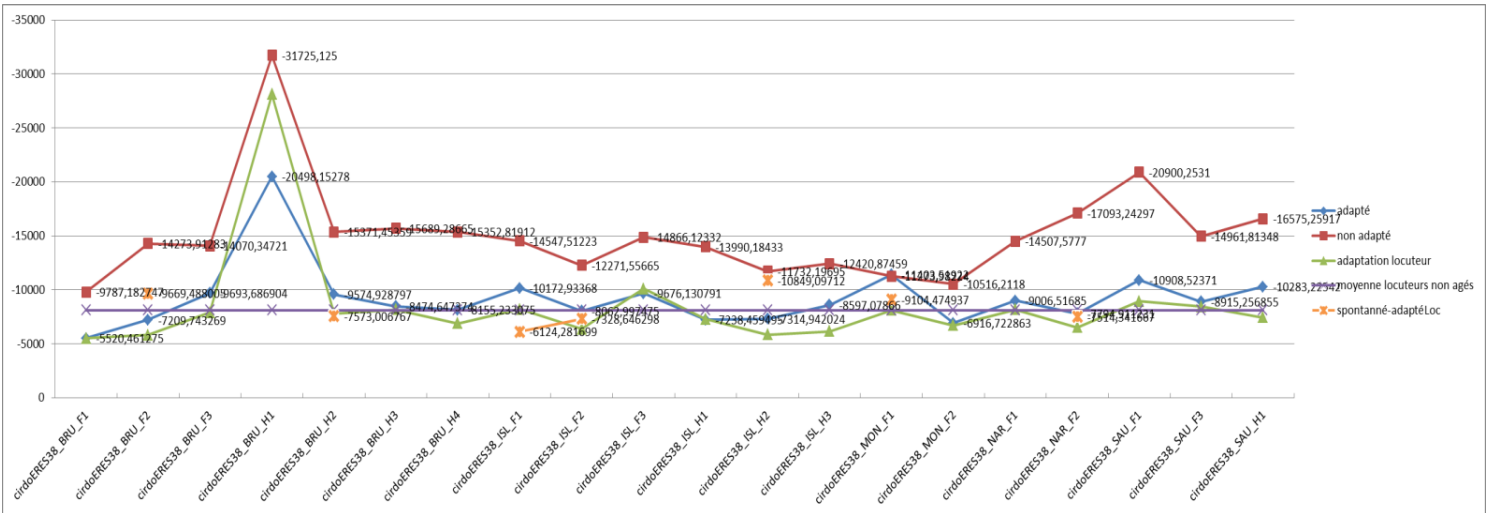
SS



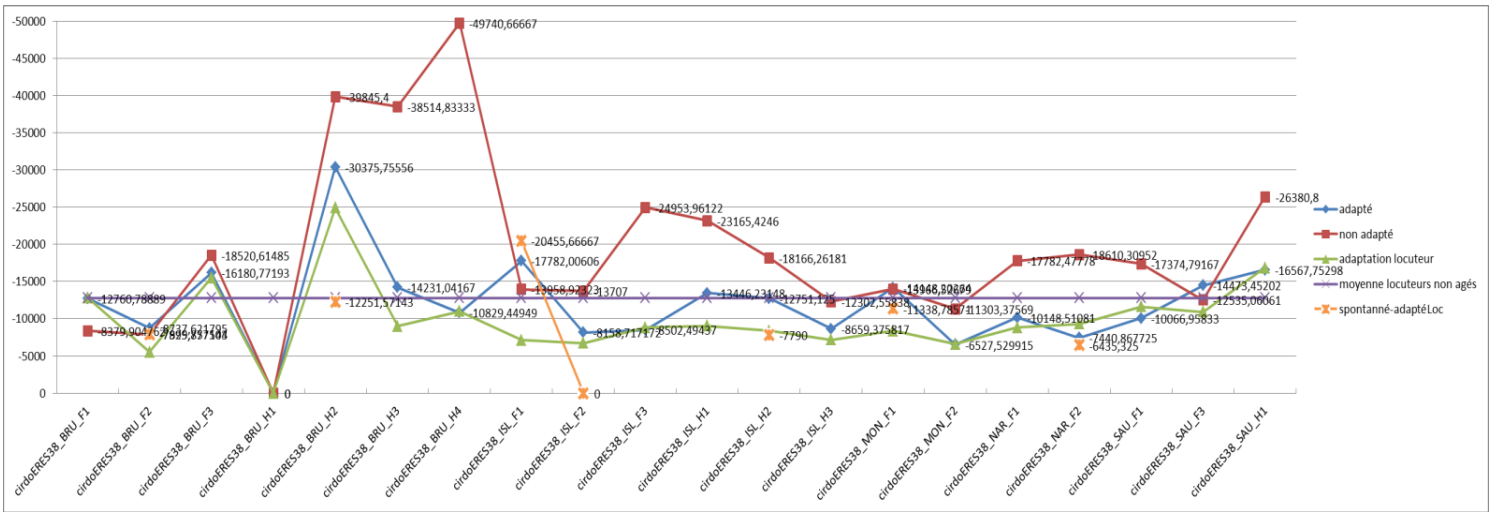
t



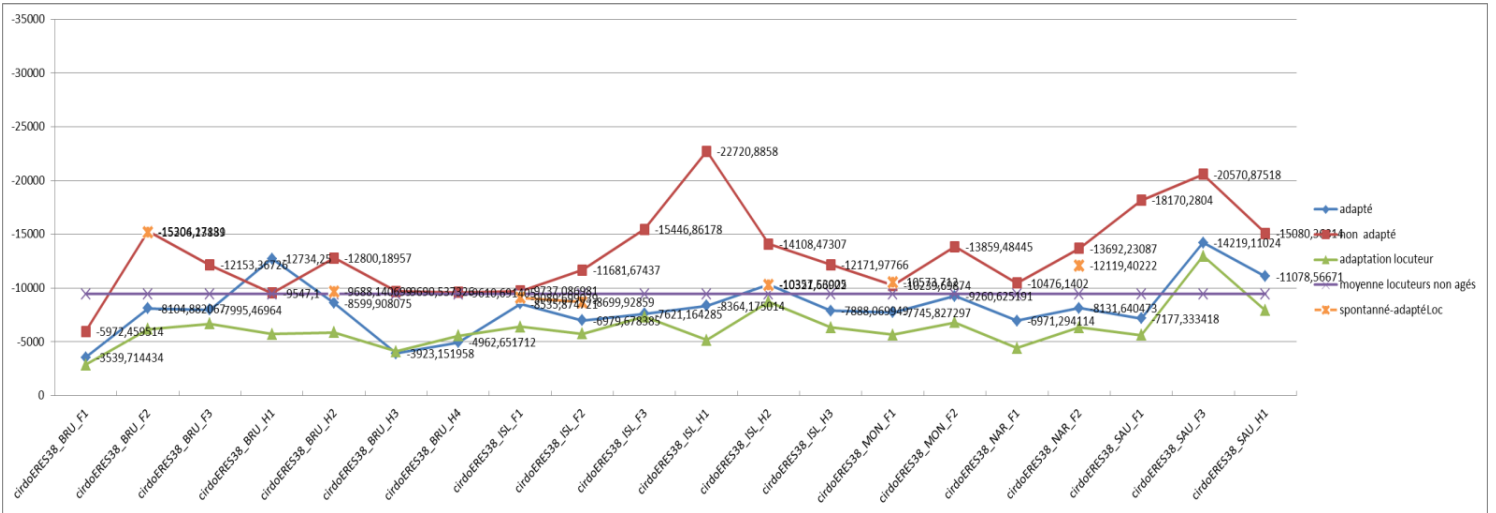
u



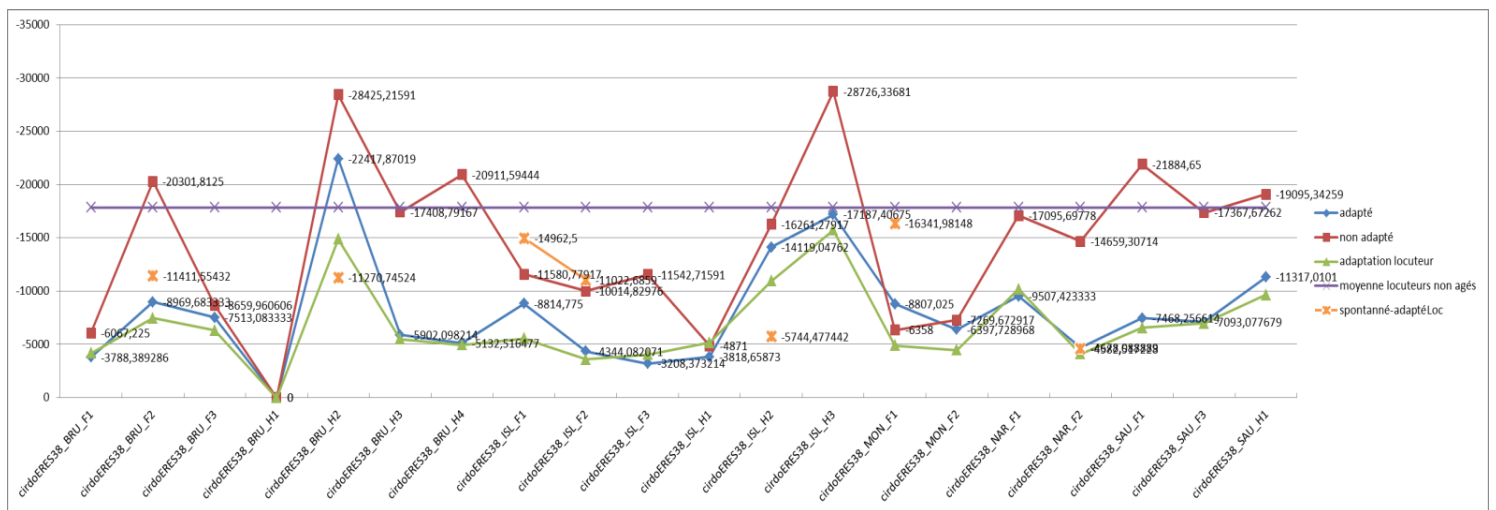
un



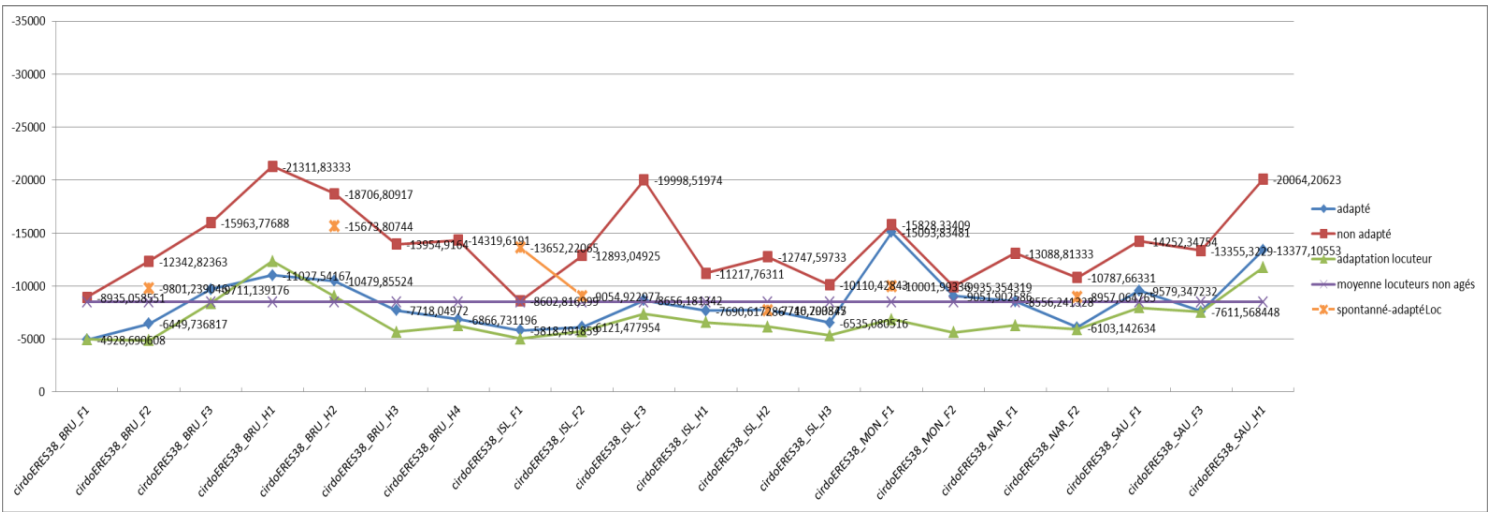
v



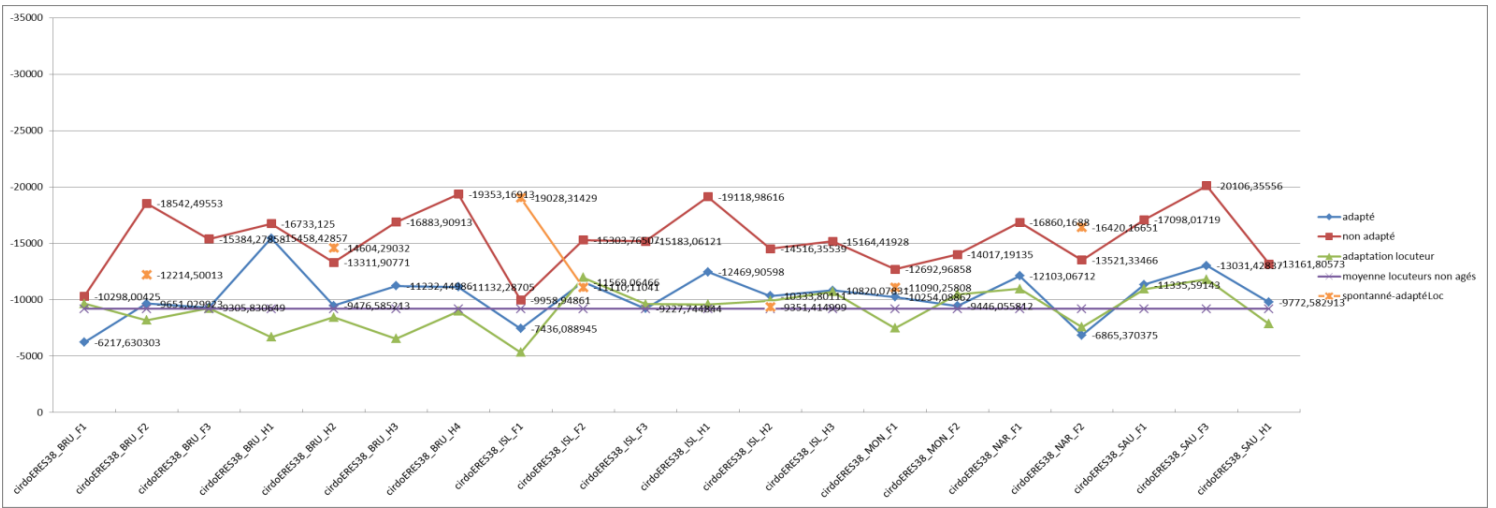
w



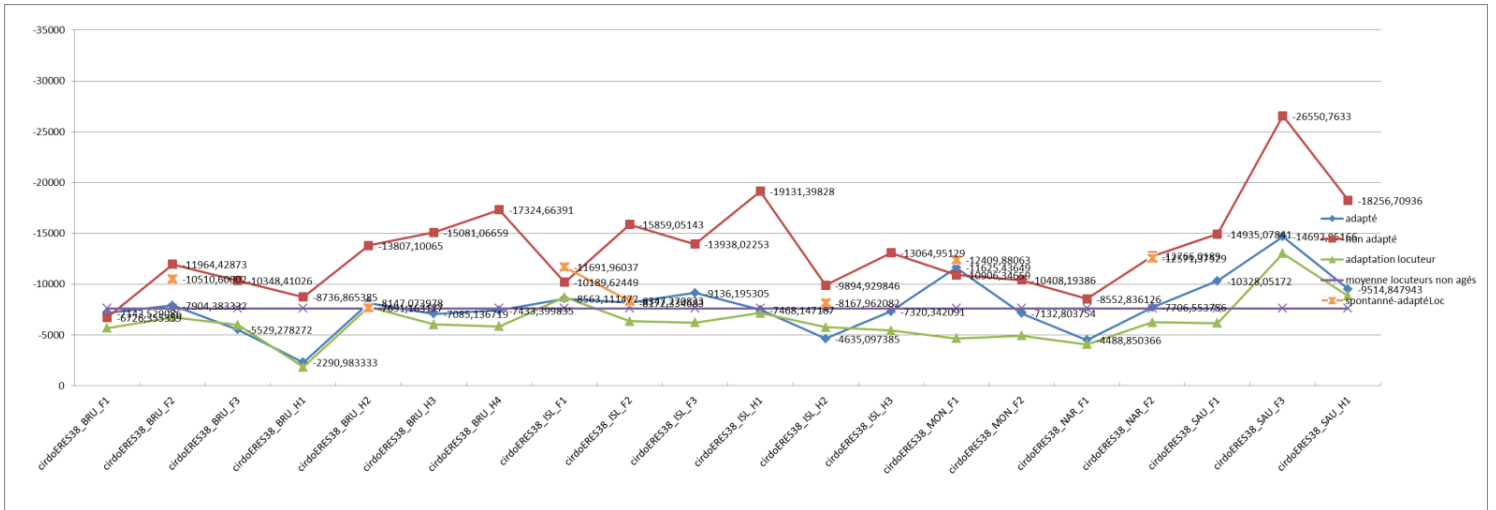
Y



Z

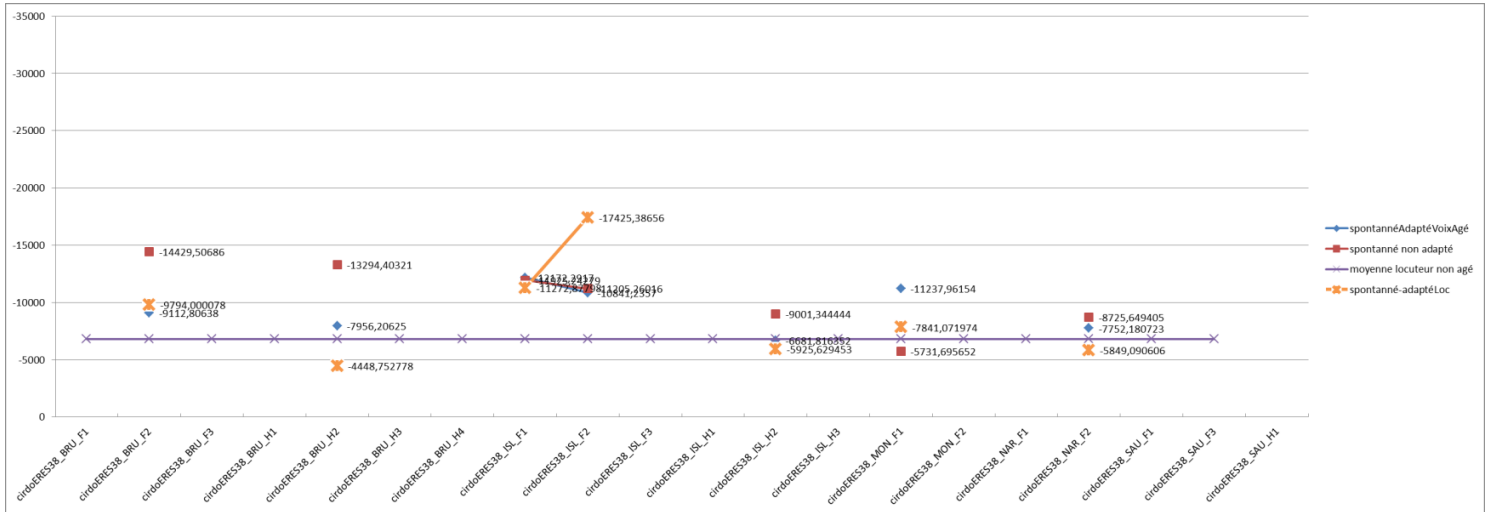


ZZ

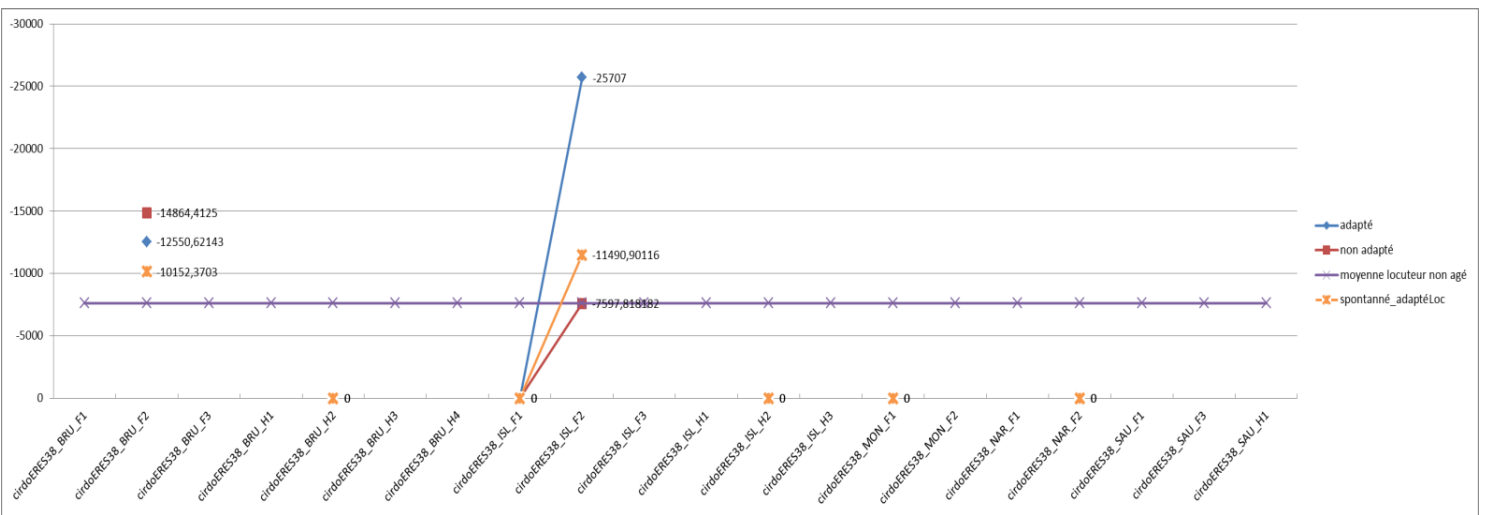


**8.8 Diagrammes représentant pour chaque phonème
et par locuteur la courbe des scores d'alignements
en fonction des modèles acoustiques
Cas de la parole spontanée**

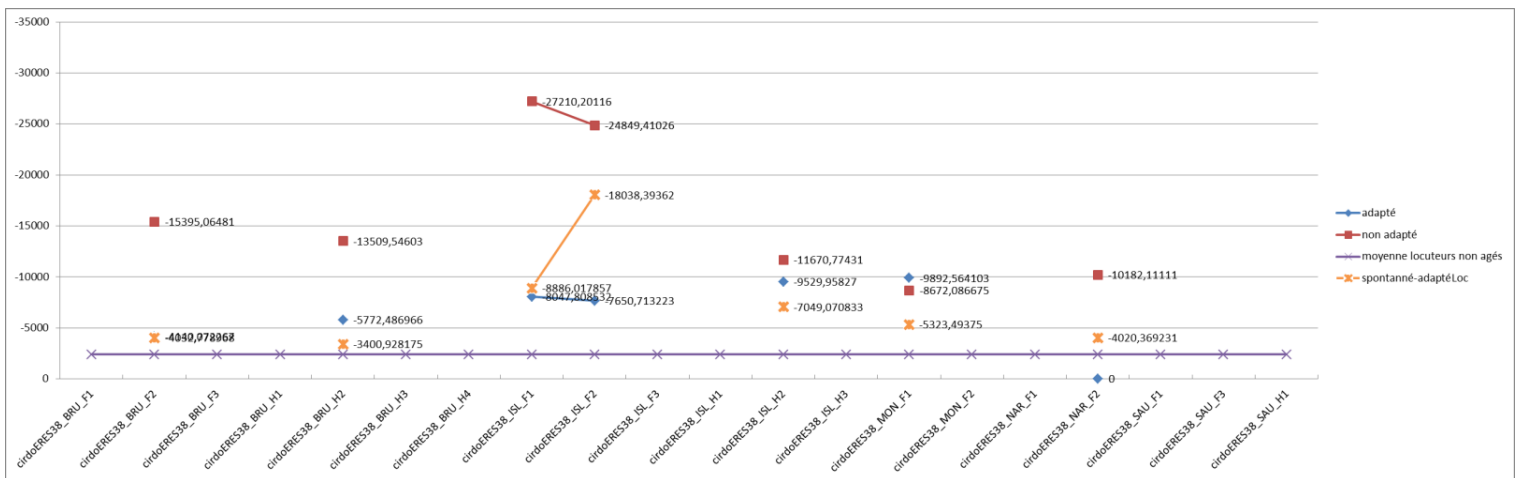
2



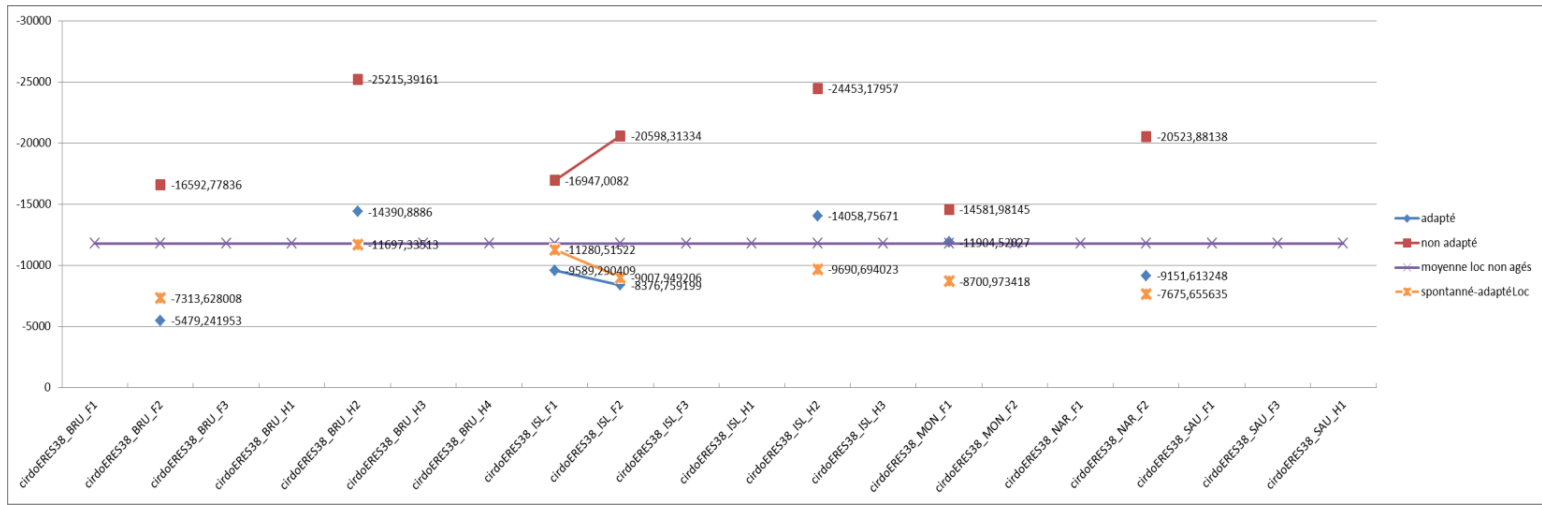
9



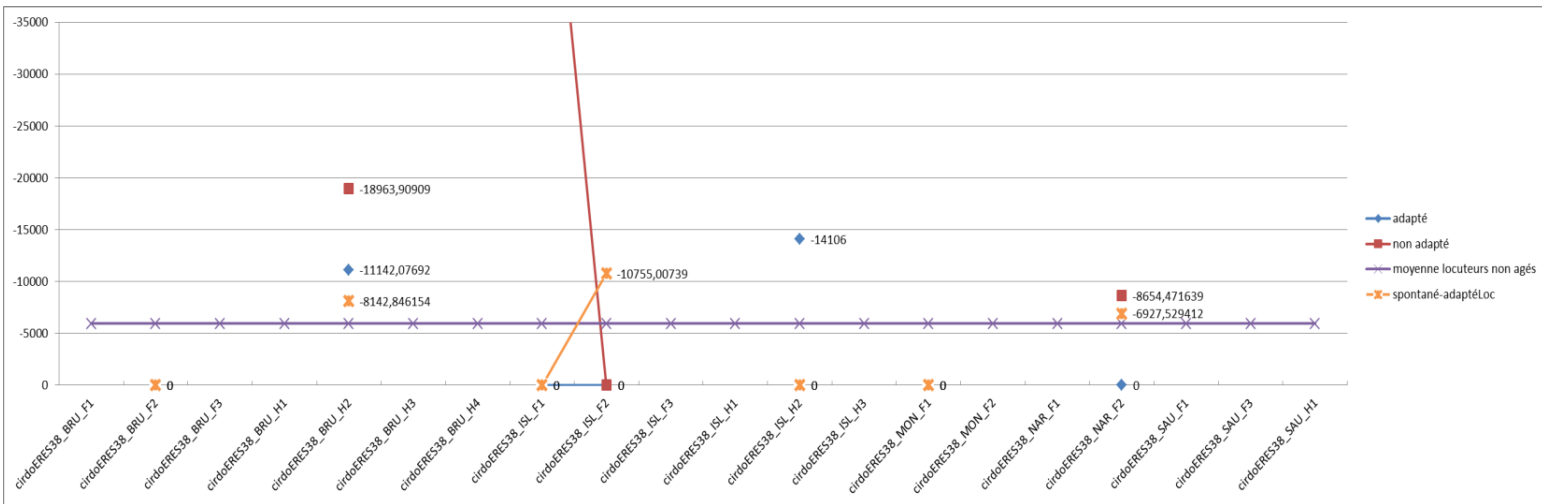
29



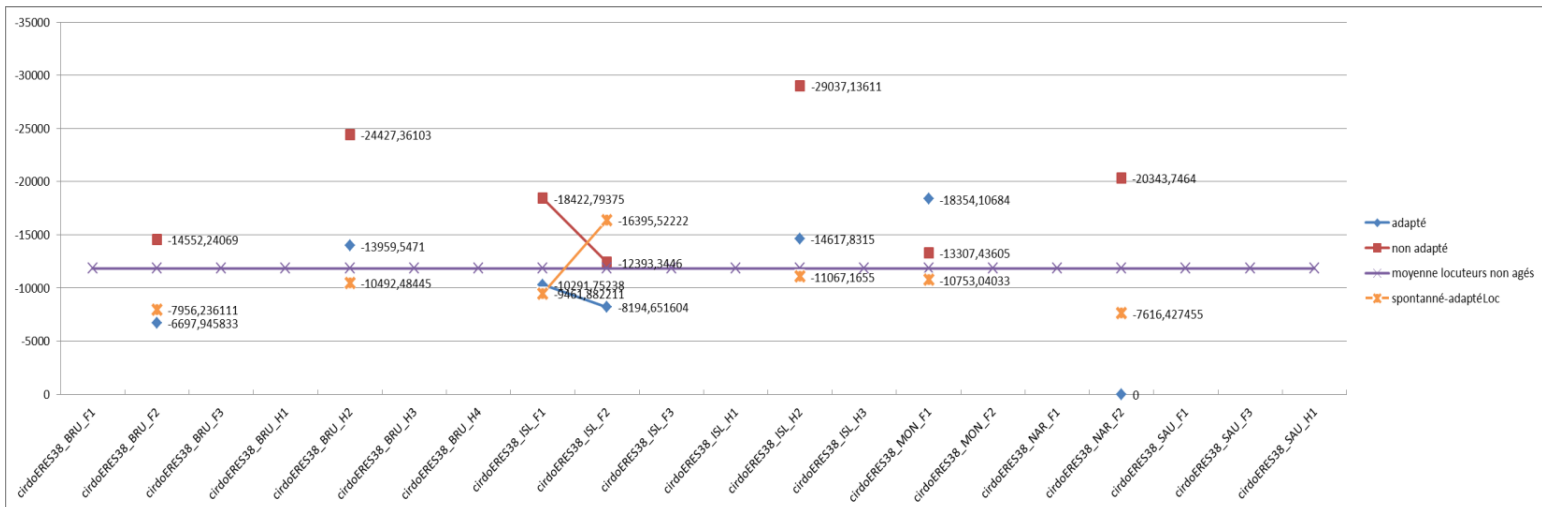
a



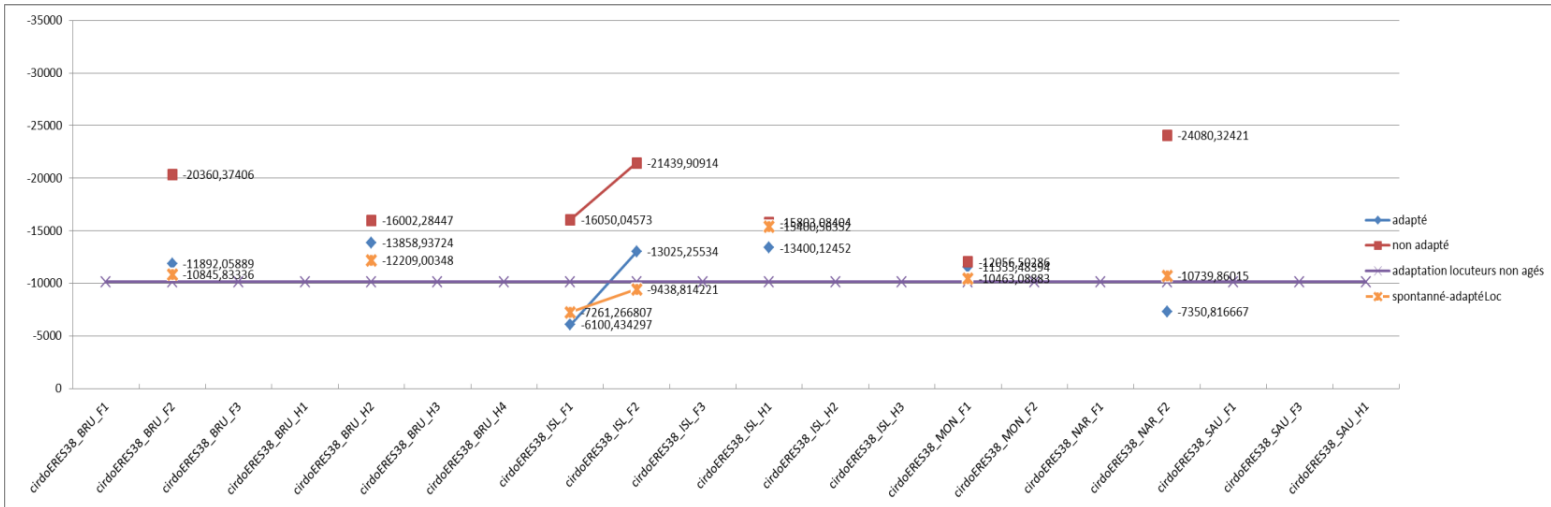
AA



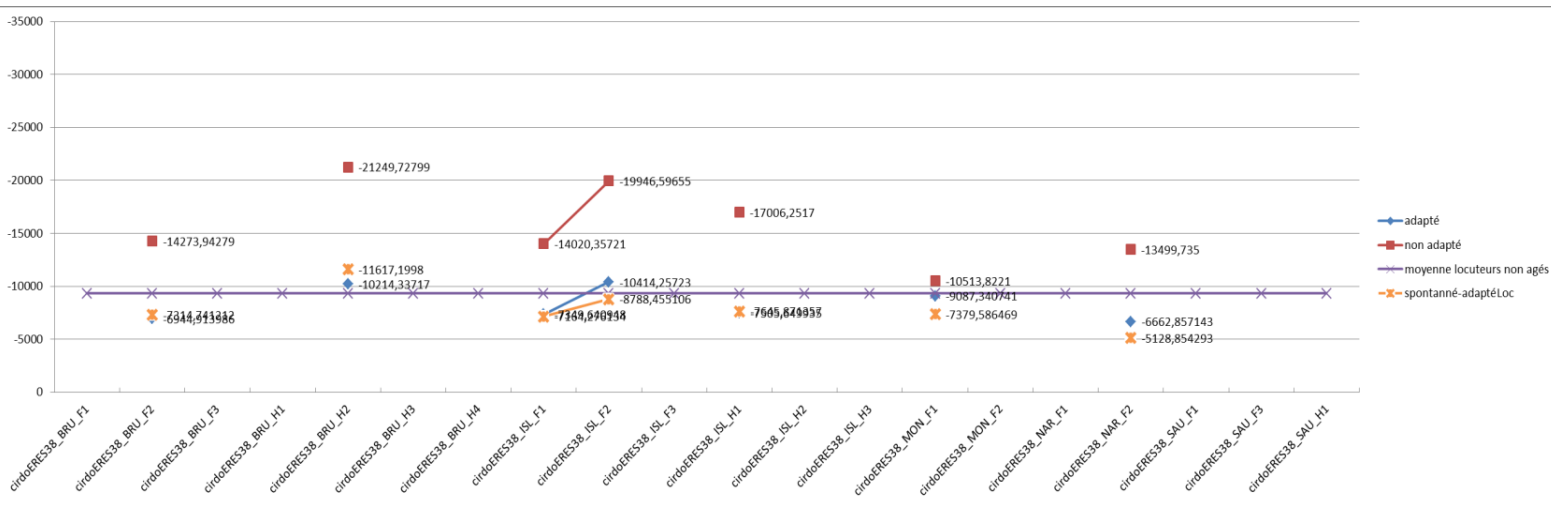
aAA



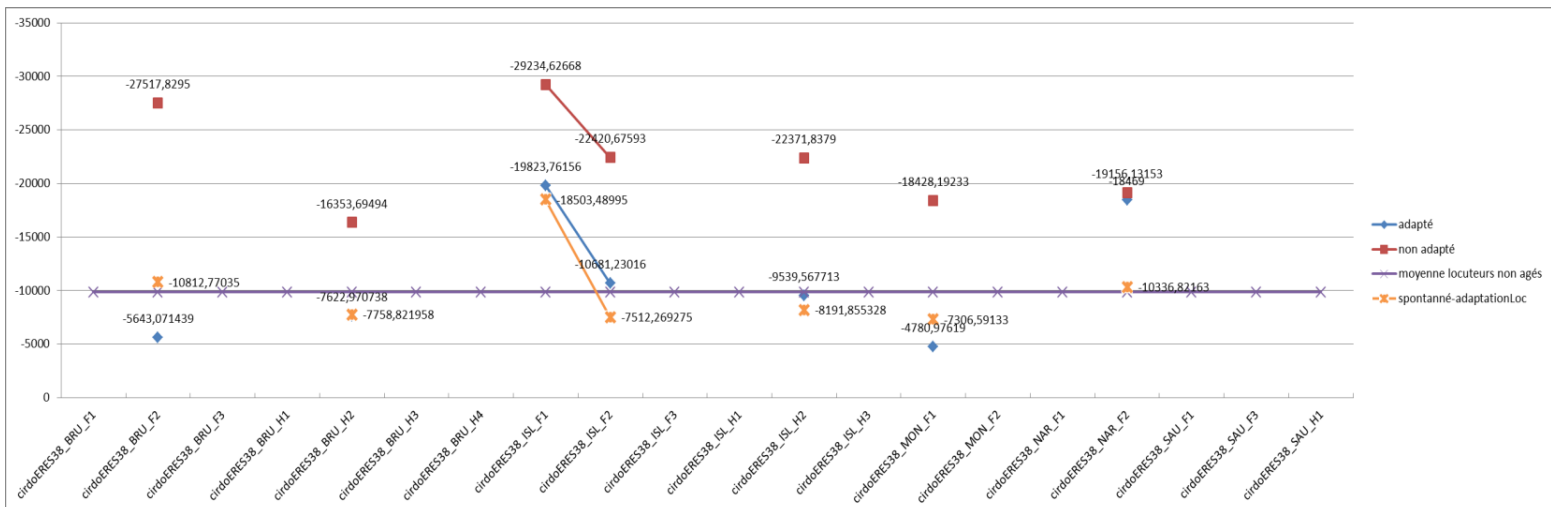
aeA



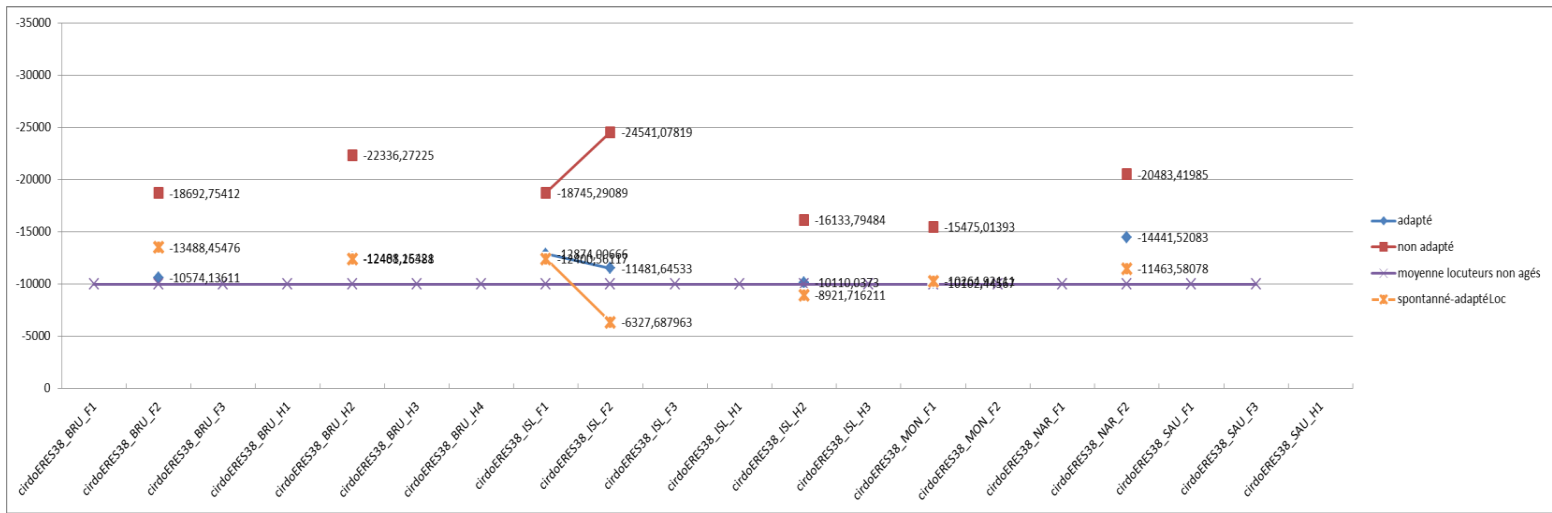
an



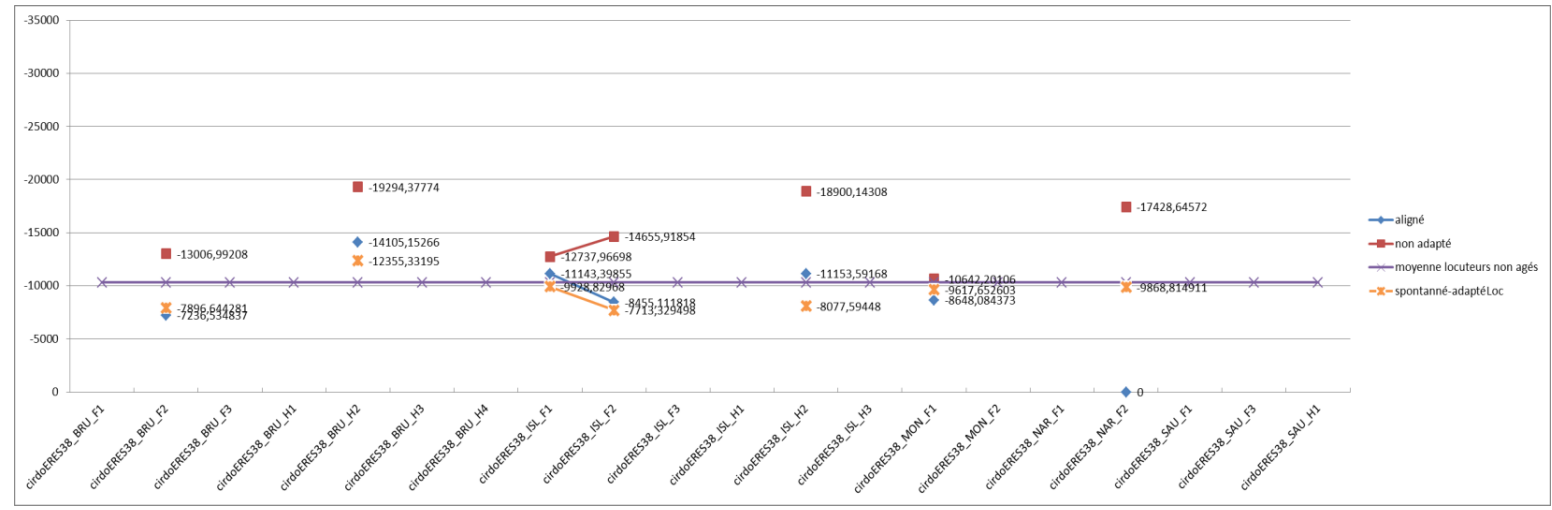
b



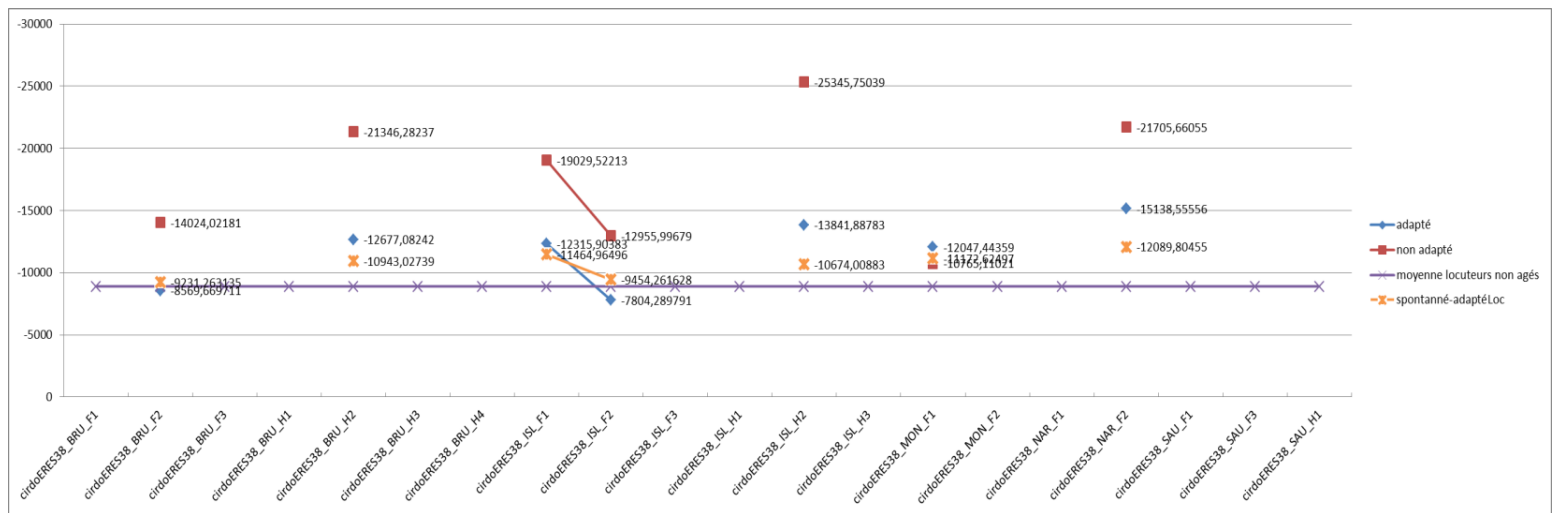
d



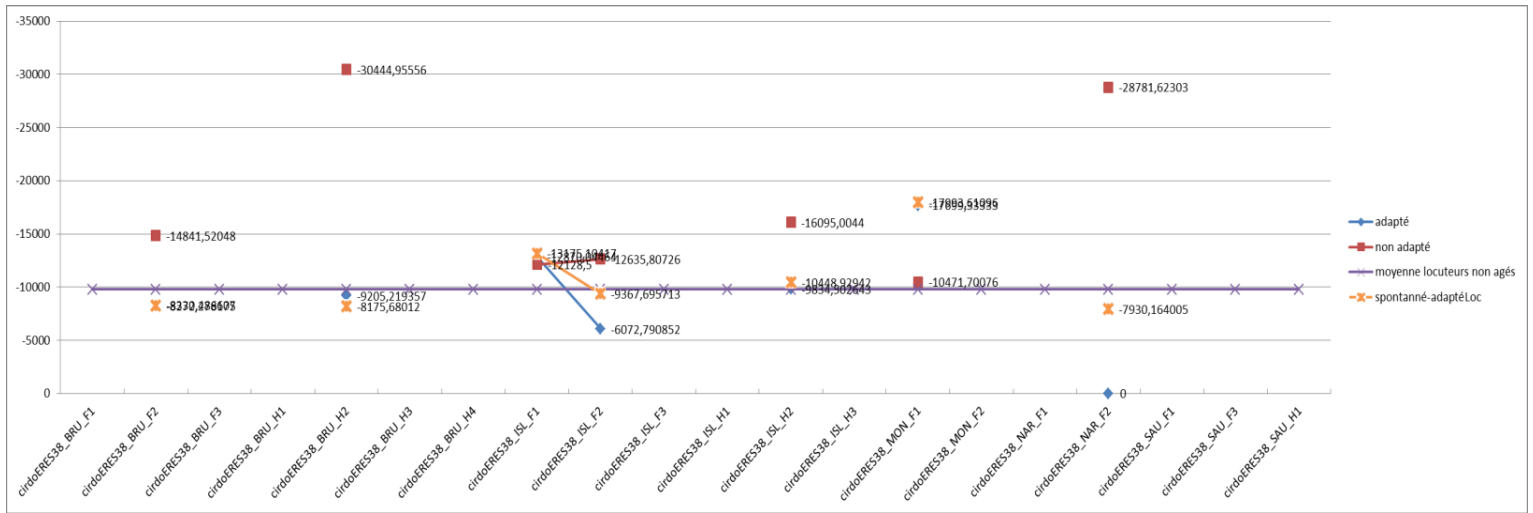
e



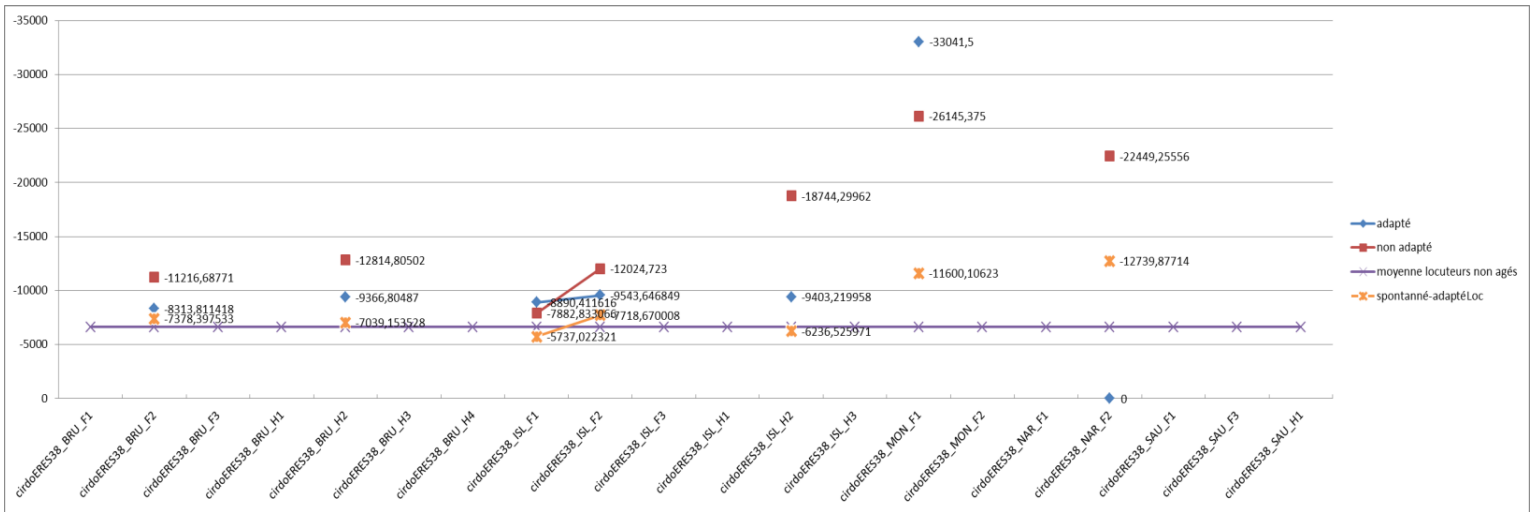
EE



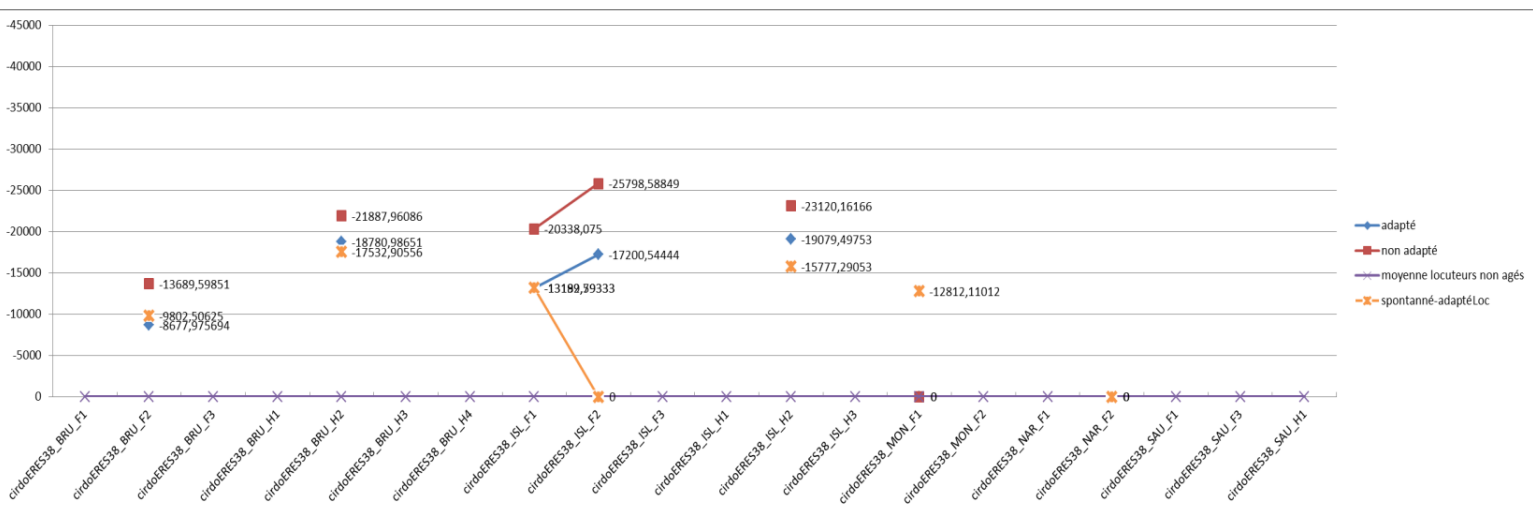
eEE



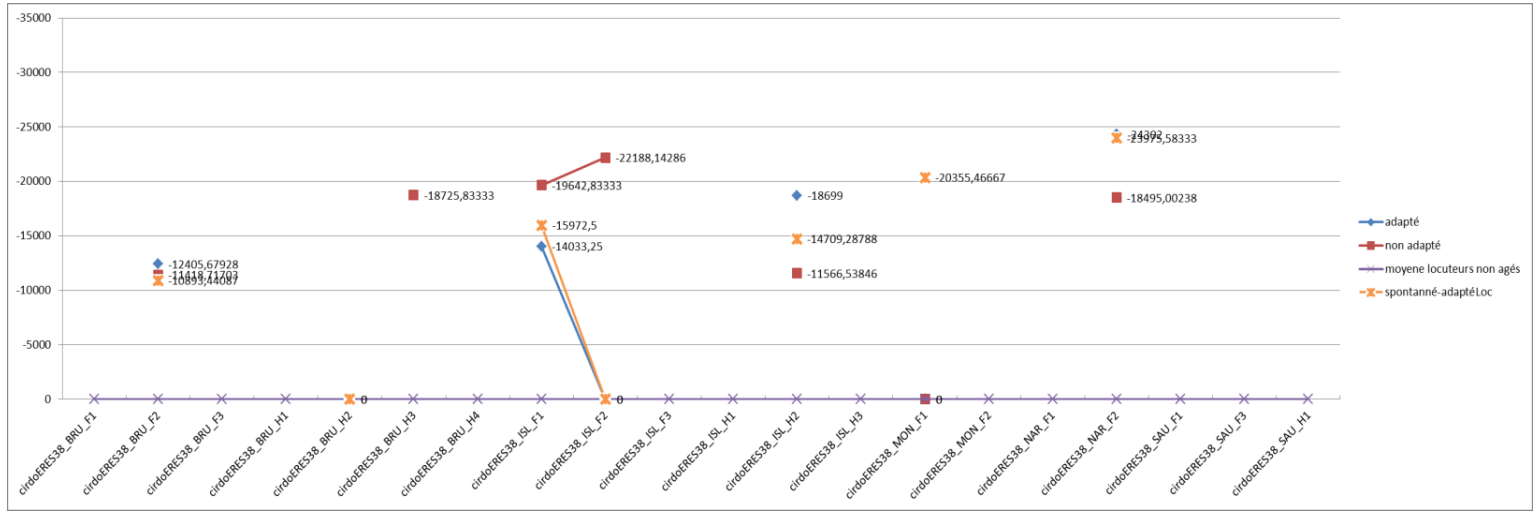
f



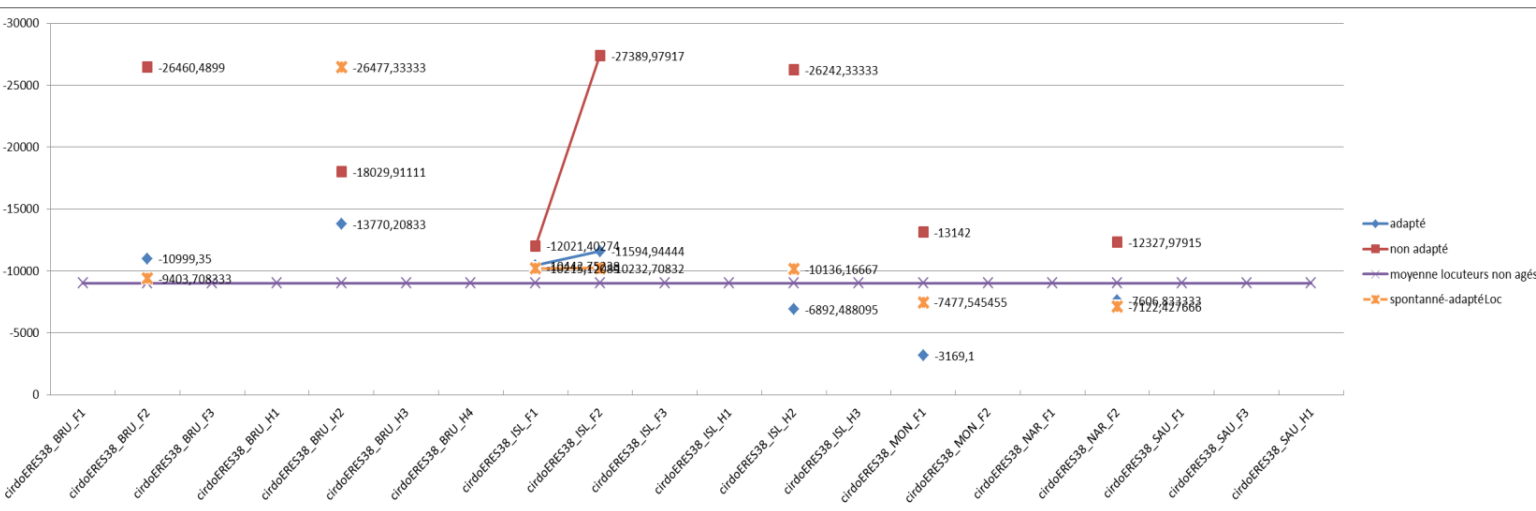
g



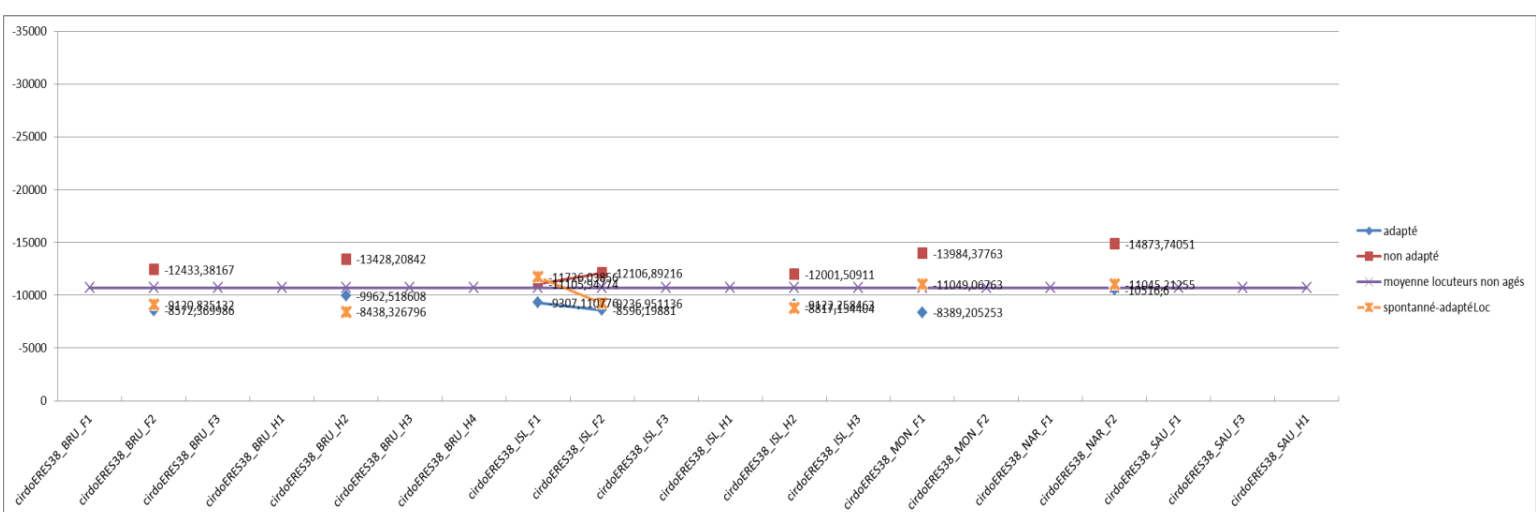
h



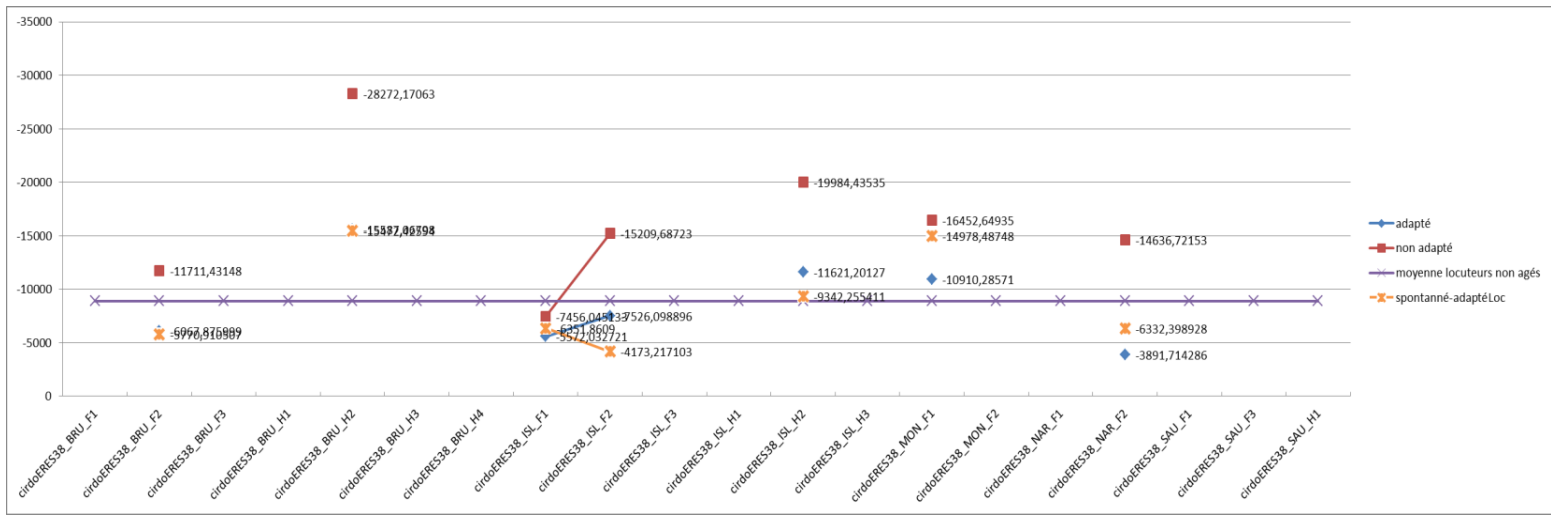
HH



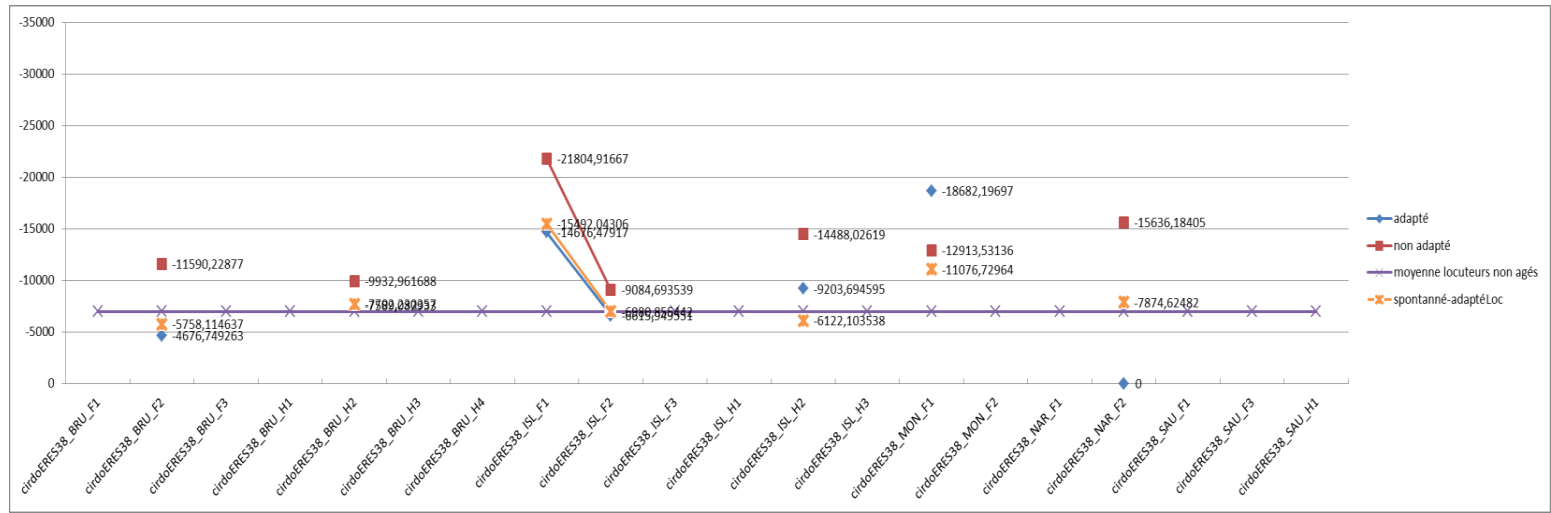
i



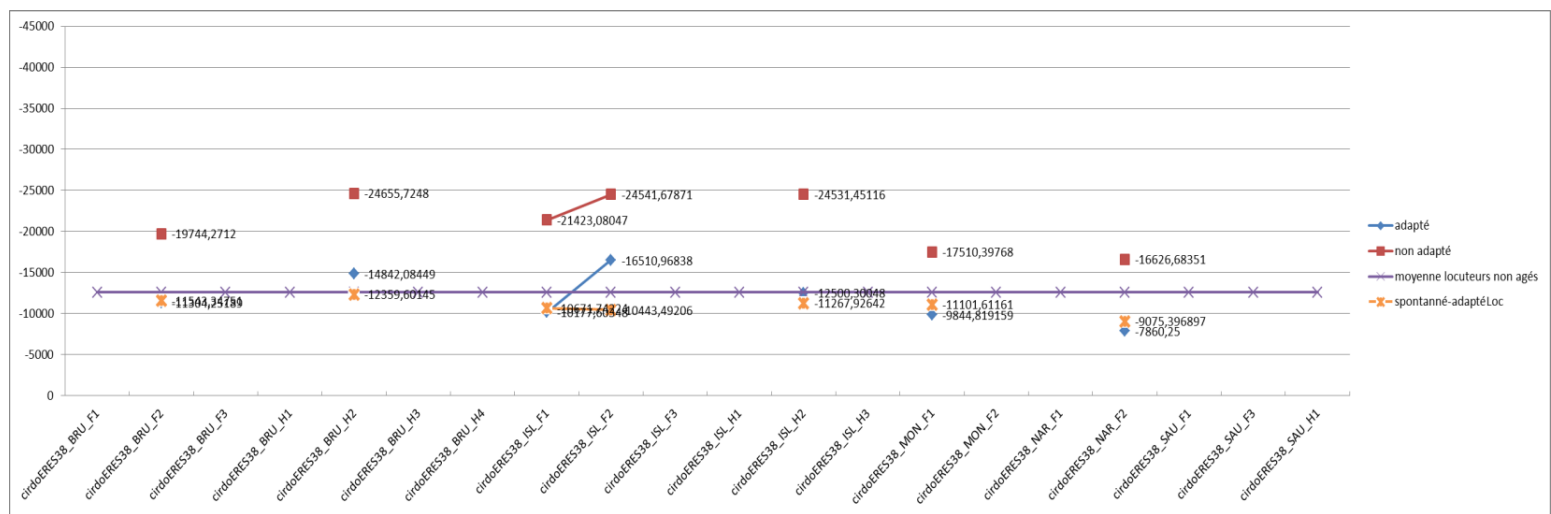
in



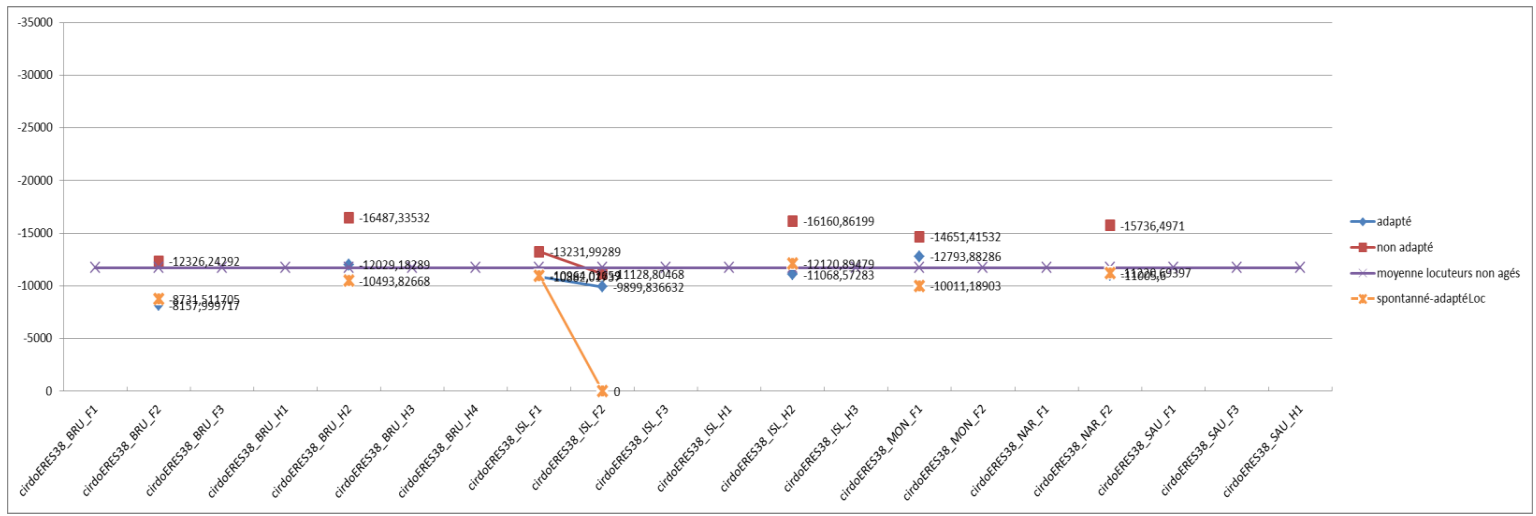
j



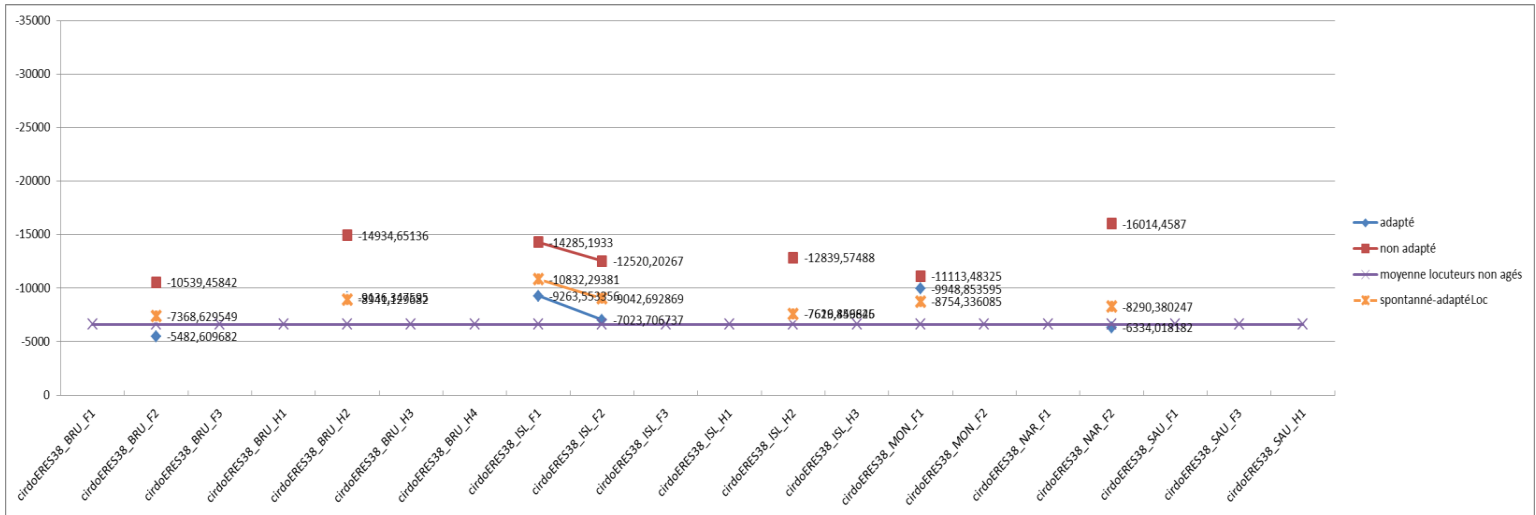
k



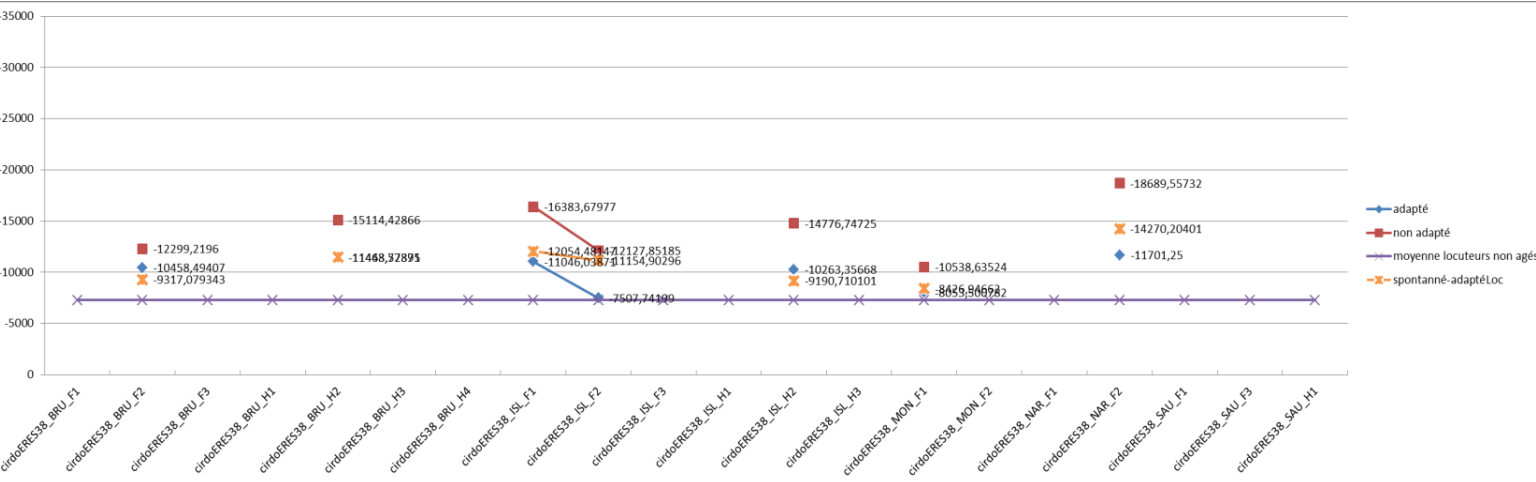
l



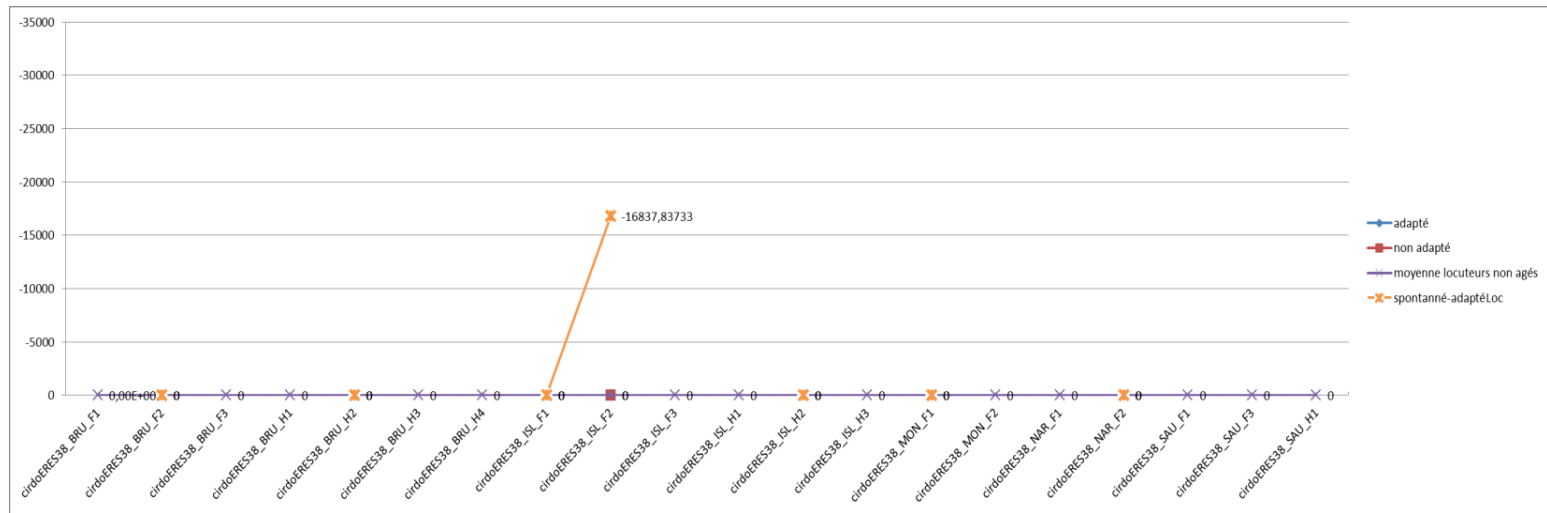
m



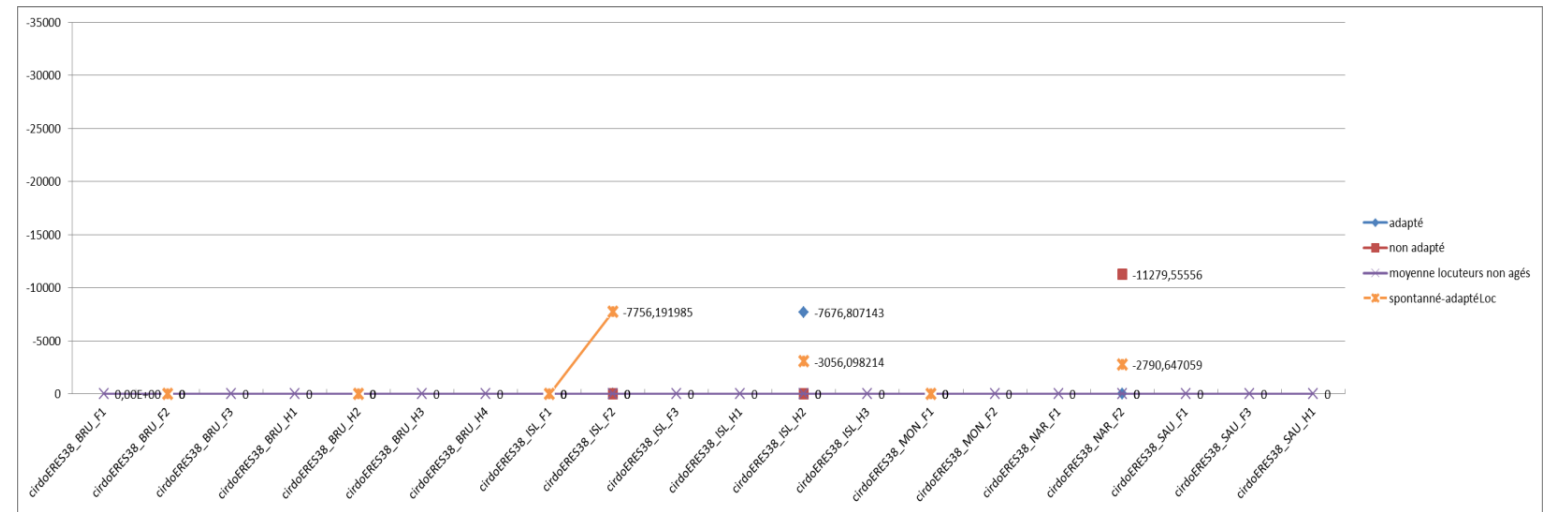
n



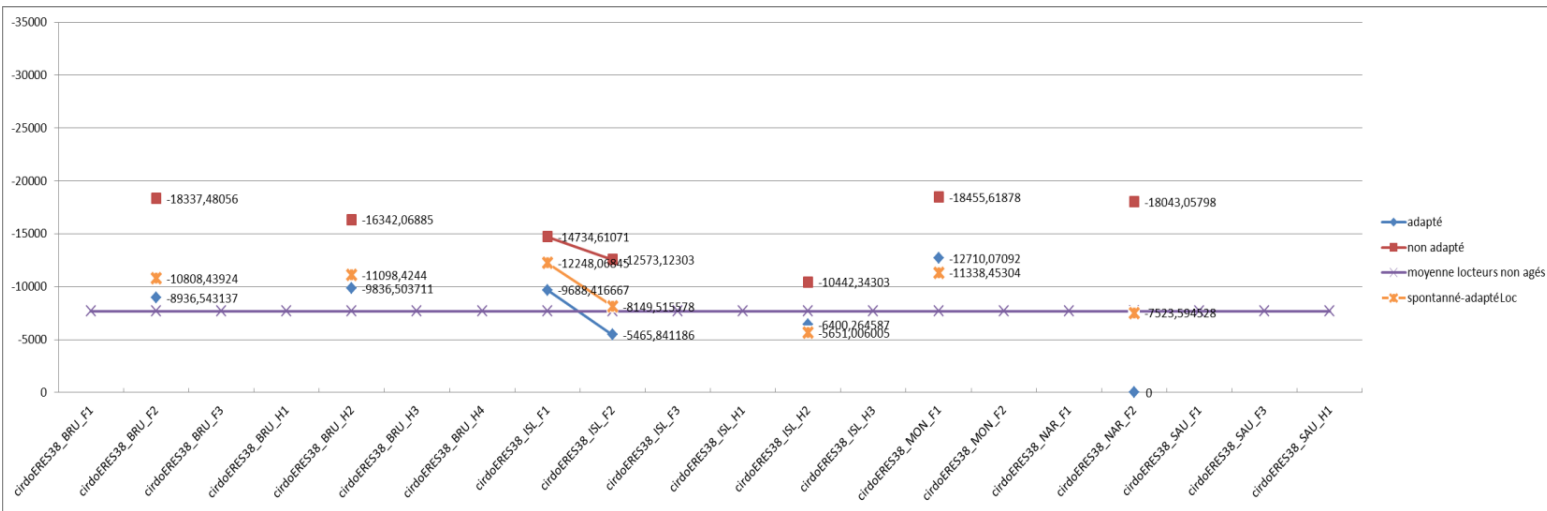
NG



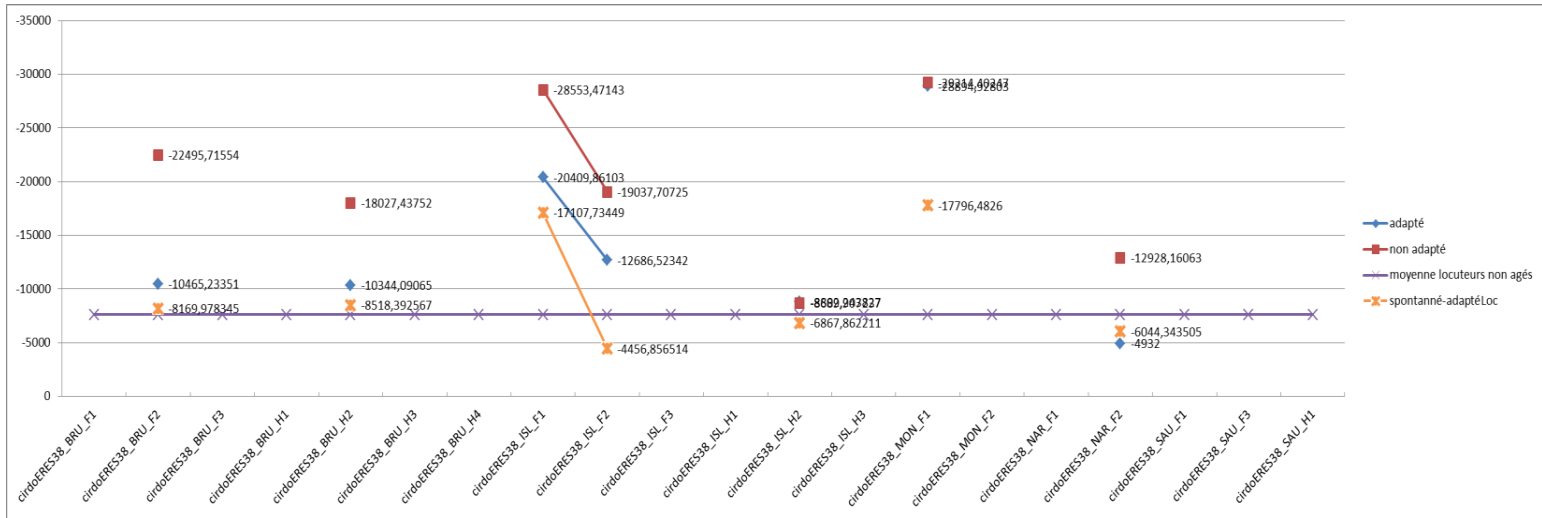
NJ



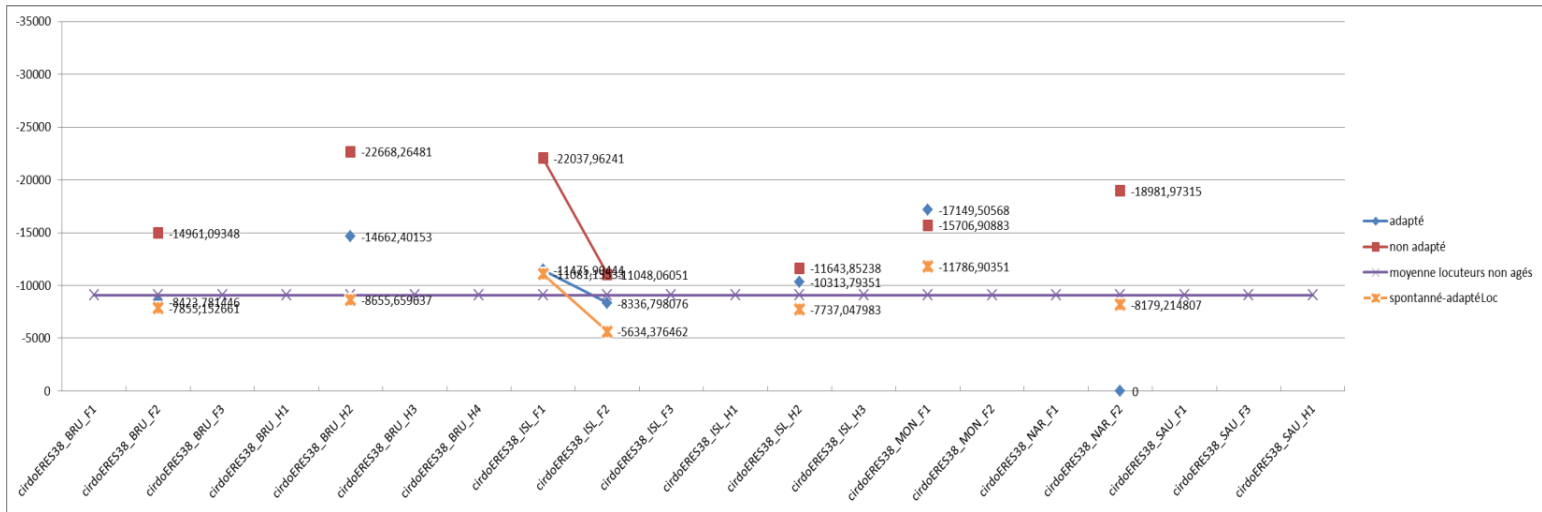
O



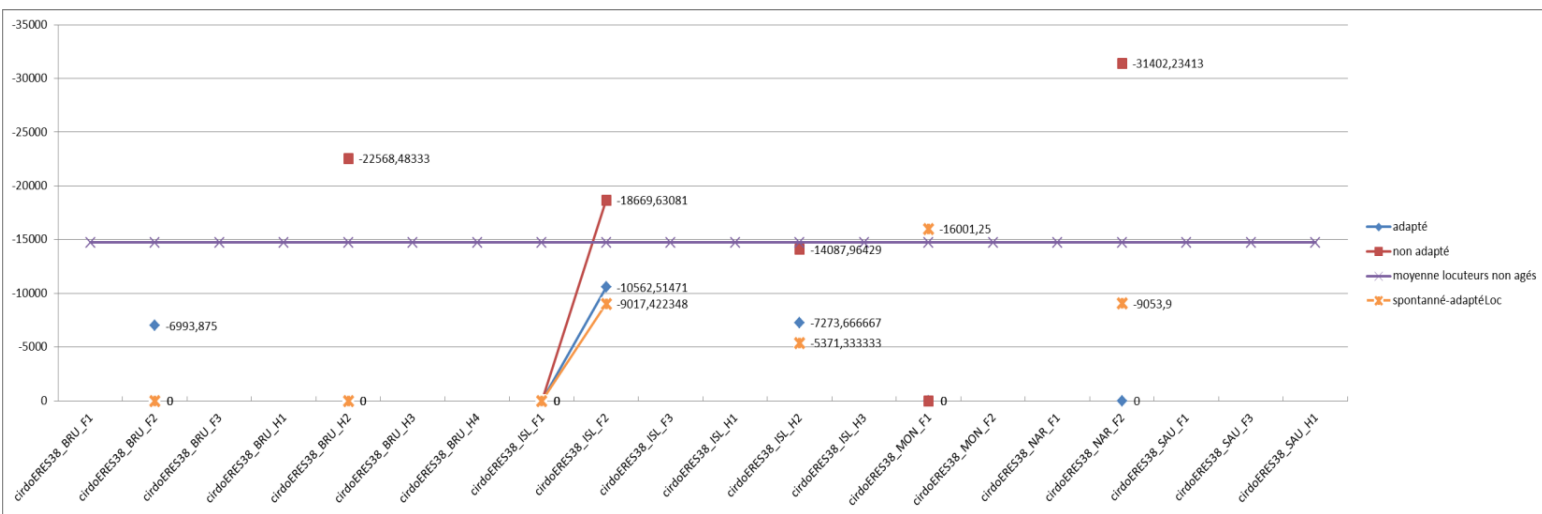
on



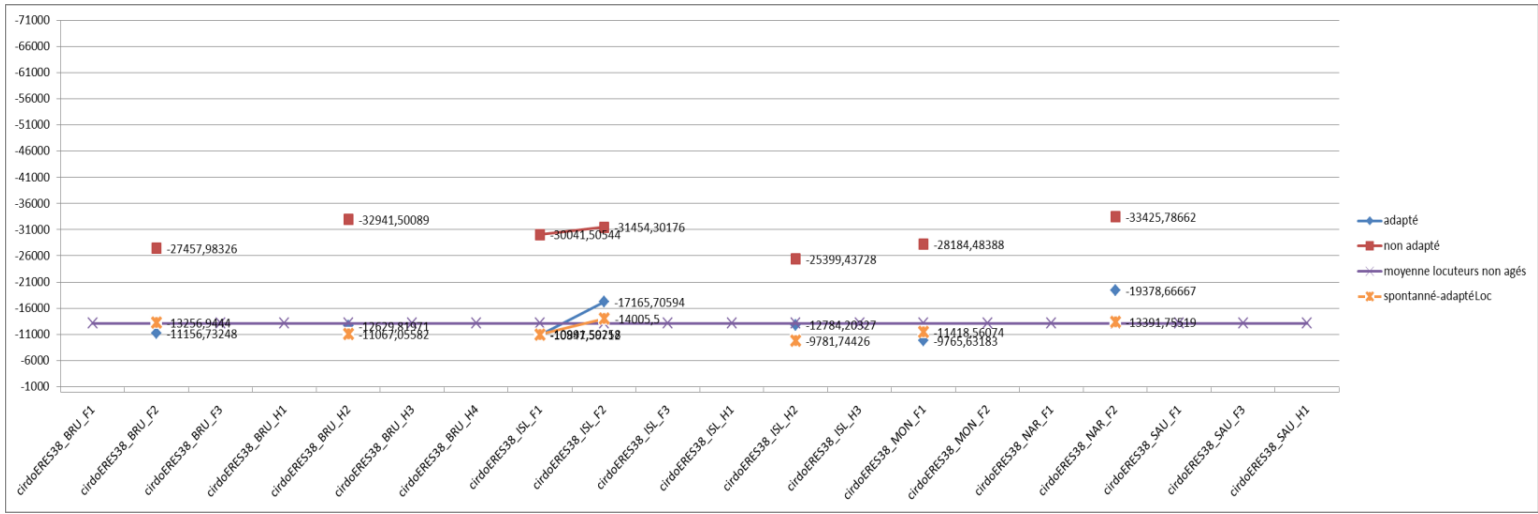
oo



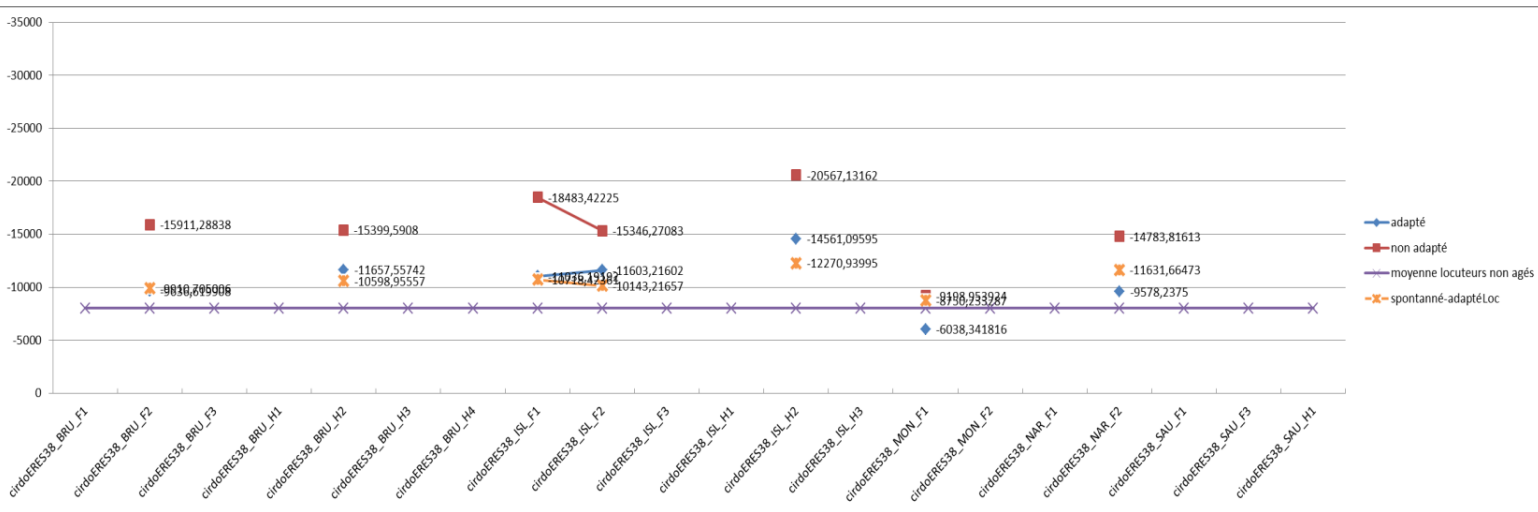
ooo



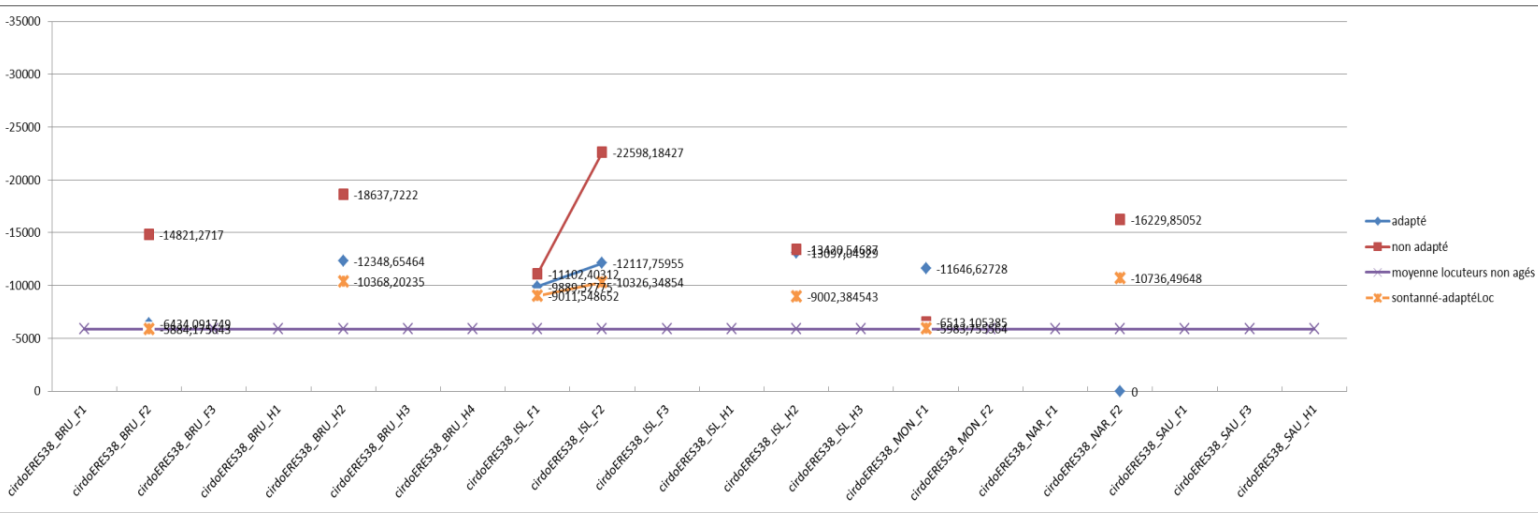
p



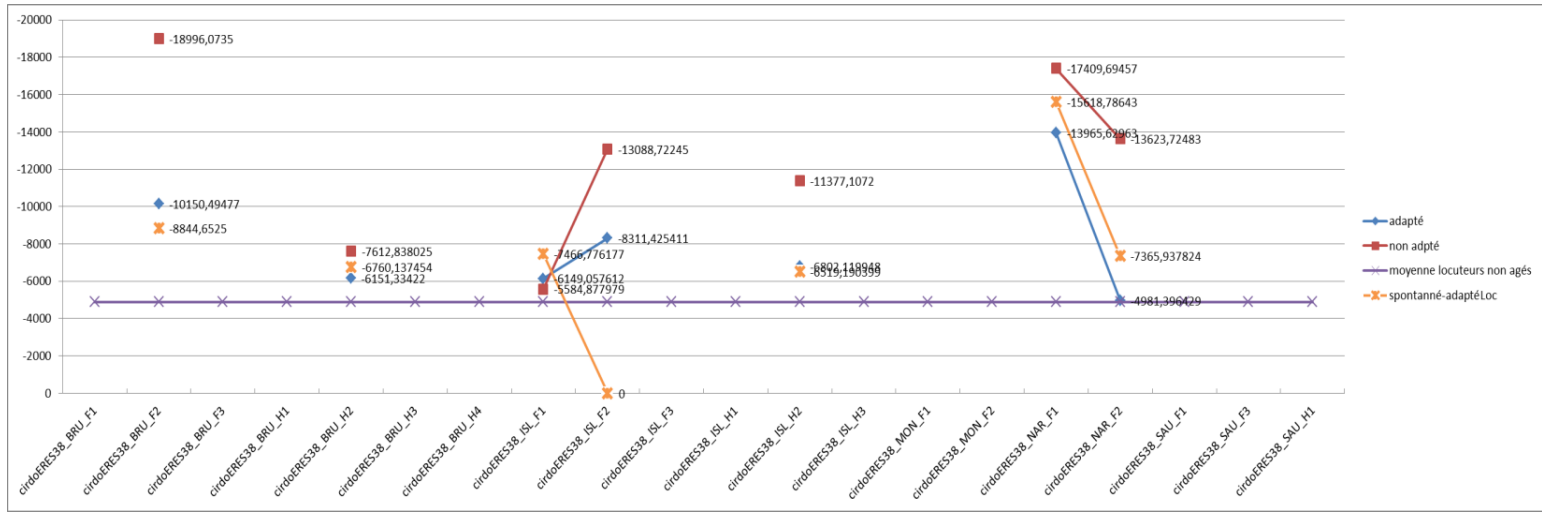
RR



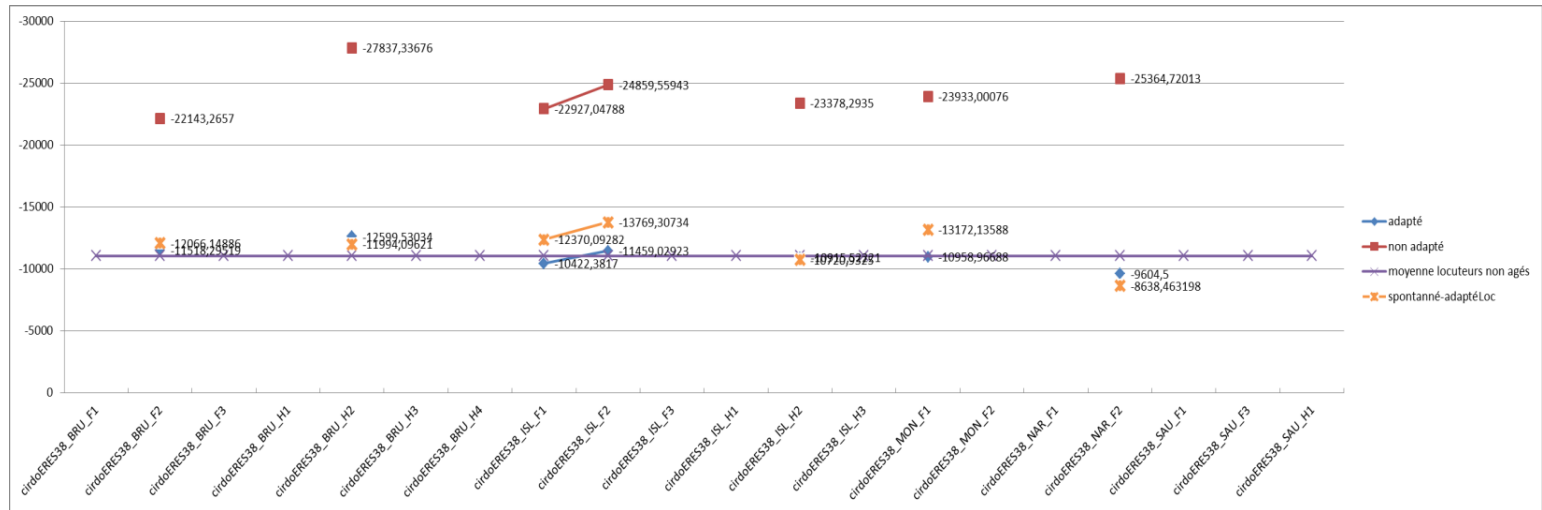
S



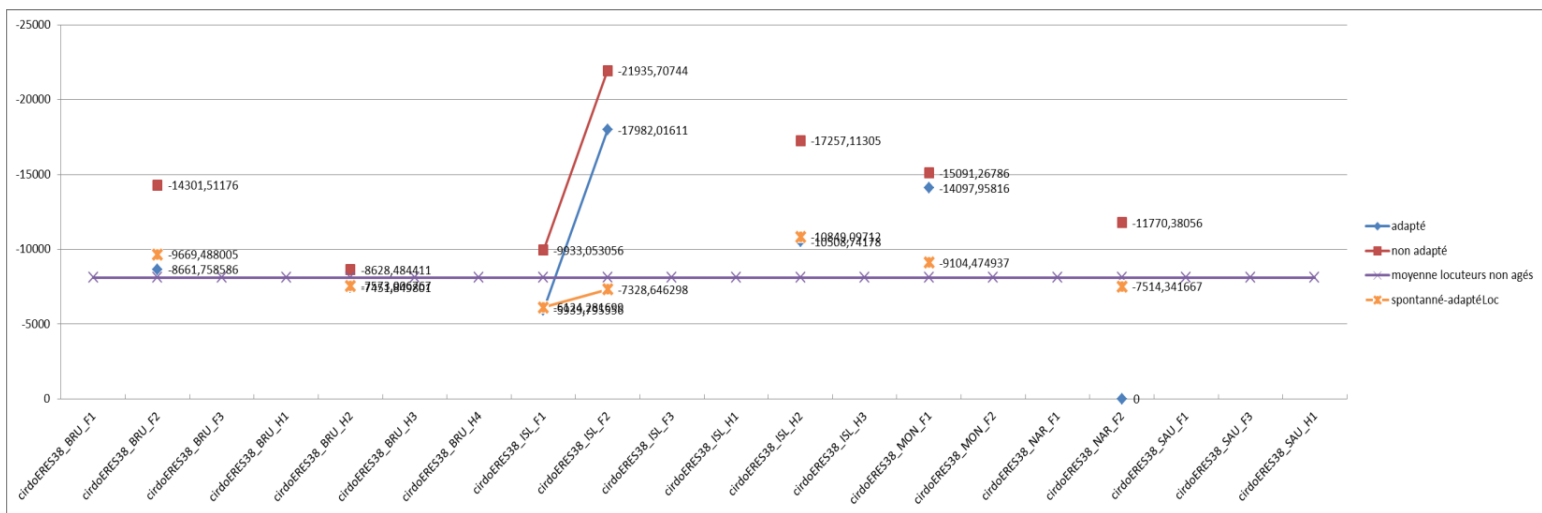
SS



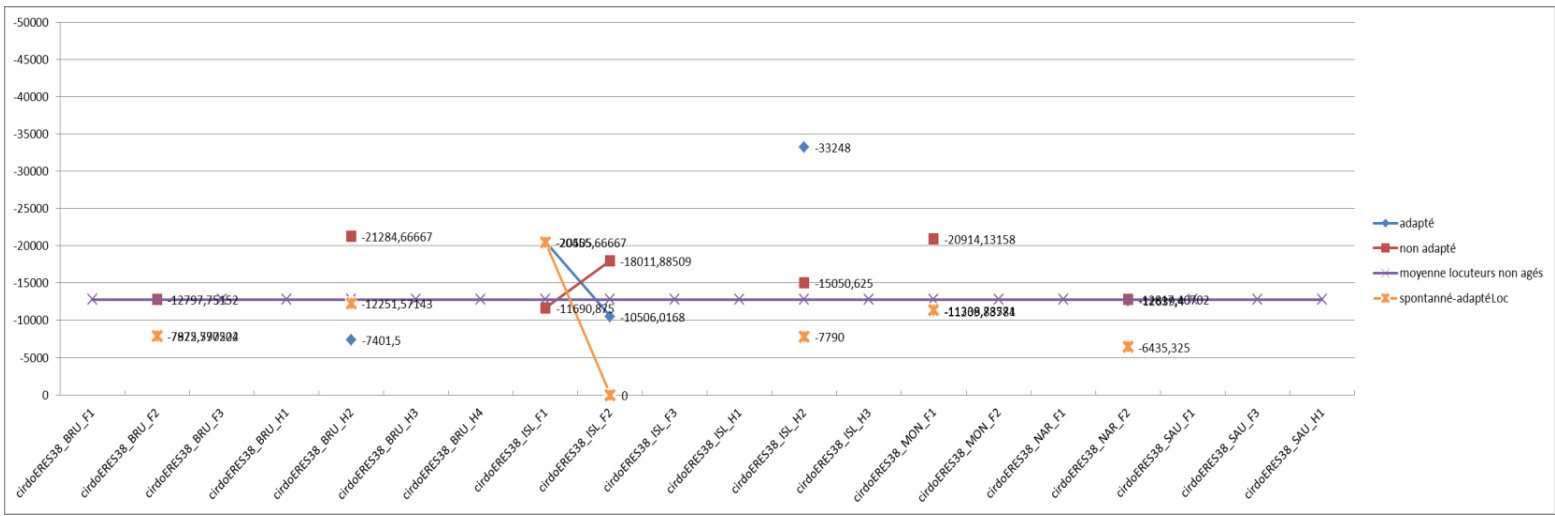
t



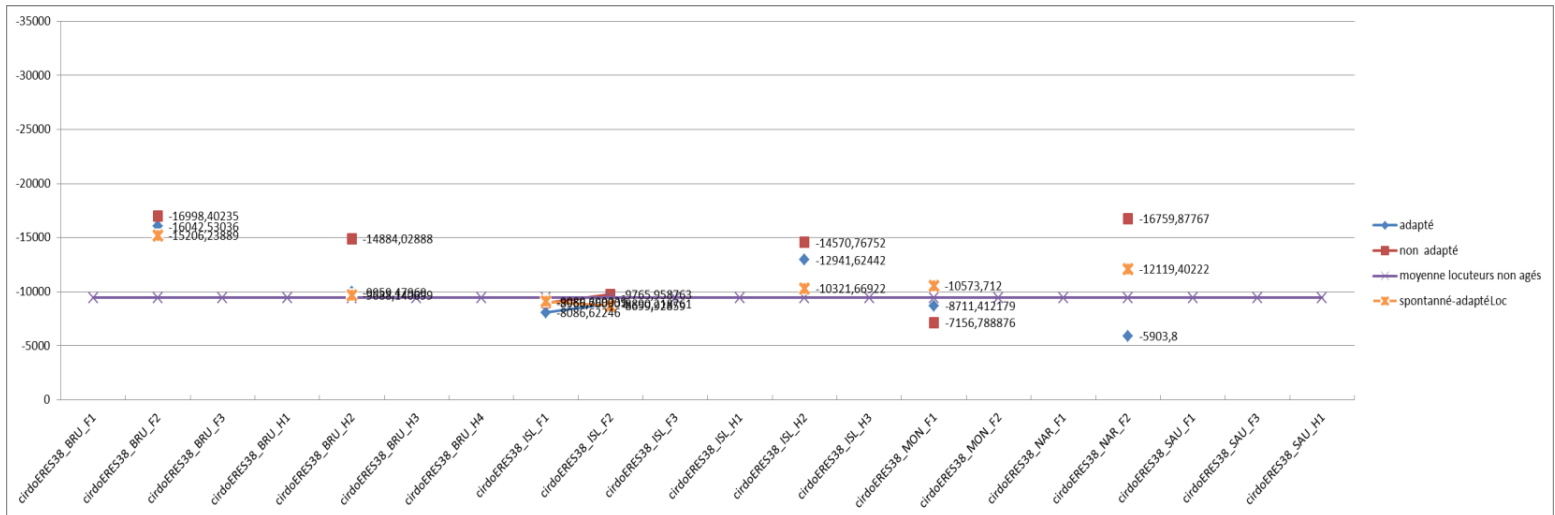
u



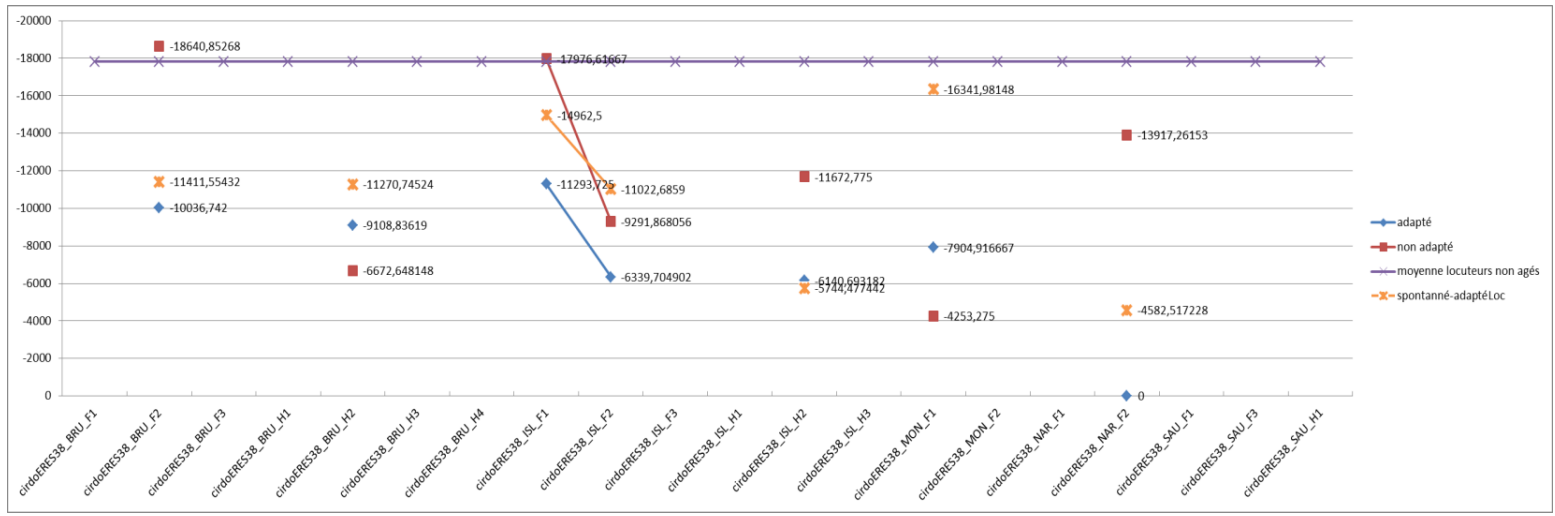
un



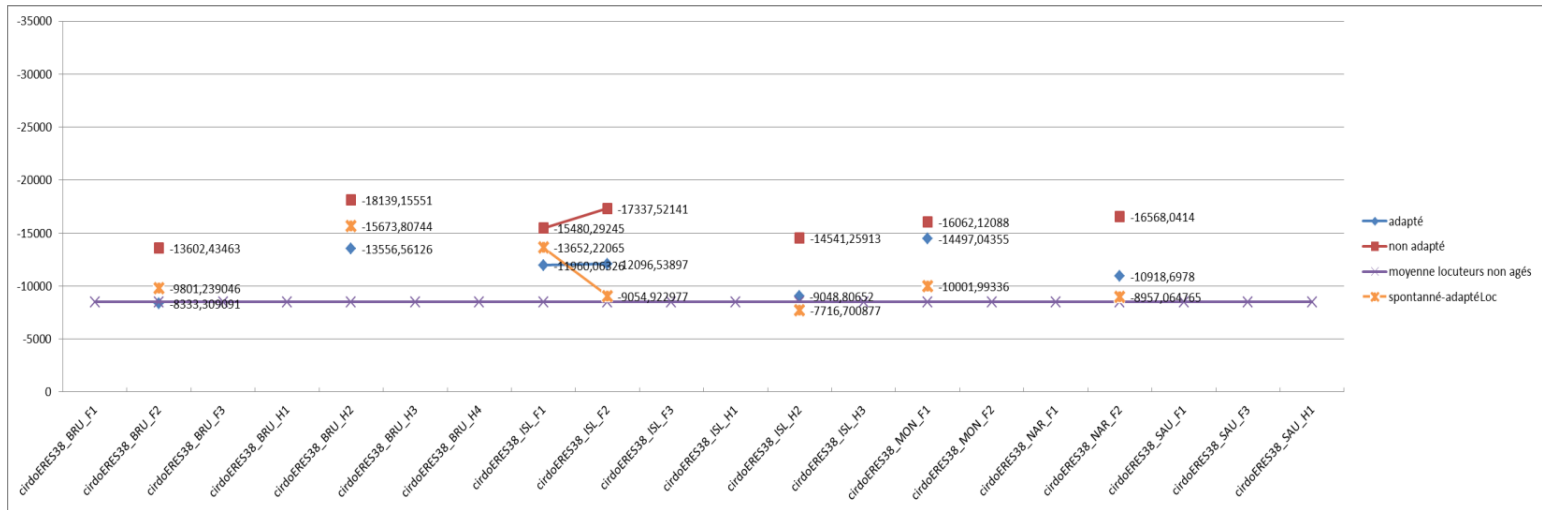
v



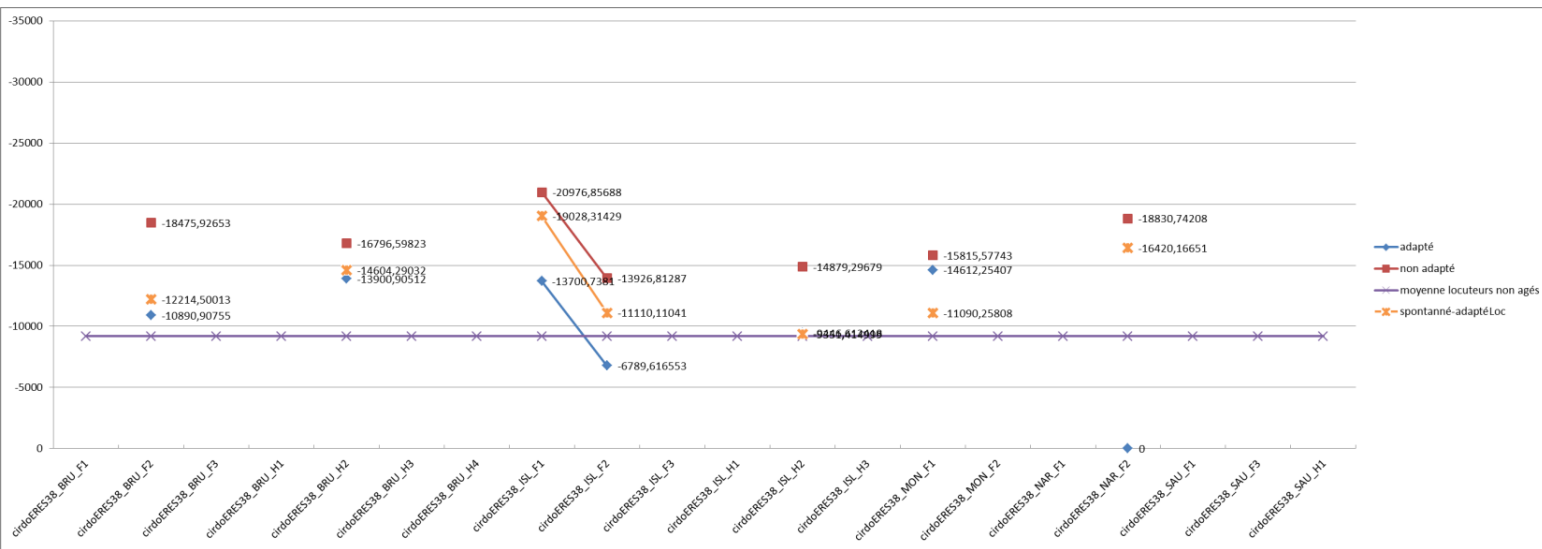
w



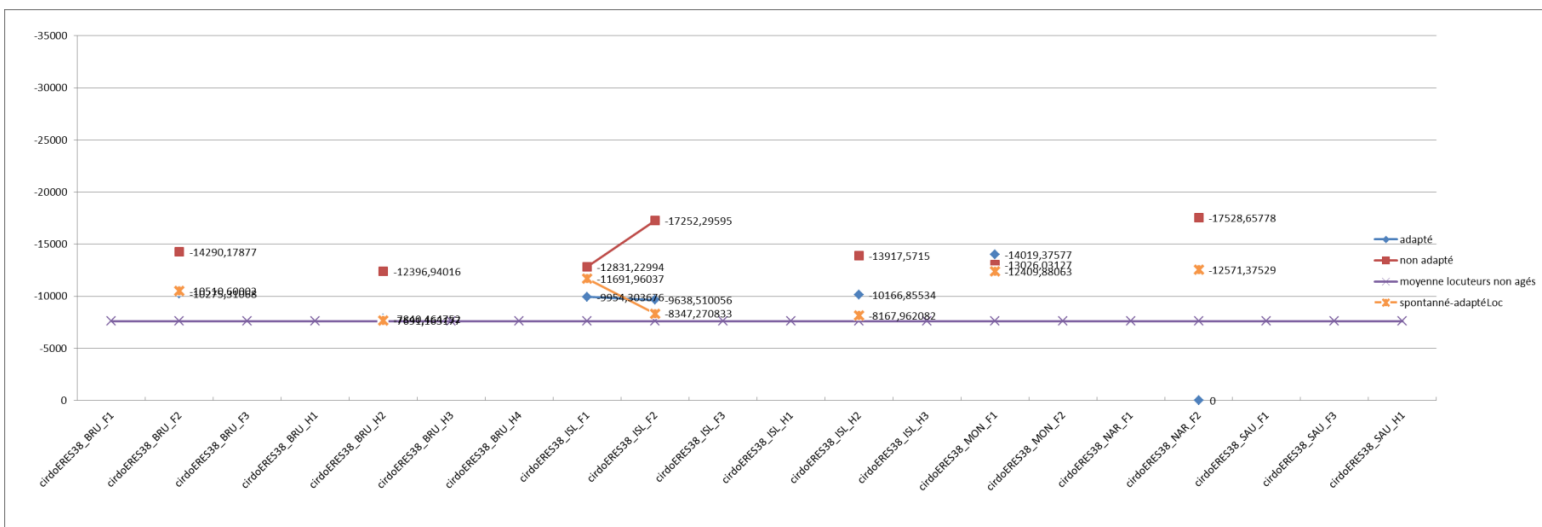
Y



Z



ZZ



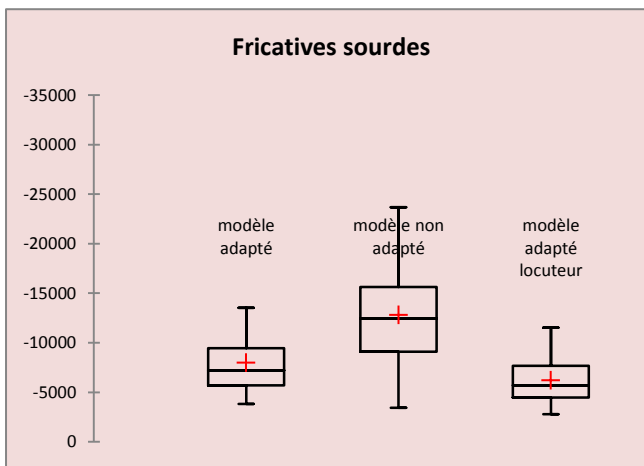
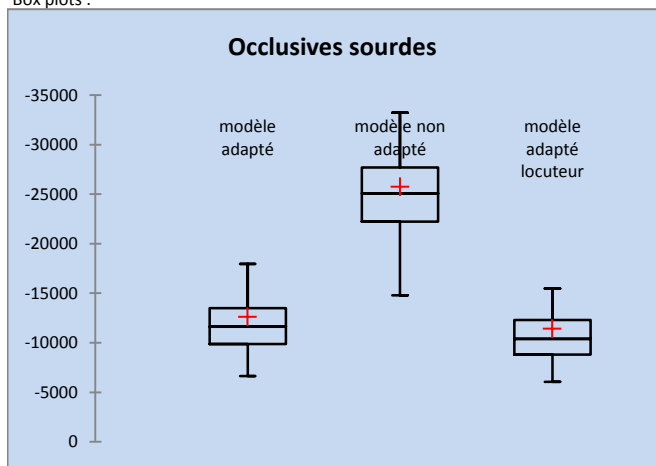
8.9 Diagrammes représentant les scores de l'alignement forcé pour la parole lue

Données quantitatives : Classeur = diagrammesStatsClassesPhon2.xlsx / Feuille = Feuil3 / Plage = Feuil3!\$B\$3:\$AH\$203 / 200 lignes et 33 colonnes

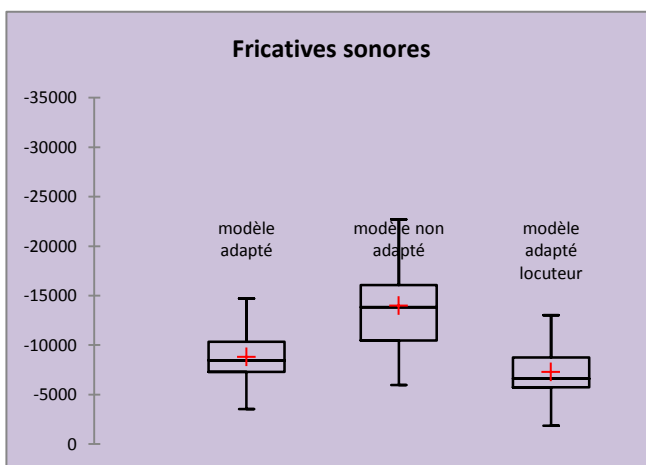
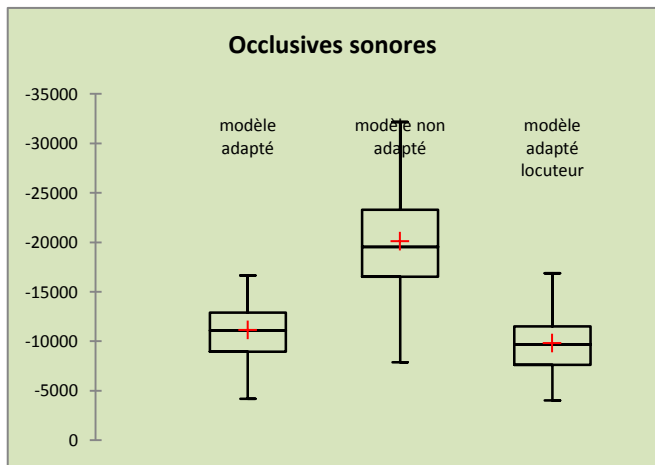
Statistiques descriptives (Données quantitatives) :

| Statistique | modèle adapté | modèle non adapté | dèle adapté locuteur | modèle adapté | modèle non adapté | dèle adapté locuteur |
|-------------------|---------------|-------------------|----------------------|---------------|-------------------|----------------------|
| Nb. d'observation | 200 | 200 | 200 | 200 | 200 | 200 |
| Minimum | -63831,667 | -70433,667 | -65666,667 | -16090,723 | -31976,849 | -14341,855 |
| Maximum | -6638,817 | -14801,237 | -6062,427 | -3805,160 | -3422,582 | -2788,658 |
| 1er Quartile | -13499,471 | -27701,289 | -12285,217 | -9460,531 | -15637,993 | -7677,568 |
| Médiane | -11625,218 | -25071,470 | -10401,728 | -7209,767 | -12423,271 | -5705,847 |
| 3ème Quartile | -9874,269 | -22217,792 | -8802,860 | -5711,050 | -9105,245 | -4480,550 |
| Moyenne | -12613,861 | -25727,544 | -11414,264 | -7961,300 | -12810,690 | -6208,211 |
| Variance (n-1) | ##### | 54550801,161 | 56225353,692 | 8828261,280 | 27837496,665 | 5446045,878 |
| Ecart-type (n-1) | 7223,063 | 7385,851 | 7498,357 | 2971,239 | 5276,125 | 2333,676 |

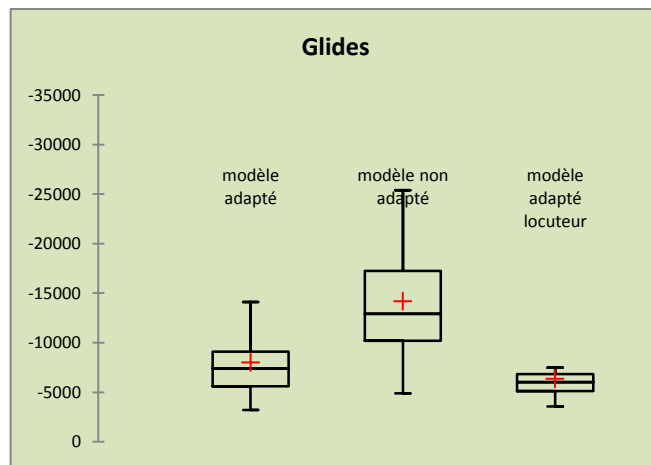
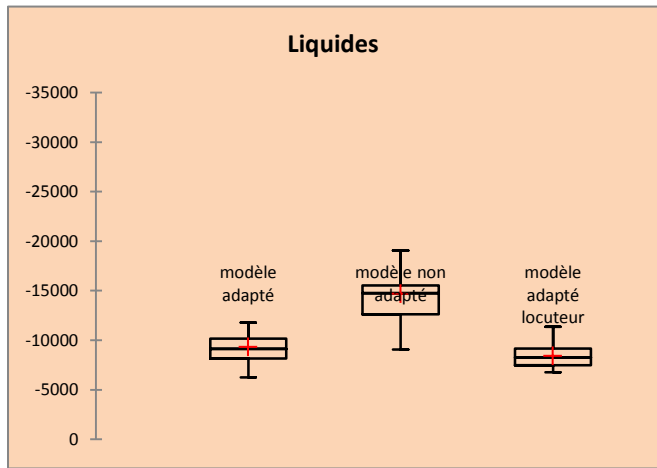
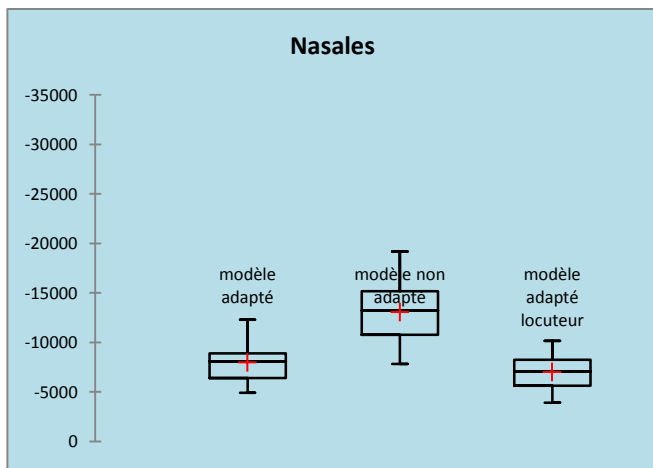
Box plots :



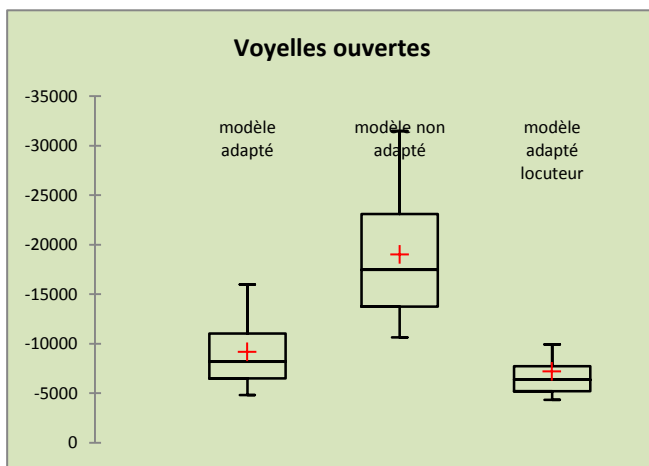
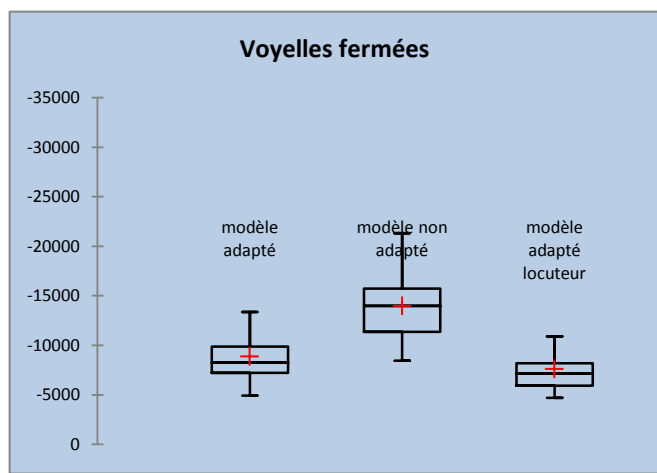
| Statistique | modèle adapté | modèle non adapté | dèle adapté locuteur | modèle adapté | modèle non adapté | dèle adapté locuteur |
|-------------------|---------------|-------------------|----------------------|---------------|-------------------|----------------------|
| Nb. d'observation | 200 | 200 | 200 | 200 | 200 | 200 |
| Minimum | -21477,497 | -42897,833 | -16861,517 | -15458,429 | -26550,763 | -13015,062 |
| Maximum | -4181,389 | -7869,904 | -4023,511 | -2290,983 | -5972,460 | -1848,830 |
| 1er Quartile | -12883,878 | -23292,296 | -11498,996 | -10339,733 | -16077,570 | -8758,270 |
| Médiane | -11088,864 | -19543,942 | -9664,733 | -8450,025 | -13833,293 | -6622,150 |
| 3ème Quartile | -8938,329 | -16528,540 | -7611,358 | -7284,590 | -10459,154 | -5741,720 |
| Moyenne | -11115,909 | -20086,540 | -9783,424 | -8794,142 | -13978,127 | -7266,993 |
| Variance (n-1) | 11604898,195 | 28935508,215 | 7750962,303 | 7136661,830 | 15436931,042 | 5650732,620 |
| Ecart-type (n-1) | 3406,596 | 5379,174 | 2784,055 | 2671,453 | 3928,986 | 2377,127 |



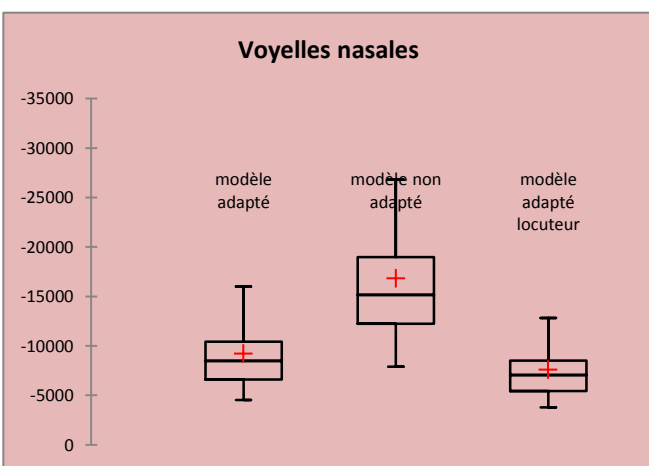
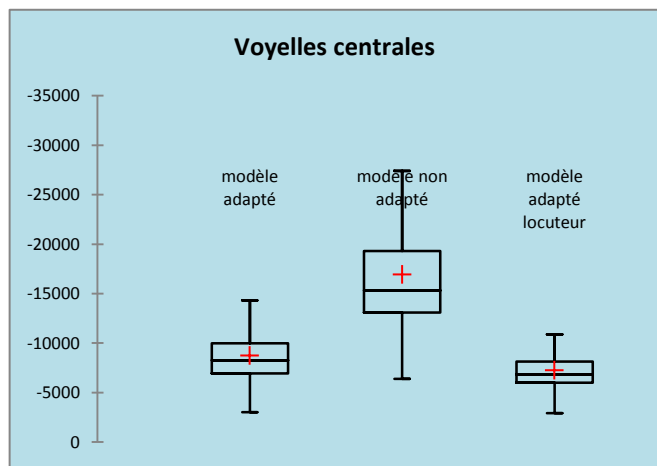
| modèle adapté | | | modèle non adapté | | | modèle adapté locuteur | | |
|---------------|-------------|-------------|-------------------|--------------|-------------|------------------------|--------------|-------------|
| 200 | 200 | 200 | 200 | 200 | 200 | 200 | 200 | 200 |
| -12289,221 | -19199,674 | -10153,201 | -16860,170 | -26512,131 | -11345,033 | -22417,870 | -30025,389 | -15708,878 |
| -4920,065 | -7814,171 | -3906,416 | -6248,956 | -9076,363 | -6756,765 | -3208,373 | -4871,000 | -3550,672 |
| -8884,637 | -15159,319 | -8260,724 | -10145,227 | -15548,570 | -9157,439 | -9101,374 | -17244,960 | -6835,526 |
| -8083,680 | -13220,007 | -7085,959 | -9141,005 | -14755,130 | -8241,052 | -7424,847 | -12915,955 | -6021,083 |
| -6403,654 | -10783,367 | -5624,864 | -8151,718 | -12634,278 | -7482,663 | -5597,665 | -10207,781 | -5128,284 |
| -7957,751 | -13060,397 | -6989,544 | -9305,678 | -14680,895 | -8401,697 | -7971,335 | -14155,082 | -6319,082 |
| 2864937,561 | 9040050,132 | 2356254,850 | 3391949,054 | 11888998,445 | 1381435,230 | 12095139,708 | 31925704,959 | 5243708,176 |
| 1692,613 | 3006,668 | 1535,010 | 1841,724 | 3448,043 | 1175,345 | 3477,807 | 5650,284 | 2289,914 |



| Statistique | modèle adapté | modèle non adapté | dèle adapté locuteur | modèle adapté | modèle non adapté | dèle adapté locuteur |
|-------------------|---------------|-------------------|----------------------|---------------|-------------------|----------------------|
| Nb. d'observation | 200 | 200 | 200 | 200 | 200 | 200 |
| Minimum | -20498,153 | -31725,125 | -28112,306 | -22139,933 | -31490,643 | -27057,333 |
| Maximum | -4928,691 | -8459,078 | -4690,410 | -4814,566 | -10641,024 | -4321,861 |
| 1er Quartile | -9877,144 | -15724,049 | -8186,204 | -11029,097 | -23103,168 | -7734,165 |
| Médiane | -8242,812 | -13972,550 | -7162,697 | -8195,275 | -17481,172 | -6355,647 |
| 3ème Quartile | -7231,280 | -11358,399 | -5911,730 | -6489,242 | -13748,334 | -5209,571 |
| Moyenne | -8845,560 | -13960,138 | -7580,977 | -9177,773 | -19004,096 | -7198,627 |
| Variance (n-1) | 6653039,981 | 15187804,023 | 10036350,072 | 13679083,236 | 37881775,386 | 17281139,554 |
| Ecart-type (n-1) | 2579,349 | 3897,153 | 3168,020 | 3698,524 | 6154,817 | 4157,059 |



| Statistique | modèle adapté | modèle non adapté | dèle adapté locuteur | modèle adapté | modèle non adapté | dèle adapté locuteur |
|-------------------|---------------|-------------------|----------------------|---------------|-------------------|----------------------|
| Nb. d'observation | 200 | 200 | 200 | 200 | 200 | 200 |
| Minimum | -30375,756 | -49740,667 | -24913,611 | -22277,412 | -45715,379 | -17089,769 |
| Maximum | -4532,134 | -7899,857 | -3790,212 | -3024,854 | -6370,700 | -2918,043 |
| 1er Quartile | -10415,681 | -18990,286 | -8498,122 | -9970,580 | -19294,262 | -8115,792 |
| Médiane | -8486,100 | -15144,062 | -7070,794 | -8249,044 | -15314,627 | -6847,202 |
| 3ème Quartile | -6591,946 | -12250,258 | -5440,624 | -6925,419 | -13093,843 | -5998,118 |
| Moyenne | -9224,268 | -16816,512 | -7578,924 | -8743,315 | -16918,879 | -7223,510 |
| Variance (n-1) | 15924844,058 | 54430275,074 | 10725571,881 | 9721513,208 | 36522980,166 | 5666063,170 |
| Ecart-type (n-1) | 3990,594 | 7377,688 | 3274,992 | 3117,934 | 6043,425 | 2380,349 |



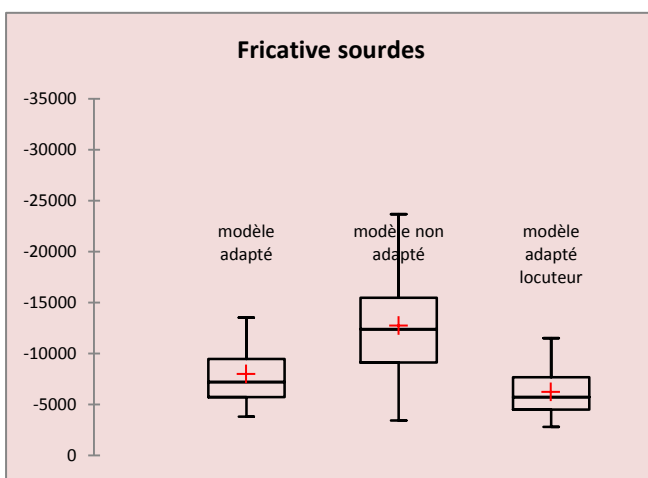
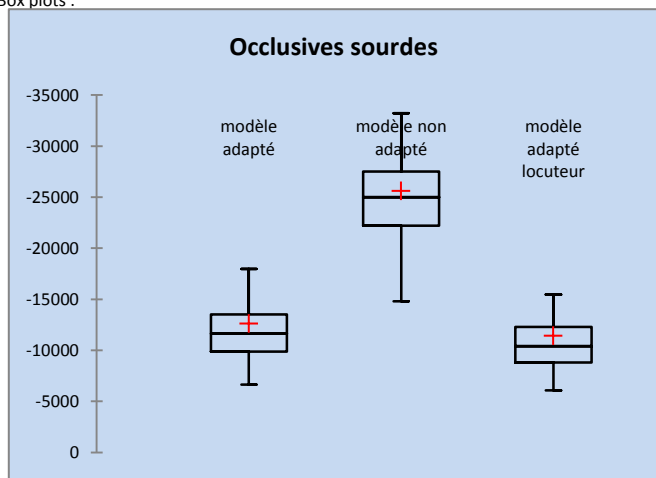
8.10 Diagrammes représentant les scores de l'alignement forcé pour la parole spontanée

Données quantitatives : Classeur = diagrammesStatsClassesPhon2.xlsx / Feuille = Feuil2 / Plage = Feuil2!\$B\$3:\$AH\$213 / 210 lignes et 33 colonnes

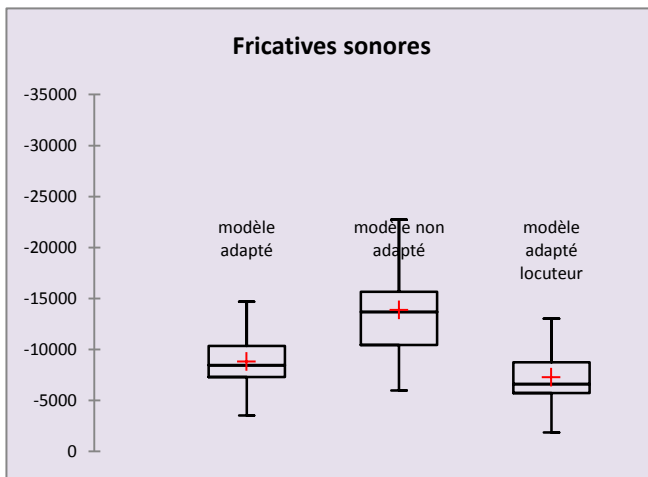
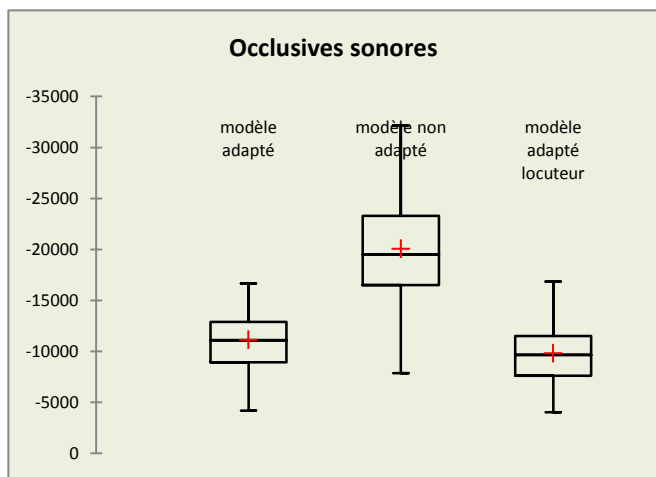
Statistiques descriptives (Données quantitatives) :

| Statistique | modèle adapté | modèle non adapté | dèle adapté locuteur | modèle adapté | modèle non adapté | dèle adapté locuteur |
|-------------------|---------------|-------------------|----------------------|---------------|-------------------|----------------------|
| Nb. d'observation | 210 | 210 | 210 | 210 | 210 | 210 |
| Minimum | -63831,667 | -70433,667 | -65666,667 | -16090,723 | -31976,849 | -14341,855 |
| Maximum | -6638,817 | -14801,237 | -6062,427 | -3805,160 | -3422,582 | -2788,658 |
| 1er Quartile | -13499,471 | -27507,704 | -12285,217 | -9460,531 | -15468,787 | -7677,568 |
| Médiane | -11625,218 | -24993,859 | -10401,728 | -7209,767 | -12390,333 | -5705,847 |
| 3ème Quartile | -9874,269 | -22215,131 | -8802,860 | -5711,050 | -9114,585 | -4480,550 |
| Moyenne | -12613,861 | -25611,579 | -11414,264 | -7961,300 | -12707,778 | -6208,211 |
| Variance (n-1) | 52172638,175 | 52256742,536 | 56225353,692 | 8828261,280 | 26738177,045 | 5446045,878 |
| Ecart-type (n-1) | 7223,063 | 7228,883 | 7498,357 | 2971,239 | 5170,897 | 2333,676 |

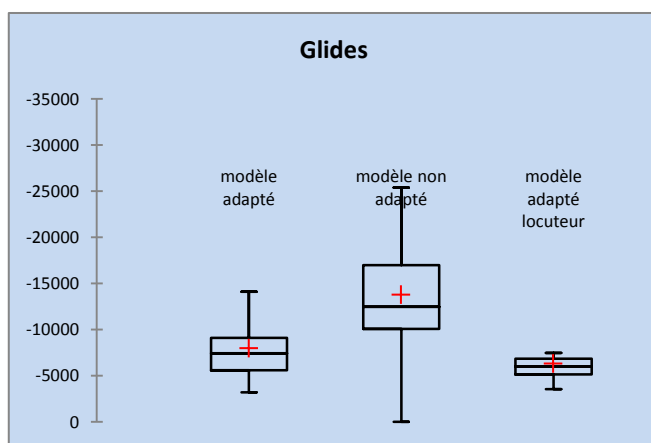
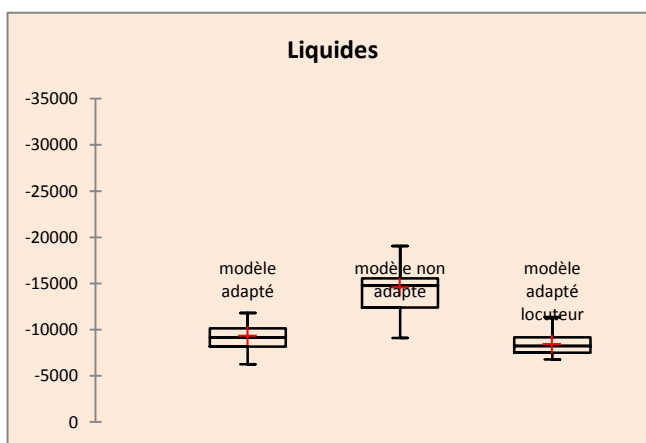
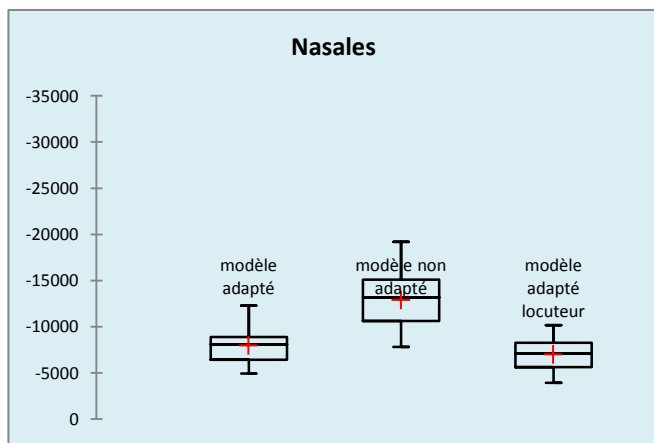
Box plots :



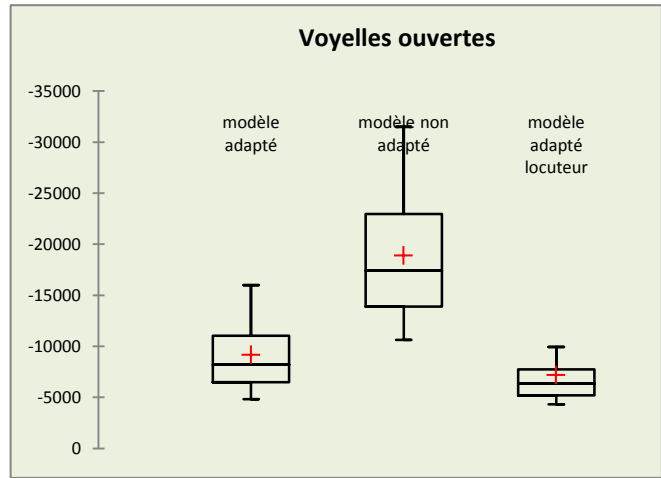
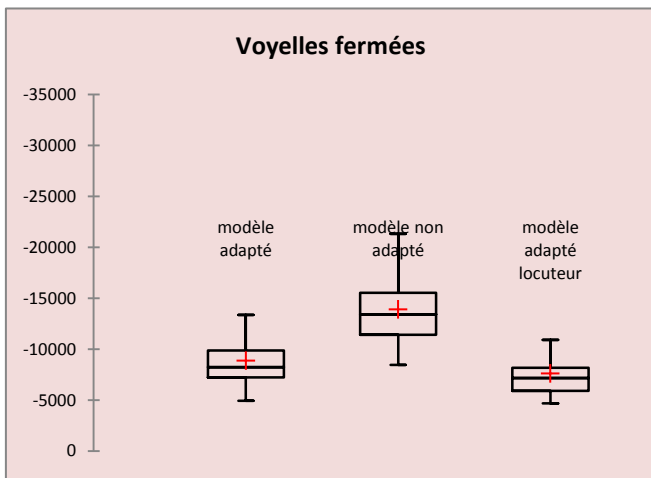
| Statistique | modèle adapté | modèle non adapté | dèle adapté locuteur | modèle adapté | modèle non adapté | dèle adapté locuteur |
|-------------------|---------------|-------------------|----------------------|---------------|-------------------|----------------------|
| Nb. d'observation | 210 | 210 | 210 | 210 | 210 | 210 |
| Minimum | -21477,497 | -42897,833 | -16861,517 | -15458,429 | -26550,763 | -13015,062 |
| Maximum | -4181,389 | -7869,904 | -4023,511 | -2290,983 | -5972,460 | -1848,830 |
| 1er Quartile | -12883,878 | -23299,159 | -11498,996 | -10339,733 | -15652,957 | -8758,270 |
| Médiane | -11088,864 | -19515,458 | -9664,733 | -8450,025 | -13692,231 | -6622,150 |
| 3ème Quartile | -8938,329 | -16496,183 | -7611,358 | -7284,590 | -10442,167 | -5741,720 |
| Moyenne | -11115,909 | -20032,089 | -9783,424 | -8794,142 | -13845,833 | -7266,993 |
| Variance (n-1) | 11604898,195 | 28144223,881 | 7750962,303 | 7136661,830 | 15109601,010 | 5650732,620 |
| Ecart-type (n-1) | 3406,596 | 5305,113 | 2784,055 | 2671,453 | 3887,107 | 2377,127 |



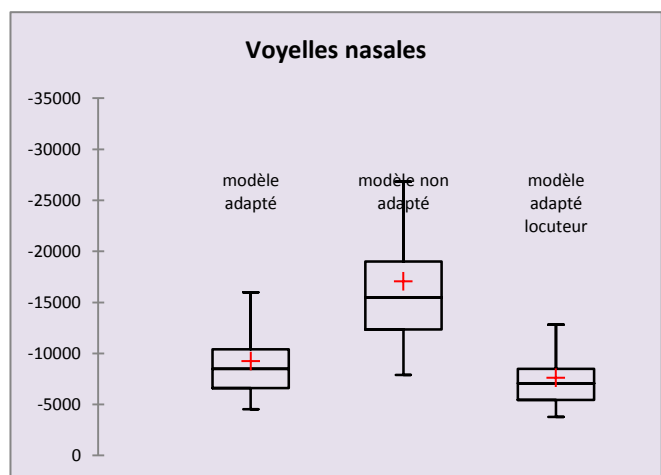
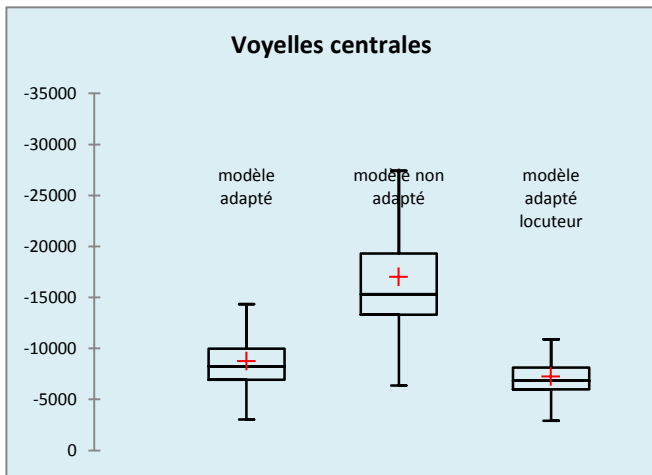
| modèle adapté | modèle non adapté | dèle adapté locut | modèle adapté | modèle non adapté | dèle adapté locut | modèle adapté | modèle non adapté | dèle adapté locut |
|---------------|-------------------|-------------------|---------------|-------------------|-------------------|---------------|-------------------|-------------------|
| 210 | 210 | 210 | 210 | 210 | 210 | 210 | 210 | 210 |
| -12289,221 | -19199,674 | -10153,201 | -16860,170 | -26512,131 | -11345,033 | -22417,870 | -30025,389 | -15708,878 |
| -4920,065 | -7814,171 | -3906,416 | -6248,956 | -9076,363 | -6756,765 | -3208,373 | 0,000 | -3550,672 |
| -8884,637 | -15080,141 | -8260,724 | -10145,227 | -15532,096 | -9157,439 | -9101,374 | -16984,940 | -6835,526 |
| -8083,680 | -13200,591 | -7085,959 | -9141,005 | -14754,465 | -8241,052 | -7424,847 | -12459,746 | -6021,083 |
| -6403,654 | -10607,898 | -5624,864 | -8151,718 | -12404,414 | -7482,663 | -5597,665 | -10061,583 | -5128,284 |
| -7957,751 | -12907,891 | -6989,544 | -9305,678 | -14568,615 | -8401,697 | -7971,335 | -13762,842 | -6319,082 |
| 2864937,561 | 9076257,321 | 2356254,850 | 3391949,054 | 11567335,846 | 1381435,230 | 12095139,708 | 33604305,140 | 5243708,176 |
| 1692,613 | 3012,683 | 1535,010 | 1841,724 | 3401,079 | 1175,345 | 3477,807 | 5796,922 | 2289,914 |



| Statistique | modèle adapté | nodèle non adapté | odèle adapté locuteur | modèle adapté | nodèle non adapté | dèle adapté locuteur |
|-------------------|---------------|-------------------|-----------------------|---------------|-------------------|----------------------|
| Nb. d'observation | 210 | 210 | 210 | 210 | 210 | 210 |
| Minimum | -20498,153 | -31725,125 | -28112,306 | -22139,933 | -31490,643 | -27057,333 |
| Maximum | -4928,691 | -8459,078 | -4690,410 | -4814,566 | -10641,024 | -4321,861 |
| 1er Quartile | -9877,144 | -15530,370 | -8186,204 | -11029,097 | -22958,020 | -7734,165 |
| Médiane | -8242,812 | -13412,646 | -7162,697 | -8195,275 | -17402,426 | -6355,647 |
| 3ème Quartile | -7231,280 | -11408,625 | -5911,730 | -6489,242 | -13901,235 | -5209,571 |
| Moyenne | -8845,560 | -13887,703 | -7580,977 | -9177,773 | -18909,157 | -7198,627 |
| Variance (n-1) | 6653039,981 | 14588642,529 | 10036350,072 | 13679083,236 | 36276258,688 | 17281139,554 |
| Ecart-type (n-1) | 2579,349 | 3819,508 | 3168,020 | 3698,524 | 6022,978 | 4157,059 |



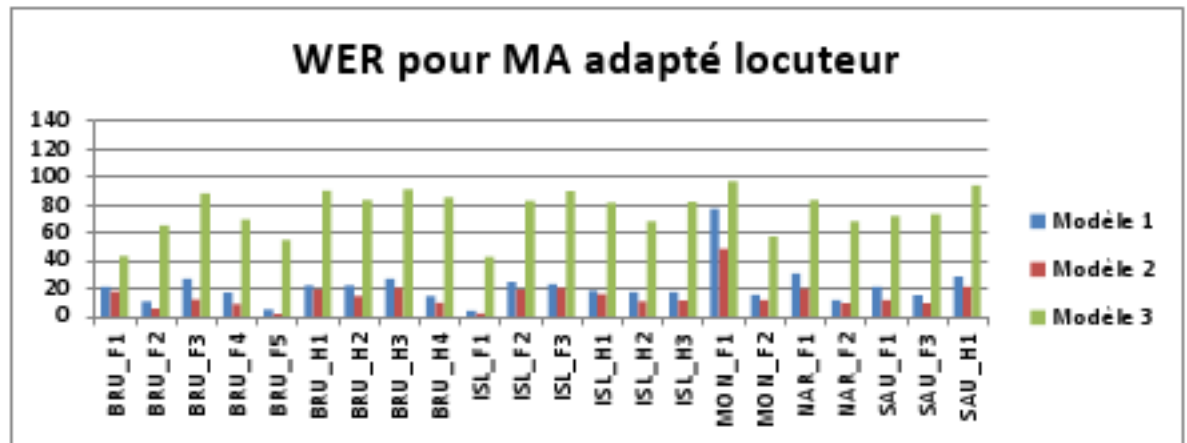
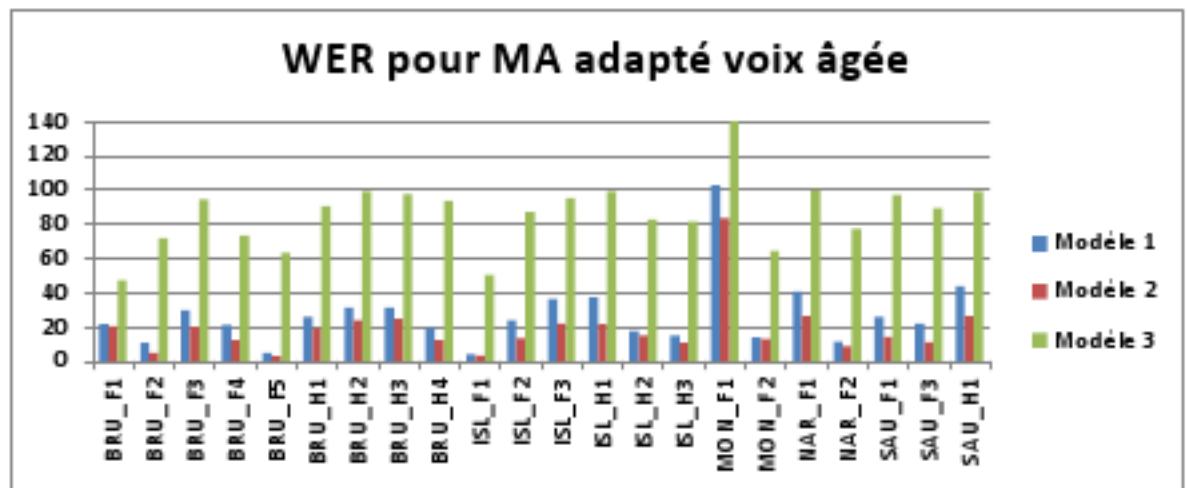
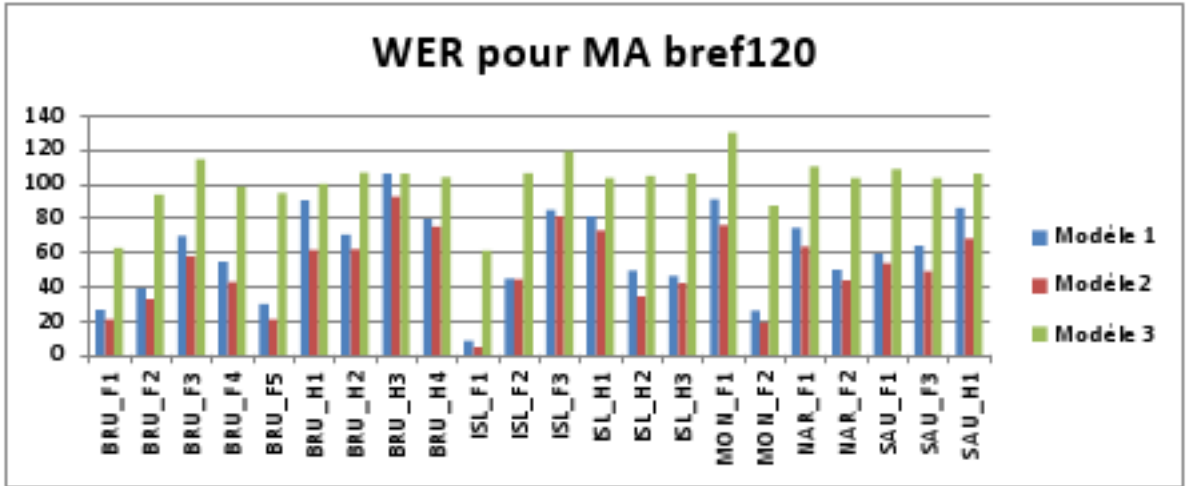
| Statistique | modèle adapté | nodèle non adapté | odèle adapté locuteur | modèle adapté | nodèle non adapté | dèle adapté locuteur |
|-------------------|---------------|-------------------|-----------------------|---------------|-------------------|----------------------|
| Nb. d'observation | 210 | 210 | 210 | 210 | 210 | 210 |
| Minimum | -30375,756 | -49740,667 | -24913,611 | -22277,412 | -45715,379 | -17089,769 |
| Maximum | -4532,134 | -7899,857 | -3790,212 | -3024,854 | -6370,700 | -2918,043 |
| 1er Quartile | -10415,681 | -18990,286 | -8498,122 | -9970,580 | -19307,819 | -8115,792 |
| Médiane | -8486,100 | -15477,577 | -7070,794 | -8249,044 | -15326,065 | -6847,202 |
| 3ème Quartile | -6591,946 | -12359,362 | -5440,624 | -6925,419 | -13306,528 | -5998,118 |
| Moyenne | -9224,268 | -17031,484 | -7578,924 | -8743,315 | -16982,677 | -7223,510 |
| Variance (n-1) | 15924844,058 | 55975442,681 | 10725571,881 | 9721513,208 | 35692760,408 | 5666063,170 |
| Ecart-type (n-1) | 3990,594 | 7481,674 | 3274,992 | 3117,934 | 5974,342 | 2380,349 |



8.11 Tableau complet ML vs MA

| locuteur | Modèle de langage 1 | | | Modèle de langage 2 | | | Modèle de langage 3 | | |
|----------|---------------------|-------------------------|------------------------|---------------------|-------------------------|------------------------|---------------------|-------------------------|------------------------|
| | Modèle bref 20 | Modèle adapté voix âgée | Modèle adapté locuteur | Modèle Bref 20 | Modèle adapté voix âgée | Modèle adapté locuteur | Modèle bref 20 | Modèle adapté voix âgée | Modèle adapté locuteur |
| | WER1 | WER2 | WER3 | WER4 | WER5 | WER6 | WER7 | WER8 | WER9 |
| BRU_F1 | 26,4 | 21,5 | 21,5 | 20,9 | 20,5 | 17,9 | 35,1 | 32 | 42,7 |
| BRU_F2 | 38,7 | 10,6 | 11 | 32,6 | 4,5 | 5,5 | 54,7 | 35 | 33,4 |
| BRU_F3 | 69,5 | 29,8 | 26,9 | 57,7 | 20 | 12,5 | 69,9 | 46,9 | 41 |
| BRU_F4 | 54,5 | 21,1 | 16,9 | 42,7 | 12,2 | 8,9 | 62,5 | 44,5 | 35,8 |
| BRU_F5 | 29,8 | 4,6 | 5,2 | 21 | 2,6 | 2 | 52,1 | 32,9 | 33,5 |
| BRU_H1 | 90,3 | 25,8 | 22,6 | 61,3 | 19,4 | 19,4 | 78,8 | 45,5 | 45,5 |
| BRU_H2 | 70,5 | 31,2 | 22,5 | 61,7 | 23,8 | 14,8 | 74,5 | 53,5 | 46,8 |
| BRU_H3 | 105,7 | 31,1 | 27,1 | 92,6 | 24,6 | 20,5 | 87,8 | 48,1 | 45 |
| BRU_H4 | 79,6 | 19,7 | 14,6 | 75,1 | 12 | 10 | 87,8 | 45,5 | 42 |
| ISL_F1 | 8,5 | 3,9 | 4,2 | 4,6 | 3,3 | 2,3 | 31,1 | 32 | 31,4 |
| ISL_F2 | 44,5 | 24 | 24,7 | 44,2 | 13,1 | 18,9 | 60,6 | 38,8 | 41,7 |
| ISL_F3 | 84,8 | 36,2 | 23,2 | 80,8 | 22 | 20,4 | 80,6 | 46,9 | 42,3 |
| ISL_H1 | 81 | 37,3 | 18,3 | 72,8 | 21,5 | 16,1 | 79,2 | 54,1 | 45,3 |
| ISL_H2 | 49,4 | 17,4 | 17,4 | 34,2 | 14,8 | 11 | 58,3 | 40,8 | 37 |
| ISL_H3 | 46,1 | 15 | 17,6 | 41,8 | 10,5 | 11,4 | 57,4 | 41,1 | 39,3 |
| MON_F1 | 91 | 102,6 | 77,5 | 75,7 | 83,5 | 48,7 | 69 | 79 | 55,9 |
| MON_F2 | 25,9 | 13,6 | 15,8 | 18,7 | 12,7 | 12 | 44,3 | 36,1 | 36,8 |
| NAR_F1 | 74,2 | 40,4 | 30,9 | 63,2 | 26,1 | 19,9 | 77,8 | 51,2 | 51 |
| NAR_F2 | 49,8 | 11,4 | 12 | 43,5 | 8,5 | 9,2 | 65,5 | 39,5 | 37,2 |
| SAU_F1 | 59,1 | 25,9 | 21,3 | 53,5 | 14,3 | 11,9 | 72,9 | 48 | 42,2 |
| SAU_F3 | 63,8 | 21,8 | 15,1 | 49 | 11,1 | 9,4 | 62,4 | 42 | 38,9 |
| SAU_H1 | 85,8 | 43,9 | 28,7 | 68,1 | 26,2 | 21,1 | 80,5 | 55,9 | 46,7 |
| TOTAL | 60,40454545 | 26,76363636 | 21,59090909 | 50,71363636 | 18,50909091 | 14,71818182 | 65,58181818 | 44,96818182 | 41,42727273 |
| Diffen % | | | | | | | | | |
| WER1-2 | 55,6926782 | | | 63,5027337 | | | 31,4319379 | | |
| WER2-3 | | 19,3274457 | | | 20,481336 | | | 7,87425452 | |
| WER1-3 | | | 64,2561517 | | | 70,9778614 | | | 36,8311616 |

8.12 Histogrammes présentant les WER en fonction des modèles acoustiques et par locuteur pour la parole lue



8.13 Détail des résultats obtenus pour la parole spontanée sur une partie du corpus ERES38

| locuteurs | modèle de langage 1 (productions pures) | | | modèle de langage 2 (productions+ESTER+PCF) | | |
|--------------------------|---|--------------|--------------|---|--------------|--------------|
| | MA bref120 | MA adaptVxAg | MA adaptLoc | MA bref120 | MA adaptVxAg | MA adaptLoc |
| BRU_F2 | 48,2 | 16,5 | 13,2 | 64 | 35,6 | 41,6 |
| BRU_H2 | 73,5 | 48,9 | 33,2 | 81 | 69,5 | 57,8 |
| ISL_F1 | 40,4 | 15,8 | 16,7 | 63,2 | 34,2 | 37,7 |
| ISL_F2 | 59,7 | 21,8 | 28,7 | 72 | 42,2 | 47,2 |
| ISL_H2 | 61,9 | 27,7 | 14,5 | 70,7 | 65,2 | 54,8 |
| MON_F1 | 36,1 | 11,4 | 16,7 | 51,9 | 43,5 | 40,1 |
| NAR_F2 | 68 | 21,1 | 24,5 | 79,6 | 53,1 | 55,8 |
| TOTAL ERES38_spc | 55,40 | 23,31 | 21,07 | 68,91 | 49,04 | 47,86 |
| TOTAL ERES38_lect | 58,95 | 23,15 | 18,93 | 65,42 | 43,35 | 40,74 |

