

Computer-Aided Pronunciation Training (CAPT) System for Non-native Learners of the Arabic Language

1- Introduction

In the following I will present the details of my research work in the field of Computer-Aided Pronunciation training (CAPT) for Arabic language. My work in the area started by supervising a master thesis titled “Automatic Pronunciation Error Detection of Nonnative Arabic Speech”. I continued the work by supervising a PhD thesis as co-supervisor with the title “Deep Learning-based Pronunciation Error Detection for Non-native Learners of the Arabic Language”. I did these two investigations as the head of the speech processing group at CCIS. Due to the group expertise in the area we were able to secure a funded project with the title “Computer-Aided Pronunciation Training System for Non-native Learners of the Arabic Language”, where I was the PI of the project.

In the following I will first present the letter of NPST about funding the project, then the abstracts of the Master thesis, PhD thesis, and the funded project, followed by a list of the papers in non-native Arabic pronunciation error detection followed by their abstract. Then I will end this report by a detailed report of our work in the funded project.

إفادة

تفيد الحطة الوطنية للعلوم والتقنية والابتكار بجامعة الملك سعود بان سعادة الدكتور / منصور بن محمد السليمان الأستاذ كلية علوم الحاسب والمعلومات، رقم وطني (١١٩٦٢٣)، هو الباحث الرئيس للمشاريع المبينة بالجدول ادناه، وهذه المشاريع ممولة من برنامج التقنيات الاستراتيجية بالحطة الوطنية للعلوم والتقنية والابتكار:

م	رقم المشروع	اسم المشروع	المدة
١	08-INF167-02	التعرف على المتحدث العربي ARABIC SPEAKER RECOGNITION	٢٠١٠ - ٢٠١٢
٢	١٢-MFD2474-02	تقييم الامراض الصوتية بالحاسب Automatic Voice Pathology Assessment	٢٠١٣ - ٢٠١٥
٣	3-17-09-001-0003	نظام حاسوبي لتعليم اللغة العربية لغير الناطقين بها Computer-Aided Pronunciation Training System for Non-native Learners of the Arabic Language	٢٠٢٠ - ٢٠٢٢
٤	٥-18-03-001-0003	نظام ترجمة محمول للغة الإشارة السعودية Saudi Sign Language Translation Companion System	٢٠٢٠ - ٢٠٢٢

وقد اعطيت له هذه الإفادة لسعدته بناء على طلبه لتتقدمها الى من يبعه الامر ودون أدنى مسؤولية على الوحدة.

مدير وحدة العلوم والتقنية والابتكار


د. أحمد بن عبد الله الحازم



2- Master Thesis

- **Title:** Automatic Pronunciation Error Detection of Non-native Arabic Speech
- **Date of defense:** 5/6/2014
- **My Role:** Supervisor, with Prof. Ghulam Muhammed and Prof. Saad Alqahtani as co-supervisors.

1- Abstract:

Computer assisted language learning (CALL) and, more specifically, computer assisted pronunciation training (CAPT) have received considerable attention in recent years. CAPT systems can provide many potential benefits to both the language learner and the teacher. They allow continuous feedback to the learner without requiring the sole attention of the teacher; they facilitate self-study and encourage interactive use of the language in preference to rote-learning. One of the important processes in CAPT system is error detection, which locates the errors in the utterance. Although Arabic is currently one of the most widely spoken languages in the world, there has been relatively little research about detection of the pronunciation error by nonnative speakers compared to the other languages. This research is concerned with detecting pronunciation errors of nonnative Arabic speakers from Pakistan and India. The sounds in this study were taken from KSU database. By analyzing the speech of the Pakistani and Indian speakers in KSU database we found that five speech sounds (Tha'a ث, Ha'a ح, Sad ص, Dad ض, Tha'a ظ) were often mispronounced by non-native speakers, hence this study will concentrate on these five pronunciation errors. The speech recognition techniques used was Hidden Markov Model (HMM). The system was built with native and non-native speakers, and tested with non-natives only. Goodness of Pronunciation (GOP) was calculated to detect if the phoneme was pronounced correctly or not. Comparison between the CAPT system judgment and the human judgment was performed. The result showed that GOP gave high accuracy, where the scoring accuracy were very good to excellent from 87% to 100% and the false rejection was from 0% to 10%.

3- PhD Thesis

- **Title:** Deep Learning-based Pronunciation Error Detection for Non-native Learners of the Arabic Language
- **Date of defense:** 26/1/2022
- **My Role:** Co-Supervisor, with Prof. Hassan mathkour as supervisor

2- Abstract:

In the recent decade, there has been great interest in computer-assisted pronunciation training (CAPT) systems. Many CAPT systems have been created for second language (L2) learners of various languages. Although Arabic is one of the most commonly spoken languages in the world, with the fifth highest number of speakers, little attention has been dedicated to computerized systems for the detection of pronunciation errors of non-Arabs. The Kingdom of Saudi Arabia is taking charge of serving the Arabic language, hence Arabic CAPT is important for the kingdom. Moreover, the CAPT system will enable the kingdom to help the large number of Muslims in the world to learn Arabic to read the Holy Quran.

Mispronunciation detection and diagnosis (MDD) module is a vital component of CAPT systems, because it will detect the mispronounced phonemes and provide different types of feedback to the learner. Compared with other languages, Arabic MDD system needs more investigating due to the scarcity of research in Arabic MDD in general and in using deep learning techniques for Arabic MDD in particular, and the lack of fully annotated non-native Arabic CAPT corpora. In this thesis, we aim to investigate different cutting edge deep learning techniques to build a high performance MDD system with feedback generation.

We tackled the research problem by several folds. In the first fold, the phoneme recognition task of MDD was formulated as an object detection task, where phonemes were considered as objects in spectral images. In the second fold, we designed a system for articulatory feature (AFs) detection by formulating the AFs detection as a multi-label detection problem. The performance of the proposed models was evaluated using Arabic corpus and benchmark English corpus. The system had excellent performance and was also light.

In the third fold, we leveraged the excellent finding of the first and second folds to develop an MDD system for non-native Arabic speech. The proposed system has the ability to detect mispronounced phonemes from the speech at the utterance level, as well as detect the AFs of each phoneme, simultaneously. Through detecting the AFs in addition to the phonemes, our proposed system can provide beneficial feedbacks to the learners, at articulatory level. Moreover, we proposed using genetic algorithm to find the best hyper-parameters of the deep neural network of the proposed models. We compared the performance of the proposed system with the state-of-the-art end-to-end MDD systems and our system had better result. In addition, we proposed using fusion between the proposed system and the end-to-end system and got better performance.

To tackle the problem of scarcity of non-native Arabic speech corpora, we investigated solving this by the use of different transfer learning techniques. We also developed a non-native Arabic speech corpus (Arabic-CAPT). Finally, we investigated using the recent neural Text to Speech (TTS) technique to develop a new synthesized non-native Arabic speech corpus.

4- **Project title:** Computer-Aided Pronunciation Training (CAPT) System for Non-native Learners of the Arabic Language

- **Funding Agency:** National Plan for Science, Technology and Innovation (MAARIFAH)
- **Period:** 2020-2022

Abstract of the Project

Computer-Aided Pronunciation Training (CAPT) System for Non-native Learners of the Arabic Language is important for the Kingdom of Saudi Arabia (KSA). This CAPT project is a joint work between two institutes of King Saud University; College of Computer and Information Sciences and the Arabic Language Institute. Both teams used their expertise in machine learning and Arabic Language to conduct research and develop an effective CAPT system. The system can be used offline or online to help learners of Arabic language correct their Arabic pronunciation.

In the first year, we designed the text for the KSU-CAPT-1 database and recorded the speech of 220 Arabic learners. The speech was verified and cleaned from extra sounds then sent to experts in Arabic language to time label it. In the second year, we cooperated with a linguistic scholar who is also an experienced instructor of Arabic as a second language to propose a new methodology to choose the text of the second database (KSU-CAPT-2) to complement the methodology that we used to select the text of the KSU-CAPT-1.

We proposed a new way of using deep learning for detection and recognition of phoneme and articulatory features (AF). In this proposed method, we treat the phonemes and AFs as objects in 3 channels spectral images of the speech. By this proposed method we were able to recognize the sequence of phoneme from the whole utterance of the non-native Arabic speakers. We used the detected phonemes for mispronunciation detection and diagnosis task and the detected AFs for feedback of error in pronunciation. This achievement was published in a two ISI papers. We did more investigation and got excellent results that are comparable or better than the state-of-the-art research and we are finalizing a new paper with these achievements.

Our proposal was based on state-of-the-art techniques at time of submission, which are based on HMM and GOP. When we started the project, the majority of the current state-of-the-art CAPT system approaches were based on End-to-End deep neural network techniques, which

implies that the system can recognize the sequence of phonemes/words from the input acoustic data using only one network. Hence our system is based on proposed End-to-End deep neural network techniques.

We aim in this project to design very accurate Arabic CAPT system that can detect the pronounced phonemes and the associated articulatory features (AFs), simultaneously, from the whole utterance using only one network. To this end, we formulated the MDD problem as a multi-label object detection problem by treating the phoneme and its AFs as objects in the three channels spectral images. We applied fast and accurate object detector as an acoustic model to recognize the sequence of phonemes and AFs from the spectral image of the whole utterance. We called the proposed system, Deep-CAPT, as shown in Figure 1. We applied the proposed system on non-native Arabic and English speech and got excellent results.

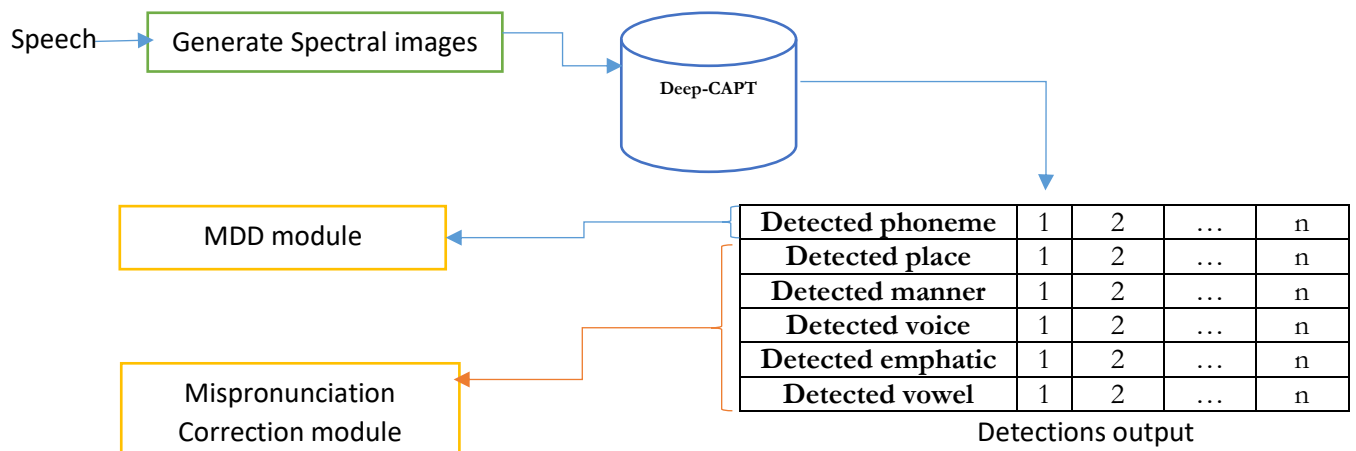


Figure 1: General diagram of the Deep-CAPT system.

5- List of papers in Computer-Aided Pronunciation Training (CAPT) and computerization of Arabic pronunciation training

ISI journals:

- 1- Algabri, Mohammed, Hassan Mathkour, Mohamed Abdelkader Bencherif, Mansour Alsulaiman, and Mohamed Amine Mekhtiche. "Towards deep object detection techniques for phoneme recognition." *IEEE Access* 8 (2020): 54663-54680.
- 2- Algabri, Mohammed, Hassan Mathkour, Mansour M. Alsulaiman, and Mohamed A. Bencherif. "Deep learning-based detection of articulatory features in arabic and english speech." *Sensors* 21, no. 4 (2021): 1205.
- 3- Algabri, Mohammed, Hassan Mathkour, Mansour M. Alsulaiman, and Mohamed A. Bencherif. "Mispronunciation error detection and diagnosis with articulatory feedback generation for non-native Arabic speech", *Mathematics* (Under review).

Conferences:

- 1- Al Hindi, Afnan, Mansour Alsulaiman, Ghulam Muhammad, and Saad Al-Kahtani. "Automatic pronunciation error detection of nonnative Arabic Speech." In 2014 IEEE/ACS 11th International Conference on Computer Systems and Applications (AICCSA), pp. 190-197. IEEE, 2014.
- 2- Alsulaiman, Mansour, Zulfiqar Ali, Ghulam Muhammad, Afnan Al Hindi, Taha Alfakih, Hussein Obeidat, and Saad Al-Kahtani. "Pronunciation errors of non-Arab learners of Arabic language." In 2014 International Conference on Computer, Communications, and Control Technology (I4CT), pp. 277-282. IEEE, 2014.

6- Papers Abstracts of the ISI journals

- 1- **Title:** Algabri, Mohammed, Hassan Mathkour, Mohamed Abdelkader Bencherif, Mansour Alsulaiman, and Mohamed Amine Mekhtiche. "Towards deep object detection techniques for phoneme recognition." *IEEE Access* 8 (2020): 54663-54680.

- **Abstract:**

The use of cutting edge object detection techniques to build an accurate phoneme sequence recognition system for English and Arabic languages is investigated in this study. Recently, numerous techniques have been proposed for object detection in daily life applications using deep learning. In this paper, we propose the use of object detection techniques in speech processing tasks. We selected two state-of-the-art object detectors, namely YOLO and CenterNet, based on a trade-off between detection accuracy and speed. We tackled the problem of phoneme sequence recognition using three systems: the domain transfer learning system (DTS) from image to speech, intra-language transfer learning system (IaTS) between speech corpora within the same language (English to English), and inter-language transfer learning system (IeTS) between speech corpora from dissimilar languages (English to Arabic). For English phoneme recognition, the Texas Instruments/Massachusetts Institute of Technology (TIMIT) corpus is used to evaluate the performance of the proposed systems. Our IaTS based on the CenterNet detector achieves the best results using the test core set of TIMIT with 15.89% phone error rate (PER). For Arabic phoneme recognition, the best performance, with 7.58% PER, was achieved using the CenterNet. These results show the effectiveness of using object detection techniques in phoneme recognition tasks. Furthermore, based on the findings of this study, speech processing tasks may be treated as object detection tasks..

- 2- **Title:** Algabri, Mohammed, Hassan Mathkour, Mansour M. Alsulaiman, and Mohamed A. Bencherif. "Deep learning-based detection of articulatory features in arabic and english speech." *Sensors* 21, no. 4 (2021): 1205.

- **Abstract:**

This study proposes using object detection techniques to recognize sequences of articulatory features (AFs) from speech utterances by treating AFs of phonemes as multi-label objects in speech spectrogram. The proposed system, called AFD-Obj, recognizes sequence of multi-label AFs in speech signal and localizes them. AFD-Obj consists of two main stages: firstly, we formulate the problem of AFs detection as an object detection problem and prepare the data to fulfill requirement of object detectors

by generating a spectral three-channel image from the speech signal and creating the corresponding annotation for each utterance. Secondly, we use annotated images to train the proposed system to detect sequences of AFs and their boundaries. We test the system by feeding spectrogram images to the system, which will recognize and localize multi-label AFs. We investigated using these AFs to detect the utterance phonemes. YOLOv3-tiny detector is selected because of its real-time property and its support for multi-label detection. We test our AFD-Obj system on Arabic and English languages using KAPD and TIMIT corpora, respectively. Additionally, we propose using YOLOv3-tiny as an Arabic phoneme detection system (i.e., PD-Obj) to recognize and localize a sequence of Arabic phonemes from whole speech utterances. The proposed AFD-Obj and PD-Obj systems achieve excellent results for Arabic corpus and comparable to the state-of-the-art method for English corpus. Moreover, we showed that using only one-scale detection is suitable for AFs detection or phoneme recognition.

Conferences:

- 1- **Title:** Al Hindi, Afnan, Mansour Alsulaiman, Ghulam Muhammad, and Saad Al-Kahtani. "Automatic pronunciation error detection of nonnative Arabic Speech." In *2014 IEEE/ACS 11th International Conference on Computer Systems and Applications (AICCSA)*, pp. 190-197. IEEE, 2014.

- **Abstract:**

Computer assisted language learning (CALL) and, more specifically, computer assisted pronunciation training (CAPT) have received considerable attention in recent years. CAPT allows continuous feedback to the learner without requiring the sole attention of the teacher; it facilitates self study and encourages interactive use of the language in preference to rote learning. One of the important processes in CAPT system is error detection, which locates the errors in the utterance. Although Arabic is currently one of the most widely spoken languages in the world, there has been relatively little research about detection of the pronunciation error by nonnative speakers compared to the other languages. This research is concerned with detecting pronunciation errors of nonnative Arabic speakers from Pakistan and India. All the sounds in this study were taken from King Saud University (KSU) Arabic Speech Database. By analyzing the speech of the Pakistani and Indian speakers in KSU database we found that five phonemes were often mispronounced by nonnative speakers, hence this research will concentrate on pronunciation errors in these five phonemes. The system was built with native and nonnative speakers, and tested with nonnative only. For each phoneme, the Goodness of Pronunciation (GOP) was calculated and compared with a threshold to decide if the phoneme was pronounced correctly or not. The result showed that GOP gave high accuracy, where the scoring accuracy was very good to excellent from 87% to 100%, and the false rejection was zero to less than 10%. This machine judgment is compared with human judgment and the comparison shows excellent agreement between them.

- 2- **Title:** Alsulaiman, Mansour, Zulfiqar Ali, Ghulam Muhammad, Afnan Al Hindi, Taha Alfakih, Hussein Obeidat, and Saad Al-Kahtani. "Pronunciation errors of non-Arab learners of Arabic language." In *2014 International Conference on Computer, Communications, and Control Technology (I4CT)*, pp. 277-282. IEEE, 2014.

- **Abstract:**

Arabic is one of the most widely spoken languages in the world, but little attention has been paid to detect pronunciation errors of non-Arabs from different nationalities. In this paper, the speech of four nationalities of Asian non-Arabs speaking different mother languages is analyzed to identify their pronunciation errors while learning Arabic as a foreign language. Two human experts have evaluated the speech of all speakers for mispronunciation, and evaluation show that the nature of errors is almost the same for all nationalities under investigation with the fact that pharyngeal, alveo-dental, and interdental sounds are difficult to pronounce by learners of Arabic. Some of the errors are due to sounds that are not present in mother languages or the learner is unable to pronounce the phoneme correctly due to similar place of production or/and manner of articulation. An interesting observation is that pronunciation errors of some learners are due to switching pairs of phonemes.

Transmittal Letter

Date: 15th/Aug/2021

Researcher name: Mansour Alsulaiman

College: Computer and Information Sciences

Department: Computer Engineering

Address: P.O. Box 51178, Riyadh 11543

E-mail: msuliman@ksu.edu.sa

Dear Prof. Ahmed Alkhazim

We are submitting to you the report, due 21/7/2021 that you requested. The report is entitled technical report for first year of the project Computer-Aided Pronunciation Training (CAPT) System for Non-native Learners of the Arabic Language. The purpose of the report is to inform you of our work in the first of the project. The content of this report concentrates on the results that we got and published for the developed CAPT system using our proposed technique. This report also discusses the database that we recorded for Non-Arab speech. If you should have any questions concerning our project, please feel free to contact Mansour Alsulaiman at 0503255927 or msuliman@ksu.edu.sa.

Sincerely,

Professor

Mansour Alsulaiman (PI)

Affiliation: College of Computer and Information Sciences, King Saud

University.

Title Page

Submitted for

National Science, Technology and Innovation Plan (NSTIP)

King Saud University

Project title

Computer-Aided Pronunciation Training (CAPT) System for Non-native Learners of the Arabic Language

Project number

3-17-09-001-0003

Project Investigator

Mansour Alsulaiman

Year

2021

Abstract

Computer-Aided Pronunciation Training (CAPT) System for Non-native Learners of the Arabic Language is important for the Kingdom of Saudi Arabia (KSA). This CAPT project is a joint work between two institutes of King Saud University; College of Computer and Information Sciences and the Arabic Language Institute. Both teams used their expertise in machine learning and Arabic Language to conduct research and develop an effective CAPT system. The system can be used offline or online to help learners of Arabic language correct their Arabic pronunciation.

In the first year, we designed the text for the project database (DB) and recorded the speech of 220 Arabic learners. The speech was verified and cleaned from extra sounds then sent to experts in Arabic language to time label it.

We proposed a new way of using deep learning for detection and recognition of phoneme and articulatory features (AF). In this proposed method, we treat the phonemes and AFs as objects in 3 channels spectral images of the speech. By this proposed method we were able to recognize the sequence of phoneme from the whole utterance of the non-native Arabic speakers. We used the detected phonemes for mispronunciation detection and diagnosis task and the detected AFs for feedback of error in pronunciation. This achievement was published in a paper in a Q2 ISI journal. We did more investigation and got excellent results that are comparable or better than the state-of-the-art research and we are finalizing a new paper with these achievements.

In next year, we will submit a paper (90% ready) on the current findings, build a new CAPT system based on the labeled DB, record and label second session of the DB, modify the CAPT system based on the second session of the DB, publish more papers based on the new researches and findings.

Acknowledgments

This work is supported by National Science, Technology and Innovation Plan (NSTIP) in King Saud University under grant number 3-17-09-001-0003. The authors are grateful for this support.

Table of Contents

Contents

Abstract.....	3
Acknowledgments.....	4
Contents	5
List of Figures	7
List of Tables	8
Report Body.....	9
1. Introduction.....	9
2. Objectives.....	11
2.1. Building the database	11
2.2. Literature Review.....	11
2.3. Building the CAPT system.....	12
2.4. Establishment of a multidisciplinary research group through the collaboration of the CCIS-team and the ALI-team Institute.....	12
2.5. Dissemination of the results and conclusions at conferences and in journals.	12
3. Development of the KSU-CAPT Non-Arabs database.....	13
3.1. Selection of text for recording of the speech corpus.....	13
3.2. Database Specifications.....	19
3.3. Establishment of the recording system.....	20
3.4. Speaker Registration	24
3.5. Database Recording (Session 1).....	27
4. Literature review of CAPT systems and of Arabic CAPT.....	29
4.1. Literature review for Arabic language CAPT in Arabic literature.....	29
4.2. Literature review of Arabic CAPT.....	36
4.3. Literature review of speech features for CAPT.	37
4.4. Classification methods for CAPT	38
4.5. Scoring in CAPT	39
4.6. Databases for CAPT.....	42
5. Design and development of a CAPT system for L2 learners of the Arabic language.	

43

5.1.	Error detection and analysis by human experts.....	44
5.2.	Labeling of speech of session 1.....	51
5.3.	Subjective evaluation of the pronunciation of L2 learners	53
5.4.	Error Analysis	55
5.5.	Implementation of speech recognition phase in the CAPT.....	55
5.6.	Investigating using the proposed technique for phoneme and AF detection in Arabic and English speech.....	57
6.	Discussion	70
7.	Future work	70
8.	References	71
9.	Publications / Presentations.....	77
10.	Appendices	78
	Appendix A - Text Selection Comparison (V1 to V3)	78
	Appendix B - Selected text for the Arabic CAPT recording system	79
	Appendix C - Tahadath App Screen Cards.....	81
	Appendix D - Durations per Speaker.....	81
	Appendix E- ELAN Annotation CAPT protocol.....	83
	Appendix F - Tahadath Application User Manual.....	84

List of Figures

<i>Figure 1, Phoneme distribution in the sentences and dual pairs of the various text selections</i>	17
<i>Figure 2, Lifecycle of the Tahadath Mobile app</i>	20
<i>Figure 3, Unity APP screenshot</i>	21
<i>Figure 4, cs2r firebase real-time database</i>	23
<i>Figure 5, Screen shot of the Google form sent to the students for the CAPT enrollment</i>	24
<i>Figure 6, Statistics of the Nationalities of the CAPT speech recording at session 1</i>	25
<i>Figure 7, CAPT Completed Recordings per University</i>	26
<i>Figure 8, Statistics of the L1 of the CAPT speech recording for the session 1 – Non-Arabs</i>	26
<i>Figure 9: General Architecture of the proposed Arabic CAPT system.</i>	44
<i>Figure 10, Screenshot of the ELAN screen work</i>	49
<i>Figure 11, Proposed System of the AFD-Obj and the PD-Obj</i>	58
<i>Figure 12, Overall Conversions of the speech signal to images (a: upper image, b: lower image)</i>	59
<i>Figure 13, Example of converting the detected AFs to the corresponding phonemes.</i>	60
<i>Figure 14, Testing example of converting the detected AFs using the YOLOv3-tiny-1S model to the corresponding phonemes and calculating the percentage of correct phonemes using the HResults tool (file “CMSSSFA”) from the KAPD corpus test set. X sign means invalid output, which occurs when the minimum hamming distance is greater than threshold (threshold = zero in case of 100% similarity).</i>	63
<i>Figure 15, Testing phase of the AFD-Obj system: calculating the frame level accuracy of the detected outputs.</i>	65

List of Tables

<i>Table 1, Sample sentences and dual phonetic words</i>	15
<i>Table 2, Benefits and drawbacks of the online recording solution</i>	16
<i>Table 3, Statistics of the CAPT-Text selections</i>	17
<i>Table 4, Comparative between the text selections of the various versions</i>	18
<i>Table 5, Sample Metadata of the CAPT recordings</i>	19
<i>Table 6, Sample screens from the Tahadath Mobile App</i>	22
<i>Table 7, Number of students that enrolled via the Google Form App</i>	25
<i>Table 8, Credentials received from /sent to the enrolled speakers (students)</i>	27
<i>Table 9: CAPT Databases' Survey.</i>	42
<i>Table 10, Details of the per-processing teams.</i>	44
<i>Table 11, Sample reports of the content checking</i>	46
<i>Table 12, Statistics about the session 1 speech recording</i>	47
<i>Table 13, Additional deep statistics about the session 1 speech recording of Non Arabs: SPW</i>	48
<i>Table 14, Additional deep statistics about the session 1 speech recording of Non Arabs: Paragraphs</i>	48
<i>Table 15, Annotator speech phoneme level segmentation sample (Empty form to be filled by annotators)</i>	51
<i>Table 16, Examples of annotating substitution and addition</i>	53
<i>Table 17, Examples of annotating deletion and other speech processes</i>	54
<i>Table 18: PER for the TIMIT test set.</i>	56
<i>Table 19: PER for non-native Arabic Speech (Small-Arabic-CAPT).</i>	57
<i>Table 20, Performance metrics of the proposed system AFD-Obj for the Arabic AFs</i>	61
<i>Table 21, PER (%) and correction rate (%) for our proposed AFD-Obj system and results of [61]</i>	64
<i>Table 22, Detection accuracy of all 28 English AFs using the proposed system AFD-Obj and state-of-the-art methods</i>	66
<i>Table 23, PER and correction rate of the Arabic phoneme recognition using the proposed models.</i>	68

Report Body

1. Introduction

Computer-Aided Pronunciation Training (CAPT) System for Non-native Learners of the Arabic Language is an important topic for the Kingdom of Saudi Arabia, as the Kingdom is taking charge of spreading knowledge about Islam and helping Muslims learn Arabic the language of the Quran. In this context, we have been working on this CAPT project which is a joint work between two institutes of King Saud University; the College of Computer and Information Sciences team (CCIS-T) and the Arabic Language Institute team (ALI-T). Both teams used their respective expertise in Machine Learning and Arabic Language to conduct research and develop an automatic solution that can detect errors of pronunciation and detect the location of the error, its type and give a feedback to correct the pronunciation. Output of the project will be a pilot system that can be used in offline or online ways, that can help learners of Arabic language, all over the world, correct their Arabic pronunciation.

Building an Arabic CAPT system requires a Non-Native Arabic Speech database that contains diverse speech and pronunciation errors. Hence this was the first objective of the project. In the proposal we aimed to record two sessions of speech of Arabic learners. The recorded database should stress on pronunciation errors, and have enough speakers with detailed phoneme annotations. ALI-T team have years of expertise in teaching the Arabic language for Non-Arabs, and conducted many research studies to enforce this expertise. Based on this expertise they were able to propose a methodology to construct the texts most suitable for the project, which took a considerable time and efforts. They proposed a text based on this methodology. The text went into many refinements from the whole project team until a set of 25 long sentences and some 61 very special pairs of words were finally selected. Due to the COVID restrictions the recording could not be done in a controlled face to face set up and instead the recording was done using an app on a mobile. Details of the DB and building it is presented in section 3.

We conducted detailed search in the literature for CAPT systems in general and CAPT for Arabic language in particular. This was second objective of the project. Details of this search is in section 4. From the search we found that some researchers are using long established techniques such as HMM, while other researchers are investigating using deep neural networks. We proposed a new technique based on deep neural networks for recognizing phonemes and AFs and built a

CAPT system using the proposed technique, hence we accomplished third objective of the project. The proposed method treated phonemes and AFs as objects in 3 channels spectral images of the speech. We published a paper on the proposed new technique. We investigated new improved models of this new technique and got excellent results that are comparable or better than state-of-the-art published research. We are finalizing a paper with these findings and accomplishments. Details of these CAPT systems are presented in 5.

2. Objectives

As we briefly presented in the introduction, the project has three technical objectives, namely:

- ❖ building a database of speech or Arabic learners
- ❖ conducting a literature review of CAPT systems in general and Arabic CAPT in particular
- ❖ Building a CAPT system using the developed speech database.

Building the database and conducting the literature review are two necessary steps in order to accomplish the main objective of the project which is building the CAPT system.

In addition to these technical objective the project has two objectives namely: establishment of a multidisciplinary research group through the collaboration of the CCIS-team and the ALI-team Institute, and dissemination of the results and conclusions at conferences and in journals.

In the following we will briefly present our accomplishments in each of these objectives.

2.1. Building the database

To build the database we had do design text best suitable for building a CAPT system. The text had to include the main or majority of pronunciation errors be Arabic learners and at the same time should be of reasonable length in order to be easily pronounced and easily recorded by non-Arabic speaking volunteers. In the initial proposal we were hoping to conduct the recording in face to face controlled sessions. Unfortunately, Covid restrictions did not allow this and we had to do the recording using a mobile app that we developed for the project. This has advantages and disadvantages as will be presented in section 3.1. Detail of the recording steps, protocols, cleaning, speech labeling are presented in sections 3 and 5.

2.2. Literature Review

We conducted a comprehensive review. This review is presented in section 4 and in the published paper. From the review we were able to propose a new technique that used object detection techniques to recognize the phonemes and AF and hence build an effective CAPT system.

2.3. Building the CAPT system

The DB that we developed is still not fully speech labeled and took long time to record, due to many reasons as will be explained in section 5, hence we used part of KSU database that is owned by the CCIS-team to build a CAPT system and got excellent results that we published in a paper. We investigated other improved models of the proposed technique to build a CAPT system and got excellent results that we will publish soon. The proposed technique and the improved models will be used to build a CAPT system using the recording of session 1. The performance of the system will be evaluated and compared to human judges. A new session of the database will be recorded. The first CAPT system will be improved based on the analysis of its performance and a new CAPT system will be built using the recording of sessions 1 and 2. The performance of the second CAPT system will be evaluated and compared to human judges. The results and analysis of the performance of the first and second CAPT system will be published in reputed journals.

2.4. Establishment of a multidisciplinary research group through the collaboration of the CCIS-team and the ALI-team Institute

The CCIS-team and ALI –team worked together and established a CAPT research group that have accomplished: Literature review of Arabic CAPT in the Arabic literature and CAPT and Arabic CAPT in the English literature, designed a text for research on CAPT, recorded speech of Arabic learners, time labeled the speech, and published a paper on novel method for building a CAPT system.

2.5. Dissemination of the results and conclusions at conferences and in journals

We published a paper on novel method for building a CAPT system. We are finalizing a new paper on improving the proposed method for building a CAPT system.

3. Development of the KSU-CAPT Non-Arabs database

Building an Arabic CAPT system requires an Arabic Speech database that contains diverse speech and pronunciation errors. At the time of writing this report, no database for Arabic L2 speakers with emphasis on pronunciation errors, and enough speakers with detailed phoneme annotations, is available.

We opted as per the proposal objectives to record a new speech dataset, having in mind the quality of the text selection and the coverage of most errors that L2 Arabic speakers might make. ALI-T teams, proposed the text for recording and had many meetings and discussions with CCIS-T, until a set of 25 long length sentences and some 61 very special pairs of words were finally selected. Details of text selection, database specifications, the recording system, registration of the speakers, speech recording are presented in sections 3.1 to 3.5 below. Details of speech labeling and analysis will be presented in section 5.

3.1. Selection of text for recording of the speech corpus

A main issue in CAPT systems is to select the optimal words and sentences that can cover the majority of errors in learning the pronunciation of L2. The selected texts must contain very specific phonemes that are difficult to pronounce by the speakers, and useful in improving the pronunciation. The text should have the following characteristics:

- ☞ Varied text containing rich diversity of phones and di-phones.
- ☞ Optimal number of sentences/words that can be pronounced by L2 Arabic learners in a minimal amount of time.

ALI-T team is well qualified for this task, as they have years of expertise in teaching the Arabic language for Non-Arabs, and conducted many research studies to enforce this expertise. Based on this expertise they were able to advice for a methodology to construct the text most suitable for the project, which took a considerable time and efforts. The proposed methodology is as below.

Methodology for the CAPT text selection

The selection of the CAPT sentences were subject to many constraints as follows:

- a) Sounds
 - Many repetitions of the same sound or phoneme are preferable.
 - Appearance of the sound at the start, middle and end of the word.
- b) Words
 - Common: Common words are preferred to specific words.
 - Diversity: words used in diverse Arabic countries are preferred.
 - Affinity: Usual and daily words are preferred.
 - Inclusion: Words used in many domains are preferred to words used in specific domains.
 - Importance: words needed by the learner are preferred.
 - Purity: Original Arabic words are preferred to Arabized Arabic words.
- c) Sentences of the text
 - Sentences must be meaningful.
 - Must have Arabic cultural aspect
 - Must be valuable.
 - Short sentences are preferred, to avoid boringness.
 - Must be consistent and clear.
 - Must be in accordance to and respect to the Kingdom's beliefs and constants.

The selection of the CAPT text passed by two main steps. The first step, contributed by ALI experts, consisted of applying the above methodology and suggesting sentences and specific words that contain phonemes that learners of Arabic have problems in pronouncing correctly, in addition to simple phonemes that can be found in other languages, such as `m`, `n`, ...etc. The second step, conducted by the project team, was refining the text to the mobile app and testing it. This second step passed by 4 stages, the first three stages were completed before starting the audio recording.

All the selected CAPT-texts were tested and adjusted over five main criteria:

- ☞ Reasonable time to read all the content.
- ☞ Complexity of the content, phoneme positions and length of words.
- ☞ Richness of the phonemic content in every sentence.
- ☞ Dual phones words must contain minimal pair diversity in phoneme pronunciations.
- ☞ It is well known that the more the sentences become long, the speaker starts damping its voice and were prone to more reading latency and stuttering, not in accordance to what the CAPT system aims to correct.

Both CCIS and ALI teams checked all the sentences and agreed to select an optimal number of 16 sentences with a minimal number of 21 words and a maximum number of words of 42 words. In addition, they selected a set of 61 minimal dual phones pairs, to target the phoneme dualities that can lead to pronunciation errors, these dual words differ by some phonemes but have a similar structure. Arabic experts from the ALI-T team stressed on the fact that these short and long dual phonetic words are very important in assessing and evaluating Non Arabs pronunciations.

Samples of two meaningful sentences and eight minimal dual phones pairs are presented in Table 1 below.

Table 1, Sample sentences and dual phonetic words

Sentence 1	
Sentence 2	
Pair of minimal dual phone words 1	/ ** / ** /
Pair of minimal dual phone words 2	/ ** / ** / ** / ** /

The initial recording setup was designed to record the CAPT students in the Arabic

Language Institute at King Saud University, in a controlled live session, face to face. If during the recording sessions, the speaker feels tired or bored, a short pause can be made, and the recording can continue after the pause session. Unfortunately, the COVID-19 restrictions imposed that students cannot come to the university, and we had to move to online recording, through a newly developed mobile app, that will be described later in the section. This online recording had some benefits and some drawbacks, as illustrated in Table 2.

Table 2, Benefits and drawbacks of the online recording solution

Benefits	Drawbacks
Any screen can be easily recorded again in a new time if any error is detected without the need to for the student to come back to the recording room (once the admin allows re-recording)	Speaker recordings had to be checked after each speaker completion. This induced EXTRA COSTS, for the listeners and Q/A checking stage.
Number of students at ALI was much lower than the number at time of submission of the proposal. Online recording allowed us to record students from inside/outside of Riyadh.	Recording solution must be compatible with diverse screens and phones, which increased the time for the testing and tuning of the App.
Diversity of the students, as students are not from one location or university.	Decrease of the total number of recorded phonemes. (see Figure 1)
Long sentences have been shortened, so the speakers could complete the recording in shorter time.	Additional explanations and discussions were necessary to make the students understand the installation and the use of the recording solution.

Due to the mobile recording constraints, the CAPT selected sentences have been additionally shortened, in order to fit into the screens of the students. Selecting the text in the app pages went into four versions. Statistics of the four versions of the texts are shown in Table 2. A detailed comparison, of the first three versions of the CAPT-Texts for each selection phase, is illustrated in Appendix A. Appendix B lists the text for version 4 (Mobile version) after further simplification and shortening so the read texts can be audio-recorded easily in the Mobile App.

Table 3, Statistics of the CAPT-Text selections

Statistics	V1	V2	V3	V4 (mobile)
Number of sentences	29	28	16	25
Total number of words (sentences)	1100	767	474	463
Maximum number of words / per sentence	72	41	16	27
Minimum number of words / per sentence	20	20	21	11
Pairs of Minimal Dual words.	113	113	113	61
Number of Screens : sentences	-	-	-	25
Number of screens: words	-	-	-	16

The numbers of phonemes, within the sentences and words of all the text selections versions, are presented in Figure 1. The phoneme distribution remained almost the same although we reduced the total number of phonemes to half from V1 to V3. V4 Mobile version is a reduced version of the V3, to fit in the screens of the Mobile, where sentences were shortened keeping the same meaning of the content, and 61 important dual pairs were kept for recording.

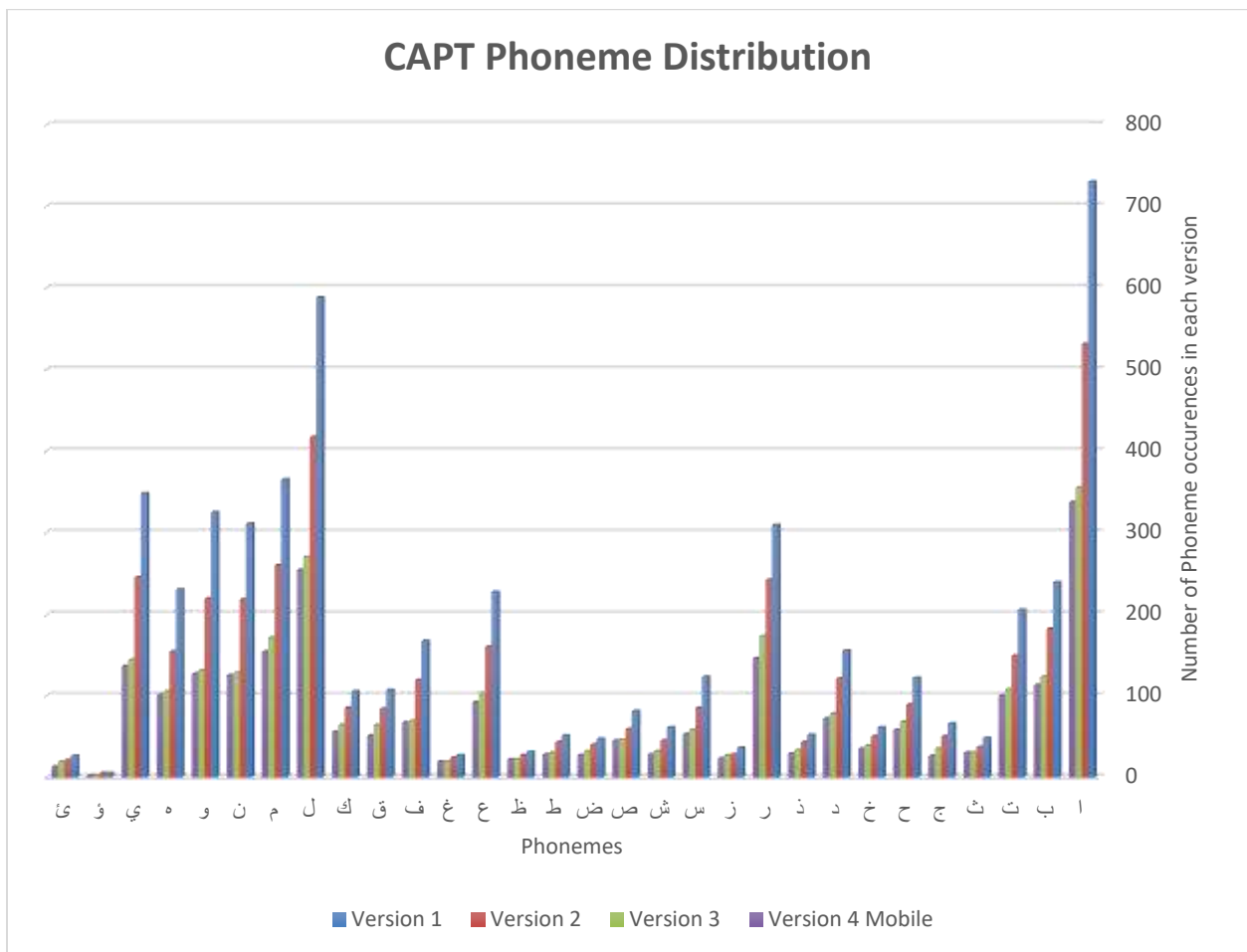


Figure 1, Phoneme distribution in the sentences and dual pairs of the various text selections

A comparative table of the phonemes of each version is detailed in Table 4.

Table 4, Comparative between the text selections of the various versions

Phoneme	Version 1	Version 2	Version 3	Version 4 Mobile
ا	730	531	355	337
ب	239	182	123	113
ت	205	149	108	100
ث	48	37	30	30
ج	66	50	35	26
ح	122	89	68	58
خ	61	50	39	35
د	155	121	78	72
ذ	52	43	33	29
ر	309	242	173	146
ز	36	28	27	23
س	123	85	58	53
ش	61	45	32	28
ص	81	59	46	45
ض	47	40	32	27
ط	51	43	30	28
ظ	31	27	22	22
ع	227	160	103	92
غ	27	24	19	19
ف	167	119	69	67
ق	107	84	64	51
ك	105	85	64	56
ل	588	417	270	254
م	365	260	171	154
ن	311	218	128	125
و	325	219	131	127
هـ	230	154	105	101
ي	348	245	144	136
و	5	5	2	2
ئ	26	21	19	13
Total Phonemes	5248	3832	2578	2369

3.2. Database Specifications

The recording step started by recruiting some sample speakers from the ALI institute, from the fourth level, in order to test the recording time and the quality of the reading. From the initial speakers' samples when recording version 3 of the text, we noticed that the duration of the recording varied between 40 and 45 minutes. The time was still long and we had to reduce the 16 long sentences (in version 3) to 25 short sentences (in version 4) with a maximum of 24 words and a minimum of 11 words per sentence. This was also a good consideration to fulfill the display constraint in reducing the displayed text on the phone screen, as mobile screens do not allow very crowded text and buttons in a convivial application. Sample screens from the Tahadath Mobile App are shown in Table 6.

Once the mobile app was developed and sent to diverse students at different geographical locations, we noticed that Non-Arabs had huge problems in reading texts without diacritics; we then updated the texts with a full diacritization. Screen shots of all the texts in the mobile app are presented in Appendix C.

The sampling rate of 8 KHz was decided upon two considerations:

1. Speech recorded at high sampling rates is in general reduced to 16 kHz or 8 kHz, for easiness of manipulation, and simplicity of use in training large models.
2. Recording from the microphone of the mobile phone allows only 8 kHz (sample metadata of the recordings are shown in Table 5).

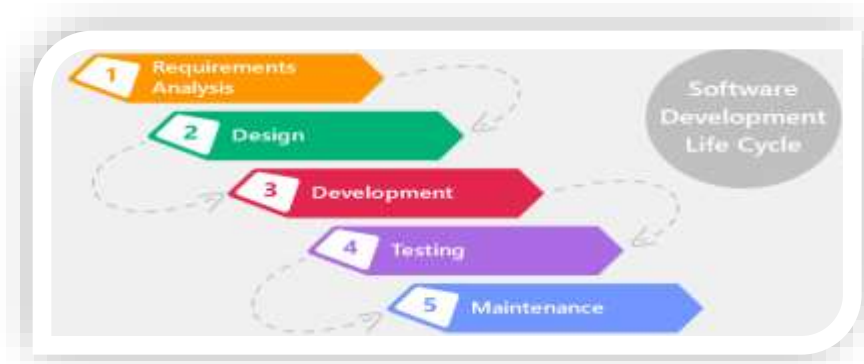
Table 5, Sample Metadata of the CAPT recordings

Input #0, ' mp3 ': <input type="checkbox"/> filetype
com.android.version: 10
Duration: 00:00:06.74, start: 0.000000, bitrate: 20 kb/s
Stream #0:0(eng): Audio: amr_nb (samr / 0x726D6173), 8000 Hz , mono, flt, 12 kb/s (default)

3.3. Establishment of the recording system

As already mentioned, we opted to develop a mobile app instead of a computerized application. We had two options, either use unity to develop the app or android using Java. The development of the recording app was subject to the known software lifecycle, presented in Figure 2.

Figure 2, Lifecycle of the Tahadath Mobile app



We have already developed, similar applications while recording speech datasets for a previous KACST funded project [1]. The requirements step was deeply discussed between the members of the team, and the first orientation was the use of a computerized application, but due to the medical restrictions against gathering of students in the institutes, due to COVID-19, we opted to use an app.

The first tentative app was developed by a specialized developer in unity, and has been tested for diverse criteria of screen sizes, colors, etc... We noticed, on the long run, that unity did not support Arabic writing in a very smooth manner, and the app developer had to load and deal with images in the application instead of writing Arabic texts directly in the app. A screenshot of the initial application is presented in Figure 3.

Figure 3, Unity APP screenshot



Unfortunately, with the numerous changes of the texts and fonts, we had to move to the development of another application in Java Android. The Java application felt more convivial to Arabic texts and font variations. The Java developer made diverse versions, as per the team requests (design-develop-test). The latest version is the 1.0.10 (10th version), in addition to a second similar application that was also developed for Arabs, as we wanted to split at the database level, Arabs from Non-Arabs recording, for a better management and checking. Some sample screens from the Android App Tahadath are presented in Table 6, a complete listing is also detailed in Appendix C.

Table 6, Sample screens from the Tahadath Mobile App



The backbone of the Android application is the Google firebase, and the application was subject to very strict access, as speakers are invited by their WhatsApp number and the access to

the application is subject to a fixed name and password, generated by our database manager. Access can be done only when a recording flag is enabled, once the speaker reads all the lists and approves its recordings, the recording flag is disabled, and no more access to the database is allowed to that speaker, unless it is reactivated by the admin for checking purposes or recording repetitions. A screenshot of the Google Firebase management platform is shown in Figure 4.

Figure 4, cs2r firebase real-time database



We can notice from Figure 4, that the control of the display font and size can be easily controlled from a centralized part of the app, in addition to the possibilities to change the display sizes at the mobile level.

Additional tests have been also made by the database team, in order to test the app in different mobiles and different versions of android. Some problems appeared in fonts and positions and were fixed as per the maintenance step of the software lifetime cycle.

3.3.1. Additional Improvements to the app

The app that has been developed is a one-way communication, i.e., the speaker records then the recordings are checked. This lead to many problems in terms of quality and recording durations, see Table 2, for benefits and drawbacks of the distant recording. When a speaker records his voice, we had to wait until he finishes his recordings to start checking because different students may be recording at the same time. Hence it was not possible to control every speaker in real time, because each speaker can record at his pace when he feels himself ready.

3.4. Speaker Registration

In the project proposal, we intended to record 200 Non-Arabs and 100 Arabs in the whole project. In order to manage such huge number of students, we followed a sample work methodology, where we start by a small number then increase to the target quota.

To ease the enrollment of the students who were mostly at distant locations, a google form, as shown in *Figure 5*, has been established and sent to the volunteers directly or to a coordinator from each institute who will send to students that he recruit at his institute. Each student needed to fill all the required fields and send it back. Once the forms are collected, the students are contacted by our team for further explanation of the recording steps or to answer any question.

Figure 5, Screen shot of the Google form sent to the students for the CAPT enrollment

نظام حاسوبي لتعليم نطق اصوات اللغة العربية للناطقين
بغيرها - مركز ابحاث الروبوتات الذكية بالاشتراك مع
معهد اللغويات العربية

* الاسم باللغة العربية Name (Arabic)
إجابتك

* الاسم باللغة الانجليزية Name (English)
إجابتك

* الجوال Mobile
ارجو ان يكون رقم الجوال سعودي حتى تتمكن من التواصل معكم
إجابتك

* WhatsApp Mobile رقم الواتس اب
إجابتك

* الجنسية Nationality
إجابتك

* المستوى الدراسي Academic level
 الاول
 الثاني
 الثالث
 الرابع

* العمر Age
إجابتك

* اللغة الام Native language
إجابتك

* الجامعة University
 جامعة الملك سعود
 الجامعة الاسلامية

* البريد الالكتروني Email
إجابتك

ملاحظات Comments
نرجو التسجيل مرة واحدة فقط للشخص الواحد وعدم تكرار التسجيل. وإذا واجهت مشكلة في التسجيل للتطبيق أو أي استفسار لتجدد النموذج يمكن
التواصل معنا على الرقم 0580874412 (واتس اب - المكالمة)

صفحة 1 من 1

إرسال

In the following part, we will present some statistics of the enrolled students. These

statistics include number of students, country of origin, language spoken, level of education, etc.

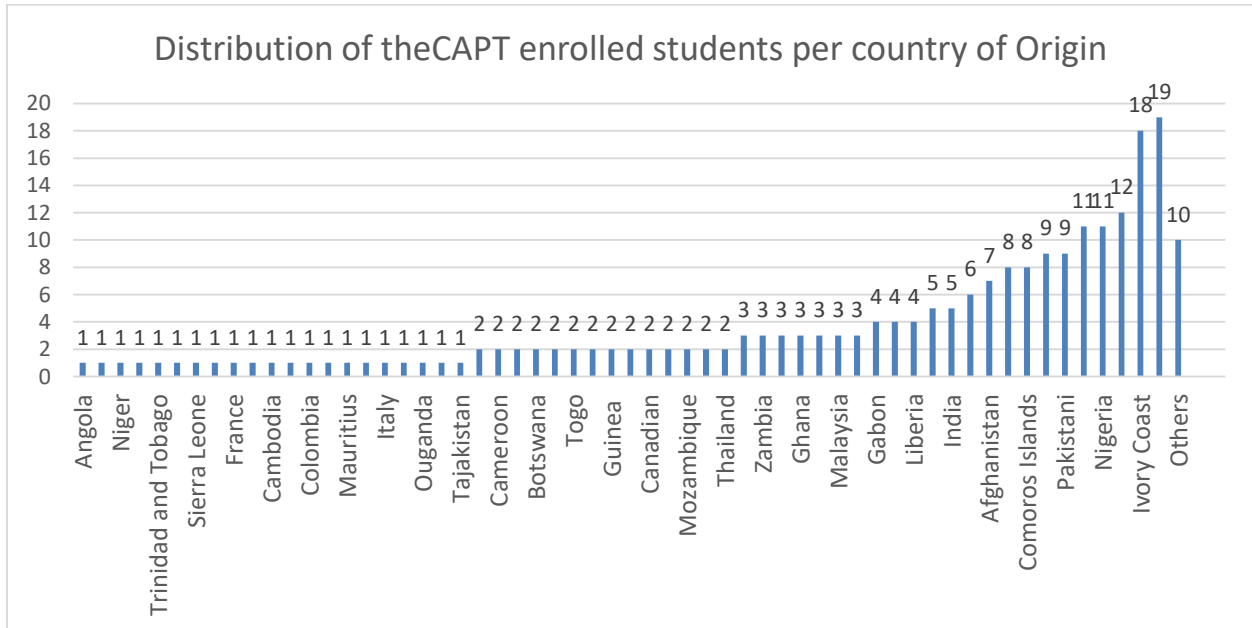
Most of the 371 collected google forms were from the Level 3 and Level 2. After many checking and controls for the validity of the pronunciation of the speakers, and testing their aptitude to pronounce Arabic even with errors, but with a minimal fluency, only 220 students had valid recordings, as shown in Table 7.

Table 7, Number of students that enrolled via the Google Form App

Academic level	Google Form Enrolled Students	Completed Valid Recordings
LEVEL 1	30	13
LEVEL 2	107	82
LEVEL 3	143	77
LEVEL 4	90	48
	370	220

In Figure 6, we present the distribution of the nationalities (more than 59 countries) of the speakers that participated to the recording of the session 1.

Figure 6, Statistics of the Nationalities of the CAPT speech recording at session 1



Statistics of the number of completed recordings per university are presented in Figure 7.

3.5 Database Recording (Session 1)

The enrolled students at session 1, were contacted via WhatsApp or by phone in order to officially make them understand the scope of the recording procedure. Each student was provided with a multimedia video tutorial, as a demo of the whole recording made by our database manager, in addition to a manual in Arabic and English, explaining all the steps of the use of the app. Each student received all the items listed in Table 8. A copy of the manual in both Arabic and English sent to every enrolled speaker is appended in Appendix F.

Table 8, Credentials received from /sent to the enrolled speakers (students)

	Item	Destination
Student ID :	Ahmed-0555555555	Received within the Google form
Username :	ahmed1	Sent to the student
Password :	123	Sent to the student
Android Application :	Apk format (through WhatsApp)	Sent to the student
Manual :	Tahadath-Manual.pdf	Sent to the student
Video :	App-Demo-Tutorial.avi	Sent to the student
Use of the recorded speech :	Consent screen in the app.	Within the Mobile App

We tried to be as clear as possible, in order to avoid any inconvenience in the use of the app.

3.5.1. Recording Constraints

- ☞ Due to the coronavirus, the decrease in the number of students at the ALI institute forced us to turn to the Islamic University of Al-Madinah, as they have more than 1500 students at their premises from more than 117 nationalities, and this helped us a lot in selecting the quality /quantity required by the project.
- ☞ The reason for selecting most of the students from outside of Riyadh, is that ALI student dropped from 300 students at the time of writing the proposal to 70 students at recording time and half of them were not physically present in Riyadh.
- ☞ The response from students at Islamic University of Al-Madinah was good at the

beginning then stalled, so we recruited students from other universities in KSA

3.5.2. Additional Remarks

- ☞ A consent text was written in Arabic in the app, the student needed to approve by clicking a check box, before starting the speech recording session.
- ☞ The recording of each speaker was accepted, when it has been double-checked, and passed the quality control criteria defined by the team.
- ☞ The student received an honorarium against his participation to the CAPT recordings.
- ☞ Many students from the level 1 could not read the texts completely, and were discarded from the recordings.

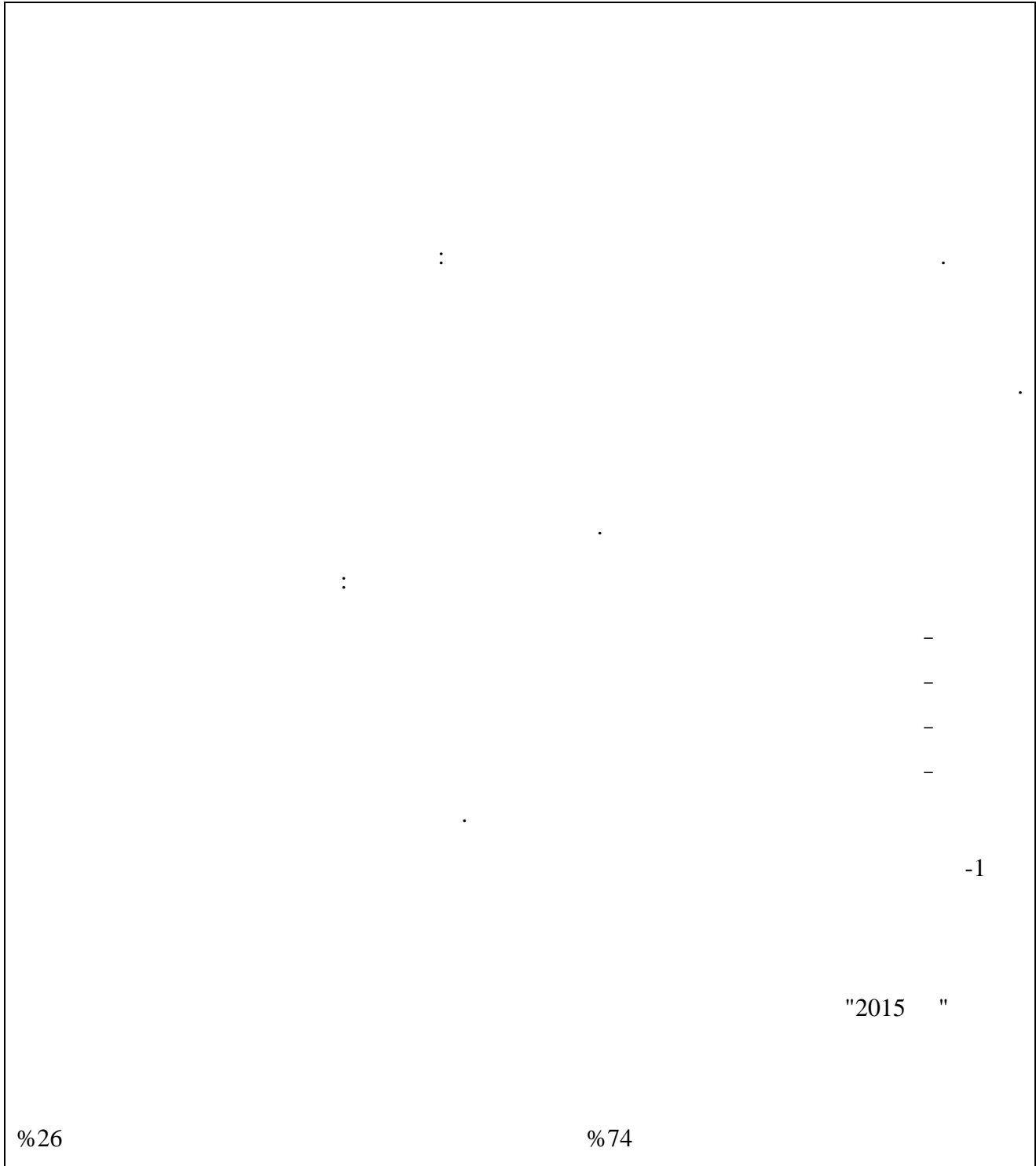
3.5.3 Recording Arab Speakers

Recording of Arab speakers started after recording of Non-Arabs. The project team tried hard to recruit Arab volunteers by personal invitation and by sending the request to participate in many WhatsApp groups. The response is very slow but the team is trying hard. The number of those who registered in the system database is 58 and among them 32 recorded their speech. The team is still trying hard to recruit. Vacation may be a major reason for the slow response.

4. Literature review of CAPT systems and of Arabic CAPT

4.1 Literature review for Arabic language CAPT in Arabic literature

(Done in Arabic by ALI from the Arabic literature)



				:	
) 2018(-
	. 2018				
:) 2012(-
			.		
:) 2016(-
			24		
:) 2005(-
.	-		3 , 32		-
) 2015(-
				:	
) 1988(-
				.29	
) 2011(-
	.117		-		
)2002(-
			17	-	
-					-
				3 23	
.					-
) 2010(-
			10	-	
2) 2014(-
:) 2008(-
				6	
					-3
	" 2010		" :		

: - - :

.

- - :

" 2016 " :

) 2012 (

) (%) (

:

:

:

/

·

) 2011 (:

·

·

:

-

· 2010 3 23

:

11 -

2020 2

:

2011

-

· 2012

-

4.2. Literature review of Arabic CAPT

Several Arabic pronunciation systems that were developed with a limited vocabulary and a small number of speakers are also available in the literature. No database is available for Arabic CAPT systems containing words or daily conversation, except those recorded using Quranic verses in [2] and [3]. In [4], a CAPT system was developed for the Arabic language. This system detects mispronunciation and assesses the pronunciation quality of learners. The ASR-based system contains forced alignment, scoring, normalization of scoring, and quantitative measurement phases. Six speakers of both genders recorded the speech database of three isolated Arabic words. One of the speakers with good pronunciation was used as a reference model, and the remaining five speakers were considered for the evaluation of the system. The reference models for the native speaker were generated by a 19-phoneme HMM. Each phoneme model contained three states and eight Gaussians per state. Thirteen MFCCs, including log energy, were extracted from each frame, and the first- and second-order derivatives of the MFCCs were calculated. Four different measurements, GLL, LLL, ROS, and ROA, were used to evaluate the quality of the phoneme-level pronunciation. A higher measurement score indicated that the quality of pronunciation was good. The GLL score provided 86.66% accuracy in the detection of mispronunciation, which is better than the accuracy of other measures. To determine the value of GLL for a larger corpus, four additional speakers with good pronunciation were added to the database. For the correct acceptance rate, LLL metric results were better than the results of other systems.

Dahan et al. designed, implemented, and evaluated an Arabic speech pronunciation scoring system in [5] for the training of Malaysian teachers of Arabic. The ASR-based system provided feedback on the pronunciation of an L2 learner of Arabic and detected mispronunciation errors. The system extracted features from a signal and fed them to the pattern recognizer. The extracted feature was the MFCC, and HMM was used as the pattern recognizer.

A computer aided language learning (CALL) system capable of providing feedback to L2 students to improve their pronunciation and evaluate learning levels was developed in [6]. The system was based on a new robust speech recognition method that was proposed in the study and implemented using Sphinx3. The method used a three-state HMM with eight Gaussians per state and 13 MFCCs with their first and second derivatives. The method provided output in a form of phonetic structure that distinguishes the proposed CALL system from other works, as presented in

[7], which suppose that learner's speech is already labeled.

In [8], a system to detect mispronunciation in Quranic recitation was implemented by Abdo et al. The system was divided into five parts: (1) segmented features were extracted by a primary feature extractor from the input speech; (2) the boundaries of the targeted segments were determined by a segment analyzer; (3) references for the correct pronunciation and errors were fed to the system through an acoustical database; (4) after the detection of the segment, discriminative features were extracted by a secondary features extractor; and (5) the distance between the test segment and the database was evaluated by a verification module. The accuracy for segmentation detection was 73%, and verification of the samples yielded a 100% accuracy. MFCC performed well as a discriminative feature and provided 95% accuracy, among other features, such as a zero crossing rate, formants, energy, local maxima of spectra, log area ratio, and linear prediction coefficients [9], [10], and [11].

Deep learning based method for pronunciation error detection for non-native Arabic speakers was proposed in [12]. The authors used a non-native Arabic database, which consisted of recording of 400 Pakistani speakers. Pronunciation error detection for non-native Arabic speakers at word level was been proposed in [13]. The authors used recording of Pakistani speakers who learned Arabic language. Deep CNN features and transfer learning parading were investigated.

4.3. Literature review of speech features for CAPT.

4.3.1. Speech features for CAPT.

Before deep learning era, hand-crafted features were playing an important role in designing the recognition systems. For example, in speech recognition systems, many features had been investigated in literature, such as MFCC, which is one of the most famous features, Linear Predictive Coding (LPC), LPCC, PLP, RASTA-PLP, etc. More information of these features and others can be found in [14][15]. With the huge improvement of deep learning technique and computation power, researchers proposed to feed the raw audio to deep learning networks in order to recognize the words/phonemes, such as in [16], [17]. On the other hand, a lot of researchers proposed to convert the audio signal to image representation such as 2 channels spectrogram / 3 channels spectrogram , and dealt with it using image based deep learning techniques, such as in [18], [19]. As will be shown latter in section 5, we investigated using deep learning techniques with 3 channels spectral images.

4.3.2. Feature reduction for CAPT (e.g., LDA, PCA)

As a middle step between feature extraction and classification steps, feature reduction techniques were used to select the discriminative features in order to optimize the recognition accuracy. In [20], authors used linear discriminative analysis (LDA) to reduce the MFCC dimension for speech recognition system. Authors in [21] proposed using principle component analysis (PCA) as a dimension reduction technique, with MFCC, as a feature extraction technique, to improve the recognition of Indonesian speech system. Details of feature reduction techniques can be found in this survey [22]. As will be shown latter in section 5, we investigated using deep learning techniques and hence did not need to use feature reduction techniques because the network structure takes care of this.

4.4. Classification methods for CAPT

A study of the speaker-independent speech recognition of non-native speakers was conducted by Alotaibi in [23]. An Arabic database from the LDC (language data consortium) was used to observe the effect of a large vocabulary for MSA (modern speech recognition). The database contained speech from 75 native speakers and 35 non-native speakers. The purpose of the study was to determine the phoneme-level differences between native and non-native speakers as well as which type of phonemes contributed to recognition among native and non-native speakers. An HMM-based system provided good results when used by non-native speakers in the training and evaluation phases, and female non-native speakers produced better results than male non-native speakers.

An automatic dialect identification system was developed in Trigui et al. [24] using GMM [25]. Nine different dialect regions (Algeria, Iraq, Morocco, Syria, Gulf countries, Tunisia, Yemen, and Lebanon) were considered in this study. Dialect varies from region to region, with gradual transitions rather than clear boundaries between them.

Another study was conducted by Trigui et al. [26] to classify the Arabic accent of non-native speakers based on a statistical HMM method. The database was recorded by male and female Arabic learners from different countries who spoke different native languages. English, German, and French accents were considered in the study. The system was divided into four components: data collection, language model construction, acoustic-phonetic decoding, and confusion. The recognition rates were 56% for French accents, 57% for English accents, and 69%

for German accents. Speech recognition is an important part of CALL systems. Speaker-independent speech recognition systems can be affected by several factors, such as accent and gender. Many studies have attempted to normalize the regional accents of speakers [27].

An ASR-based CAPT system for L2 learning was implemented in [28] by using prosodic information and the hidden Markov model. For this system, the target language was Indonesian, and the native languages of the learners were Japanese, Peruvian, and Vietnamese. Six graduate students of both genders recorded eighteen target words with all Indonesian phonemes. The pronunciation errors made by the learners depended on their native and target languages. Seven types of pronunciation errors were identified in the study for the Indonesian language learners and were considered during the system development.

End-to-End (E2E) systems based deep learning techniques have obtained promising results in an automatic speech recognition systems, such as: Deep speech [29], Deep speech 2 [30], wav2letter [31], EESSEN [32] and end-to-end attention model [33]. Likewise, E2E pronunciation error detection for CAPT systems have been proposed in the last few years and outperform the convention methods. Authors in [34] proposed E2E mispronunciation detection system based on CTC-Attention model. Another E2E system based on CTC-Attention model, for mandarin annotated spoken corpus, were proposed in [35]. Authors in [36] proposed E2E mispronunciation detection and diagnosis system for non-native English speakers. They used TIMIT corpus for native speech and L2-ARCTIC corpus for non-native speech. In our work we also using deep neural to make E2E system but in a new novel technique where we treat the phonemes or the AFs as objects in 3 channels spectral images.

4.5. Scoring in CAPT

ASR-based pronunciation training systems using different speech features and scoring measurement techniques have been developed for pronunciation training in other languages, such as Dutch, Mandarin, and Indonesian. A system to detect errors among Dutch language learners was developed by Doremalen et al. in [37]. Eleven different short and long vowels that are commonly mispronounced by Dutch learners were highlighted in this study. The experiments were conducted using the Spoken Dutch Corpus (CGN). This database contains nine million words and many speakers belonging to different age, sex, and regional groups [38].

An ASR-based system was implemented in [39] by extracting various phonetic features:

the mean pitch and intensity of the segment, and three formants with $F2 - F1$ calculated at three different locations of the sample. To reduce the variability of the measured parameters between speakers, normalization was performed using Lobanov's Z-score. After normalization, the normalized features are called spectral. In addition to the eight duration features, a zero-coefficient plus twelve MFCCs with their first- and second-order temporal derivatives were measured. The SVM with linear kernel was used for the classification, and the results were provided using different performance parameters, such as EER (equal error rate) and a 95% confidence interval. To obtain the baseline results, segment- and state-based confidence measures (CM) were calculated, and HMM models were trained by SPRAAK for phone alignment. The best EER of 15% was obtained with MFCC, and a result of 12.3% was achieved when MFCC was used with CM and duration features. The MFCC approach yielded better performance than the phonetic features, CM, and the duration features.

In the study [40] by Troung, acoustic-phonetic-based classifiers were developed by implementing linear discriminant analysis [41] and decision trees [42] to discriminate between the correct sounds of native speakers and the incorrect sounds of non-native speakers. The classifiers were designed to detect mispronunciation errors of three phonemes, /a/, /Y/, and /X/, which are frequently made by L2 learners of Dutch. A CAPT system using a confidence measure score, as presented in [43], to predict pronunciation errors did not have a high correlation between the human score and the machine-calculated score, perhaps because of the use of the same type of features for all sounds without considering acoustic characteristics. Therefore, each classifier was developed for specific errors by analyzing the acoustic differences of each sound in [40]. Moreover, gender-dependent models were used in the classifiers to optimize the performance of the developed pronunciation error detection system. Two experiments were conducted with 20 native Dutch speakers and 60 non-native speakers for both classifiers. In the first type of experiment, the classifiers were tested and trained in two ways: training and testing by native speakers and training and testing by non-native speakers. In the second type of experiment, the classifiers were trained by native speakers and tested by non-native speakers. The accuracy of the experiments for the decision tree and LDA ranged from 75% to 91.7% and from 87% to 95%, respectively.

An automatic pronunciation error detection system for L2 learners of Mandarin was proposed by Xu et al. in [44]. Prior linguistic information was used to improve the performance of

the developed CALL system. Methods based on LPP (log posterior probability) and RLPP (revised log-posterior probability) were implemented. The second method used linguistic knowledge, whereas the first did not. The results show that the performance of the method using linguistic knowledge was better than the other method. In another approach, RLPP was used to construct the restricted pronunciation space (RPS) for each phone to observe its pronunciation variation. A database containing 1,585 words pronounced by non-native Mandarin speakers was used to evaluate the system performance, and training was performed by native Mandarin speakers. The accuracy for the detection of pronunciation errors showed that the RPS-based system produced the best performance. Another CAPT system for the Mandarin language was proposed by Liang et al. [45]. The system was divided into two tasks: sentence verification and syllable identification. For sentence verification, acoustic models for the tri-phone, garbage model, pronunciation manner cluster (PMC), and anti-PMC were developed by using HMM. Forty-eight MFCCs with four energy coefficients were extracted from the Formosa Speech Database (ForSDAT) [46] to train the models. Five sentences containing all Mandarin phonemes were recorded for the testing of the system. Syllables were divided into two categories (out-of-task and confusion), and 160 sentences were developed by combining them. The sentence verification task obtained the highest F-measure [47] with the tri-phone vs. anti-PCM acoustic model, at 91%. The output of this model was fed to the syllable identification task. To extend the pronunciation lexicon, pronunciation variation rules were used. The best F-measure for syllable identification was 77.2%.

4.6. Databases for CAPT

Several databases were used in the literature for CAPT system. Some of them are publicly available and some of them are private. Table 9 shows a summary of some CAPT corpora for Arabic and non-Arabic languages.

Table 9: CAPT Databases' Survey.

Database Name	Year	Language	No. of recorded speakers	Total duration (hours)	Recording environment
L2-ARCTIC [48]	2017	English	24	11.2	quiet room
CU-CHLOE Corpus [49]	2015	English from native speakers of Chinese that learn English	211		sound-dampened room
iCALL Corpus [50]	2016	Mandarin	305	142	quiet office rooms,
CSLU [51]	2005 and 2007 update	English	90	30	digital telephone lines
West Point Arabic Speech [52]	2002	Arabic	110	11.42	
The CrossTowns Corpus [53]	2006	German, English, French, Italian, Netherlands	161	16	noise-controlled cabin and small room
Speech accent archive [54]	2016	English	646	N/A	Online
IDEA. The International Dialects of English Archive [55]	1998 -2020	English	N/A	170	Online

5. Design and development of a CAPT system for L2 learners of the Arabic language.

In this section we will show the achievements that we did for building the CAPT system which is the project main objective. The section will also present the materials and methods and the results.

When we submitted the proposal, we proposed to use the conventional way in speech recognition systems, which consists of the following steps: feature extraction, feature reduction, and classification. With the great advancement of the deep learning and the huge improvement in computation power, many of state-of-the-art system in the current time are using deep learning to construct end-to-end high performance speech recognition systems. Hence in our work in the first year we investigated using deep learning networks for phoneme recognition system and for mispronunciation detection and diagnosis system (MDD). We proposed a new way of using deep learning for detection and recognition of phoneme and articulatory features (AF). In the new proposed method, we treat the phonemes and ATFs as objects in 3 channels spectral images of the speech. By this proposed method we were able to recognizing the sequence of phoneme from the whole utterance of the non-native Arabic speakers. Then we used the detected phonemes for mispronunciation detection and diagnosis task.

Providing feedback to non-native learners is very important especially at articulation level. Hence, by the proposed method we detect and recognize the AFs as objects in the 3 channels spectral images, then we use the detected AFs for mispronunciation correction and providing feedback at articulatory level. Figure 9 shows the general architectures of the proposed Arabic CAPT system. The database recording phase was described in details in section 3 of this report. In the following sections, we explain the details of the proposed Arabic CAPT system.

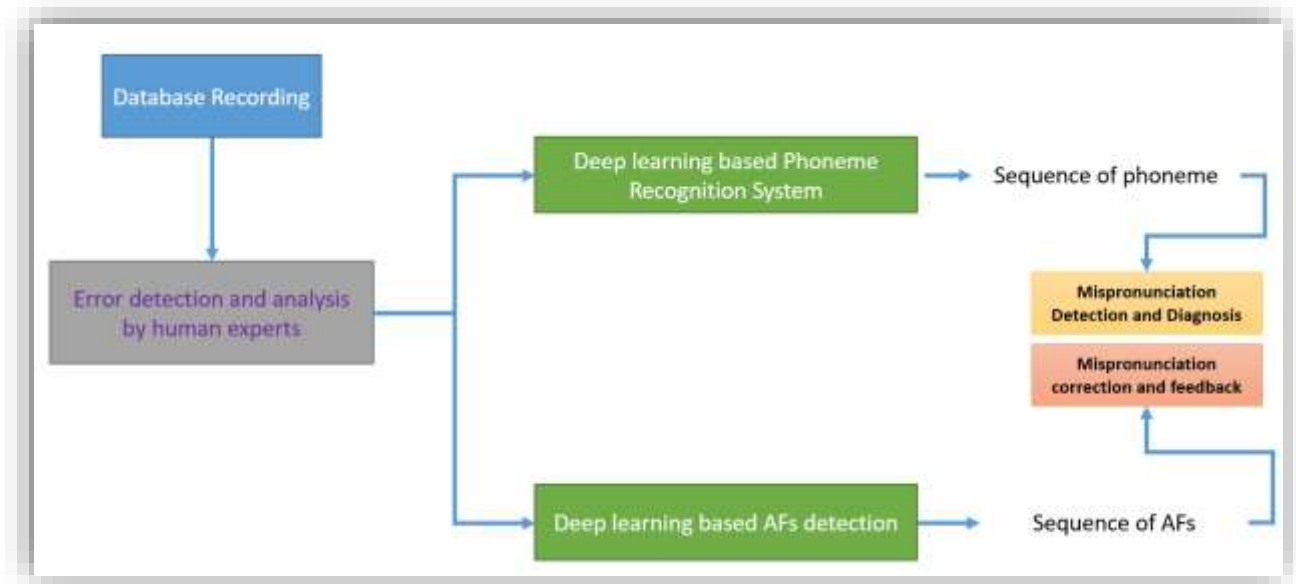


Figure 9: General Architecture of the proposed Arabic CAPT system.

5.1. Error detection and analysis by human experts

Before sending the recorded wave files to the human annotators for speech labelling, subjective evaluation and error analysis, many steps of cleaning and checking have to be realized. These steps range from checking the data consistency to extra sounds removal. This step is a mandatory pre-processing step, and will be described in details throughout the next paragraphs.

5.1.1. Pre-processing of the recorded wave files

Two specific tasks have been assigned to two different per-processing teams that were required to check the content and mention any errors in separate files, and never alter the original wave files by any manner. Details of the tasks of team A & B are described in Table 10.

Table 10, Details of the per-processing teams.

Team ID	Task	Details of Tasks	Number of Checkers per team	Software
A	Content checking:	Check the content of each wave file at the content conformity or consistency level.	2	Audio reader

B	Detailed content analysis: Mark Additional words or sounds	Check for any additional words or sounds that are not part of the CAPT content, using the ELAN software. (Marking temporally the outlier segments)	3	ELAN software
---	------------------------------------------------------------	----------------------------------------------------------------------------------------------------------------------------------------------------	---	---------------

The first **Team A**, had a huge load of daily work, as they were checking each day all the recorded speech of all the app enrolled students and were giving quantitative reports of the correctness and quality of each recorded content. The task was very tedious and took a lot of efforts and time. Detailed samples of the daily reports of team A is presented in Table 11. Where “Pxx” means paragraph number xx (one to three sentences grouped all together), and “SPWyy” means list yy of Sequence of Pairs of Words in the text to be recorded, (each list contains 2 to three pairs of words).

Table 11, Sample reports of the content checking

NAME	DATE	RECORDING	VERIFICATION	COMMENTS	
		STATUS		(Coding Used : Pxx : Paragraph xx, SPWyy : List of Sequence of Pairs of Words yy)	
M	17/11/2020	COMPLETE	Done	SPW6, SPW8, SPW13, and SPW15 are not complete (2 words missed at the end of each record)	
	17/11/2020	COMPLETE	Done		
	17/11/2020	COMPLETE	Done	SPW6, SPW8, SPW13, and SPW15 are not complete (2 words missed at the end of each record)	
	17/11/2020	COMPLETE	Done	Lot of stutter & repetitions in P1.P2.P3.P4.P5.P10.P17.P18.P23.P25 SPW6, SPW8, SPW13, and SPW15 are not complete (2 words missed at the end of each record)	
	17/11/2020	INCOMPLETE			
	17/11/2020	COMPLETE	Done	Small noise in SPW14	
	17/11/2020	COMPLETE	Done	Small noise in P7	
	17/11/2020	INCOMPLETE			
	17/11/2020	COMPLETE	Done	SPW6, SPW8, SPW9, SPW10, SPW11, SPW12, SPW14, SPW16 missed the last 2 words in SPW6, SPW8, SPW13, SPW15	
	17/11/2020	INCOMPLETE			
	17/11/2020	COMPLETE	Done		
	17/11/2020	INCOMPLETE			
	17/11/2020	COMPLETE	Done	repetitions & noise in P1.P9/ SPW6, SPW8, SPW13, and SPW15 are not complete (2 words missed at the end of each record)	
	17/11/2020	INCOMPLETE			
	17/11/2020	INCOMPLETE			
	N	17/11/2020	COMPLETE	Done	repetition in P5 repetition & noise in P13.P14.P19 noise in P20 P24 missed the first word SPW6, SPW8, SPW13, and SPW15 are not complete (2 words missed at the end of each record)
		18/11/2020	INCOMPLETE		
18/11/2020		INCOMPLETE			
18/11/2020		COMPLETE	Done	P1, P4 wrong record stutter & repetitions P7.P8.P19.P21.P23.P24	
18/11/2020		INCOMPLETE			
18/11/2020		INCOMPLETE			
18/11/2020		INCOMPLETE			
18/11/2020		INCOMPLETE			
18/11/2020		INCOMPLETE			
18/11/2020		COMPLETE	Done	bad stutter & repetitions P2.P3.P4.P5.P6.P12.P13.P14.P15.P16.P18.P19.P22.P23.P24.P25.SPW1.SPW14.SPW15 stutter & repetitions & noise P7.P8.P9.P10.P11.P17.P20.P21 SPW6, SPW8, SPW13, and SPW15 are not complete (2 words missed at the end of each record)	
19/11/2020		COMPLETE	Done	stutter & repetitions P3.P20.P22.P23.P14.P16.P25.SPW7 SPW6, SPW8, SPW13, and SPW15 are not complete (2 words missed at the end of each record)	
19/11/2020		INCOMPLETE			
21/11/2020		COMPLETE	Done	SPW6, SPW8, SPW13, and SPW15 are not complete (2 words missed at the end of each record)	
21/11/2020		COMPLETE	Done	stutter, repetitions & noise P1.P2.P3.P4.P5.P6.P7.P10.P11.P13.P14.P16.P17.P20.P23.P24.SPW9.SPW15 SPW6, SPW8, SPW13, and SPW15 are not complete 2 words missed at the end of each	
21/11/2020		COMPLETE	Done	stutter repetition P11.P13.P14.P19.P24.SPW1.SPW13	
21/11/2020		COMPLETE	Done	stutter in P4.P10 SPWwrong sentences: P25 SPW6, SPW8, SPW13, and SPW15 are not complete (2 words missed at the end of each record)	
17/11/2020		INCOMPLETE			
ames	21/11/2020	COMPLETE	Done	Incomplete records P2.P3.P4.P5.P7.P11.P12.P13.P14.P15.SPW5.SPW16 wrong sentences: P8 noise in: P21.P24.P25 empty record: SPW8 SPW6, SPW7, SPW9, SPW13, and SPW15 are not complete (2 words missed at the end of each record)	
	21/11/2020	COMPLETE	Done	stutter & repetition in P3.P10.P11.P21.P22.P25.SPW3 SPW6, SPW8, SPW13, and SPW15 are not complete 2 words missed at the end of each	
	22/11/2020	COMPLETE	Done	empty record: P1 stutter & repetition in P2.P4.P5.P6.P7.P8.P9.P10.P11.P14.P15.P16.P17.P18.P19.P20.P21.P22.P23.SPW9.SPW14 wrong sentence: P24.P25 SPW6, SPW8, SPW13, and SPW15 are not complete 2 words missed at the end of each	
	21/11/2020	COMPLETE	Done	P1 wrong record stutter & repetition in P3.P4.P10.P13.P14.P16.P19.P22.SPW7.SPW9.SPW15 noise in P5.P11.P23 SPW6, SPW8, SPW13, and SPW15 are not complete 2 words missed at the end of each record	
	23/11/2020	INCOMPLETE		3 files missed while according to the app he's finished	
	23/11/2020	COMPLETE	Done	stutter repetition P7.P12.P16.P18.P20.P22 SPW6, SPW8, SPW13, and SPW15 are not complete 2 words missed at the end of each record	
	22/11/2020	COMPLETE	Done	stutter repetition P3 SPW6, SPW8, SPW13, and SPW15 are not complete 2 words missed at the end of each record	
	23/11/2020	COMPLETE	Done	stutter repetition P2.P3.P9.P14.P15.P17.P20.P25 noise P4.P6.P19 SPW6, SPW8, SPW13, and SPW15 are not complete 2 words missed at the end of each record	

At this stage, the recordings have been checked for correctness of global content, i.e.: the content of each speaker directory should contain 41 wave files, as described in Appendix B. Once Team A, has completed its task, a second Team B will continue the verification process by marking each speaker content for extra sounds that have been introduced during the App recording. Examples of these sounds includes stuttering, door knocking, baby in background, etc..., as many students were recording from home or from the university housing. The main objective of this step in the pre-processing was to be sure that what was written in the App cards, regardless if it was correctly pronounced or not, is the only speech or sounds within the wave files.

We opted to use the ELAN software due to its capabilities in annotation and export possibilities to text grids, further used by our other software applications. In order for the Team B to work fluently and at fast pace, the technical team concatenated the wave files, to ease the step of listening to the speech and marking the segments of text or special sounds that are not within the original text.

This task was executed by the three checkers in team B and required a lot concentration and repeated listening and required long time. For example, a checker in team B spent one hour to process the speech of one speaker that was recorded in 10 min speech.

Table 12 to Table 13 , present some statistics about the minimum, average, and maximum speech duration for the 220 Non-Arabs recorded speakers. A very detailed speaker duration is listed in Appendix D.

Table 12, Statistics about the session 1 speech recording

Item	Details
Minimum time	6.37 min
Average time	12.06 min
Maximum time	37.4 min
Speakers <10 min	66 speakers : (Avg time: 8.98min)
Speaker Rec. >=10 min and <15min	128 speakers (Avg time : 12.13min)
Speaker Rec. >=15 min and <20min	19 speakers (Avg time : 16.96min)
Speaker Rec. >=20 min	7 speakers (Avg time : 26.43min)
Short Paragraphs duration	31.29hours
Sequence of pair of words duration	12.94hours
Total recorded time for all the speakers	44.23Hours
Number of volunteers who registered in the system	371
Number of Speakers who installed the app and started recording	235
Number of Speakers who completed Recording	224
Number of Speakers with valid recordings	220
Number of Speakers with completed transcription checking	132 (60%)
Number of Speakers with completed phoneme level labeling	40 (18%)

Table 13, Additional deep statistics about the session 1 speech recording of Non Arabs: SPW

SPW ID	Avg time (sec)	Min time (sec)	Max time (sec)	Total time for the SPW(Sec)
SPW 1	13.95	3.76	27.04	3068
SPW 2	8	3.02	22.02	1761
SPW 3	11.7	5.6	29.2	2574
SPW 4	10.29	5.22	20.86	2265
SPW 5	7.39	3.36	15.94	1626
SPW 6	16.02	6.92	30.56	3524
SPW 7	14.61	6.82	31.08	3214
SPW 8	15.99	6.46	30.16	3518
SPW 9	14.3	5.7	41.6	3146
SPW 10	14.06	1.28	23.68	3092
SPW 11	10.9	5.12	27	2399
SPW 12	10.01	3.7	23.8	2202
SPW 13	15.07	5.92	67.48	3315
SPW 14	15.38	6.5	57.3	3383
SPW 15	20.05	6.16	48.34	4411
SPW 16	14.08	2.7	38.08	3098

Table 14, Additional deep statistics about the session 1 speech recording of Non Arabs: Paragraphs

Short	Avg. time (sec)	Min time (sec)	Max time (sec)	Total time per type of
P1	21.91	11.08	61.28	4820
P2	25.74	12.66	423.66	5664
P3	27.02	9.86	121.4	5944
P4	22.27	10.38	84.46	4899
P5	18.55	9.24	67	4080
P6	15.82	7.96	111.76	3480
P7	18.68	10.32	74.64	4109
P8	16.52	8.78	66.14	3635
P9	16.64	7.4	71.32	3661
P10	15.73	7.82	78.9	3462
P11	22.89	11.4	78.5	5035
P12	21.94	10.5	71.52	4827
P13	27.36	8.22	78.26	6020
P14	25.61	3.36	72.46	5633
P15	15.88	7.62	89.72	3495
P16	14.76	7.64	50.38	3248
P17	19.8	7.36	89.62	4355
P18	14.23	7.16	61.56	3130
P19	26.8	10.6	206.24	5896
P20	17.8	7.34	76.34	3915
P21	24.8	12.8	77.26	5455
P22	22.43	9.48	89.04	4934
P23	26.79	7.34	87.94	5894
P24	19.26	9.1	100.08	4237
P25	12.81	5.06	82.68	2818

In order to check the conformity of the written transcription (i.e. displayed text in the mobile screens) to the content of the wave file, we opted to use the ELAN annotation tool, which is known for its labeling capability for both audio signals and video sequences, in a single or multi-tiers labeling. Hence we used ELAN software as an aiding tool in speech segmentation. Using ELAN Team B listened carefully to the wave files, and marked any additional noise, repeated word or background sound, by putting time boundaries around them. This process kept the original files as is, and generated a PRAAT TextGrid file, that was used later, by our technical team, for segmenting the wave files, where only the valuable parts were kept.

In order to ensure that Team B, did a job without apparent errors, an additional effort of rechecking 5% of the speakers is performed by the project technical team to recheck the conformity of the content. If no error is detected, no further action is made, else more samples will be rechecked.

All intermediate work and results have a backup copy, in case of need for traceability of errors at the labeling level.

Team B, followed a very strict protocol, from the wave file opening in ELAN to the text-grid generated at the end of the checking session. A copy of the protocol is presented in Appendix E. In Figure 10, we present a screenshot of the ELAN software work for sentence P10.

Figure 10, Screenshot of the ELAN screen work

The screenshot displays the ELAN software interface. At the top, there is a menu bar with options: File, Edit, Annotation, Tier, Type, Search, View, Options, Window, Help. Below the menu bar is a toolbar with various icons for file operations and editing. The main window is divided into several sections:

- Sentences List:** A table with columns for 'Nr', 'Annotation', 'Begin Time', 'End Time', and 'Duration'. Sentence P10 is highlighted in red.

Nr	Annotation	Begin Time	End Time	Duration
8 P8	أفمن الله بعد الله ثم أمتنا الأضال القاصلة، والمنتارا الصبيحة، الذين يزعمون إلى الخير رخصاً :	00:03:0	00:03:2	00:00:2
9 P9	وأمتنا الأضال القاصلة، يعيرون الضعيف والمريض، ويأهضون الضلال والإفراط أملا في مراضة الله :	00:03:2	00:03:4	00:00:1
10 P10	أثارت الثمن وأرسلت أمتها، فتعوت بالعبادة والمزور، وترعت أمير نخر الشيطان :	00:03:4	00:04:0	00:00:2
11 P11	تأخرت على سوية منبه، دعاني صاحبها، فركبت معه، ثم استألف أعمل لصيد الأسماك، وكان بعضها يلين الشباك ولا يدخل فيها :	00:04:0	00:04:3	00:00:2
12 P12	أمتنا صديقي تمنعني نصيحة عظيمة، فترك عليه المومن التي تدعو إلى العسر، ففسد على ما أصابه، فأصبح الله حوله، وأصبح قوة للمنازين :	00:04:3	00:04:5	00:00:2
13 P13	عزم زيد على زيارة أحد الأعمى الثنايين بعداً عن القران فقرأه ساعة ثم أتته، فقال له: يا بني: طهر قلبك بالآية، وإلا فإنا عند الله، وعن عزير النفس :	00:04:5	00:05:3	00:00:3
- Audio Waveform:** A visual representation of the audio signal for sentence P10, showing amplitude over time. A red vertical line indicates the current selection point at 00:00:00.000.
- Transcription Tiers:** Below the waveform, there are several tiers for text transcription. The 'default' tier contains the transcription for sentence P10: "أثارت الثمن وأرسلت أمتها، فتعوت بالعبادة والمزور، وترعت أمير نخر الشيطان :". Below it, there are tiers for 'Sentences' and 'To_Remove', with 'P10' and 'R-P10' visible.

From the different reports submitted by teams A and B, we collected some remarks that will be taken into account while recording the speakers in the next session.

Conclusions from the written feedback of the team B.

- ∞ The voices of some speakers were not very clear
 - ∞ Some speakers had difficulties the reading and read in very low speed.
 - ∞ Many stuttering and repetitions.
 - ∞ Difficulties to read some sentences or words.
-

Speech annotation experts will listen to each wave file and check the phonemes of each word:

- If a word is well pronounced, no remark is written.
- If the word is not well pronounced, by either addition, deletion or substitution of any phoneme, the annotator will mark the location of the error compared to the reference original text.

The final output of this step will generate excel files containing specific information of each wave file, such as the phonemes present in the text (reference text), phonemes correctly pronounced (correct phonemes), phonemes that were not pronounced (deleted), the phonemes badly pronounced (substituted), or added phonemes (added).

After initial work by the annotators we realized there are more in speech labeling and subjective evaluation than the information to put in Table 15. Hence with consultation of the Arabic language experts we came a new format as in the next section.

5.3. Subjective evaluation of the pronunciation of L2 learners

As we mentioned in the previous section in addition to the difference between the canonical text and the pronounced speech due to substitution, insertion, and deletion, which are errors, there are other ways change that might be correct or wrong changes. Hence we had to add notation for these changes. Moreover, to make annotation faster we proposed a new table that will cover speech labeling as well as subjective evaluation. For substitution and addition no symbol is used and the annotators write the canonical and pronounced text as in Table 16, the other two columns are for explanation for the report only.

Table 16, Examples of annotating substitution and addition

طريقة تمثيل الاستبدال والإضافة... (بدون رموز)

توضيح Explanation	النص المنطوق Pronounced text	النص النموذج Canonical text	نوع الخطأ Type of error
			الاستبدال Substitution
			الإضافة Addition
()			

For deletion and other operations, we used the symbols as in Table 17 below. The annotators write the canonical and pronounced text as in Table 17, the other three columns are for explanation for the report only.

Table 17, Examples of annotating deletion and other speech processes

طريقة تمثيل الحذف والظواهر الصوتية الأخرى التي قد يستفاد منها في رصد التباين الصوتي بين مختلف المتحدثين أثناء

ملاحظة الأخطاء... (باستخدام الرموز)

توضيح	Example		Process	Symbol
	النص المنطوق	النص النموذج		
—	#		التفخيم	#
.	#			
.	@			@
) (× *			*
) (&			&
) (>			>
(:) () .) (× ×			×
.	^			^
	~			~
	/			/
	/			/

5.4. Error Analysis

Until finishing the first year report, our Arabic experts annotated only the speech of 40 speakers from the 220 recorded speakers. Hence, we cannot do a complete analysis in order to present the most common errors for non-native Arabic speakers and the total number of mispronounced phonemes. Once, the annotators finish the evaluation process, we will make the agreement between the annotators, when we find any disagreement between the two annotators, we will send these utterances to another Arabic expert to evaluate it, who is a CO-PI in this project. We expect to finish the annotation process of session 1 as soon as possible, and we will present the deep analysis in the next report.

5.5. Implementation of speech recognition phase in the CAPT

The proposal was based on state of the art techniques at time of submitting the proposal, and we mentioned deep neural networks (DNN) as something to investigate. When we started working in the project DNN was widely used in speech recognition, hence we focused on DNN and proposed a new method of using object detection techniques, based on DNN, to detect the phonemes which resulted in publication [56]. We continued in this direction and published the new investigation details and results in the project first paper [57], and we are continuing this in a new paper that we are finalizing. This shift meant that we will not investigate many hand crafted speech features, as we did on some of our previous researches, as the DNN is powerful and has the ability to even work in raw data. This also resulted in combining 3.2.2 into 3.2.3 since the DNN will perform the reduction in its initial layers.

We started our investigation of using DNN in our paper [56], which was supported and funded by the Deanship of Scientific Research at King Saud University. This showed the great potential for our proposed method. We continued in this direction and in [57], which is the first paper out of this project, and investigated our proposed method to detect the phonemes, hence accomplish speech recognition, as well detect and locate the AFs. By detecting the AFs, we will enable our CAPT system to give feedback to the speaker about the error and may suggest a way to correct it. This is an important addition to the project that was not in initial proposal but we were able to do. We are continuing enhancing and modifying our proposed method to be able to recognize the phoneme and AF in one system. We reached excellent results that is comparable or better than state of the art methods and we are finalizing a new paper based on these results.

Before presenting the accomplished results and systems of this project and to put our work in the project in perspective we will very briefly summarize our work and results in [56].

5.5.1 Summary of our work in [56].

In [56], we proposed the use of object detection techniques for recognizing sequence of phonemes from the whole spectrogram. We converted the utterance of speech to a three channel spectral image then we used a deep object detection models for detecting the phonemes from the spectral image. The novelty of the proposed system is represented by treating of phonemes of utterance as objects in spectral images. We chose YOLO and CenterNet, two cutting-edge object detectors, based on a trade-off between detection accuracy and speed. Our study is the first study in literature, to the best of our knowledge that used object detection for phoneme recognition system. We evaluated the proposed system using native English and non-native Arabic speech corpora. For English phoneme recognition, we used the TIMIT dataset which is a well-known English speech corpus. For non-native Arabic phoneme recognition, we used a small part of the KSU speech corpus [1].

Due to the small size of the corpora, we investigated using different types of transfer learning techniques as follow:

- Transfer learning from image to speech databases (DTS)
- Transfer learning between speech corpora within the same language (IaTS)
- Transfer learning between speech corpora within the different language (IeTS)

Table 18 shows the result of the two proposed systems DTS and IaTS for the test set of the TIMIT corpus. We can see that the performance of the transfer learning improved the results. We achieved the best PER using the IaTS using CenterNet detector with DLA backbone network which was 15.89%.

Table 18: PER for the TIMIT test set.

<i>System</i>	<i>Object Detector</i>	<i>Model</i>	<i>PER</i>
<i>DTS</i>	YOLO	YOLOv3-tiny	28.25
<i>DTS</i>	YOLO	YOLOv3	20.2
<i>DTS</i>	CenterNet	ResNet	21.09
<i>DTS</i>	CenterNet	DLA	19.06

<i>IaTS</i>	YOLO	YOLOv3-tiny	25.57
<i>IaTS</i>	YOLO	YOLOv3	16.34
<i>IaTS</i>	CenterNet	ResNet	17.16
<i>IaTS</i>	CenterNet	DLA	15.89

Table 19 presents the performance of the proposed system *IeTS* for non-native Arabic phoneme recognition. In this experiments, we used speech of (15 non-native speakers and 5 native speakers) for training and 11 non-native speakers for testing. The total number of phonemes in training and testing utterances is 14413. Using the YOLO detector tiny version, we achieved 10.15% PER and using the CenterNet detector, we achieved a 7.58% PER.

Table 19: PER for non-native Arabic Speech (Small-Arabic-CAPT).

<i>System</i>	<i>Object Detector</i>	<i>Model</i>	<i>PER</i>
<i>IeTS</i>	YOLO	YOLOv3-tiny	10.15
<i>IeTS</i>	CenterNet	DLA	7.58

The results in [56] are comparable or better than state of the art published researches.

5.6. Investigating using the proposed technique for phoneme and AF detection in Arabic and English speech

In this part we will present the main results of this task out of the project within the first year. The details of the work and the results were published in [57]. In the following we will try to highlight the main work, findings and results in the paper. The paper title is “Deep learning-based detection of articulatory features in Arabic and English speech” but it also investigated detection of the phonemes in two ways. It detected the phonemes directly from the spectral images or based on the detected AFs. To detect the AFs, we proposed using object detection techniques to recognize sequence of AFs from speech utterances by treating AFs of phonemes as multi-label objects in spectral images, while for phoneme recognition we treated the phonemes as single label objects. We tested the proposed system on English corpus, TIMIT, and on Arabic speech corpus, KAPD [59]. Figure 11 sows the general overview of the proposed systems, where the system to detect the AFs is called AFD-Obj and the system to detect the phonemes is called PD-Obj. By detecting the AFs, we can provide feedback to the non-native Arabic learners at articulatory level. This is a new important and beneficial feature that we aim to include in this Arabic-CAPT system, though it was

not in the proposal. Moreover, we studied the effects of the number of detection levels of YOLOv3-tiny detector.

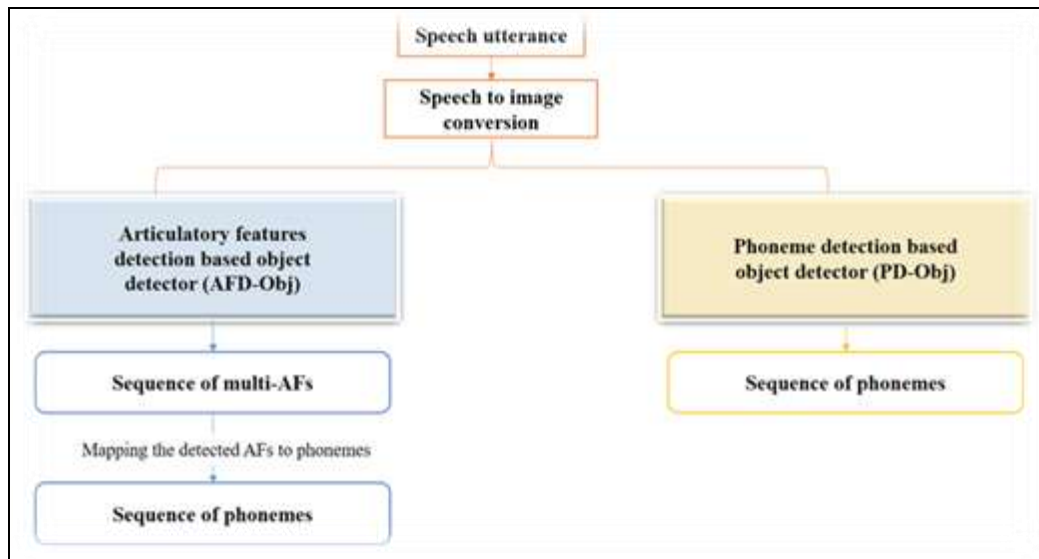


Figure 11, Proposed System of the AFD-Obj and the PD-Obj

5.6.1. Feature extraction

We used the speech-to-image transformation presented in detail in our previous study [56]. We concatenated the power Mel-spectrogram and the first and second derivatives to generate a three-channel image. Then, using the time boundaries, we calculate the bounding box of each object. Figure 12 shows an overall picture of converting the speech to image to be able to recognize the phonemes in Figure 12-a and the AFs in Figure 12-b.

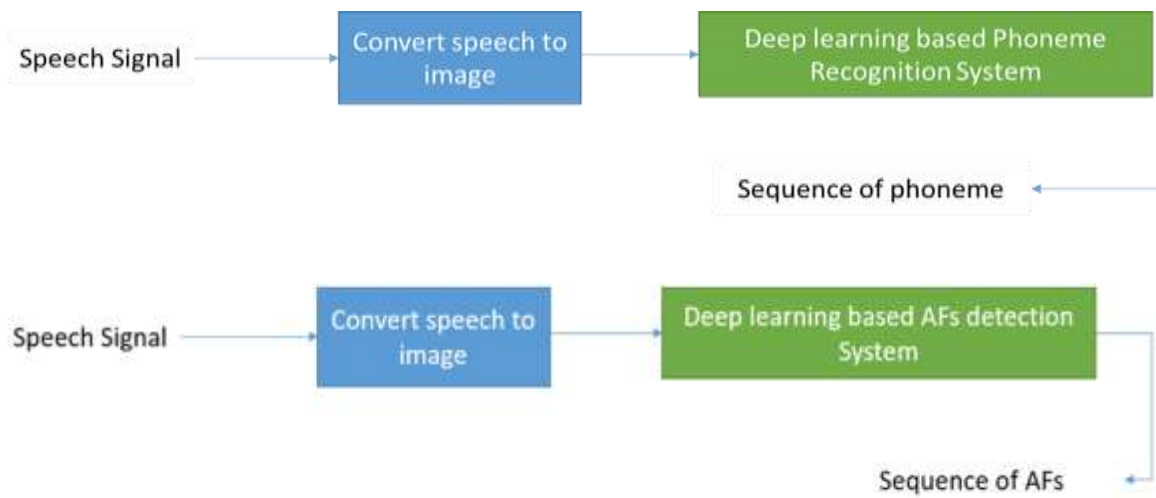


Figure 12, Overall Conversions of the speech signal to images (a: upper image, b: lower image)

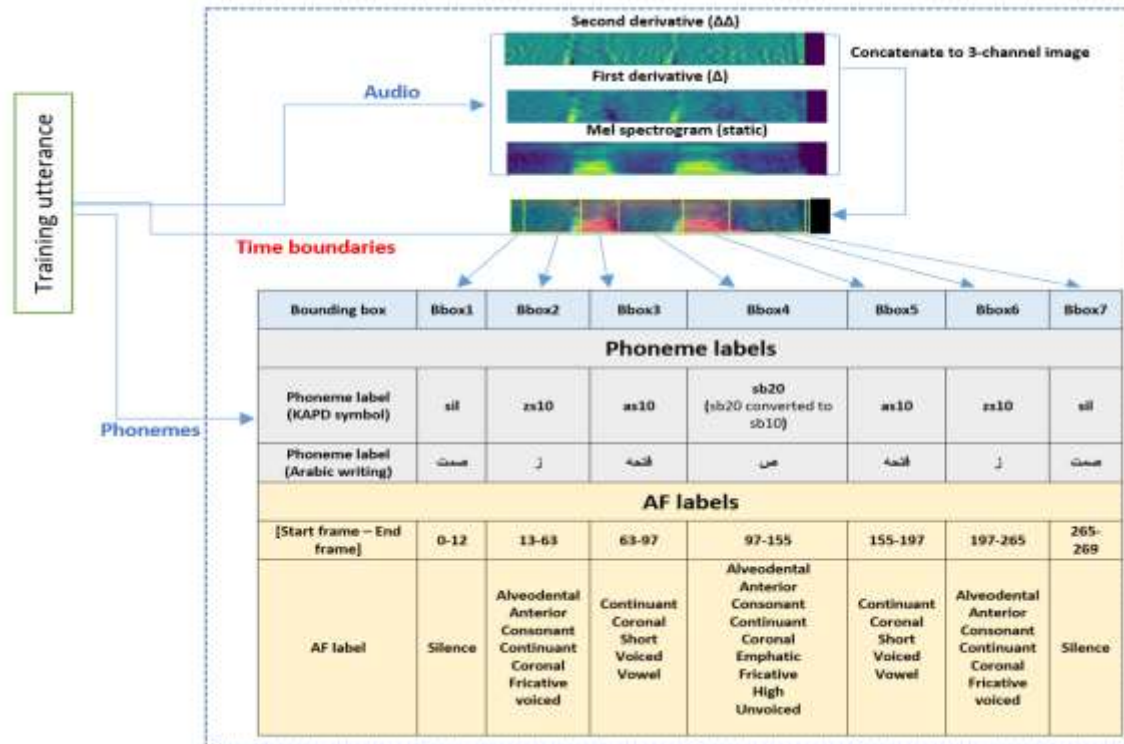


Figure 13, Example of converting the detected AFs to the corresponding phonemes.

Figure 13 shows a detailed example of the process of creating spectral images with annotations. It shows generating a spectral three-channel image from the speech and creating the associated bounding boxes for the utterance (GHSBGMA) from the KAPD training set [59].

5.6.2. Deep learning based AFs detection system (AFD-Obj)

We selected the yolov3-tiny detector for this investigation because its real time property and its support of multi-label detection. The real time property will allow our system to be used online on mobile devices. The detector consists of two main parts: backbone network and the detection layers. We started by training the backbone network of yolov3-tiny detector, which is called darknet-reference, for speech command classification task using Google speech command corpus (V2) [58]. Then, we used the weights of the backbone network to initialize the weights of yolov3-tiny detector for AFs detection task. We trained the proposed system AFD-Obj for AFs detection for the Arabic corpus and for AFs detection for English corpus. For each task, we investigated different models of AFD-Obj system by changing the number detection's scale which are: YOLOv3-tiny-1S, YOLOv3-tiny-2S, and YOLOv3-tiny-3S for one scale, two scale, and three scale of detection, respectively.

5.6.3. Results of AFD-Obj for Detecting AFs in Arabic Corpus

We used the KAPD corpus in this section to detect Arabic AFs and recognize Arabic phonemes from the detected AFs. KAPD was developed by King Abdul-Aziz City for Science and Technology at 2003 [59]. In our experiments, we used the latest version of KAPD corpus, which was developed by [60]. The total number of phonemes is 20283 for training and 8138 for testing. for mapping phoneme to AFs, we used the mapping table of [61]. The geometric mean (GM) and F-measure for 31 Arabic AFs are tabulated in Table 20 below.

Table 20, Performance metrics of the proposed system AFD-Obj for the Arabic AFs.

	YOLOv3-tiny-1S		YOLOv3-tiny-2S		YOLOv3-tiny-3S	
	GM	F-measure	GM	F-measure	GM	F-measure
affricative	0.929	0.927	0.931	0.929	0.931	0.929
alveodental	0.988	0.982	0.989	0.986	0.992	0.989
alveopalatal	0.938	0.936	0.927	0.925	0.945	0.943
anterior	0.980	0.982	0.985	0.986	0.989	0.990
aspirated	0.988	0.907	0.978	0.918	0.994	0.941
bilabial	0.954	0.876	0.930	0.868	0.940	0.908
consonant	0.998	0.998	0.997	0.997	0.999	0.998
continuant	0.992	0.993	0.994	0.994	0.993	0.994
coronal	0.977	0.975	0.980	0.978	0.984	0.983
emphatic	0.904	0.891	0.912	0.900	0.913	0.904
fricative	0.992	0.990	0.993	0.991	0.990	0.990
glottal	0.968	0.903	0.984	0.915	0.975	0.933
high	0.932	0.918	0.939	0.923	0.927	0.920
interdental	0.856	0.775	0.865	0.811	0.879	0.833
labiodental	0.795	0.721	0.838	0.776	0.803	0.729
labiovelar	1.000	0.967	0.988	0.953	0.988	0.966
lateral	0.960	0.922	0.960	0.897	0.969	0.873
nasal	0.979	0.963	0.951	0.929	0.978	0.973
palatal	0.978	0.967	0.978	0.977	0.967	0.945
pharyngeal	0.984	0.984	0.966	0.960	0.967	0.961
plosive	0.960	0.913	0.961	0.926	0.965	0.936
rounded	0.982	0.940	0.987	0.949	0.990	0.966
semivowel	0.989	0.967	0.994	0.983	0.989	0.972
short	0.997	0.995	0.999	0.998	0.999	0.998
silence	0.999	0.999	0.998	0.998	1.000	0.999
trill	0.955	0.933	0.954	0.932	0.919	0.874
unvoiced	0.985	0.964	0.983	0.972	0.982	0.971
uvular	0.960	0.917	0.953	0.932	0.938	0.926
velar	0.989	0.958	0.968	0.928	0.989	0.948
voiced	0.995	0.996	0.996	0.996	0.996	0.997
vowel	0.999	0.999	1.000	0.999	0.999	0.999
Average	0.965	0.941	0.964	0.943	0.964	0.945

For all AFs, the systems achieved a GM greater than 80%, except for labiodental. For the F-measure of all AFs, the systems achieved accuracies greater than 80%, except for labiodental, which had an F-measure of 72.1%, 77.6%, and 72.9% using YOLOv3-tiny-1S, YOLOv3-tiny-2S, and YOLOv3-tiny-3S, respectively, and interdental, which had an F-measure of 77.5% using YOLOv3-tiny-1S. In general, we achieved GM and F-measure average accuracies of 96.5% and 94.1% for the YOLOv3-tiny-1S model, 96.4% and 94.3% for the YOLOv3-tiny-2S model, and 96.4% and 94.5% for the YOLOv3-tiny-3S model. These results are better than those of state-of-the-art results [61], where approximately 45% of the AFs obtained less than 80% for GM and approximately 61% obtained less than 80% for the F-measure using their best model (i.e., DBN–DNN). We achieved our results using a single network for all AFs, while Ref. [61] used a different network for each AF. Moreover, our testing input is a whole utterance without time boundary information, while their testing input was speech phonemes. We also detected the time boundaries of each AF; therefore, we can calculate the accuracy at the frame level.

(i.e., DBN–DNN) by almost 40% at 100% similarity and outperformed the matching rate of their classifier by approximately 4% at 90% similarity [61]. Using 100% similarity, we achieved correction rates of 86.04%, 88.06%, and 89.35% for YOLOv3-tiny-3S, YOLOv3-tiny-2S, and YOLOv3-tiny-1S, respectively, compared to 64% (matching rate) for the model in Ref. [61]. These values increased to 91.16%, 92.38%, and 92.59%, respectively, when using 90% similarity for all three models compared to 89% for that in Ref. [61]. This increase can be attributed to the fact that the correction rate measure ignored the insertion errors; hence, we ignored many insertion errors when using only 90% similarity.

For 100% similarity, our models obtained PERs of 14.13%, 12.09%, and 10.84%, respectively, which increased to 20.1%, 15.53%, and 12.57%, respectively, for 90% similarity. Ref. [61] did not provide the PER result. Another point to highlight is that these observations confirmed our postulation for not needing the second and third scales of the YOLO detector in the AF detection and phoneme recognition. The PER results also illustrate that using 90% similarity during AF matching to generate the corresponding phonemes is not acceptable because more wrong phonemes can be recognized as correct, as shown in Table 21.

Table 21, PER (%) and correction rate (%) for our proposed AFD-Obj system and results of [61].

Matching rate (# bits)	Model	PER (%)	Correction rate (%)
100% (0 bit)	YOLOv3-tiny-3S	14.13	86.04
	YOLOv3-tiny-2S	12.09	88.06
	YOLOv3-tiny-1S	10.84	89.35
	DBN-DNN [61]	-	64.00 (Exact matching rate)
	YOLOv3-tiny-3S	20.1	91.16
90% (3 bits)	YOLOv3-tiny-2S	15.53	92.38
	YOLOv3-tiny-1S	12.57	92.59
	DBN-DNN [61]	-	89.00 (Matching rate)

5.6.5. Results of the AFD-Obj System for Detecting AFs in the English Corpus

This section presents the results of applying our proposed system for detecting the English AFs using the TIMIT corpus. We used accuracy at the frame level, by considering the bounding box coordinates as the start and end frames, as our evaluation metric, which is calculated as shown in Figure 15.

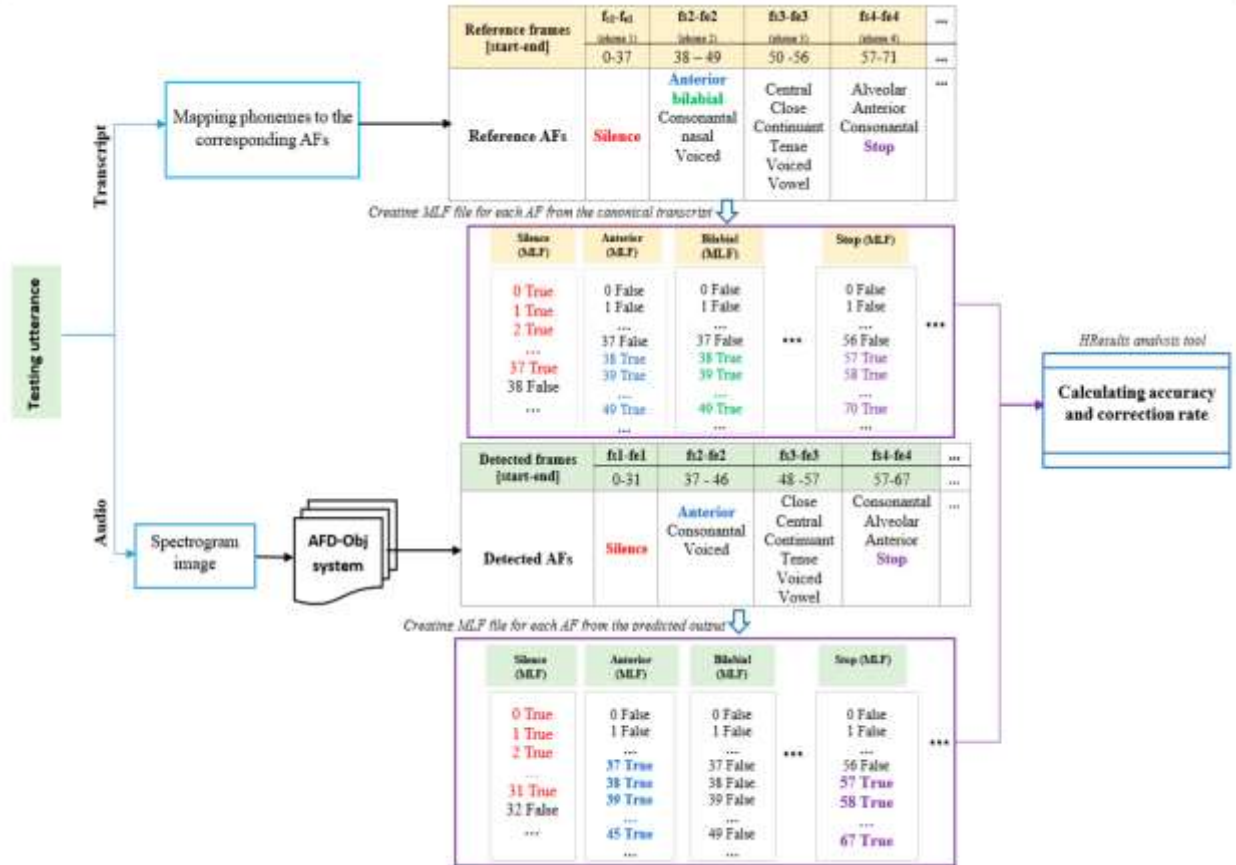


Figure 15, Testing phase of the AFD-Obj system: calculating the frame level accuracy of the detected outputs.

We then compared the results of the proposed system with that of the state-of-art published work in AFs detection using TIMIT [62], called LAS-MTL-M. We considered for our comparison the results of LAS-MTL-M which were reported at frame level. Authors of [62] used the TIMIT segments markup (time boundaries) to calculate the accuracies of the column “markup frames” and the DTW algorithm to convert soft attention to hard attention to calculate the accuracies of the “frames” column. In both cases, they dealt with the different number of predicted and target frames by taking the minimum length of target and prediction, as shown in the code provided. For better comparison, we calculated the accuracy of our proposed system using the coordinates of the detected bounding boxes as markup frames, after taking the minimum length of the predicted and

target frames. This result is presented in the column “bounding box coord.” of Table 22. We also used another measure to deal with the problem of the difference in length between the predicted and target frames, where we used HResults analysis tool to align the predicted and target frames, then we calculate the accuracy. Moreover, HResults accuracy is more precise because it considers the insertion errors. This accuracy is presented in column “HResult align.” of Table 22 . We also show the results with those in [63], which worked on detecting only some AFs in TIMIT.

Table 22 presents the result of our proposed system AFD-Obj with the three models. The table shows that our system achieved results comparable to the published result for all TIMIT AFs. Our models had an average accuracy (with bounding box coord.) of 94.29%, 95.04%, and 95.13%, and average accuracy (using HResult) of 93.23%, 93.47%, and 93.66% for YOLOv3-tiny-3S, YOLOv3-tiny-2S, and YOLOv3-tiny-1S, respectively, while the phones-las-frames model had an average detection accuracy of 95.5% using markup-frames. Since the test results of the “markup-frames” of [62] depend on segmenting the speech into markup frames, while our system doesn’t depend on any segmentation, hence we think fair comparison should be with the result of the “frames” column of [62].

Table 22, Detection accuracy of all 28 English AFs using the proposed system AFD-Obj and state-of-the-art methods.

Articulatory features	AFD-Obj system						LAS-MTL-M markup-frames [62]	LAS-MTL-M frames [62]	KT [63]
	YOLOv3-tiny-1S		YOLOv3-tiny-2S		YOLOv3-tiny-3S				
	Bounding box coord.	HResult align.	Bounding box coord.	HResult align.	Bounding box coord.	HResult align.			
Alveolar	91.05	90.22	90.92	90.01	89.31	88.96	95	77	
Anterior	89.69	89.34	89.55	89.02	87.92	88.08	90	69	90
Approximant	97.12	95.39	97.17	95.32	96.87	95.39	98	94	68
Bilabial	97.70	95.89	97.53	95.57	97.30	95.78	98	93	
Central	93.73	92.31	93.73	92.18	93.36	92.27	99	91	
Close	94.13	92.65	94.02	92.46	93.36	92.33	97	88	86
Consonantal	88.97	88.75	88.75	88.42	87.32	87.64	88	64	90
Continuant	91.37	90.46	90.88	90.04	88.60	88.38	89	68	86
Fricative	96.03	94.56	95.73	94.21	95.04	94.06	95	83	88
Front	93.33	91.96	93.42	91.92	92.06	91.12	95	89	84
Glottal	98.67	96.69	98.62	96.48	98.42	96.82	99	98	
labiodental	98.88	96.89	98.80	96.71	98.57	96.94	99	96	
Lateral approximant	98.21	96.34	98.11	96.07	97.88	96.31	99	96	
Mid	90.28	89.09	90.04	88.77	88.87	88.3	97	82	
Nasal	97.59	95.95	97.55	95.72	97.15	95.74	99	93	84

Non sibilant fricative	97.60	95.8	97.50	95.58	97.22	95.74	97	94	
Open	96.09	94.31	95.83	93.98	95.63	94.23	98	91	93
palatal	99.60	97.54	99.63	97.4	99.57	97.81	99	99	
postalveolar	99.18	97.12	99.17	96.94	98.96	97.21	99	97	
Round	94.99	93.36	94.70	92.97	94.30	93.04	98	91	92
Sibilant affricate	99.50	97.41	99.51	97.29	99.38	97.64	99	99	
Sibilant fricative	97.97	96.1	97.81	95.95	97.37	96	98	90	
Silence	96.79	95.21	97.05	95.29	96.68	95.35	80	63	89
Stop	95.03	93.74	95.05	93.61	94.46	93.53	97	85	96
Tense	89.63	88.46	89.92	88.73	88.65	87.94	97	81	87
Velar	98.37	96.47	98.31	96.25	98.01	96.37	99	95	
Voiced	90.86	89.9	90.71	89.69	88.62	88.19	84	72	93
vowel	91.31	90.69	91.19	90.57	89.29	89.26	92	70	92
<u>Average</u>	95.13	93.66	95.04	93.47	94.29	93.23	95.5	86	

An important observation from Table 22 is that our models detected silence within the utterance with a high accuracy compared to the phone-las models [62], which achieved only 63% and 80% for frames and markup frames, respectively. This high performance in detecting silence in continuous speech is very promising and can be looked at as an important achievement by itself. YOLOv3-tiny-1S had the best average detection accuracy (95.13%), while our model YOLOv3-tiny-2S had almost the same average detection accuracy (95.04%). The average detection accuracy of YOLOv3-tiny-3S was 94.29%. These results reinforce our previous assumption for not needing the second and third scales of YOLO detection for our specific application.

5.6.6. Results of PD-Obj for Detecting phonemes in Arabic Corpus

Table 23 presents the results of investigating the proposed PD-Obj system for phoneme recognition using the KAPD corpus. Our proposed system using the YOLOv3-tiny-2S model achieved the lowest PER of 5.63%, while the YOLOv3-tiny-1S and YOLOv3-tiny-3S models achieved 5.79% and 6.29% PER, respectively. These results are remarkable and show that our proposed system has an excellent potential compared to the recent state-of-the-art systems on this corpus [64]. Ref. [64] used the HMM for the Arabic phoneme recognition using the DPF elements. The results also reinforced our previous assumption for not needing the second and third scales of the YOLO detection for our specific application.

Table 23, PER and correction rate of the Arabic phoneme recognition using the proposed models.

Model	PER (%)	Correction rate (%)
PD-Obj (YOLOv3-tiny-3S)	6.29	93.94
PD-Obj (YOLOv3-tiny-2S)	5.63	94.56
PD-Obj (YOLOv3-tiny-1S)	5.79	94.34
AFD-Obj (YOLOv3-tiny-1S)	10.85	89.33
PDF-HMM [64]	39.57	70.68

We observe from Table 23 that the PD-Obj system obtained better results than the system based on the AFD-Obj. However, detecting the AFs is important in many applications such as, pronunciation error correction and diagnosis. And the AFs are universal between many languages. An interesting point for future work is to see how to improve the accuracy of the system that performs phoneme recognition based on the detected AFs.

We also calculated the correction rate of each phoneme for the YOLOv3-tiny-1S model. We found that 79% of the Arabic phonemes had a correction rate greater than 80%, while 44% had a correction rate greater than 90%.

5.6.7. Implementation of the score measurement in the CAPT

The analysis of the errors in pronunciation by the language experts is ongoing and not complete yet because it requires long time since the text that we designed for the CAPT is very long. Nonetheless we are using part of publically available KSU speech database to investigate the scoring measures available and investigating new proposed scores that may be more suitable for our proposed method of treating the phonemes and the AFs as objects in 3 channels spectral images and for our end-to-end recognition systems.

6. Discussion

We successfully applied object detection techniques for phoneme and AFs detection using Arabic and English speech. The achieved results showed the effectiveness of the proposed method and encouraged us to investigate using object detection for mispronunciation error detection task and pronunciation scoring for non-native Arabic speech. The investigation showed excellent results that is comparable or better than state of the art research in mispronunciation error detection task and pronunciation scoring in CAPT systems. We are finalizing a paper with this investigation and its results.

7. Future work

The proposed technique and the improved models will be used to build a CAPT system using the recording of session 1. The performance of the system will be evaluated and compared to human judges. A new session of the database will be recorded. The first CAPT system will be improved based on the analysis of its performance and a new CAPT system will be built using the recording of sessions 1 and 2. The performance of the second CAPT system will be evaluated and compared to human judges. The results and analysis of the performance of the first and second CAPT system will be published in reputed journals.

8. References

- [1] M. Alsulaiman, Z. Ali, G. Muhammed, M. Bencherif, and A. Mahmood, "KSU speech database: text selection, recording and verification," in *2013 European Modelling Symposium*, 2013, pp. 237–242.
 - [2] S. E. Hamid, O. Abdel-Hamid, and M. Rashwan, "Performance Tuning and System Evaluation for Computer Aided Pronunciation Learning," in *Proceedings of NEMLAR International Conference on Arabic Language Resources and Tools*, 2009, pp. 140–143.
 - [3] S. M. Abdou *et al.*, "Computer aided pronunciation learning system using speech recognition techniques," in *Ninth International Conference on Spoken Language Processing*, 2006.
 - [4] K. Necibi and H. Bahi, "An arabic mispronunciation detection system by means of automatic speech recognition technology," in *The 13th International Arab Conference on Information Technology Proceedings*, 2012, pp. 303–308.
 - [5] H. Dahan, A. Hussin, Z. Razak, and M. Odelha, "Automatic arabic pronunciation scoring for language instruction," 2011.
 - [6] M. Belgacem, A. Maatallaoui, and M. Zrigui, "Arabic language learning assistance based on automatic speech recognition system," in *Proceedings of the International Conference on e-Learning, e-Business, Enterprise Information Systems, and e-Government (EEE)*, 2011, p. 1.
 - [7] M. S. El-Kasasy, "An Automatic Speech Verification System," Ph. D. Thesis, Cairo University, Faculty of Engineering, Department of~..., 1992.
 - [8] M. S. Abdo, A. H. Kandil, A. M. El-Bialy, and S. A. Fawzy, "Automatic detection for some common pronunciation mistakes applied to chosen Quran sounds," in *2010 5th Cairo International Biomedical Engineering Conference*, 2010, pp. 219–222.
 - [9] L. R. Rabiner and B.-H. Juang, *Fundamentals of speech recognition*. Tsinghua University Press, 1999.
 - [10] R. M. Hegde, "Fourier transform phase-based features for speech recognition," *Indian Inst. Technol. Madras*, 2005.
-

-
- [11] B. S. Atal, "Effectiveness of linear prediction characteristics of the speech wave for automatic speaker identification and verification," *J. Acoust. Soc. Am.*, vol. 55, no. 6, pp. 1304–1312, 1974.
- [12] F. Nazir, M. N. Majeed, M. A. Ghazanfar, and M. Maqsood, "Mispronunciation detection using deep convolutional neural network features and transfer learning-based model for Arabic phonemes," *IEEE Access*, vol. 7, pp. 52589–52608, 2019.
- [13] S. Akhtar *et al.*, "Improving Mispronunciation Detection of Arabic Words for Non-Native Learners Using Deep Convolutional Neural Network Features," *Electronics*, vol. 9, no. 6, p. 963, 2020.
- [14] U. Shrawankar and V. M. Thakare, "Techniques for feature extraction in speech recognition system: A comparative study," *arXiv Prepr. arXiv1305.1145*, 2013.
- [15] S. K. Gaikwad, B. W. Gawali, and P. Yannawar, "A review on speech recognition technique," *Int. J. Comput. Appl.*, vol. 10, no. 3, pp. 16–24, 2010.
- [16] D. Palaz, R. Collobert, and others, "Analysis of cnn-based speech recognition system using raw speech as input," 2015.
- [17] J. Lee, T. Kim, J. Park, and J. Nam, "Raw waveform-based audio classification using sample-level CNN architectures," *arXiv Prepr. arXiv1712.00866*, 2017.
- [18] G. Kovács, L. Tóth, D. Van Compernelle, and S. Ganapathy, "Increasing the robustness of CNN acoustic models using autoregressive moving average spectrogram features and channel dropout," *Pattern Recognit. Lett.*, vol. 100, pp. 44–50, 2017.
- [19] H. Meng, T. Yan, F. Yuan, and H. Wei, "Speech emotion recognition from 3D log-mel spectrograms with deep learning network," *IEEE Access*, vol. 7, pp. 125868–125881, 2019.
- [20] N. Souissi and A. Cherif, "Speech recognition system based on short-term cepstral parameters, feature reduction method and artificial neural networks," in *2016 2nd international conference on advanced technologies for signal and image processing (ATSIP)*, 2016, pp. 667–671.
- [21] A. Winursito, R. Hidayat, and A. Bejo, "Improvement of MFCC feature extraction accuracy using PCA in Indonesian speech recognition," in *2018 International Conference on Information and Communications Technology (ICOIACT)*, 2018, pp. 379–383.
-

-
- [22] C. O. S. Sorzano, J. Vargas, and A. P. Montano, "A survey of dimensionality reduction techniques," *arXiv Prepr. arXiv1403.2877*, 2014.
- [23] Y. Alotaibi, S.-A. Selouani, and D. O'shaughnessy, "Experiments on automatic recognition of nonnative Arabic speech," *EURASIP J. Audio, Speech, Music Process.*, vol. 2008, no. 1, p. 679831, 2008.
- [24] M. Belgacem and M. Zrigui, "Automatic Identification System of Arabic Dialects," in *IPCV 2010: proceedings of the 2010 international conference on image processing, computer vision, & pattern recognition (Las Vegas NV, July 12-15, 2010)*, 2010, pp. 740–749.
- [25] W. J. J. Roberts and J. P. Willmore, "Automatic speaker recognition using Gaussian mixture models," in *1999 Information, Decision and Control. Data and Information Fusion Symposium, Signal Processing and Communications Symposium and Decision and Control Symposium. Proceedings (Cat. No. 99EX251)*, 1999, pp. 465–470.
- [26] A. Trigui, A. Mars, M. A. Ben Jannet, M. Maraoui, and M. Zrigui, "Foreign accent classification for Arabic speech learning," in *Proceedings of the International Conference on Image Processing, Computer Vision, and Pattern Recognition (IPCV)*, 2011, p. 1.
- [27] W. J. Barry, C. E. Hoequist, and F. J. Nolan, "An approach to the problem of regional accent in automatic speech recognition," *Comput. Speech Lang.*, vol. 3, no. 4, pp. 355–366, 1989.
- [28] L. Indrayanti, T. Usagawa, Y. Chisaki, and T. Dutono, "Evaluation of pronunciation by means of automatic speech recognition system for computer aided Indonesian language learning," in *2006 7th International Conference on Information Technology Based Higher Education and Training*, 2006, pp. 553–556.
- [29] A. Hannun *et al.*, "Deep speech: Scaling up end-to-end speech recognition," *arXiv Prepr. arXiv1412.5567*, 2014.
- [30] D. Amodei *et al.*, "Deep speech 2: End-to-end speech recognition in english and mandarin," in *International conference on machine learning*, 2016, pp. 173–182.
- [31] R. Collobert, C. Puhersch, and G. Synnaeve, "Wav2letter: an end-to-end convnet-based speech recognition system," *arXiv Prepr. arXiv1609.03193*, 2016.
- [32] Y. Miao, M. Gowayyed, and F. Metze, "EESSEN: End-to-end speech recognition using deep RNN models and WFST-based decoding," in *2015 IEEE Workshop on Automatic Speech*
-

-
- Recognition and Understanding (ASRU)*, 2015, pp. 167–174.
- [33] A. Zeyer, K. Irie, R. Schlüter, and H. Ney, “Improved training of end-to-end attention models for speech recognition,” *arXiv Prepr. arXiv1805.03294*, 2018.
- [34] L. Zhang *et al.*, “End-to-End Automatic Pronunciation Error Detection Based on Improved Hybrid CTC/Attention Architecture,” *Sensors*, vol. 20, no. 7, p. 1809, 2020.
- [35] T.-H. Lo, S.-Y. Weng, H.-J. Chang, and B. Chen, “An Effective End-to-End Modeling Approach for Mispronunciation Detection,” *arXiv Prepr. arXiv2005.08440*, 2020.
- [36] Y. Feng, G. Fu, Q. Chen, and K. Chen, “SED-MDD: Towards Sentence Dependent End-To-End Mispronunciation Detection and Diagnosis,” in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2020, pp. 3492–3496.
- [37] J. Van Doremalen, C. Cucchiari, and H. Strik, “Automatic detection of vowel pronunciation errors using multiple information sources,” in *2009 IEEE Workshop on Automatic Speech Recognition & Understanding*, 2009, pp. 580–585.
- [38] N. Oostdijk, “The Spoken Dutch Corpus. Overview and First Evaluation.,” in *LREC*, 2000, pp. 887–894.
- [39] H. Franco, L. Neumeyer, V. Digalakis, and O. Ronen, “Combination of machine scores for automatic grading of pronunciation quality,” *Speech Commun.*, vol. 30, no. 2–3, pp. 121–130, 2000.
- [40] K. P. Truong, A. Neri, F. de Wet, C. Cucchiari, and H. Strik, “Automatic detection of frequent pronunciation errors made by L2-learners,” in *Ninth European Conference on Speech Communication and Technology*, 2005.
- [41] K. Fukunaga, *Introduction to statistical pattern recognition*. Elsevier, 2013.
- [42] L. F. Weigelt, S. J. Sadoff, and J. D. Miller, “Plosive/fricative distinction: The voiceless case,” *J. Acoust. Soc. Am.*, vol. 87, no. 6, pp. 2729–2737, 1990.
- [43] Y. Kim, H. Franco, and L. Neumeyer, “Automatic pronunciation scoring of specific phone segments for language instruction,” in *Fifth European Conference on Speech Communication and Technology*, 1997.
-

-
- [44] S. Xu, J. Jiang, Z. Chen, and B. Xu, "Automatic pronunciation error detection based on linguistic knowledge and pronunciation space," in *2009 IEEE International Conference on Acoustics, Speech and Signal Processing*, 2009, pp. 4841–4844.
- [45] M.-S. Liang, J.-Y. Hung, R.-Y. Lyu, and Y.-C. Chiang, "Pronunciation error detection for computer assisted pronunciation teaching in mandarin," in *2008 6th International Symposium on Chinese Spoken Language Processing*, 2008, pp. 1–4.
- [46] R. Lyu, M. Liang, and Y. Chiang, "Toward constructing a multilingual speech corpus for Taiwanese (Min-nan), Hakka, and Mandarin," in *International Journal of Computational Linguistics & Chinese Language Processing, Volume 9, Number 2, August 2004: Special Issue on New Trends of Speech and Language Processing*, 2004, pp. 1–12.
- [47] K. Dembczynski, W. Waegeman, W. Cheng, and E. Hüllermeier, "An exact algorithm for F-measure maximization," *Adv. Neural Inf. Process. Syst.*, vol. 24, pp. 1404–1412, 2011.
- [48] G. Zhao *et al.*, "L2-ARCTIC: A non-native English speech corpus," *Percept. Sens. Instrum. Lab*, 2018.
- [49] K. Li, X. Qian, and H. Meng, "Mispronunciation detection and diagnosis in 12 english speech using multidistribution deep neural networks," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 25, no. 1, pp. 193–207, 2016.
- [50] N. F. Chen, R. Tong, D. Wee, P. Lee, B. Ma, and H. Li, "iCALL corpus: Mandarin Chinese spoken by non-native speakers of European descent," in *Sixteenth Annual Conference of the International Speech Communication Association*, 2015.
- [51] T. Lander, "CSLU: Foreign Accented English." Release, 2007.
- [52] S. Col, A. LaRocca, and R. Chouairi, "West Point Arabic Speech," *LDC Cat. LDC2002S02*, 2002.
- [53] S. Schaden and U. Jekosch, "'Casselberveetovallarga' and other Unpronounceable Places: The CrossTowns Corpus.," in *LREC*, 2006, pp. 993–998.
- [54] S. Weinberger, "Speech accent archive," *Georg. Mason Univ.*, 2015.
- [55] P. Meier, "International dialects of English archive," *IDEA-The Int. Dialects English Arch.*, 1997.
-

- [56] M. Algabri, H. Mathkour, M. A. Bencherif, M. Alsulaiman, and M. A. Mekhtiche, “Towards Deep Object Detection Techniques for Phoneme Recognition,” *IEEE Access*, vol. 8, pp. 54663–54680, 2020.
 - [57] M. Algabri, H. Mathkour, M. M. Alsulaiman, and M. A. Bencherif, “Deep learning-based detection of articulatory features in arabic and english speech,” *Sensors (Switzerland)*, vol. 21, no. 4, 2021, doi: 10.3390/s21041205.
 - [58] P. Warden, “Speech commands: A dataset for limited-vocabulary speech recognition,” *arXiv Prepr. arXiv1804.03209*, 2018.
 - [59] M. Alghmadi, “KACST arabic phonetic database,” in *the Fifteenth International Congress of Phonetics Science, Barcelona*, 2003, pp. 3109–3112.
 - [60] Y. Seddiq, A. Meftah, M. Alghamdi, and Y. Alotaibi, “Reintroducing KAPD as a Dataset for Machine Learning and Data Mining Applications,” in *2016 European Modelling Symposium (EMS)*, 2016, pp. 70–74.
 - [61] Y. Seddiq, Y. A. Alotaibi, S.-A. Selouani, and A. H. Meftah, “Distinctive Phonetic Features Modeling and Extraction Using Deep Neural Networks,” *IEEE Access*, vol. 7, pp. 81382–81396, 2019.
 - [62] I. Karaulov and D. Tkanov, “Attention model for articulatory features detection,” *arXiv Prepr. arXiv1907.01914*, 2019.
 - [63] S. King and P. Taylor, “Detection of phonological features in continuous speech using neural networks,” 2000.
 - [64] Alotaibi, Y.A.; Selouani, S.-A.; Yakoub, M.S.; Seddiq, Y.M.; Meftah, A. A canonicalization of distinctive phonetic features to improve arabic speech recognition. *Acta Acust. United Acust.* 2019, 105, 1269–1277.
-

9. Publications / Presentations

We published the following paper:

M. Algabri, H. Mathkour, M. M. Alsulaiman, and M. A. Bencherif, “Deep learning-based detection of articulatory features in arabic and english speech,” *Sensors (Switzerland)*, vol. 21, no. 4, 2021, doi: 10.3390/s21041205.

This paper is within the thesis of Mohammad Aljabri (1st author). PI Mansour Alsulaiman is the co-advisor for this thesis.

10. Appendices

Appendix A - Text Selection Comparison (V1 to V3)

The table has been trimmed because the work is under publication

V3	V2	V1
حَرْفُ الْيَاءِ:	حَرْفُ الْيَاءِ:	حَرْفُ الْيَاءِ:
	-	-
حَرْفُ الدَّالِ:		
	حَرْفُ الدَّالِ:	
		حَرْفُ الدَّالِ:

نصوص مشروع الأصوات وتدريباته

- أولاً: نصوص الأصوات.
- ثانياً: تدريبات التمييز الصوتي.

أولاً: نصوص الأصوات

The item has been deleted because the work is under publication

ثانياً: تدريبات التمييز الصوتي

The item has been deleted because the work is under publication

Appendix C - Tahadath App Screen Cards



Appendix D - Durations per Speaker

Paragraphs and SPW for 44 first Speakers

Appendix E- ELAN Annotation CAPT protocol

In order to unify the checking procedure, Each ELAN-Checker was asked to follow the following protocol: (including an explanatory video in Appendix F)

STEP 1:

1. Open a new eaf file :
2. Include the wavefile of the speakers
3. Import the TextGrid of the same speakers
4. Save the work an speaker EAF file, SAME DIRECTORY

STEP 2:

1. Remove default tier
2. Activate the To_Remove tier



Figure E.1: Elan Tier Selection

1. Listen to the whole segment of the file.
2. The Arabic text is inside the Sentences tier.
3. If any repeated speech is heard:
 - a. Try to locate the left and right boundaries of that segment.
 - b. Select the tier To_Remove,
 - c. Insert the left boundary split
 - d. Insert the left boundary split
4. Move to next Speech sentence.

N.B :

Please mark whenever you can:

- Any word not in the text
- Any additional text, at start or end of file.
- Any extra sound (baby, cat, door opening/closing, phone notification,)
- If any speaker is not worth hearing, or understanding, do not continue on its speech, but put a remark in an additional file. (problems.docx).

Appendix F - Tahadath Application User Manual

Tahadth Application User Manual

دليل المستخدم لتطبيق تحدث



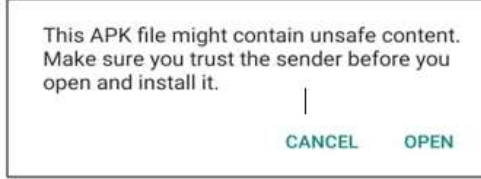
نظام حاسوبي لتعليم نطق أصوات اللغة العربية
للساطقين بغيرها

مركز أبحاث الروبوتات الذكية معهد اللغويات العربية
كلية علوم الحاسب والمعلومات

جامعة الملك سعود

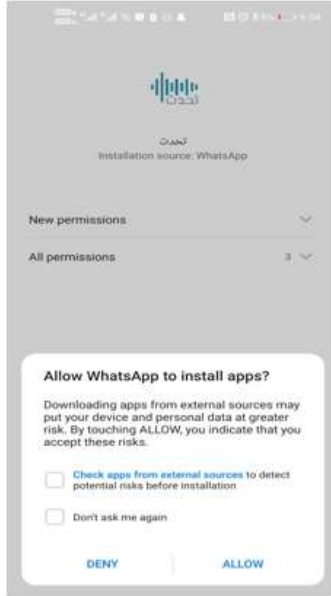
The application "Tahadth" is used to record some Arabic texts.

- To install the "Tahadath" application,
- Click on the application, the following screen will appear:



Ignore the message content, this application is safe.

- Click on the "OPEN" button, the following screen will appear:



- Click on the "ALLOW" button, then the following screen will appear:



إن تطبيق "تحدث" Tahadth يستخدم لتسجيل بعض النصوص العربية.

- لتحميل تطبيق "تحدث Tahadth" اضغط على التطبيق ، ومثل اي تطبيق ستظهر شاشة تحمل المحتوى التالي:



تجاهل محتوى الرسالة ، حيث أن التطبيق آمن طالما تثقيت هذا التطبيق من جهة معروفة.

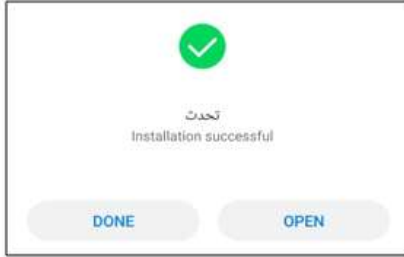
- اضغط على زر "فتح" و بعدها تظهر الشاشة التالية:



- اضغط على زر "سماح" ومن ثم تظهر الشاشة التالية ليده تثبيت التطبيق:



- Click on the "INSTALL" button to install the application, after that the following screen appears to open the application:



- To open the application, click on the "Open" button, the application has been installed on the device as "Tahadth", the following screen appears:



- Wait until the following login screen appears:



- اضغط على زر "تثبيت" و عند الانتهاء من تثبيت التطبيق تظهر الشاشة التالية لفتح التطبيق:



- لفتح التطبيق اضغط على زر "فتح" ، حيث أن التطبيق كذلك تم تثبيته على الجهاز باسم "تحدث" ، ومن ثم تظهر الشاشة التالية:



- إنتظر حتى تظهر شاشة تسجيل الدخول التالية:



- Enter the user name and password, usually the username is in lowercase letters. (credentials are sent privately to you)
- After the login, the following screen will appear, click on the check box of the phrase "I agree to allow the use of my recordings for research purposes", which confirms your approval to use your recordings for research studies, after that click on "انقر للبدء" button:



- أدخل اسم المستخدم وكلمة المرور اللذين تم ارسالهما على الخاص وعادةً ما يكون اسم المستخدم بالأحرف الانجليزية الصغيرة.
- بعد عملية تسجيل الدخول ستظهر الشاشة التالية ، قم بالتأشير على عبارة " أقر بالسماح باستخدام التسجيلات لأغراض بحثية" التي تؤكد إقرارك و موافقتك على استخدام تسجيلاتك لأغراض بحثية ، بعدها اضغط على زر "انقر للبدء ":



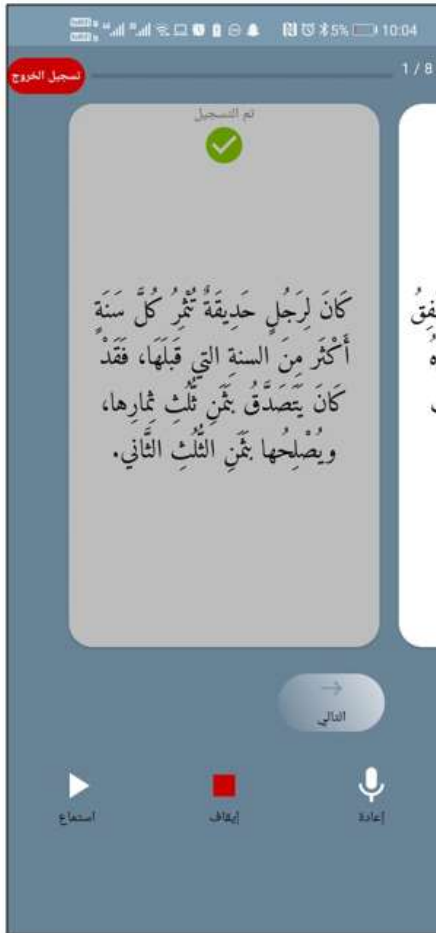
- The first screen of the application appears, where the text to be recorded appears in a card:



- تظهر الشاشة الأولى من التطبيق، التي يظهر فيها النص المراد تسجيله:



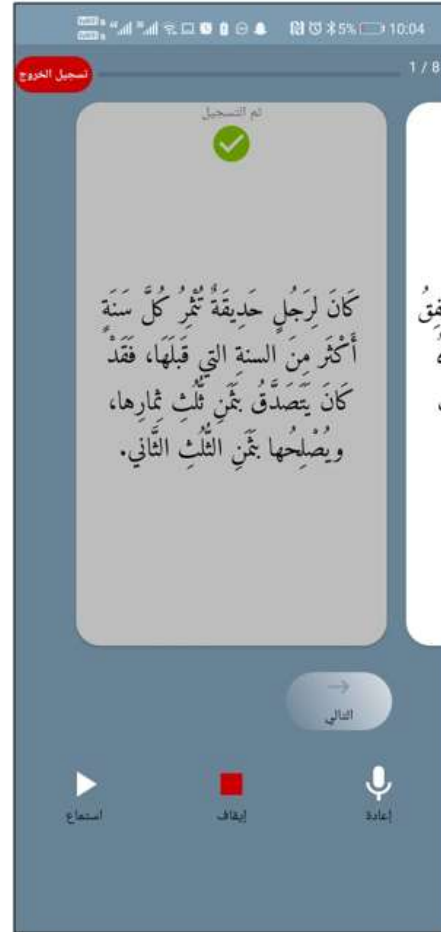
- Two buttons on the recording screen, a microphone-recording button "تسجيل" and a button to stop recording as a red square "إيقاف". To start the recording process, click on the microphone button and read the text shown on the screen, after you finish click on the stop recording button.
- The recorded sound is automatically approved, and the sign end of recording appears at the top of this screen "تم التسجيل", as shown in the following screen:



- To move to the next text, click on the "التالي" button, or by swiping the screen to the left. You can also listen to your recorded voice by clicking on the "استماع" button.

• في شاشة التسجيل يظهر زر تسجيل بشكل ميكروفون و زر "إيقاف" التسجيل كمرجع احمر ، لبدء عملية التسجيل ، اضغط على زر الميكروفون و قراءة النص الظاهر على الشاشة بصوت واضح وعند الانتهاء من قراءة النص اضغط على زر "إيقاف" التسجيل، (يفضل استخدام ميكروفون سماعة خارجية).

• يتم اعتماد الصوت المسجل آليا و تظهر علامة انتهاء التسجيل لهذه الشاشة و ذلك بظهور رسالة تأكيد إتمام التسجيل في اعلى الشاشة "تم التسجيل" كما هو موضح بالشاشة التالية:



- إنتقل للنص التالي وذلك بالضغط على زر التالي او بسحب الشاشة لليسار. يمكنك قبل الانتقال الاستماع لصوتك المسجل بالضغط على زر استماع.

- Once all the recordings are complete, the following message appears:



- Click on "نعم" button. If you are sure that all screens have been recorded, please click on the approval button "اعتماد" , which appears in the following screen:



- وهكذا يتم تسجيل نصوص كل الشاشات التالية.

- بعد الانتهاء من تسجيل نص آخر شاشة و في حالة تأكد النظام من تسجيلك لجميع الشاشات ستظهر لك الرسالة المبنية التالية:

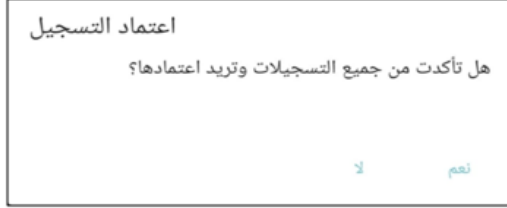


- اضغط على زر "نعم".

- لاعتماد التسجيل اضغط على زر اعتماد الظاهر في اسفل الشاشة والموضح في الصورة التالية:



- When you click on the approval button "اعتماد", a confirmation screen appears as follows:



- The previous screen represents the recording approval.
- Press the "yes نعم" button, you confirm that you have completed the recording, then you are automatically logged out of the application and back to the login screen.
- Click on the "لا No" option, if you want to listen to the previous recordings. Press the previous buttons to scroll back over the recordings.
- **Note: If you don't press "نعم yes", the recordings are not approved.**
- When you press the "Approval اعتماد" button, and you have some screens that have not been recorded, a pop-up message will appear to inform you of the number of screens that you have not recorded, as in the following example::



- Press on the "Yes نعم" button, as you will be kept on the last screen. You can return to screens that you did not record by using the "السابق Previous" button to moves between screens
- Once you finish recording all the required texts, go to the last screen and press the approve button "إعتماد" approved to your recordings.

- عند الضغط على زر اعتماد تظهر الشاشة المبتدئة التالية وذلك للتأكد من رغبتك في اعتماد تسجيلاتك .

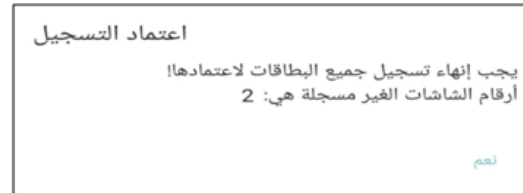


- اضغط على زر "نعم" لاعتماد تسجيلك و موافقتك على رفع ما قمت بتسجيله ، ويتم بعدها ليا تسجيل خروجك من التطبيق و العودة إلى شاشة الدخول، ولن تستطيع لاحقاً الدخول للنظام .

- اضغط على زر "لا" في حالة ما اذا كنت تريد سماع او إعادة تسجيل اي نص، يمكنك بعد الضغط الرجوع للشاشات السابقة وذلك باستخدام زر "السابق" للتنقل بين النصوص.

- **ملاحظة:** لن يتم اعتماد تسجيلاتك الا بالضغط على زر "نعم" في الرسالة المبتدئة عند ضغطك على زر اعتماد .

- في حالة ما اذا تم الضغط على زر "اعتماد" وكان هناك شاشات لم يتم تسجيلها ستظهر لك رسالة مبنقة لإبلاغك بأرقام الشاشات التي لم يتم بتسجيلها كما في المثال التالي:



- اضغط على زر "نعم" حيث سيتم إبقاءك في الشاشة الأخيرة. يمكنك الرجوع للشاشات التي لم يتم بتسجيلها وذلك باستخدام زر "السابق" للتنقل بين النصوص.

- بعد الانتهاء من تسجيل نصوص الشاشات المطلوبة انتقل إلى الشاشة الأخيرة واضغط على زر "اعتماد" ، لاعتماد تسجيلاتك.