

# A Multi-Purpose Audio-Visual Corpus for Multi-Modal Persian Speech Recognition: the Arman-AV Dataset

Javad Peymanfard<sup>\*</sup>, Samin Heydarian<sup>\*§</sup>, Ali Lashini<sup>\*§</sup>, Hossein Zeinali<sup>†</sup>,  
Mohammad Reza Mohammadi<sup>\*</sup>, and Nasser Mozayani<sup>\*</sup>

<sup>\*</sup> School of Computer Engineering

Iran University of Science and Technology, Tehran, Iran

Email: javad\_peymanfard@comp.iust.ac.ir, samin\_heydarian@comp.iust.ac.ir  
a\_lashini@comp.iust.ac.ir, mrmohammadi@iust.ac.ir  
mozayani@iust.ac.ir

<sup>†</sup> Department of Computer Engineering

Amirkabir University of Technology, Tehran, Iran

Email: hzeinali@aut.ac.ir

**Abstract**—In recent years, significant progress has been made in automatic lip reading. But these methods require large-scale datasets that do not exist for many low-resource languages. In this paper, we have presented a new multipurpose audio-visual dataset for Persian. This dataset consists of almost 220 hours of videos with 1760 corresponding speakers. In addition to lip reading, the dataset is suitable for automatic speech recognition, audio-visual speech recognition, and speaker recognition. Also, it is the first large-scale lip reading dataset in Persian. A baseline method was provided for each mentioned task. In addition, we have proposed a technique to detect visemes (a visual equivalent of a phoneme) in Persian. The visemes obtained by this method increase the accuracy of the lip reading task by 7% relatively compared to the previously proposed visemes, which can be applied to other languages as well.

**Index Terms**—persian dataset, audio-visual speech recognition, lip reading, viseme

## I. INTRODUCTION

Automatic speech recognition (ASR) is a task to understand speech from audio signals. This task has been developed over the years and got mature enough to be used on any device. But the trained models, apart from the high ability of speech recognition, have weaknesses in special conditions like environments with loud noises. That's where audio-visual speech recognition (AVSR) comes in to overcome this limitation. AVSR uses visual information alongside audio in order to decode speech more effectively in noisy environments like inside cars. AVSR is more like human comprehension, which uses visual and audio perception to understand each other. Lip reading is another task which is only consuming visual information.

The traditional approaches [1]–[3] used a two-stage algorithm to deal with such problems. In the first stage, a hand-crafted feature extractor, extract useful features from lip movements

alongside another feature extractor that extracts feature from audio signals and then fuse them. In the second stage, a classifier such as the hidden Markov model or artificial neural networks is used to classify digits, characters, etc. However, in the last few years, the availability of large public datasets and emerge of deep neural networks had a substantial impact on this field. Deep learning approaches usually consist of two parts, which are front-end and back-end, like traditional ones except that here they are end-to-end trainable. For front-end part usually use convolution neural networks to extract visual and audio features and temporal networks like RNNs, attention model, and transformers on the back-end side to model temporal information. Recently, methods such as [4]–[6] used knowledge distillation to train lip reading and AVSR models. To do this, usually use the ASR model as the teacher and the lip reading model as the student.

The primary dataset was collected under laboratory conditions [7], [8] - Usually, a person stands in front of the camera and read some words or sentences in a quiet place at normal speed. Emerge of deep learning and automatic pipeline led to building datasets in the "wild" condition [9]–[11] which is larger and also more challenging and close to the real condition than before. Unlike the old datasets which were restricted to digit, character, or phrase classification, nowadays, datasets are collected for word-level classification and sentence-level AVSR. Unfortunately, most of these datasets are in English and there are a few datasets which come in other languages [12]. In this paper, we will more concentrate on sentence-level AVSR, but we did not limit our dataset to just speech recognition tasks. We make our dataset in a way that can be used for different tasks. In the following sections, we will explain more about it.

Currently, LRS2 [13] is the most widely used audio-visual dataset. The dataset contains more than 240 hours of videos and about 118,000 utterances. Also, they used BBC news

<sup>§</sup>These authors contributed equally.

and talk shows as sources for their dataset. The dataset is publicly available to the research community. LSVSR [11] is the largest dataset with over 3800 hours of data collected from YouTube-uploaded videos. This dataset was collected by Google DeepMind and, unfortunately, they did not make it publicly available.

In this paper we propose a novel large-scale dataset that is collected in the "wild" condition from reviews, movies, etc. on the Aparat website. The dataset contains over 220 hours of videos from 1760 Persian celebrities. Also, we provide the name of the celebrity in each sample as labels to be used for different purposes, such as speaker recognition. To the best of our knowledge, this is the largest audio-visual dataset in the Persian language.

Since this is an audio-visual dataset, it could be used for a number of different applications such as automatic speech recognition, lip reading, speaker recognition, audio-visual speech synthesis, etc.

The organization of this paper is as follows. In section 2, we look at some of the most recent approaches and datasets. In section 3, we discuss the statistic of the data and the pipeline that we build for collecting data. In section 4, we will use a base model to evaluate our dataset. Finally, conclude the paper in section 5.

## II. RELATED WORKS

In this section, we will review the two tasks of speech recognition and speaker recognition. Since we have introduced a new dataset and tested the performance of the baseline models on it, we will explore the related works from the perspective of the proposed methods and datasets.

### A. Methods

1) *Speech Recognition: Lip reading.* There is a large body of research on automated lip reading. The advent of deep learning and the availability of large-scale datasets cause massive progress in the automation of this field. Here we discuss the works which focused on the sentence-level task in lip reading as an open set character. For sentence-level recognition, the recent works can be divided into approaches. The first approach utilizes Connectionist Temporal Classification (CTC). The model predicts frame-wise characters and tries to maximize and marginalizes all possible paths to find optimal alignment between prediction and ground truth. An example based on this approach is LipNet [28]. LipNet was the first sentence-level lip reading model that used spatiotemporal CNNs as the front-end, followed by Bidirectional Gated Recurrent Unit (Bi-GRU) as the back-end of architecture, and employed CTC loss to train the network. Shillingford *et al.* [11] proposed the V2P model closest to previous work. The model outputs a sequence of phoneme distributions and uses CTC loss at train time. At inference time, a decoder based on finite-state transducers (FSTs) maps a sequence of phoneme distributions to a word sequence.

The second approach follows the sequence-to-sequence model, in which output characters are conditioned based on

each other, unlike CTC-based models. WAS [9] was the first sequence-to-sequence model. The model consists of two key components called the image encoder Watch and the character decoder Spell and has a unique dual attention mechanism. Chung *et al.* [15] extended the WAS model to the MV-WAS model that can decode visual sequences across all poses and show that it is possible to read lips in profile, but the standard is inferior to reading frontal faces. Peymanfard *et al.* [29] suggested external viseme decoding, which divides the sequence-to-sequence model into two stages, video to viseme and viseme to character, respectively. The network outputs character sequence given viseme through external text data. Afouras *et al.* [30] proposed a visual transformer pooling (VTP) module that learns to track and aggregate the lip movements representation. This work also utilizes a wordPiece tokenizer (one of the subword algorithms) to learn the language easily, resulting in time and memory efficiency

Petridis *et al.* [31] presented a hybrid architecture and used a joint decoder including RNN, attention, and CTC. Afouras *et al.* [32] took a hybrid approach as well but suggested replacing RNN with a transformer. This work proposed a sequence-to-sequence and CTC on top of the transformer self-attention as the back-end. Both models use spatiotemporal CNNs, followed by ResNet as a common front-end. The extension of [31] is [33] that employed the conformer variant of the transformer to model both local and global dependencies of an image or audio sequence.

There is a trend in lip reading that benefits from the knowledge distillation area to improve performance by extracting information from the audio-only counterpart of datasets. In this approach, the Automatic Speech Recognition (ASR) model plays a role as the teacher for the Visual Speech Recognition (VSR) model as the student. The examples of this trend are [4], [5] and [6].

**Automatic Speech Recognition.** In ASR, self-supervision plays the main role in the state-of-the-art models. This approach has two steps, pre-training and fine-tuning. In the pre-training step, the model learns powerful speech representation from large amounts of unlabeled data through a pretext task. Then, the output of the previous step (speech representation) is fed to an acoustic model to be fine-tuned for a downstream task. In wav2vec [34], the pre-training step has two networks named encoder and context. The encoder network converts raw audio input to feature representations, and the context network converts feature representations to contextualized features. Both networks are convolutional networks. The objective of this model is defined by contrastive loss between future steps and distractors. The vq-wav2vec [35] is the improved version of the previous work, which benefited from the BERT model. This work is also composed of two stages. In the first stage, the continuous audio feature representations are discretized via Gumbel-Softmax or k-means clustering quantization methods. Thanks to this, the speech data takes a structure similar to language, and the BERT model can be applied. So in the second stage, the discretized representations are fed to a BERT model for contextualized representations. The wav2vec 2.0 [36]

TABLE I: Statistical comparison between well-known datasets.

Category	Dataset	General info							Transcription info		Video info				Speaker info		
		Year	Lang.	Task	Classes	Utter.	Availability	Source	Vocab.	Words	Num.	Duration	Resolution	FPS		View(°)	Num.
Audio-Visual Speech Recognition	SFAVD [14]	2013	FA	Sent.	-	600	-	Lab environment	1000	-	-	-	131×105	30	Frontal	1	
	LRW [9]	2016	ENG	Words	500	400K	Avail.* <sup>1</sup>	BBC Programmes	500	400K	-	-	256×256	25	-30 30	1000+	
	LRS2 [13]	2017	ENG	Sent.	-	118K	Avail.*	BBC Programmes	174K	807K	-	246h	160×160	25	-30 30	1000+	
	MV-LRS [15]	2017	ENG	Sent.	-	74K	Avail.*	Programs	14K	-	-	165h	160×160	25	-90 90	1000+	
	VLRf [16]	2017	SPA	Sent.	-	600	Public	Lab environment	1374	10K	-	180min	1280×720	50	Frontal	24	
	LRS3-TED [10]	2018	ENG	Sent.	-	165K	Avail.*	TED Talks	57K	-	-	475h	224×224	25	-90 90	1000+	
	LSVSR [11]	2018	ENG	Sent.	-	2,934K	Private	YouTube	127K	-	-	3,886h	128×128	30	-30 30	1000+	
	LRW-1000 [17]	2019	CHI	Words	1000	718K	Avail.*	Broadcast news	1K	-	-	-	Distributed	25	-90 90	2000+	
	CMLR [12]	2019	CHI	Words	-	102K	Public	News Broadcast	-	-	-	-	-	-	-	-	11
	LRWR [18]	2021	RUS	Words	235	-	-	YouTube	235	117K	-	-	112×112	25	0 20	135	
	GLips [19]	2022	GER	Words	500	-	Public	Hessian Parliament	500	250K	250K	-	256×256	25	-	100	
	RUSAVIC [20]	2022	RUS	-	62	62	Avail.*	Vehicle environment	-	-	200	-	1920×1080	60	-30 30	20	
	PLRW [21]	2022	FA	Words	500	244K	Public	Aparat	500	-	244K	30h	224×224	25	-	1800	
	Our Dataset	2022	FA	Sent.	-	89K	Public	Aparat	42K	2.5M	89K	220h	224×224	25	-	1760	
Active Speaker Detection	AVA [22]	2019	-	-	3	-	Public	YouTube	-	-	188	38.5h	128×128	-	-	-	
	ASW [23]	2021	-	-	2	-	-	VoxConverse dataset	-	-	212	30.9h	-	-	-	-	
Speaker Recognition	Voxceleb [24]	2017	ENG	-	1251	153K	Public	YouTube	-	-	22K	352h	-	-	-	1,251	
	Voxceleb2 [25]	2018	ENG	-	6112	1,128K	Public	YouTube	-	-	150K	2442h	-	-	-	6,112	
	DeepMine [26]	2019	FA&ENG	-	1850	544K	Avail.*	Crowdsourcing	-	-	540	480h	-	-	-	1,850	
	CN-Celeb [27]	2019	CHI	-	1000	130K	Public	bilibili.com	-	-	-	273.73h	-	-	-	1,000	

<sup>1</sup> Available by contact.

showed better results by joint learning of the two-stages pre-training pipeline mentioned above. The model is pre-trained by new contrastive loss between contextualized output and quantized representation. The HuBERT [37] is very similar to wav2vec 2.0 in model architecture but different in the training process. First, HuBERT uses the cross-entropy loss for model optimization. Second, Data discretization is done through a k-means algorithm instead of a quantization process. Third, The training process alternates between the hidden units discovery and target prediction. The model re-uses embeddings from the BERT encoder in the clustering step.

**Audio-Visual Speech Recognition.** The task of audio-visual speech recognition is lip reading with the presence of audio. Shi *et al.* [38] presented the audio-visual counterpart of HuBERT model as an AV-HuBERT. Also, this work reported the result of visual HuBERT performance for the lip reading task. The works like [11], [31] and [32] addressed both the VSR and the AVSR problems, which are later explained in the lip reading section.

2) *Speaker Recognition:* Since this is a classification problem, the general models consist of two modules. The first is a feature extractor module and the second is a classification module. Depending on the input data, the feature extractor could be a convolutional network like VGG-M or ResNet [25] or any other neural network that could be used for sequential data. For classification, a fully connected network is the leading choice, and the output size is the number of speakers in the dataset. Since the input data is an audio file, input data to the network could be a spectrogram [25] of data or the output of a feature extraction algorithm like MFCC, etc., which usually is a 1D feature vector of data.

## B. Datasets

1) *Speech Recognition:* As we know, the progress in deep learning is due to large datasets. The performance of lip reading systems is affected by position, illumination and speaker diversity, so the quality and quantity of the dataset are important. In recent years some datasets like LRS2 [13], LRS3-TED [10] and LSVSR [11] have tried to increase the number of utterances and video hours and MV-LRS [15] has tried to cover more

face angles. Table I shows three categories of datasets named Audio-Visual Speech Recognition, Active Speaker Detection and Speaker recognition. The table has general, transcription, video and speaker information columns for each dataset. We can see useful statistics in this information, such as the type of task run on the dataset, the source of the collected videos and the number of speakers in the videos. The language of most of the datasets in this table is English, but fortunately, other languages have been published in recent years, like CMLR [12], LRWR [18] and GLips [19]. Since 2013, no Persian dataset has been published for lip reading until now; we released Audio-Visual Speech Recognition in Persian (Arman\_AV) with 89k utterances, 220 hours and 1760 speakers.

2) *Speaker Recognition:* Most of the datasets in this field are collected under laboratory conditions and are not available to be used freely for research purposes. One of the first free datasets which were collected under the "wild" conditions is the Speakers in the Wild (SITW) [39] dataset. This dataset contains multimedia content of 299 speakers with hand-annotated speech samples. VoxCeleb [24] is a large-scale dataset that contains speech samples of over 1,000 celebrities with more than 100,000 utterances, extracted from YouTube videos; this dataset is gender-balanced (45% of the speakers are female), along with speakers with various accents, racial backgrounds, and ages. The dataset's creators included information about each speaker's gender and country of origin from Wikipedia. Videos contained in the dataset are recorded in a huge number of challenging environments such as quiet and noisy studio interviews, open-air stadiums, etc., and all are debased with real-world commotion, consisting of giggling, covering discourse, foundation chatter, etc. Also, another goal of the authors of the dataset was to propose a pipeline to create fully automated datasets. VoxCeleb 2 [25] which came shortly after VoxCeleb (and is very similar but much larger) is the largest dataset in this field which includes more than six thousand speakers that speak more than a million utterances. The three mentioned datasets are in English. CN-Celeb [27] is another large-scale dataset which contains over 130,000 utterances and is very similar to the VoxCeleb dataset except in the three following aspects: first, CN-Celeb covers 11 genres of speech such as singing, entertainment, etc.,

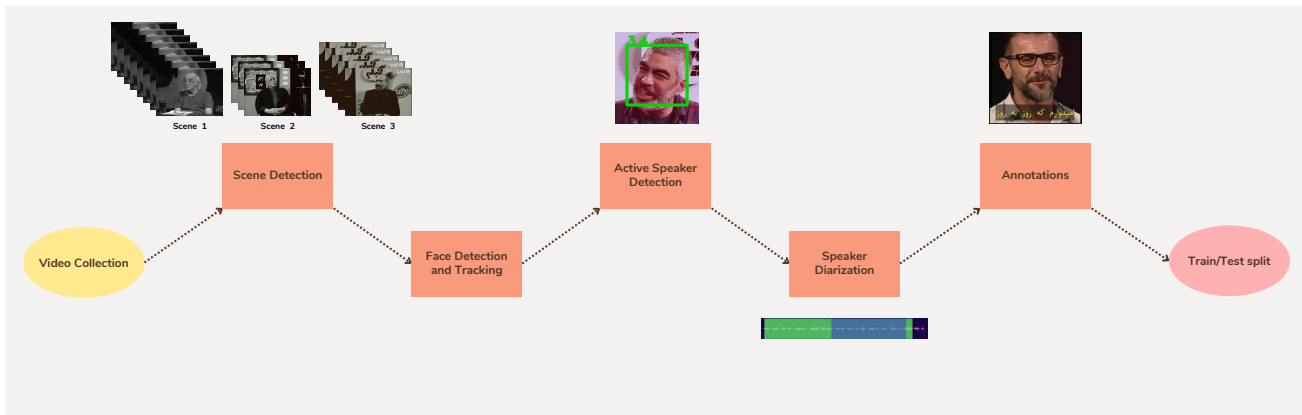


Fig. 1: Data collection pipeline.



Fig. 2: Our dataset samples.

which is more than VoxCeleb. Second, CN-Celeb concentrated on Chinese celebrities, containing videos of 1,000 celebrities. And last but not least is that the dataset is not automated completely, and they considered human supervision for the dataset. To the best of our knowledge, there is no suitable Persian dataset in this field.

### III. DATASET

#### A. Automated Pipeline

In this subsection, we see the steps of data collection. Figure 1 shows the overview of these steps.

**Step 1: Video Collection.** We pick three kinds of videos for our dataset described in the following segments.

- 1) **Interviews.** Interviews and biographies are well-suited for audio-visual datasets. The main goal of these videos is to speak to artists, politicians or famous people in general. These kinds of videos are divided into two groups. In the first group, the shows include a host and narrator (which is not desirable in this case), and the second group consists of those that don't have them.
- 2) **Series and Movies.** These kinds of videos can be used if suitable preprocessing is applied to them. The main challenges are dealing with different directions and angles of the camera (When a speaker is talking or is not in the shot all of the time, unlike in the interviews). Also, the

entire video has a significant amount of silence or music. Another challenge is multi-speaker simultaneous speech. In general, this type of video is not ideal for collecting data because the percentage of redundant data in them is relatively high. In the investigation we did, less than 10% of the input data is considered appropriate.

- 3) **Vods.** One of the other ideas to collect more data is different keyword searching in ugs. Then we review the resulting videos and choose the top keywords for the dataset.

To obtain data, multiple search terms are used

**Step 2: Scene Detection.** In this stage, we want to detect each scene. For this purpose, the PySceneDetect<sup>1</sup> is utilized. In this algorithm, we subtract the pixel values of two consecutive frames to detect different scenes. This difference value is calculated in HSV colour space.

**Step 3: Face Detection and Face Tracking.** Now each video is divided into smaller pieces called a scene. Up to this point, the video had a temporal clipping. Note that we need frames in which we see faces. So the S<sup>3</sup>FD model [40] and algorithm based on IoU [41] are used for face detection and face tracking respectively. Depending on the number of faces found in each part of the video, each video may be split into one or more videos or removed if there is no face.

<sup>1</sup><https://github.com/Breakthrough/PySceneDetect>

#### Step 4: Active Speaker Detection.

In this step, we automatically found the parts of the video that contain the speaker’s face.

**Step 5: Speaker Diarization.** As stated, up to stage 4, the video parts including a speaker with a constant shooting angle have been chosen. But there is another potential challenge in interviews, a situation where the host and the guest talk at the same time. To address this problem, we can use Speaker Diarization<sup>2</sup> based on UIS-RNN.

**Step 6: Annotations.** As we mentioned before, the source of our dataset videos is Aparat. The videos on this video hosting website lack subtitles, so we utilize the commercial Aipaa’s<sup>3</sup> ASR service to construct rough sentence-level transcripts of the videos.

**Step 7: Face Recognition and Dataset Split.** We use ArcFace [42] and produce face feature embeddings for face recognition. This stage aims to create a dataset with a speaker-independent property. A dataset like this is beneficial for applications like lip reading. It should be noted that the results of the face recognition algorithm were checked manually, and its errors were corrected.

#### B. statistics

First of all, the dataset will be available as mp4 face-cropped videos with 224\*224 pixels resolution and a frame rate of 25 FPS. We leveraged Aparat (a Persian video-sharing website like YouTube) as a source for our dataset. The collected videos came from different channels and programs, which means our dataset contains various kinds of challenging environments and circumstances. The outcome of our work comprises 220 hours of video, which contains over 89 thousand utterances and 2.5 million words of 1760 celebrities. We split our dataset into train/ validation, and test sections which each contain 211, 9 hours of data. Figure 2 shows some examples of the video samples. In table I, in addition to the statistics of our dataset, we have also presented the statistics of other recent datasets.

### IV. VISEME ANALYSIS

The viseme is a significant challenge in the lip reading problem. Visemes consist of a set of phonemes spoken by the same form of lips and aren’t exactly known for any languages. This mapping may even be different in various speakers. However, some research has addressed the issue in languages such as English and Persian. Some proposed models for said languages. In this study, we proposed a method to automatically identify the Persian visemes on large-scale data that one can apply to other languages by only having the appropriate dataset of the interested language. Some works have addressed the visemes for the Persian language [43]. Collecting a set of laboratory datasets and considering the similarities in speaking the various phonemes, the authors in [43] proposed a categorization method for phonemes to identify the visemes in the language. Combining the deep learning technique to automatically extract the features and clustering in this study,

<sup>2</sup><https://github.com/taylorlu/Speaker-Diarization>

<sup>3</sup><https://aipaa.ir/>

we introduced a method to categorize the phonemes. Then, we compared a viseme-based model for lip reading with these visemes and the visemes addressed in previous works. Employing the Kaldi tool, we first determined the phoneme level transcription for each sample in this experiment. Each phoneme spoken by the speaker is known every 30 milliseconds in this transcription. Then, we employed the av-Hubert pre-trained model to obtain the visual embedding for each phoneme, which was used as a feature for our clustering using the k-means algorithm. In this step, we expected the phonemes having the same visual features, like the same lip movements and forms, to be grouped in the same category. We used the method in [29] to train our model. Employing the two models of the video-to-viseme and the viseme-to-character models and finally merging the two, we obtained the lip reading model. One benefit of such a method is that one can use the textual data to improve the quality of the lip reading model. However, the model requires a proper grapheme to phoneme (G2P) model to convert the Persian text into the phonemes and then convert it into the visemes, for which we used an appropriate G2P. For the video-to-viseme model dataset, we used the phoneme level transcription to convert the phonemes into visemes with the help of its corresponding mapping. We also employed about 1 million utterances for the visemes into character conversion datasets. In this regard, we first employed a G2P model to obtain the phoneme sequence for each utterance and then used the phoneme-to-viseme mapping to extract the viseme sequence for each one. To train the viseme-to-character model, we used the viseme sequence as the input, whereas the Persian text was considered the desired output. Employing the attention mechanism to implement the video-to-viseme model, we used a seq2seq one having a two-layer GRU. Furthermore, we employed the same architecture with fewer parameters and a different input type to train our viseme-to-character model. As seen in table II, we obtained a higher precision in the lip reading problem by employing the new proposed model to identify the Persian visemes for which we believe this phoneme-to-viseme mapping to be more appropriate. Furthermore, one can use the method to identify visemes in any given language. All the data and their related transcriptions are publicly available.

### V. EXPERIMENTS

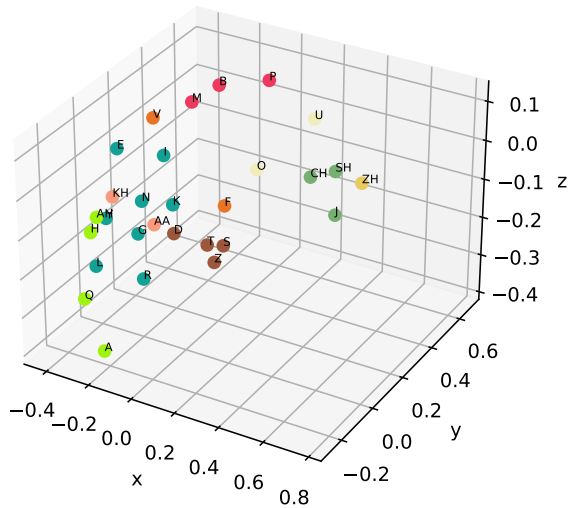
Here, the baseline methods are explained. We demonstrate how well they work with our dataset and provide information on training setup. The image sequence and audio sequence inputs are denoted as  $x_{1:T}^v$  and  $x_{1:T}^a$  in all baselines, respectively.

#### A. Automatic Speech Recognition (ASR)

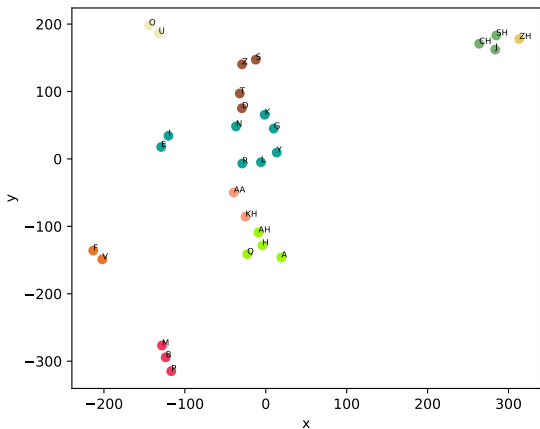
The present section reviews the experiments conducted on speech recognition using the produced datasets. Audio data alone was used in the first experiment to train Persian speech recognition. To this end, the pre-trained AV-HuBERT [44] model uses  $x_{1:T}^a$  to extract a 768-dimensional feature vector  $e_{1:T}^a$  for each 40ms of audio input. In this experiment, we used the pre-trained AV-HuBERT base model for feature extraction on the English language data. The output embedding in this

TABLE II: Lip reading results on our dataset

Viseme Mapping	CER (Greedy decoding)	WER (Greedy decoding)	CER (Beam search decoding)	WER (Beam search decoding)
Traditional [43]	%54.04	<b>%76.32</b>	%52.19	%73.24
AV-HuBERT + K-means (Ours)	<b>%51.24</b>	%76.71	<b>%48.54</b>	<b>%72.83</b>



(a)



(b)

Fig. 3: 2D and 3D projection of AV-HuBERT for Persian phonemes using our dataset. a) 3D projection of embedding of phonemes using PCA. b) 2D projection of embedding of phonemes using tSNE.

TABLE III: Persian visemes using clustering and AV-HuBERT

cluster id	phonemes
1	/F/ /V/
2	/AA/ /KH/
3	/B/ /P/ /M/
4	/ZH/
5	/CH/ /JH/ /SH/
6	/O/ /U/
7	/S/ /Z/ /D/ /T/
8	/A/ /' / /H/ /GH/
9	/G/ /K/ /E/ /L/ /I/ /Y/ /N/ /R/
10	/sil/

network for each window is a representation of that audio segment, including speaker and speech features. We expect the model to use speech-related features during the learning process. As such, we expect features corresponding to the speaker to be implicitly ignored, as they are irrelevant to our purpose. One further consideration during the experiment is whether or not a model pre-trained via the English language data can be used for the Persian language.

In the next step, the  $e_{1:T}^a$  vector was fed to a two-stacked Bidirectional LSTM, resulting in  $o_{1:T}^a$ . This network was trained using CTC. We assume  $y = (y_1, y_2, \dots, y_n)$  is a transcription and  $\pi$  are paths that function B will map  $\pi$  to  $y$ . We calculated the  $p_t^{CTC}$  probability over 52 characters by applying Softmax on  $o_{1:T}^a$ . As our CTC loss  $L_{CTC}$  decreases, the network begins to learn.

As evident in table IV, using the proposed architecture, we achieved a %84.66 character accuracy rate (CAR) using only 211 hours of Persian speech data, which is a suitable value, considering the data volume. Moreover, the experiment has demonstrated that the AV-HuBERT model, trained on data in English, could extract a proper embedding for speech data in Persian.

### B. Audio Visual Speech Recognition (AVSR)

During the second experiment, we trained an audio-visual speech recognition model. Image pre-processing was first used in this model to crop only the mouth ROI. Then the AV-HuBERT model was used per cropped video frame to obtain a 768-dimensional vector. In this section, an architecture similar to the previous stage was used, but we also gave the visual features as input to the model. Moreover, we used less output for visual features as they presumably contain less information than speech information. This model is presented in Figure

4. The input image sequence  $x_{1:T}^v$  is fed to the pre-trained Single-modal visual HuBERT to obtain a 768-dimensional feature vector  $e_{1:T}^v$ . The 768-dimensional feature vector  $e_{1:T}^a$  is obtained from input audio sequence  $x_{1:T}^a$  through the pre-trained Single-modal audio HuBERT concurrently. In the next step, the separated two-stacked Bi-LSTM are used to map  $e_{1:T}^v$  and  $e_{1:T}^a$  to  $o_{1:T}^v$  and  $o_{1:T}^a$  respectively. Then the two 1-dimensional convolutional filters are applied on  $c_t$  as a result of the concatenation of  $o_t^v$  and  $o_t^a$ . The output of this step is  $f_t$ . As mentioned above, CTC loss is used to train the model. First, we calculated the  $p_t^{CTC}$  probability. This probability is obtained by applying Softmax and linear projection on  $f_t$  (Where  $W \in \mathbb{R}^{52 \times (256+1024)}$  and  $b \in \mathbb{R}^{52}$ ). Then, the model will be optimized through  $L_{CTC}$  minimizing the value.

This model was trained with the exact mechanism as the previous one. As visible in table IV, through this model, which provides both audio and visual features, we have achieved higher accuracy than the ASR model. As visible in this experiment, embeddings obtained in the pre-trained AV-HuBERT model, which are trained on the English language data, can be used in the Farsi language, and relatively desirable features of lip movement can be extracted from it.

### C. Speaker recognition

Another problem we tested using the generated data was speaker recognition. For this, we first separated the people with more than five samples. Then we randomly selected five samples for each speaker. For each speaker, four samples were considered as a reference, and only one of the samples was used for testing. Then, using the TitaNet [45] pre-trained model, which was trained on 5 datasets (Voxceleb [24] and Voxceleb2 [25], NIST SRE portion of datasets from 2004-2008 (LDC2009E100), Switchboard-Cellular1 and Switchboard-Cellular2 [46], Fisher [47] and Librispeech [48]). We fed  $x_{1:T}^{a_t}$  as a test and  $x_{1:T}^{a_0}, \dots, x_{1:T}^{a_3}$  as references to the TitaNet [45] model. The model produces the 192-dimensional embeddings of each audio sample, for example,  $v_{1:T}^{a_t}$ . In the next step, the test embedding vector is compared with every reference of every person through cosine similarity (Here  $i \in 0, \dots, 3$ ). After that, we calculate the average with the previous step's results. Then, each sample in the test was assigned to the person who had the most similarity with that in terms of the average cosine similarity criterion. Using this method, we achieved an accuracy of 84.5.

## VI. CONCLUSION

Arman-AV is a large-scale multipurpose dataset for Persian that can be used in various tasks such as lip reading, automatic speech recognition, audio-visual speech recognition, and speaker recognition, and it is publicly available. The dataset consists of almost 220 hours of video data (about 89,000 samples) from 1760 speakers. Although the main goal of collecting this dataset was lip reading, the results of a baseline method were reported for each of the mentioned tasks. According to the obtained results, the character error rate of speech recognition was reduced by 14.7% relatively

using visual features. We also proposed a method for obtaining Persian visemes. By using these visemes, higher accuracy in lip reading was achieved. In addition, we have provided various analyses for the dataset such as the age and gender of speakers, and the synchronization of audio and video for each segment. All of this metadata is available along with the dataset. The dataset can be used for other tasks such as face recognition, face verification, and audio-visual speaker recognition, which we did not cover in this paper and can be considered for future work on the dataset.

## VII. ACKNOWLEDGEMENT

The project was mainly supported by Arman Rayan Sharif company, an AI company in Iran.

## REFERENCES

- [1] I. Matthews, T. F. Cootes, J. A. Bangham, S. Cox, and R. Harvey, "Extraction of visual features for lipreading," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 24, no. 2, pp. 198–213, 2002.
- [2] S. S. Morade and S. Patnaik, "A novel lip reading algorithm by using localized acm and hmm: Tested for digit recognition," *Optik*, vol. 125, no. 18, pp. 5181–5186, 2014. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0030402614005786>
- [3] Z. Zhou, G. Zhao, and M. Pietikäinen, "Towards a practical lipreading system," in *CVPR 2011*. IEEE, 2011, pp. 137–144.
- [4] Y. Zhao, R. Xu, X. Wang, P. Hou, H. Tang, and M. Song, "Hearing lips: Improving lip reading by distilling speech recognizers," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 04, 2020, pp. 6917–6924.
- [5] T. Afouras, J. S. Chung, and A. Zisserman, "Asr is all you need: Cross-modal distillation for lip reading," in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 2143–2147.
- [6] S. Ren, Y. Du, J. Lv, G. Han, and S. He, "Learning from the master: Distilling cross-modal advanced knowledge for lip reading," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 13 325–13 333.
- [7] M. Cooke, J. Barker, S. Cunningham, and X. Shao, "An audio-visual corpus for speech perception and automatic speech recognition," *The Journal of the Acoustical Society of America*, vol. 120, no. 5, pp. 2421–2424, 2006.
- [8] I. Anina, Z. Zhou, G. Zhao, and M. Pietikäinen, "Ouluvs2: A multi-view audiovisual database for non-rigid mouth motion analysis," in *2015 11th IEEE International Conference and Workshops on Automatic Face and Gesture Recognition (FG)*, vol. 1. IEEE, 2015, pp. 1–5.
- [9] J. S. Chung and A. Zisserman, "Lip reading in the wild," in *Asian conference on computer vision*. Springer, 2017, pp. 87–103.
- [10] T. Afouras, J. S. Chung, and A. Zisserman, "Lrs3-ted: a large-scale dataset for visual speech recognition," *arXiv preprint arXiv:1809.00496*, 2018.
- [11] B. Shillingford, Y. Assael, M. W. Hoffman, T. Paine, C. Hughes, U. Prabhu, H. Liao, H. Sak, K. Rao, L. Bennett, M. Mulville, M. Denil, B. Coppin, B. Laurie, A. Senior, and N. de Freitas, "Large-Scale Visual Speech Recognition," in *Proc. Interspeech 2019*, 2019, pp. 4135–4139.
- [12] Y. Zhao, R. Xu, and M. Song, "A cascade sequence-to-sequence model for chinese mandarin lip reading," in *Proceedings of the ACM Multimedia Asia*, 2019, pp. 1–6.
- [13] J. Son Chung, A. Senior, O. Vinyals, and A. Zisserman, "Lip reading sentences in the wild," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 6447–6456.
- [14] Z. Naraghi and M. Jamzad, "Sfavn: Sharif farsi audio visual database," in *The 5th Conference on Information and Knowledge Technology*. IEEE, 2013, pp. 417–421.
- [15] J. S. Chung and A. Zisserman, "Lip reading in profile," 2017.
- [16] A. Fernandez-Lopez, O. Martinez, and F. M. Sukno, "Towards estimating the upper bound of visual-speech recognition: The visual lip-reading feasibility database," in *2017 12th IEEE international conference on automatic face & gesture recognition (FG 2017)*. IEEE, 2017, pp. 208–215.

TABLE IV: Speech recognition results on our dataset

Method	CER (Greedy decoding)	WER (Greedy decoding)	CER (Beam search decoding)	WER (Beam search decoding)
ASR(CTC)	%15.77	%46.43	%15.34	%45.86
AVSR(CTC)	<b>%13.76</b>	<b>%44.98</b>	<b>%13.08</b>	<b>%43.43</b>

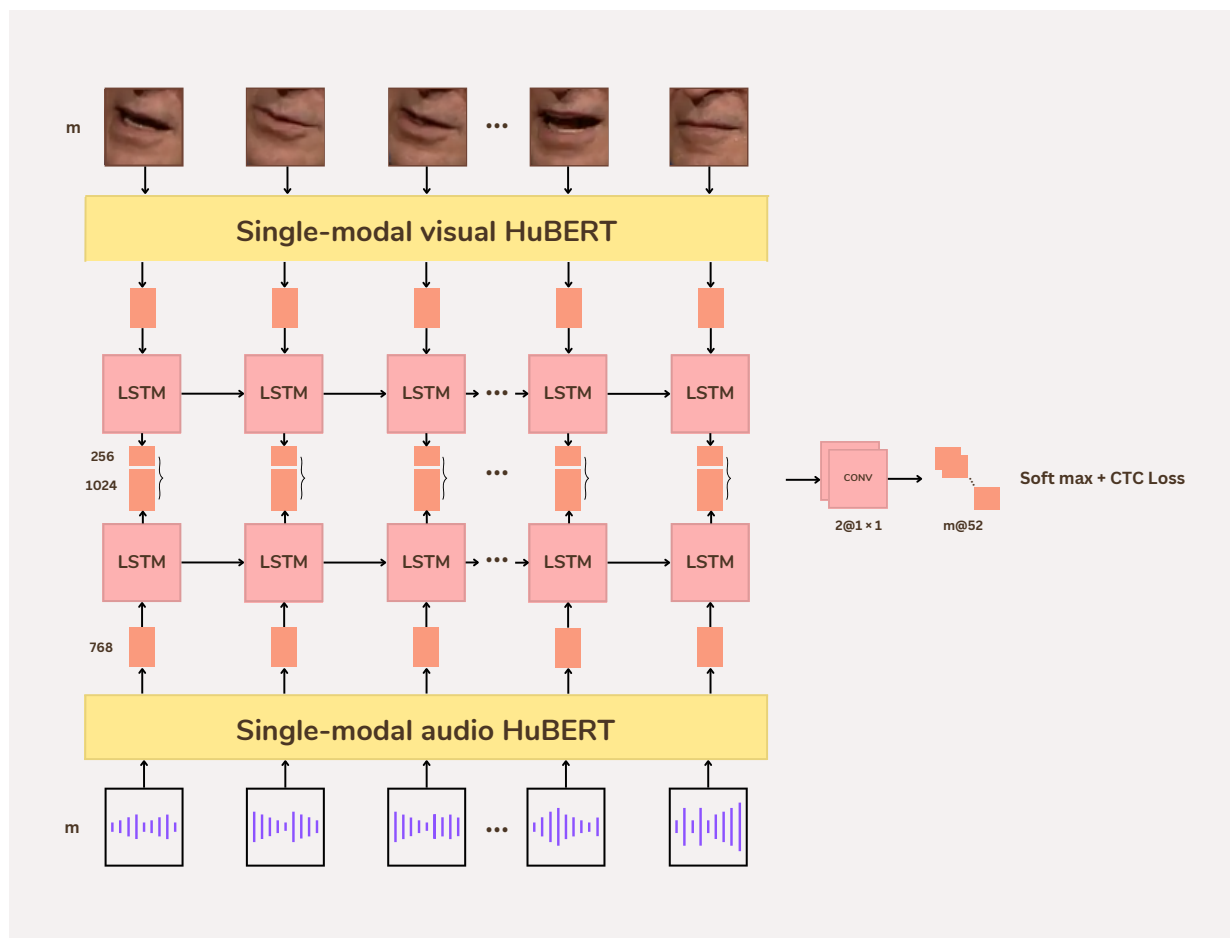


Fig. 4: The audio-visual architecture.

- [17] S. Yang, Y. Zhang, D. Feng, M. Yang, C. Wang, J. Xiao, K. Long, S. Shan, and X. Chen, "Lrw-1000: A naturally-distributed large-scale benchmark for lip reading in the wild," in *2019 14th IEEE international conference on automatic face & gesture recognition (FG 2019)*. IEEE, 2019, pp. 1–8.
- [18] E. Egorov, V. Kostyumov, M. Konyk, and S. Kolesnikov, "Lrwr: Large-scale benchmark for lip reading in russian language," *arXiv preprint arXiv:2109.06692*, 2021.
- [19] G. Schwiebert, C. Weber, L. Qu, H. Siqueira, and S. Wermter, "A multimodal german dataset for automatic lip reading systems and transfer learning," *arXiv preprint arXiv:2202.13403*, 2022.
- [20] D. Ivanko, A. Axyonov, D. Ryumin, A. Kashevnik, and A. Karpov, "Rusavic corpus: Russian audio-visual speech in cars," in *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, 2022, pp. 1555–1559.
- [21] J. Peymanfard, A. Lashini, S. Heydarian, H. Zeinali, and N. Mozayani, "Word-level persian lipreading dataset," in *2022 12th International Conference on Computer and Knowledge Engineering (ICCKE)*. IEEE, 2022, pp. 225–230.
- [22] J. Roth, S. Chaudhuri, O. Klejch, R. Marvin, A. Gallagher, L. Kaver, S. Ramaswamy, A. Stopczynski, C. Schmid, Z. Xi *et al.*, "Ava active speaker: An audio-visual dataset for active speaker detection," in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 4492–4496.
- [23] Y. J. Kim, H.-S. Heo, S. Choe, S.-W. Chung, Y. Kwon, B.-J. Lee, Y. Kwon, and J. S. Chung, "Look who's talking: Active speaker detection in the wild," *arXiv preprint arXiv:2108.07640*, 2021.
- [24] A. Nagrani, J. S. Chung, and A. Zisserman, "Voxceleb: a large-scale speaker identification dataset," in *INTERSPEECH*, 2017.
- [25] J. S. Chung, A. Nagrani, and A. Zisserman, "Voxceleb2: Deep speaker recognition," in *INTERSPEECH*, 2018.
- [26] H. Zeinali, H. Sameti, and T. Stafylakis, "Deepmine speech processing database: Text-dependent and independent speaker verification and speech recognition in persian and english," in *Odyssey*, 2018, pp. 386–392.
- [27] Y. Fan, J. Kang, L. Li, K. Li, H. Chen, S. Cheng, P. Zhang, Z. Zhou, Y. Cai, and D. Wang, "Cn-celeb: a challenging chinese speaker recognition dataset," in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 7604–7608.
- [28] Y. M. Assael, B. Shillingford, S. Whiteson, and N. de Freitas, "Lipnet: Sentence-level lipreading," *CoRR*, vol. abs/1611.01599, 2016. [Online]. Available: <http://arxiv.org/abs/1611.01599>
- [29] J. Peymanfard, M. R. Mohammadi, H. Zeinali, and N. Mozayani, "Lip reading using external viseme decoding," in *2022 International*



- Conference on Machine Vision and Image Processing (MVIP)*. IEEE, 2022, pp. 1–5.
- [30] T. Afouras, A. Zisserman *et al.*, “Sub-word level lip reading with visual attention,” *arXiv preprint arXiv:2110.07603*, 2021.
- [31] S. Petridis, T. Stafylakis, P. Ma, G. Tzimiropoulos, and M. Pantic, “Audio-visual speech recognition with a hybrid ctc/attention architecture,” in *2018 IEEE Spoken Language Technology Workshop (SLT)*. IEEE, 2018, pp. 513–520.
- [32] T. Afouras, J. S. Chung, A. Senior, O. Vinyals, and A. Zisserman, “Deep audio-visual speech recognition,” *IEEE transactions on pattern analysis and machine intelligence*, 2018.
- [33] P. Ma, S. Petridis, and M. Pantic, “End-to-end audio-visual speech recognition with conformers,” in *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2021, pp. 7613–7617.
- [34] S. Schneider, A. Baevski, R. Collobert, and M. Auli, “wav2vec: Unsupervised pre-training for speech recognition,” *arXiv preprint arXiv:1904.05862*, 2019.
- [35] A. Baevski, S. Schneider, and M. Auli, “vq-wav2vec: Self-supervised learning of discrete speech representations,” *arXiv preprint arXiv:1910.05453*, 2019.
- [36] A. Baevski, Y. Zhou, A. Mohamed, and M. Auli, “wav2vec 2.0: A framework for self-supervised learning of speech representations,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 12 449–12 460, 2020.
- [37] W.-N. Hsu, B. Bolte, Y.-H. H. Tsai, K. Lakhotia, R. Salakhutdinov, and A. Mohamed, “Hubert: Self-supervised speech representation learning by masked prediction of hidden units,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 29, pp. 3451–3460, 2021.
- [38] B. Shi, W.-N. Hsu, K. Lakhotia, and A. Mohamed, “Learning audio-visual speech representation by masked multimodal cluster prediction,” *arXiv preprint arXiv:2201.02184*, 2022.
- [39] M. McLaren, L. Ferrer, D. Castán, and A. D. Lawson, “The speakers in the wild (sitw) speaker recognition database,” in *INTERSPEECH*, 2016.
- [40] S. Zhang, X. Zhu, Z. Lei, H. Shi, X. Wang, and S. Z. Li, “S3fd: Single shot scale-invariant face detector,” in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 192–201.
- [41] R. Tao, Z. Pan, R. K. Das, X. Qian, M. Z. Shou, and H. Li, “Is someone speaking? exploring long-term temporal features for audio-visual active speaker detection,” in *Proceedings of the 29th ACM International Conference on Multimedia*, 2021, pp. 3927–3935.
- [42] J. Deng, J. Guo, N. Xue, and S. Zafeiriou, “Arcface: Additive angular margin loss for deep face recognition,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 4690–4699.
- [43] M. Aghaahmadi, M. M. Dehshibi, A. Bastanfard, and M. Fazlali, “Clustering persian viseme using phoneme subspace for developing visual speech application,” *Multimedia tools and applications*, vol. 65, no. 3, pp. 521–541, 2013.
- [44] B. Shi, W.-N. Hsu, K. Lakhotia, and A. Mohamed, “Learning audio-visual speech representation by masked multimodal cluster prediction,” *arXiv preprint arXiv:2201.02184*, 2022.
- [45] N. R. Koluguri, T. Park, and B. Ginsburg, “Titanet: Neural model for speaker representation with 1d depth-wise separable convolutions and global context,” in *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2022, pp. 8102–8106.
- [46] J. Godfrey and E. Holliman, “Switchboard-1 release 2 ldc97s62,” *Linguistic Data Consortium*, 1993.
- [47] C. Cieri, D. Graff, O. Kimball, D. Miller, and K. Walker, “Fisher english training speech part 1 transcripts,” *Philadelphia: Linguistic Data Consortium*, 2004.
- [48] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, “Librispeech: an asr corpus based on public domain audio books,” in *2015 IEEE international conference on acoustics, speech and signal processing (ICASSP)*. IEEE, 2015, pp. 5206–5210.