

# AKVSR: Audio Knowledge Empowered Visual Speech Recognition by Compressing Audio Knowledge of a Pretrained Model

Jeong Hun Yeo, Minsu Kim, Jeongsoo Choi, Dae Hoe Kim, and Yong Man Ro, *Senior Member, IEEE*

**Abstract**—Visual Speech Recognition (VSR) is the task of predicting spoken words from silent lip movements. VSR is regarded as a challenging task because of the insufficient information on lip movements. In this paper, we propose an Audio Knowledge empowered Visual Speech Recognition framework (AKVSR) to complement the insufficient speech information of visual modality by using audio modality. Different from the previous methods, the proposed AKVSR 1) utilizes rich audio knowledge encoded by a large-scale pretrained audio model, 2) saves the linguistic information of audio knowledge in compact audio memory by discarding the non-linguistic information from the audio through quantization, and 3) includes Audio Bridging Module which can find the best-matched audio features from the compact audio memory, which makes our training possible without audio inputs, once after the compact audio memory is composed. We validate the effectiveness of the proposed method through extensive experiments, and achieve new state-of-the-art performances on the widely-used LRS3 dataset.

**Index Terms**—Audio Knowledge via memory, Audio Knowledge Quantization, Audio Empowered Visual Speech Recognition, Audio Pretrained Model, VSR

## I. INTRODUCTION

VISUAL Speech Recognition (VSR) is a task of predicting speech content from lip movement without sound. VSR has received a lot of attention due to its practical applications. It can be used as a subtitling tool for silent movies, an auxiliary tool for speech recognition in noisy environments, and a conversation tool for the hearing impaired.

VSR has significantly improved in its performance along with the development of Deep Learning [1]–[9]. Many efforts have been made to improve the network architecture of the VSR systems. A visual encoder based on the combination of a 3D convolution layer and a 2D Convolutional Neural Network (CNN) is suggested by [10] to encode spatio-temporal visual features from lip movements. To capture the context information from the encoded visual features, prior works [10]–[12] adopted Recurrent Neural Network (RNN) [13] after the visual encoder. Recently, inspired by the success of the Transformer [2] in Natural Language Processing (NLP), the VSR model

augmented with the Transformer achieved significant speech recognition performances [14]–[17]. Unlike the RNN, the self-attention mechanism of the Transformer enables it to capture dependencies between any two positions in the lip sequence, facilitating a more comprehensive understanding of linguistic content from lip movement. Despite the development of VSR architectures, VSR is still regarded as a challenging task due to the characteristics of visual speech. Different from audio speech, visual speech inherently contains insufficient information to fully represent speech content, as speech is not only produced with the parts that are visible (*i.e.*, lips) but also with diverse internal human organs [18]. Hence, another research stream focuses on complementing insufficient visual information by augmenting the VSR model with additional information.

To complement the insufficient visual information, several prior works proposed to provide audio knowledge into the VSR model. Knowledge Distillation (KD) [19] is one of the most popular schemes for transferring superior knowledge of a teacher model to a student model. [20]–[25] tried to transfer the audio knowledge of the teacher model into the visual student model. These approaches supervised the student model to follow the soft-label or audio features generated from the teacher model. However, because of the differences in inherent properties between audio and visual modalities called heterogeneity gap [26]–[28], some knowledge can be discarded during knowledge distillation [21]. To bypass the heterogeneity gap, [29]–[31] proposed audio-visual multimodal bridging frameworks based on a memory network [32]. They built a visual-to-audio mapping function using a visual key memory and an audio value memory. Through the learned mapping function, the VSR model can utilize the saved audio knowledge. All of the aforementioned methods showed that the VSR systems can better model speech by complementing visual information with audio knowledge. Nevertheless these successes, the previous methods utilizing audio [20]–[25], [29]–[31] do not focus only on transferring linguistic information of audio. For example, prior works [20]–[25] utilized Knowledge Distillation (KD) to make the visual feature to be close to the audio feature without considering the characteristics of audio. An aim of the other works based on memory [29]–[31] also save audio features and reconstruct audio features using visual features without focusing on linguistic information. However, the audio contains not only linguistic information but also contains diverse information such as speaker characteristics, background noises, etc. If we do not consider these diverse

J. H. Yeo, M. Kim, J. Choi, and Y. M. Ro are with the Image and Video Systems Laboratory, School of Electrical Engineering, Korea Advanced Institute of Science and Technology (KAIST), Republic of Korea (e-mail: sedne246@kaist.ac.kr; ms.k@kaist.ac.kr; jeongsoo.choi@kaist.ac.kr; ymro@kaist.ac.kr)

D. H. Kim is with the Visual Intelligence Research Section, Superintelligence Creative Research Laboratory, Electronics and Telecommunications Research Institute (ETRI), Republic of Korea (e-mail: dhkim19@etri.re.kr). Corresponding author: Y. M. Ro (fax: 82-42-350-5494)

factors of audio when using it for VSR training, the complementary effects of audio can be degraded.

Recently, Large-scale pretrained models from raw audio data such as wav2vec2.0 [33] and Hidden-Unit BERT (HuBERT) [34] achieved significant performance improvement in Audio-based Automatic Speech Recognition (ASR) [35]–[39]. Especially, HuBERT which is pretrained by masked prediction like BERT [40], [41] to capture context information from unmasked audio features achieved state-of-the-art performance in ASR after finetuned on paired audio-text dataset. Motivated by the recent success of 1) complementing visual modality with audio modality through audio memory in VSR and 2) self-supervised pretraining in ASR, we try to empower VSR models with audio knowledge extracted from a pretrained model.

In this paper, we propose a novel Audio Knowledge empowered Visual Speech Recognition framework (AKVSR), where the audio knowledge of a large-scale audio pretrained model is extracted with compact representation discarding non-linguistic factors like speaker and noise, and utilized to empower the VSR model. Different from previous approaches complementing visual modality with audio modality [20]–[25], [29]–[31], the proposed method is the first work to adopt the large-scale pretrained audio model in VSR and transfer the linguistic information of audio by considering the properties of audio modality. We would like to highlight that directly utilizing the large-scale pretrained audio model for VSR without considering the characteristics lying in audio such as speaker, and noise may reduce the beneficial effects of using the large-scale pretrained model in VSR. Therefore, we complement the insufficient visual information by using only linguistic information of audio except for other characteristics of audio (i.e., speaker characteristics and noise). To achieve this, the audio knowledge of pretrained HuBERT is vector quantized [42]–[44] with a fixed size of clusters which are learned to contain only linguistic information through ASR.R. Therefore, the most representative knowledge in predicting speech can be extracted. The extracted audio knowledge of the pretrained audio model composes the knowledge in compact audio memory, which will be incorporated in VSR models. To employ the audio knowledge in training VSR models without input audio, we build an Audio Bridging Module (ABM). ABM is for finding the best-matched audio knowledge with input visual representation from the memory through cross-modal attention. Then, the best-matched audio knowledge is added with encoded visual features to empower speech representations.

Our proposed method has three differences compared to the existing methods based on memory networks [29]–[31]. 1) We utilize rich audio knowledge encoded by a large-scale pretrained audio model and transform the audio knowledge into a compact representation to store only the linguistic information (e.g., phoneme content). 2) The proposed method does not require additional audio input and audio model to supplement insufficient visual representation due to the existence of the ABM when we train the VSR model in contrast to the existing VSR methods. 3) The compact audio memory can be utilized to furnish audio information to any VSR model

as the representative knowledge needed for speech prediction, such as phoneme information, does not vary between datasets in speech recognition tasks.

In summary, our key contributions are as follows:

- We introduce a novel Audio Knowledge empowered Visual Speech Recognition (AKVSR) framework. To the best of our knowledge, this is the first work to refine non-linguistic factors of audio and transfer the knowledge of a large-scale pretrained audio model into the VSR model.
- We do not need an additional audio model when we train the VSR model utilizing audio knowledge.
- We validate through ASR that the representative knowledge of predicting speech is stored in compact audio memory. Moreover, we verify that the compact audio memory can be adapted to any VSR model.
- The proposed AKVSR outperforms the current state-of-the-art VSR model on the most popular sentence-level LRS3 dataset.

## II. RELATED WORKS

### A. Visual Speech Recognition

With the great development of deep learning, many research contributions have been made to VSR, especially in terms of architecture and data. Chung *et al.* [45] proposed an English word-level VSR data, LRW, and proposed a VGG-based VSR model. Stafylakis *et al.* [10] improved the architecture of the VSR model by using ResNet-34 [46] with one 3D convolution layer and Bi-LSTM. Some works [11], [47] proposed an end-to-end Audio-Visual Speech Recognition (AVSR) model, and Petridis *et al.* [11] set a strong baseline in word-level VSR. Some works [48], [49] tried to capture the lip movements in detail by using two-stream networks which utilize both RGB frames and optical flows. Zhao *et al.* [12] introduced mutual information maximization-based method to enhance the relations of the features with the speech content. Zhang *et al.* [50] proved that using face region instead of using lip region only, is beneficial to VSR. Martinez *et al.* [51] improved the temporal encoding of the back-end by proposing Multi-Scale Temporal Convolutional Network (MS-TCN). Ma *et al.* [52] proposed a distillation-based method of [53] in VSR. They repeatedly trained new models through born-again distillation, where the trained model becomes the new teacher. With the distillation, the VSR model can be lightened without loss of performance. Kim *et al.* [54] explored speaker dependency of pretrained VSR models and proposed a speaker adaptation method. For the sentence-level VSR, Assael *et al.* [55] proposed an end-to-end VSR framework using Connectionist Temporal Classification (CTC) [56]. Chung *et al.* [57] improved it to unconstrained sentence-level VSR by proposing LRS2 dataset and sequence-to-sequence architecture [3]. Recently, Transformer-based [2] architectures became the basics for visual speech modeling as they achieved significant VSR performances [14]–[16], [58], [59]. The transformer-based encoder-decoder structure [2] that utilizes attention mechanisms enables the handling of input and output sequences with varying lengths. This adaptable characteristic empowers the model to accurately transcribe and generate variable-length

outputs in the domain of Visual Speech Recognition (VSR). The proposed method also exploits this useful characteristic to effectively predict different output lengths regardless of the input lengths.

In this paper, we try to improve VSR systems by complementing the limited information of lip movements by proposing a compact audio memory, instead of improving the network architecture. In the next section, we will delve into the recent advancements in incorporating audio information to enhance visual information in the field of VSR.

### B. Complementing Visual using Audio in VSR

There are other efforts trying to augment the VSR model with audio modal knowledge. Afouras *et al.* [22] proposed a method of utilizing a large number of unlabeled audio data. By distilling the predicted logits of a pretrained ASR model into the VSR model, the VSR model can be learned from large-scale unlabelled audio-visual data. On a similar line, [20], [21] proposed knowledge distillation methods [19] by using pretrained ASR models from large-scale audio corpus datasets. By guiding the VSR model to follow the audio features at different levels encoded from the ASR model, the VSR model is expected to extract more discriminative visual features. Another research stream is utilizing memory network [32], [60] for saving audio knowledge, Kim *et al.* [29], [61] proposed Visual-Audio Memory which can save the audio features during training and read the saved audio knowledge from the learned memory with just visual inputs during inference. They improved the memory network to be able to consider the one-to-many mapping of viseme-to-phoneme by proposing multi-head memory architectures [30].

Existing methods for VSR utilize the audio modality to complement the visual modality. In this perspective, the KD [20]–[22] and memory-based [29], [31], [61] approaches improve the VSR system through multi-modal learning. However, the audio modality contains many characteristics such as speaker, noise, and linguistic content. Transferring audio knowledge without considering non-linguistic factors such as speaker and noise may reduce the complementing effect of the VSR system. Different from the previous approaches, this is the first work to utilize a large-scale pretrained audio model and transfer the audio knowledge to a VSR system considering the linguistic factor of audio modality. Namely, we aim to transfer the audio knowledge focused on only the linguistic factor.

### C. Pretraining on Large-scale Databases

Pretraining neural networks (*e.g.*, BERT [40]) on large-scale datasets achieved significant performances when the pretrained model is adapted to the downstream tasks, in diverse research areas [62]–[69].

It has also achieved promising results in Automatic Speech Recognition (ASR). Prior works, wav2vec2.0 [33] and HuBERT [34], proposed to learn the speech representation from raw audio in a self-supervised manner. Since the methods do not need text annotations, they can be trained via large-scale audio databases. In visual speech modeling, [70], [71]

proposed self-supervised pretraining methods using audio-visual correspondences. They showed that finetuning the pre-trained model on the VSR task can achieve better performance than learning the VSR model from scratch. Recently, AV-HuBERT [17] which proposed to pretrain the model with masked predictions using audio-visual databases achieved state-of-the-art performance and showed the powerful speech representation power of the model. Moreover, Zhang *et al.* [24] produces online target features by self-distillation during masked prediction training and then reduces the training cost of self-supervised speech representation learning.

The advantage of utilizing a large-scale pretrained audio model for complementing a visual using audio is that we can acquire improved quality audio knowledge. However, audio knowledge generated by a large-scale pretrained audio model contains a wide range of information, including speech content, speaker characteristics, and noise, as stated in [44], [72]. In this paper, our focus is to store only linguistic audio knowledge in the compact audio memory. Moreover, we inject the audio knowledge from the memory into the visual modality through ABM when we train the VSR model without the audio model.

## III. METHODS

### A. Overview

The Audio Knowledge empowered Visual Speech Recognition (AKVSR) framework is proposed to improve the complementary effect of audio information in visual speech recognition. Unlike traditional methods, AKVSR obtains improved audio knowledge by adopting a large-scale pretrained audio model and extracts linguistic information by eliminating non-linguistic factors, such as speaker information, through vector quantization. We then store the obtained linguistic information in compact audio memory. This allows us to provide linguistic information of audio knowledge to the VSR model without the additional audio model and inputs. The Audio Bridging Module (ABM) is employed to find the most appropriate audio knowledge from the compact audio memory matched with the visual feature. Therefore, it is possible to complement insufficient visual information by injecting the found linguistic information of the audio knowledge into the VSR model. In the following subsections, we will provide details of each proposed method.

### B. Compressing Linguistic Audio knowledge into Compact Audio Memory

Recent studies [44], [73], [74] report that the vector quantized self-supervised speech representation can disentangle linguistic content from speaker characteristics and noise. They show that the same linguistic content can be obtained at least the content is the same, even if the speaker is changed. The rationale behind this lies in the application of vector quantization to speech features at each time step. This process compels the resulting quantized speech feature to exhibit the most discriminative representations of self-supervised speech feature. Here, the most discriminative representation of self-supervised speech feature indicates the linguistic information

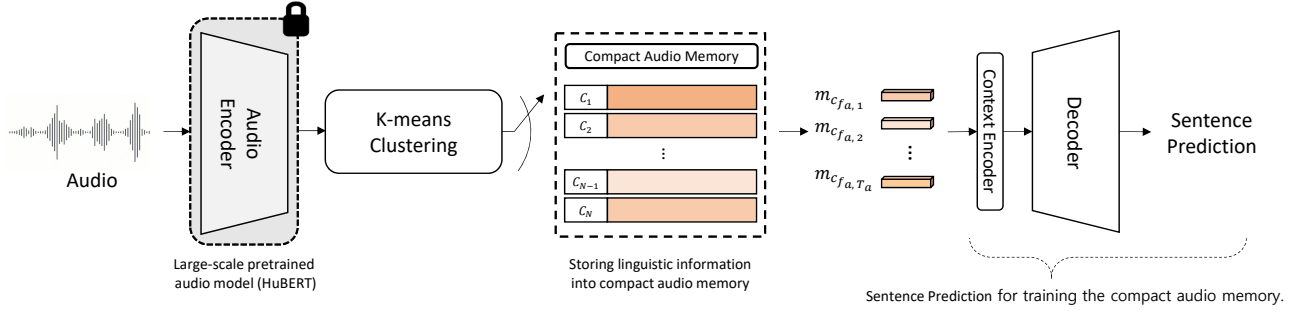


Fig. 1. Overview of building the compact audio memory to store linguistic information of large-scale pretrained audio model. The audio features are generated by the large-scale pretrained audio model, and the features are transformed into discrete representations in compact audio memory. Sentence Prediction is conducted so as to store linguistic information in the compact audio memory. Note that trained compact audio memory is used in the proposed AKVSR.

which is the phonetic information, which is described in [75]. Motivated by evidence of the prior works, we separately train the K-means clustering and compact audio memory module, and extract the linguistic information from the self-supervised speech representation model through vector quantization.

To this end, as shown in Fig 1, from a given speech utterance  $x_a \in \mathbb{R}^T$ , where  $T$  is the length of the speech utterance, we encode the audio features  $f_a = \{f_{a,i}\}_{i=1}^{T_a} \in \mathbb{R}^{T_a \times d_a}$  using a large-scale pretrained audio model  $E_a$ , where  $T_a$  is the frame length of the audio feature, and  $d_a$  is the embedding dimension of the audio feature. This process can be expressed as  $f_a = E_a(x_a)$ . Since we utilize pretrained audio model on large audio data, the extracted features can be regarded as containing rich speech knowledge.

Next, we aim to store linguistic information at the phoneme level by refining the audio features. To achieve this, we first create a k-means clustering model using an audio corpus dataset consisting of a single speaker to minimize non-linguistic factors about speakers inspired by [44]. Then, we cluster the audio features  $f_a$  into  $N$  clusters,  $C = [1, 2, \dots, N]$ , through the k-means clustering model. The cluster labels of each audio feature are determined by the frame-wise quantization,  $\mathbf{q}(\cdot)$ , as follows:

$$c_{f_a} = \mathbf{q}(f_a) = \{c_{f_{a,i}}\}_{i=1}^{T_a} \quad (1)$$

where  $c_{f_{a,i}} \in C$  is the cluster label of the each audio feature.

After the clustering step, we introduce a trainable compact audio memory to store linguistic information (*i.e.*, representative speech feature) for each cluster group. The compact audio memory  $M_a = \{m_n\}_{n=1}^N$  is comprised of  $N$  discrete representations equal to the number of cluster groups. Each representation denoted  $m_n$  has an embedding dimension of  $d$ . By utilizing the cluster labels generated from the audio features, we are able to access each slot of the compact audio memory. The process of accessing the memory can be generalized as follows:

$$m_{c_{f_a}} = \mathbf{M}(c_{f_a}) = \{m_{c_{f_{a,i}}}\}_{i=1}^{T_a} \quad (2)$$

where  $\mathbf{M}$  is the frame-wise memory accessing function. For instance,  $\mathbf{M}(c_{f_{a,i}} = 1) = m_1$  represents that the discrete representation has been extracted from the first slot of the compact audio memory.

Finally, to store representative knowledge (*i.e.*, linguistic in-

formation only) in predicting speech at compact audio memory while discarding non-linguistic information, we perform ASR using the memory through paired audio-text data. To conduct ASR, we employ a context encoder  $E_c$  and a decoder  $D$ . When the discrete representations are extracted from the compact audio memory, the context encoder captures the context information from these representations. The contextualized representations are then used to predict the speech content  $\hat{y}$  through the decoder as follows:  $\hat{y} = D(E_c(m_{c_{f_a}}))$ . To train all components of our model, we use a hybrid CTC/attention loss [76], which is a commonly used loss function for the speech recognition task. Further details about the losses are provided in subsection III-D. By performing ASR with the compact audio memory, we can extract and store the linguistic knowledge of pretrained audio model in compact audio memory while disentangling the non-linguistic information which is not important in predicting speech.

### C. Injecting Audio Knowledge into VSR with Audio Bridging Module

Our prior discussion addressed the methods for storing linguistic information within a compact audio memory. This section describes how the saved linguistic audio knowledge can be employed for VSR. To this end, we propose Audio Bridging Module (ABM), which aims to inject the best-matched audio knowledge saved in the memory with visual features into the VSR model to complement insufficient visual modality.

The ABM is trained to identify the most relevant audio knowledge stored in the compact audio memory. Once training is completed, the most suitable audio can be injected into the VSR model through the ABM to enhance the limited lip movement information during inference, even without the presence of audio input data. Therefore, the proposed method uses only visual input for extracting audio knowledge from the compact audio memory and training the VSR model. This is different from the existing VSR methods that require both modal inputs during training for Knowledge Distillation (KD) or contrastive learning. The entire pipeline for finding audio knowledge and complementing the visual information is illustrated in Fig. 2.

Given the lip video  $x_v \in \mathbb{R}^{T_v \times C \times H \times W}$ , where  $T_v$  is the number of frames of video,  $C$  is the channel dimension,

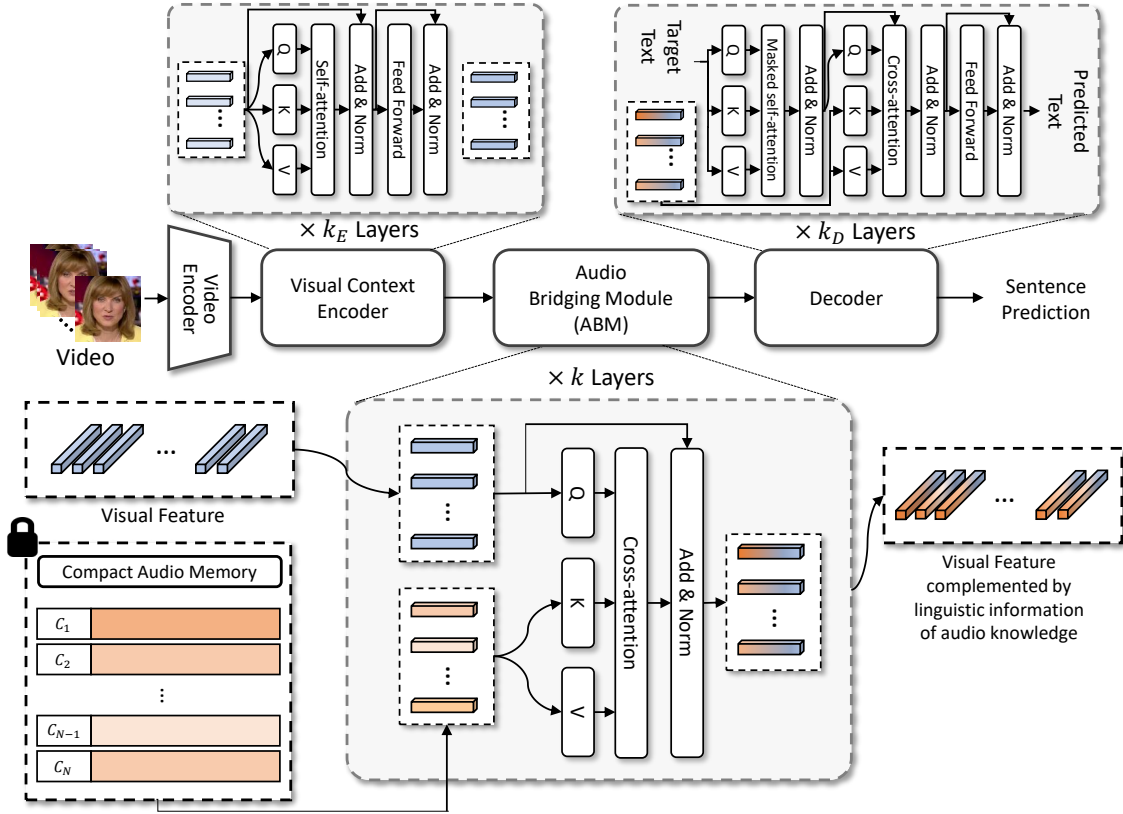


Fig. 2. The overall framework of a proposed AKVSR for complementing visual modality with audio modality. The AKVSR mainly consists of 2 parts: 1) The compact audio memory provides linguistic information from audio knowledge generated by a large-scale pretrained audio model. The meaning of N in compact audio memory is the number of discrete representations in the memory. Moreover, the number of discrete representations is the same as the number of clustering groups of audio features. 2) The proposed ABM finds best-matched information in compact audio memory and injects the linguistic information into the visual feature to complement the insufficient information of lip movements.

$H$  is the height of lip video, and  $W$  is the width of lip video. We encode the lip video to visual features through the video encoder and visual context encoder. The video encoder captures compact spatio-temporal information from lip movements [1]. Next, the visual context encoder can consider the context information of neighboring words via a self-attention mechanism [2]. The encoding process can be formulated as:  $f_v = E_v(x_v) = \{f_{v,i}\}_{i=1}^{T_v}$  where  $E_v$  indicate both the video encoder and visual context encoder,  $f_{v,i} \in \mathbb{R}^d$  is visual feature, and  $d$  represents the embedding dimension.

We construct the ABM to find the best-matched audio knowledge with visual features  $f_v$  from the compact audio memory utilizing a cross-attention mechanism. The primary motivation behind employing cross-attention in our model is to receive the linguistic information from unaligned compact audio memory for complementing visual features. While compact audio memory is comprised of a fixed number of N discrete representations, the number of visual features varies when the input video changes. Recent study [77] shows that cross-attention can be used for aligning multi-modal language sequences, and receiving information from another modality to supplement one modality. Motivated by the success of these works, we utilize attention scores based on cross-attention to complement the visual features through compact audio memory. Moreover, through this approach, the correlation between unaligned visual features and discrete representations of compact audio memory can be calculated. In cross-attention,

the attention score measures how much the given visual feature correlates to the discrete representation of compact audio memory.

To calculate the attention scores, we define the visual query of  $i$ -th visual feature as  $Q_{f_{v,i}} = f_{v,i}W_{Q_v}$ , and audio keys of  $j$ -th linguistic information in compact audio memory as  $K_{a,j} = m_jW_{K_a}$ , where  $W_{Q_v} \in \mathbb{R}^{d \times d_k}$ , and  $W_{K_a} \in \mathbb{R}^{d \times d_k}$  are weight matrices. Then, to find appropriate linguistic information through visual features, we calculate the attention score  $A_{i,j}^{(0)}$  between a  $i$ -th visual feature  $f_{v,i}$  and each audio feature stored in compact audio memory as follows:

$$A_{i,j}^{(0)} = \text{Softmax}\left(\frac{Q_{f_{v,i}}K_{a,j}^T}{\tau}\right) = \frac{\exp(f_{v,i}W_{Q_v}W_{K_a}^Tm_j^T/\tau)}{\sum_{n=1}^N \exp(f_{v,i}W_{Q_v}W_{K_a}^Tm_n^T/\tau)}, \quad (3)$$

where the  $\tau$  is a scaling parameter. Intuitively, the attention score provides how much linguistic information in the  $j$ -th slot in compact audio memory is related to  $i$ -th visual feature. For example, the higher attention score can load more audio information from one of the N discrete representations in compact audio memory. Moreover, we would like to emphasize that the proposed method uses only visual input to recall the audio information inputs because N discrete units of compact audio memory can represent all audio features.

The linguistic audio knowledge in the compact audio mem-

ory for complementing  $i$ -th visual feature can be reconstructed by using the attention score, which can be denoted as:

$$m'_{f_v,i}(0) = \sum_{j=1}^N A_{i,j}(0) (m_j W_{V_a}), \quad (4)$$

where  $W_{V_a} \in \mathbb{R}^{d \times d_v}$  is a weight matrix. Then, we employ a weight matrix  $W_o \in \mathbb{R}^{d_v \times d}$  to match the dimensions of the reconstructed audio knowledge and visual feature. Moreover, since one linguistic audio knowledge is reconstructed for each visual feature, an alignment process between visual and audio features is not required. Therefore, we utilize the reconstructed audio knowledge for the  $i$ -th visual feature as follows:

$$f_{v,i}^{(1)} = LN(f_{v,i} + m'_{f_v,i}(0) W_o), \quad (5)$$

where  $f_{v,i}^{(1)} \in \mathbb{R}^{d_v}$  is the visual feature complemented once by the information of compact audio memory, and the  $LN$  denotes the Layer Normalization [78].

We can generalize the above processes as consisting of two parts. 1) We bring the linguistic information from compact audio memory via a cross-attention mechanism, and 2) we inject the information into the visual feature sequence. To this end, we define the visual feature sequence complemented  $k-1$  times as  $f_v^{(k-1)} \in \mathbb{R}^{T_v \times d}$  and denote the query of visual feature sequence  $Q_{f_v}^{(k-1)} = f_v^{(k-1)} W_{Q_v}^{(k-1)}$ , the entire audio keys  $K_a^{(k-1)} = M_a W_{K_a}^{(k-1)}$ , and the audio values of entire compact audio memory as  $V_a^{(k-1)} = M_a W_{V_a}^{(k-1)}$ . Then, we can formulate the first process of getting linguistic audio knowledge from compact audio memory as follows:

$$\begin{aligned} m'_{f_v}{}^{(k-1)} &= A^{(k-1)} V_a^{(k-1)} \\ &= \text{Softmax}\left(\frac{Q_{f_v}^{(k-1)} K_a^{(k-1)T}}{\tau}\right) V_a^{(k-1)} \end{aligned} \quad (6)$$

where  $A^{(k-1)} = \{A_i^{(k-1)}\}_{i=1}^{T_v}$  is attention scores between visual features sequence and every slot in compact audio memory. We then complement the visual feature sequence with knowledge brought by attention scores from compact audio memory as follows:  $f_v^{(k)} = LN(f_v^{(k-1)} + m'_{f_v}{}^{(k-1)} W_o^{(k-1)})$

#### D. Lip-To-Text Translation

In the previous section, the proposed ABM complements the visual feature through the linguistic information stored in compact audio memory. In this section, we introduce the process of lip-to-text translation through the complemented visual feature.

For the decoder, we use transformer, following the previous methods, to predict sentences. Different from previous works, our proposed method can provide additional linguistic audio knowledge to the visual feature sequence. Therefore, the complemented visual feature is fed into the decoder as input, and the decoder predicted  $L$  subwords through the features. We then employ the hybrid CTC/attention [76] loss to supervise the proposed model through the predicted sentence and target sentence.

CTC [56] loss is widely used to guide the VSR model through frame-wise prediction based on conditional independence. The frame-wise posterior distribution  $p(s_t|\mathbf{x})$ , where

the  $s_t$  is the target subword corresponding to  $t$ -th frame and  $\mathbf{x}$  is a visual input. The CTC probability and CTC loss can be formulated as follows:

$$p_{ctc}(\mathbf{s}|\mathbf{x}) \approx \sum_s \prod_{t=1}^T p(s_t|\mathbf{x}) \quad (7)$$

$$\mathcal{L}_{ctc} = \log p_{ctc}(\mathbf{s}|\mathbf{x}) \quad (8)$$

where the  $\mathbf{s}$  is a target subwords containing blank symbols. Attention loss is employed to learn an implicit language model. The decoder infers the next target subword conditioned on the previous prediction. The attention loss can be formulated as:

$$p_{att}(\mathbf{s}|\mathbf{x}) = \prod_{l=1}^L p(s_l|s_1, \dots, s_{l-1}, \mathbf{x}) \quad (9)$$

$$\mathcal{L}_{att} = \log p_{att}(\mathbf{s}|\mathbf{x}) \quad (10)$$

where  $s_l$  is the predicted subword at time step  $l$ . Finally, the hybrid CTC/Attention loss can be formulated by weighted summation of the two loss functions:

$$\mathcal{L}_{tot} = (1 - \lambda)\mathcal{L}_{att} + \lambda\mathcal{L}_{ctc}, \quad (11)$$

where  $\lambda$  is the balancing weight.

## IV. EXPERIMENTAL SETUP

### A. Dataset.

**LRS2 & LRS3** are two of the most popular publicly available sentence-level VSR datasets. LRS2 [57] and LRS3 [79] datasets are extracted from BBC television and TED & TEDx talks, respectively. The LRS2 consists of 28 hours of video for training and 195 hours of video for pretraining. The difference between the training dataset and the pretraining dataset is the duration of each video clip. The duration of pretraining videos is longer than the trainset. In the same way, the LRS3 comprised 30 hours of video for training and 403 hours of video for pretraining. We divide the LRS3 dataset into 30 hours of video for low-resource settings and 433 hours of video for high-resource settings to verify the effectiveness according to the amount the labeled data like AV-HuBERT [17]. In addition, The LRS2 dataset is divided into 28 hours of video for low-resource settings and 223 hours of video for high-resource settings.

We follow the preprocessing pipeline of AV-HuBERT [17]. Firstly, we extract the landmarks from each video clip using dlib [80], and employ affine transformation to align each frame to the reference face frame. Next, we crop the video into  $96 \times 96$  corresponding to the lip region. For data augmentation at train time, the random crop and random horizontal flip are used.

### B. Implementation Details.

For the VSR model, we employ a state-of-the-art architecture based on transformer encoder-decoder models detailed in [17]. The video encoder network consists of a 3D CNN, ResNet-18, and the visual context encoder based on the transformer encoder. We employ a Transformer BASE (12 layers) and a Transformer LARGE (24 layers) same as AV-HuBERT

TABLE I

VSR PERFORMANCE COMPARISON BETWEEN THE PROPOSED METHODS AND PREVIOUS MODELS ON THE LRS3 DATASET. † INDICATES THAT THE NON-PUBLIC VIDEO-TEXT DATA IS USED FOR TRAINING THE VSR MODEL. ‡ INDICATES THAT THE SYNTHVSR USES AN ADDITIONAL 3,652 HOURS OF SYNTHETIC VIDEO-TEXT DATA TO TRAIN THE CONFORMER-BASE MODEL. (+ $\alpha$ ) INDICATES THE AMOUNT OF PSEUDO-TEXT LABELS GENERATED BY THE PRETRAINED ASR MODEL.

Method	Backbone	Labeled word data (hrs)	Labeled sentence data (hrs)	Unlabeled sentence data (hrs)	WER(%)
Afouras et al. [22]	CNN	157	433	-	68.8
Zhang et al. [81]	CNN	157	698	-	60.1
Afouras et al. [79]	Transformer	157	1362	-	58.9
Xu et al. [82]	RNN	157	433	-	57.8
Shillingford et al. [83]	RNN	-	3.886	-	55.1
Ma et al. [15]	Conformer	-	433	-	46.9
Ma et al. [15]	Conformer	157	433	-	43.3
Prajwal et al. [84]	Transformer	-	698	-	40.6
Ma et al. [16]	Conformer	-	433	-	37.9
Ma et al. [16]	Conformer	157	433	-	35.1
Makino et al.† [85]	RNN	-	31,000	-	33.6
Serdyuk et al.† [86]	Conformer	-	90,000	-	17.0
SynthVSR [87]	Conformer-BASE	-	30	3652‡	43.3
	Conformer-BASE	-	433	3652‡	27.9
AV-HuBERT [17]	Transformer-BASE	-	30	1759	46.1
	Transformer-BASE	-	433	1759	34.8
	Transformer-LARGE	-	30	1759	32.5
	Transformer-LARGE	-	433	1759	28.6
	Transformer-LARGE	-	433(+1,326)	1759	26.8
Lohrenz. et al. [88]	Transformer-LARGE	-	30	1326	44.0
	Transformer-LARGE	-	433	1326	28.8
	Transformer-LARGE	-	433(+1,326)	1326	26.3
RAVE <sub>n</sub> [89]	Transformer-LARGE	-	30	1759	32.5
	Transformer-LARGE	-	433	1759	27.8
	Transformer-LARGE	-	433(+1,326)	1759	24.4
<b>Proposed Method</b>	Transformer-BASE	-	30	1759	<b>41.6</b>
	Transformer-BASE	-	433	1759	<b>34.2</b>
	Transformer-LARGE	-	30	1759	<b>29.1</b>
	Transformer-LARGE	-	433	1759	<b>27.6</b>
	Transformer-LARGE	-	433(+1,326)	1759	<b>23.6</b>

[17], which models have 103M and 325M of the number of parameters, respectively. We initialize the parameters of both the video encoder and visual context encoder from the BASE model and the LARGE model of AV-HuBERT, respectively. Similar to the transformer encoder, there are two versions of the decoder. transformer decoders have 6 layers for the BASE and 9 layers for the LARGE. We then utilize a unigram-based subword unit [90] composed of 1000 subwords like Av-HuBERT to decode the features produced by the encoder into subword units. We initialize the parameters of the transformer decoders fine-tuned on the audio corpus of LRS2 and LRS3, respectively. The proposed compact audio memory is constructed of 200 discrete units. The embedding dimension of the memory for BASE and LARGE are 768 and 1024, and the number of parameters are 0.15M and 0.2M, respectively. To train the compact audio memory with both context encoder and decoder through the ASR task, we use four transformer layers for context encoder, and transformer decoders (BASE)

such as the VSR model decoder. We then use 2 cross-attention layers with 8 multi-head attention mechanisms for ABM. The ABM has 4.7M parameters. The balancing weight 0.1 is used at the proposed method.

## V. EXPERIMENTAL RESULTS

In this section, we validate the effectiveness of the AKVSR framework utilizing both compact audio memory and ABM through a comparison of the proposed method with the previous VSR methods. Then, to demonstrate the availability of each component, we provide various experimental results. Firstly, we examine how the linguistic information constructed by different pretrained audio model affects the VSR models when we inject the information into the visual modality. Next, we investigate the impact on VSR according to the amount of injected audio information through the proposed ABM. Finally, we demonstrate that the AKVSR can complement visual information of the other VSR models.

TABLE II  
VSR PERFORMANCE COMPARISON BETWEEN THE PROPOSED METHODS AND PREVIOUS MODELS ON LRS2 DATASET.

Method	Backbone	Labeled word data (hrs)	Labeled sentence data (hrs)	Unlabeled sentence data (hrs)	WER(%)
Afouras et al. [22]	CNN	157	223	-	58.5
Zhang et al. [81]	CNN	157	698	-	51.7
Afouras et al. [79]	Transformer	157	1362	-	48.3
Kim et al. [29]	Transformer	157	656	-	46.2
Kim et al. [30]	Transformer	157	656	-	44.5
Ma et al. [15]	Conformer	-	223	-	39.1
Ma et al. [15]	Conformer	157	223	-	37.9
Ma et al. [16]	Conformer	-	223	-	32.9
K R Prajwal et al. [84]	Transformer	-	698	-	28.9
Ma et al. [16]	Conformer	157	223	-	28.7
AV-HuBERT [17]	Transformer-BASE	-	28	1759	43.3
	Transformer-BASE	-	223	1759	31.2
	Transformer-LARGE	-	28	1759	32.2
	Transformer-LARGE	-	223	1759	25.5
<b>Proposed Method</b>	Transformer-BASE	-	28	1759	<b>40.5</b>
	Transformer-BASE	-	223	1759	<b>30.1</b>
	Transformer-LARGE	-	28	1759	<b>28.7</b>
	Transformer-LARGE	-	223	1759	<b>24.1</b>

### A. Comparisons with the State-of-the-art

To verify the effectiveness of the AKVSR, we first compared it to AV-HuBERT on the LRS3 dataset. We adopt the AKVSR to both the Transformer BASE and Transformer LARGE models. To compare performance based on the amount of video-text label data, we perform fine-tuning using both low-resource (30 hours) and high-resource (433 hours) settings. Thus, we conduct experiments with four different settings on the LRS3 dataset.

Table I presents the results of the AVKSR method on the LRS3 dataset. In the [86], the state-of-the-art model using a conformer encoder as a visual front-end shows a best WER of 17.0%. However, we would like to emphasize that the state-of-the-art performance model is trained on 90,000 hours of non-public video-text data. On the other hand, we use publicly available video-text data for training. Therefore, we would like to note that a direct comparison between the method [86] (trained on 90,000 hours of data) and the proposed method is not fair in terms of state-of-the-art (SOTA) comparison. Instead, to expand the experimental dataset, we utilize an additional VoxCeleb2 dataset, which datasets comprise 1,326 hours of video. Similar to recent works [17], [89], we generate pseudo-text labels from the VoxCeleb2 Dataset.

In the LRS3 dataset, the detailed comparison of our proposed methods' performances with recent state-of-the-art performances is shown in Table I. We would like to emphasize the proposed methods consistently outperform the AV-HuBERT model across various training data sizes. When trained on a 30-hour dataset, the proposed model achieves a WER of 29.1%, which is notably lower than the 32.5% WER of AV-HuBERT. Similarly, with a larger dataset of 433 hours, the proposed method shows an improved WER of 27.6% compared to

AV-HuBERT's 28.6%. This trend of enhanced performance continues with the most extensive dataset of 433(+1,326) hours, where the proposed method attains a WER of 23.6%, surpassing AV-HuBERT's 26.8%. Moreover, we compare this performance of the proposed model with the recent state-of-the-art methods [88], [89]. Our methods achieve a WER of 23.6%, surpassing the recent Raven [89] method, which records a WER of 24.4%. Additionally, it outperforms the previous approach [88] relaxed attention, which has 26.3% WER. Please note that these reported WERs of [17], [87]–[89] are based on not using the language model during the inference stage for fair comparison.

Regarding correlated results to the model parameters, the LARGE model surpasses 12.5% WER compared to the BASE model since the LARGE model has 122M more parameters. On the other hand, the baseline method achieves a WER of 46.1% and 32.5%. We would like to emphasize that the compact audio memory and ABM utilize only around 5M additional parameters compared to the baseline.

To validate the effectiveness of the proposed method in another dataset, we additionally conduct experiments on the LRS2 dataset. The AV-HuBERT method does not provide results for the LRS2 dataset, because the video data of LRS2 is open for academic research. In our comparative analysis showcased in Table II, the performance of the AV-HuBERT model and our proposed method across different configurations is detailed. For the Transformer-BASE configuration, AV-HuBERT achieves a WER of 43.3% with 28 hours of training data and 31.2% with 223 hours, while our proposed method shows improved results with a WER of 40.5% and 30.1%. In addition, the proposed methods also outperform the AV-HuBERT by achieving a WER of 24.1%, in the Transformer-



TABLE III  
VSR PERFORMANCE OF THE PROPOSED AKVSR WHEN DIFFERENT PRETRAINED AUDIO MODELS (CPC, WAV2VEC2.0, AND HUBERT) ARE UTILIZED TO CONSTRUCT COMPACT AUDIO MEMORY

Methods	WER(%)
Baseline	46.1
<b>AKVSR (CPC [91])</b>	42.7
<b>AKVSR (Wav2vec2.0 [33])</b>	42.4
<b>AKVSR (HuBERT [34])</b>	<b>41.6</b>

TABLE IV  
IMPACT OF REFINING NON-LINGUISTIC FACTORS SUCH AS SPEAKERS AND NOISE OF AUDIO KNOWLEDGE WHEN COMPLEMENTING VISUAL MODALITY WITH AUDIO MODALITY

Method	WER(%)
AV-HuBERT [17] + <b>KD [19]</b>	<b>43.6</b>
AV-HuBERT [17] + <b>Auxiliary task [16]</b>	<b>43.4</b>
AV-HuBERT [17] + <b>AKVSR</b>	<b>41.6</b>

LARGE configuration.

#### B. Effect of different pretrained audio models on constructing compact audio memory

In the previous section, the proposed AKVSR injects linguistic information into the VSR model and shows promising results outperforming the current state-of-the-art VSR model. In this section, to assess the effectiveness of different large-scale pretrained audio models, we create compact audio memory using CPC, Wav2vec2.0, and HuBERT and apply the compact audio memory to the VSR model, respectively. The results are shown in Table III. While the baseline model (*i.e.* without compact audio memory and ABM) achieves 46.1% WER, our proposed method adopted by CPC, Wav2vec2.0, and HuBERT accomplishes a WER of 42.7%, 42.4%, and 41.6%, respectively. The results show that the proposed compact audio memory can store the linguistic information of any large-scale pretrained audio model and improves the performance of the existing VSR methods by supplementing visual modality with the information. We utilize the HuBERT audio model in other experiments as it achieves the best performance.

#### C. Effect of refining non-linguistic factors when transferring the audio knowledge.

In this section, we verify the effect of refining non-linguistic factors. To verify this, we experiment with two methods [16], [19] utilizing KD and compare them to the proposed method. The results are shown in Table IV. Following the method [19], we use pretrained ASR model as a teacher network and utilize the AV-HuBERT as a student network. From guiding the pretrained ASR model, we achieve a WER of 43.6%. The other experiment is transferring the knowledge between intermediate layers. The effectiveness of this method called auxiliary task recently is more effective than [19] and verified in multiple languages VSR in [16]. Through this approach, we

TABLE V  
VSR PERFORMANCES ACCORDING TO THE NUMBER OF CROSS-ATTENTION LAYERS IN ABM

# Cross-attention layer	WER(%)
1	41.8
2	<b>41.6</b>
3	41.9
4	41.8

TABLE VI  
IMPROVING PERFORMANCES OF EXISTING VSR METHODS BY APPLYING THE PROPOSED AKVSR

Method	WER(%)
TM-seq2seq [14]	59.9
TM-seq2seq [14] + <b>AKVSR</b>	<b>54.5</b>
AV-HuBERT [17]	46.1
AV-HuBERT [17] + <b>AKVSR</b>	<b>41.6</b>

accomplish 43.4% WER. At this time, our proposed method achieve a 41.6% WER by using the same training data by refining the non-linguistic factors such as speakers and noise of audio knowledge and transferring via ABM.

#### D. Effect of according to different number of layers in Audio Bridging Module

The proposed ABM can be composed of different numbers of cross-attention layers. In order to evaluate the effect of the number of layers, we build 4 variants of ABM by differing the number of layers from 1 to 4 and perform VSR on LRS3. The ablation result is shown in Table V. By using ABM consisting of a single cross-attention layer, we can achieve a significant performance improvement, 4.3% WER from the baseline. This result confirms that utilizing compact audio knowledge through the proposed memory is effective in VSR by complementing insufficient visual information with audio information provided by a large-scale pretrained audio model. The best result is obtained when 2 cross-attention layers are utilized for ABM and we found that increasing the number of layers to more than 2 does not give more performance gain. Therefore, we employ 2 cross-attention layers for ABM in other experiments. Moreover, by comparing the number of parameters of ABM with the baseline model, adding one cross-attention layer increases 2.4M parameters, which is 1.5% of that of the baseline model. Therefore, the best-performed model (*i.e.*, 2 cross-attention layers) just requires only 3% additional parameters of the baseline while improving 9.76% relative performance from the baseline.

#### E. Effect of AKVSR on different VSR methods

We conducted experiments to demonstrate that our proposed method, which incorporates compact audio memory and ABM, can be applied to other VSR methods. To validate this, we apply the proposed method to another popular VSR method, TM-seq2seq [14]. This model is originally designed to apply

TABLE VII

IMPACT OF THE NUMBER OF CLUSTERS ON THE RESULTS OF SENTENCE PREDICTION. "A" DENOTES THE ASR PERFORMANCE OF USING ONLY DISCRETE UNITS OF THE MEMORY. "V" REPRESENTS THE VSR PERFORMANCE OF THE BASELINE MODEL WHEN IT ADAPTS COMPACT AUDIO MEMORY USING THE ABM.

Methods	# of clusters	WER(%)	
		A	V
AKVSR (HuBERT [34])	200	9.4	<b>41.6</b>
AKVSR (HuBERT [34])	500	8.1	42.0
AKVSR (HuBERT [34])	1000	<b>6.2</b>	41.8

VSR in wild videos. We re-implemented the TM-seq2seq and trained it using the curriculum learning approach detailed in [14]. The results are shown in VI. We can achieve a WER of 54.1% by applying the proposed AKVSR at TM-seq2seq. In addition, when we adopt our proposed method to the state-of-the-art AV-HuBERT VSR model. Our implementation resulted in a 4.5% WER improvement. As a result, the proposed method injects the linguistic information of a large-scale pretrained audio model into visual modality at the feature level and can bring improvements to various VSR methods.

#### F. Effect of varying the number of clusters

We evaluate how storing linguistic information in compact audio memory is affected by varying the number of distinct clusters used to organize the information. The results are shown in Table VII. We conduct three experiments using 200, 500, and 1000 clusters referring to that utilizing audio feature clustering [17], [44]. In these experiments, we employ HuBERT to extract and cluster the audio features. After storing the linguistic information in trainable compact audio memory via ASR task on the LRS3 dataset, we achieve a WER of 9.4% a WER of 8.1% WER, and a WER of 6.2%, respectively. The results show that if the number of clusters increases, the WER of ASR decrease. However, when we apply the compact audio memories to the baseline models, we observe that the memory using 200 clusters obtains the most 41.6% WER of performance. We analyze these results as that using the 200 discrete units in the compact audio memory is more appropriate to find the best-matched audio representation via ABM from the visual feature than employing other compact audio memories consisting of 500 and 1000 discrete units.

#### G. Effect of varying dimension of compact audio memory

We evaluate how the discrete representation dimension in compact audio memory affects the VSR performance of the proposed method. To this end, we conduct three experiments, varying the dimensions of the compact audio memory to 512, 768, and 1024, when constructing the memory to store linguistic information from audio. Then, these memories are applied to the baseline VSR model. The results are shown in Table VIII. When the compact audio memory of 768 dimensions is used, we obtain the best performance, a WER of 41.5%. The other settings using 512 and 1024 dimensions achieve 42.0% and 41.7% WERs, respectively. Please note that

TABLE VIII

ABLATION STUDY ON THE PROPOSED METHOD PERFORMANCE WITH VARYING DIMENSION OF COMPACT AUDIO MEMORY.

Dimension	WER(%)
512	41.7
768	<b>41.5</b>
1024	42.0

TABLE IX

ABLATION STUDIES ON THE EFFECTS OF VARYING TRAINING DATASETS IN THE CONSTRUCTION OF COMPACT AUDIO MEMORY

Training Datasets	Duration (hrs)	WER(%)
LRS2	223	44.7
LRS3	433	43.3
LRS2, LRS3	656	<b>41.5</b>

all experiments improve the performance of the baseline VSR model achieving a WER of 46.1%.

#### H. Effect of varying training dataset when constructing the compact audio memory

We verify how the varying datasets for training the compact audio memory affect the performance of the baseline. For this, LRS2, LRS3, and merging LRS2 with LRS3 are used for building compact audio memories. The results are shown in Table IX. In these settings, we obtain 44.7%, 43.3%, and 41.5% WERs, respectively. Based on these results, the compact audio memory trained on merging LRS2 with LRS3 datasets is applied to the VSR model in other experiments.

## VI. CONCLUSION

This paper has presented the Audio Knowledge empowered Visual Speech Recognition (AKVSR) framework to enhance the complementary effect of audio knowledge for visual modality, by removing the non-linguistic factors. This proposed framework consists of three components: (1) Unlike previous methods, the proposed approach utilizes a large-scale pretrained audio model to acquire audio knowledge, and the non-linguistic factors, such as speaker and noise, are then removed through vector quantization from the audio knowledge. (2) The resulting audio knowledge, focusing on linguistic information, has been stored in a compact audio memory to allow for the utilization of audio information. (3) The Audio Bridging Module has been devised to match the best audio knowledge with the visual features in the compact audio memory to complement the insufficient visual information and inject linguistic information into the VSR model. We have demonstrated the effectiveness of the proposed method by outperforming the state-of-the-art VSR models on the LRS3 dataset.

## VII. DISCUSSION

Although the proposed AKVSR can improve the VSR systems using the audio knowledge of large-scale pretrained audio models, there is some limitation and future research: (1) In this paper, the compact audio memory needs to be constructed

before training the VSR model. It requires additional time consumption. In future research, we are going to further investigate a time-efficient training method. (3) Compared with the ASR systems, the performance of the VSR system (including the proposed method) is still inferior. To bridge the thin gap between VSR and ASR, in future research, how to utilize the audio model needs to be continually investigated.

## VIII. ACKNOWLEDGMENT

This work was supported in part by the IITP grant funded by the Korea government (MSIT) (No.2020-0-00004, Development of Previsional Intelligence based on Long-Term Visual Memory Network), in part by the National Research Foundation of Korea (NRF) grant funded by the Korea government (MSIT) (No.NRF-2022R1A2C2005529), and in part by the BK21 FOUR(Connected AI Education & Research Program for Industry and Society Innovation, KAIST EE, No. 4120200113769).

## REFERENCES

- [1] S. Ji, W. Xu, M. Yang, and K. Yu, "3d convolutional neural networks for human action recognition," *IEEE transactions on pattern analysis and machine intelligence*, vol. 35, no. 1, pp. 221–231, 2012.
- [2] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," pp. 5998–6008, 2017.
- [3] K. Cho, B. Van Merriënboer, D. Bahdanau, and Y. Bengio, "On the properties of neural machine translation: Encoder-decoder approaches," *arXiv preprint arXiv:1409.1259*, 2014.
- [4] D. Ivanko, D. Ryumin, A. Kashevnik, A. Axyonov, and A. Karnov, "Visual speech recognition in a driver assistance system," in *2022 30th European Signal Processing Conference (EUSIPCO)*. IEEE, 2022, pp. 1131–1135.
- [5] C. Sheng, L. Liu, W. Deng, L. Bai, Z. Liu, S. Lao, G. Kuang, and M. Pietikäinen, "Importance-aware information bottleneck learning paradigm for lip reading," *IEEE Transactions on Multimedia*, 2022.
- [6] C. Sheng, X. Zhu, H. Xu, M. Pietikäinen, and L. Liu, "Adaptive semantic-spatio-temporal graph convolutional network for lip reading," *IEEE Transactions on Multimedia*, vol. 24, pp. 3545–3557, 2021.
- [7] T. Saitoh, Z. Zhou, G. Zhao, and M. Pietikäinen, "Concatenated frame image based cnn for visual speech recognition," in *Computer Vision–ACCV 2016 Workshops: ACCV 2016 International Workshops, Taipei, Taiwan, November 20–24, 2016, Revised Selected Papers, Part II 13*. Springer, 2017, pp. 277–289.
- [8] M. Sadeghi, S. Leglaive, X. Alameda-Pineda, L. Girin, and R. Horaud, "Audio-visual speech enhancement using conditional variational auto-encoders," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 28, pp. 1788–1800, 2020.
- [9] H. Akbari, H. Arora, L. Cao, and N. Mesgarani, "Lip2audspec: Speech reconstruction from silent lip movements video," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018, pp. 2516–2520.
- [10] T. Stafylakis and G. Tzimiropoulos, "Combining residual networks with lstms for lipreading," *arXiv preprint arXiv:1703.04105*, 2017.
- [11] S. Petridis, T. Stafylakis, P. Ma, F. Cai, G. Tzimiropoulos, and M. Pantic, "End-to-end audiovisual speech recognition," in *2018 IEEE international conference on acoustics, speech and signal processing (ICASSP)*. IEEE, 2018, pp. 6548–6552.
- [12] X. Zhao, S. Yang, S. Shan, and X. Chen, "Mutual information maximization for effective lip reading," in *2020 15th IEEE International Conference on Automatic Face and Gesture Recognition (FG 2020)*. IEEE, 2020, pp. 420–427.
- [13] T. Mikolov, M. Karafiát, L. Burget, J. Cernocký, and S. Khudanpur, "Recurrent neural network based language model," in *Interspeech*, vol. 2, no. 3. Makuhari, 2010, pp. 1045–1048.
- [14] 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.
- [15] 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.
- [16] M. P. Pingchuan Ma, Stavros Petridis, "Visual speech recognition for multiple languages in the wild," *Nature Machine Intelligence*, pp. 1–10, 2022.
- [17] B. Shi, W.-N. Hsu, K. Lakhota, and A. Mohamed, "Learning audio-visual speech representation by masked multimodal cluster prediction," *arXiv preprint arXiv:2201.02184*, 2022.
- [18] R. T. Sataloff, "The human voice," *Scientific American*, vol. 267, no. 6, pp. 108–115, 1992.
- [19] G. Hinton, O. Vinyals, J. Dean *et al.*, "Distilling the knowledge in a neural network," *arXiv preprint arXiv:1503.02531*, vol. 2, no. 7, 2015.
- [20] 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.
- [21] 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.
- [22] 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.
- [23] H. Mabrouk, O. Abugabal, N. Sakr, and H. M. Eraqi, "Lip-listening: Mixing senses to understand lips using cross modality knowledge distillation for word-based models," *arXiv preprint arXiv:2207.05692*, 2022.
- [24] J.-X. Zhang, G. Wan, Z.-H. Ling, J. Pan, J. Gao, and C. Liu, "Self-supervised audio-visual speech representations learning by multimodal self-distillation," in *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2023, pp. 1–5.
- [25] S. Elashmawy, M. Ramsis, H. M. Eraqi, F. Eldeshnawy, H. Mabrouk, O. Abugabal, and N. Sakr, "Spatio-temporal attention mechanism and knowledge distillation for lip reading," *arXiv preprint arXiv:2108.03543*, 2021.
- [26] X. Huang and Y. Peng, "Deep cross-media knowledge transfer," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 8837–8846.
- [27] Y. Peng and J. Qi, "Cm-gans: Cross-modal generative adversarial networks for common representation learning," *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, vol. 15, no. 1, pp. 1–24, 2019.
- [28] K. Lin, X. Xu, L. Gao, Z. Wang, and H. T. Shen, "Learning cross-aligned latent embeddings for zero-shot cross-modal retrieval," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 07, 2020, pp. 11 515–11 522.
- [29] M. Kim, J. Hong, S. J. Park, and Y. M. Ro, "Cromm-vsr: Cross-modal memory augmented visual speech recognition," *IEEE Transactions on Multimedia*, 2021.
- [30] M. Kim, J. Yeo, and Y. M. Ro, "Distinguishing homophenes using multi-head visual-audio memory," in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2022.
- [31] J. H. Yeo, M. Kim, and Y. M. Ro, "Multi-temporal lip-audio memory for visual speech recognition," in *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2023, pp. 1–5.
- [32] J. Weston, S. Chopra, and A. Bordes, "Memory networks," *arXiv preprint arXiv:1410.3916*, 2014.
- [33] 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.
- [34] W.-N. Hsu, B. Bolte, Y.-H. H. Tsai, K. Lakhota, 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.
- [35] J.-S. Lee and C. H. Park, "Robust audio-visual speech recognition based on late integration," *IEEE Transactions on Multimedia*, vol. 10, no. 5, pp. 767–779, 2008.
- [36] S. Dupont and J. Luetttin, "Audio-visual speech modeling for continuous speech recognition," *IEEE transactions on multimedia*, vol. 2, no. 3, pp. 141–151, 2000.

- [37] C. Peláez-Moreno, A. Gallardo-Antolín, and F. Díaz-de María, “Recognizing voice over ip: A robust front-end for speech recognition on the world wide web,” *IEEE Transactions on Multimedia*, vol. 3, no. 2, pp. 209–218, 2001.
- [38] C. C. Chibelushi, F. Deravi, and J. S. Mason, “A review of speech-based bimodal recognition,” *IEEE transactions on multimedia*, vol. 4, no. 1, pp. 23–37, 2002.
- [39] G. Zheng, Y. Xiao, K. Gong, P. Zhou, X. Liang, and L. Lin, “Wavbert: Cooperative acoustic and linguistic representation learning for low-resource speech recognition,” *arXiv preprint arXiv:2109.09161*, 2021.
- [40] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “Bert: Pre-training of deep bidirectional transformers for language understanding,” *arXiv preprint arXiv:1810.04805*, 2018.
- [41] Y. Zhao, G. Wang, C. Luo, W. Zeng, and Z.-J. Zha, “Self-supervised visual representations learning by contrastive mask prediction,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 10 160–10 169.
- [42] A. Van Den Oord, O. Vinyals *et al.*, “Neural discrete representation learning,” *Advances in neural information processing systems*, vol. 30, 2017.
- [43] P. Esser, R. Rombach, and B. Ommer, “Taming transformers for high-resolution image synthesis,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 12 873–12 883.
- [44] K. Lakhotia, E. Kharitonov, W.-N. Hsu, Y. Adi, A. Polyak, B. Bolte, T.-A. Nguyen, J. Copet, A. Baevski, A. Mohamed *et al.*, “On generative spoken language modeling from raw audio,” *Transactions of the Association for Computational Linguistics*, vol. 9, pp. 1336–1354, 2021.
- [45] J. S. Chung and A. Zisserman, “Lip reading in the wild,” in *Asian conference on computer vision*. Springer, 2016, pp. 87–103.
- [46] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [47] F. Tao and C. Busso, “End-to-end audiovisual speech recognition system with multitask learning,” *IEEE Transactions on Multimedia*, vol. 23, pp. 1–11, 2020.
- [48] X. Weng and K. Kitani, “Learning spatio-temporal features with two-stream deep 3d cnns for lipreading,” *arXiv preprint arXiv:1905.02540*, 2019.
- [49] J. Xiao, S. Yang, Y. Zhang, S. Shan, and X. Chen, “Deformation flow based two-stream network for lip reading,” in *2020 15th IEEE International Conference on Automatic Face and Gesture Recognition (FG 2020)*. IEEE, 2020, pp. 364–370.
- [50] Y. Zhang, S. Yang, J. Xiao, S. Shan, and X. Chen, “Can we read speech beyond the lips? rethinking roi selection for deep visual speech recognition,” in *2020 15th IEEE International Conference on Automatic Face and Gesture Recognition (FG 2020)*. IEEE, 2020, pp. 356–363.
- [51] B. Martinez, P. Ma, S. Petridis, and M. Pantic, “Lipreading using temporal convolutional networks,” in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 6319–6323.
- [52] P. Ma, B. Martinez, S. Petridis, and M. Pantic, “Towards practical lipreading with distilled and efficient models,” in *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2021, pp. 7608–7612.
- [53] T. Furlanello, Z. Lipton, M. Tschannen, L. Itti, and A. Anandkumar, “Born again neural networks,” in *International Conference on Machine Learning*. PMLR, 2018, pp. 1607–1616.
- [54] M. Kim, H. Kim, and Y. M. Ro, “Speaker-adaptive lip reading with user-dependent padding,” in *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXVI*. Springer, 2022, pp. 576–593.
- [55] Y. M. Assael, B. Shillingford, S. Whiteson, and N. De Freitas, “Lipnet: End-to-end sentence-level lipreading,” *arXiv preprint arXiv:1611.01599*, 2016.
- [56] A. Graves, S. Fernández, F. Gomez, and J. Schmidhuber, “Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks,” in *Proceedings of the 23rd international conference on Machine learning*, 2006, pp. 369–376.
- [57] J. S. Chung, A. Senior, O. Vinyals, and A. Zisserman, “Lip reading sentences in the wild,” in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2017, pp. 3444–3453.
- [58] J. Hong, M. Kim, and Y. M. Ro, “Visual context-driven audio feature enhancement for robust end-to-end audio-visual speech recognition,” in *23rd Annual Conference of the International Speech Communication Association, INTERSPEECH 2022*. International Speech Communication Association, 2022, pp. 2838–2842.
- [59] A. Koumparoulis and G. Potamianos, “Accurate and resource-efficient lipreading with efficientnetv2 and transformers,” in *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2022, pp. 8467–8471.
- [60] S. Lee, H. G. Kim, D. H. Choi, H.-I. Kim, and Y. M. Ro, “Video prediction recalling long-term motion context via memory alignment learning,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 3054–3063.
- [61] M. Kim, J. Hong, S. J. Park, and Y. M. Ro, “Multi-modality associative bridging through memory: Speech sound recollected from face video,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 296–306.
- [62] L. H. Li, M. Yatskar, D. Yin, C.-J. Hsieh, and K.-W. Chang, “Visualbert: A simple and performant baseline for vision and language,” *arXiv preprint arXiv:1908.03557*, 2019.
- [63] Y.-C. Chen, L. Li, L. Yu, A. El Kholy, F. Ahmed, Z. Gan, Y. Cheng, and J. Liu, “Uniter: Universal image-text representation learning,” in *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXX*. Springer, 2020, pp. 104–120.
- [64] L. Zhou, H. Palangi, L. Zhang, H. Hu, J. Corso, and J. Gao, “Unified vision-language pre-training for image captioning and vqa,” in *Proceedings of the AAAI conference on artificial intelligence*, vol. 34, no. 07, 2020, pp. 13 041–13 049.
- [65] S. Zhang, T. Jiang, T. Wang, K. Kuang, Z. Zhao, J. Zhu, J. Yu, H. Yang, and F. Wu, “Devlibert: Learning deconfounded visio-linguistic representations,” in *Proceedings of the 28th ACM International Conference on Multimedia*, 2020, pp. 4373–4382.
- [66] Z.-Y. Dou, Y. Xu, Z. Gan, J. Wang, S. Wang, L. Wang, C. Zhu, P. Zhang, L. Yuan, N. Peng *et al.*, “An empirical study of training end-to-end vision-and-language transformers,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 18 166–18 176.
- [67] C. Sun, A. Myers, C. Vondrick, K. Murphy, and C. Schmid, “Videobert: A joint model for video and language representation learning,” in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 7464–7473.
- [68] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, “Learning transferable visual models from natural language supervision,” in *International conference on machine learning*. PMLR, 2021, pp. 8748–8763.
- [69] J. Zhao and W.-Q. Zhang, “Improving automatic speech recognition performance for low-resource languages with self-supervised models,” *IEEE Journal of Selected Topics in Signal Processing*, vol. 16, no. 6, pp. 1227–1241, 2022.
- [70] J. S. Chung and A. Zisserman, “Out of time: automated lip sync in the wild,” in *Asian conference on computer vision*. Springer, 2016, pp. 251–263.
- [71] P. Ma, R. Mira, S. Petridis, B. W. Schuller, and M. Pantic, “Lira: Learning visual speech representations from audio through self-supervision,” *arXiv preprint arXiv:2106.09171*, 2021.
- [72] Y. Wang, A. Boumadane, and A. Heba, “A fine-tuned wav2vec 2.0/hubert benchmark for speech emotion recognition, speaker verification and spoken language understanding,” *arXiv preprint arXiv:2111.02735*, 2021.
- [73] X. Chang, B. Yan, Y. Fujita, T. Maekaku, and S. Watanabe, “Exploration of efficient end-to-end asr using discretized input from self-supervised learning,” *arXiv preprint arXiv:2305.18108*, 2023.
- [74] A. Lee, P.-J. Chen, C. Wang, J. Gu, S. Popuri, X. Ma, A. Polyak, Y. Adi, Q. He, Y. Tang *et al.*, “Direct speech-to-speech translation with discrete units,” *arXiv preprint arXiv:2107.05604*, 2021.
- [75] A. Pasad, B. Shi, and K. Livescu, “Comparative layer-wise analysis of self-supervised speech models,” in *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2023, pp. 1–5.
- [76] S. Watanabe, T. Hori, S. Kim, J. R. Hershey, and T. Hayashi, “Hybrid ctc/attention architecture for end-to-end speech recognition,” *IEEE Journal of Selected Topics in Signal Processing*, vol. 11, no. 8, pp. 1240–1253, 2017.
- [77] Y.-H. H. Tsai, S. Bai, P. P. Liang, J. Z. Kolter, L.-P. Morency, and R. Salakhutdinov, “Multimodal transformer for unaligned multimodal language sequences,” in *Proceedings of the conference. Association for Computational Linguistics. Meeting*, vol. 2019. NIH Public Access, 2019, p. 6558.
- [78] J. L. Ba, J. R. Kiros, and G. E. Hinton, “Layer normalization,” *arXiv preprint arXiv:1607.06450*, 2016.

- [79] T. Afouras, J. S. Chung, and A. Zisserman, "Lrs3-ted: a large-scale dataset for visual speech recognition," *arXiv preprint arXiv:1809.00496*, 2018.
- [80] D. E. King, "Dlib-ml: A machine learning toolkit," *The Journal of Machine Learning Research*, vol. 10, pp. 1755–1758, 2009.
- [81] X. Zhang, F. Cheng, and S. Wang, "Spatio-temporal fusion based convolutional sequence learning for lip reading," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 713–722.
- [82] B. Xu, C. Lu, Y. Guo, and J. Wang, "Discriminative multi-modality speech recognition," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 14 433–14 442.
- [83] B. Shillingford, Y. Assael, M. W. Hoffman, T. Paine, C. Hughes, U. Prabhu, H. Liao, H. Sak, K. Rao, L. Bennett *et al.*, "Large-scale visual speech recognition," *arXiv preprint arXiv:1807.05162*, 2018.
- [84] K. Prajwal, T. Afouras, and A. Zisserman, "Sub-word level lip reading with visual attention," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 5162–5172.
- [85] T. Makino, H. Liao, Y. Assael, B. Shillingford, B. Garcia, O. Braga, and O. Siohan, "Recurrent neural network transducer for audio-visual speech recognition," in *2019 IEEE automatic speech recognition and understanding workshop (ASRU)*. IEEE, 2019, pp. 905–912.
- [86] D. Serdyuk, O. Braga, and O. Siohan, "Transformer-based video front-ends for audio-visual speech recognition for single and multi-person video," *arXiv preprint arXiv:2201.10439*, 2022.
- [87] X. Liu, E. Lakomkin, K. Vougioukas, P. Ma, H. Chen, R. Xie, M. Doulaty, N. Moritz, J. Kolar, S. Petridis *et al.*, "Synthvsr: Scaling up visual speech recognition with synthetic supervision," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 18 806–18 815.
- [88] T. Lohrenz, B. Möller, Z. Li, and T. Fingscheidt, "Relaxed attention for transformer models," in *2023 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2023, pp. 1–10.
- [89] A. Haliassos, P. Ma, R. Mira, S. Petridis, and M. Pantic, "Jointly learning visual and auditory speech representations from raw data," *arXiv preprint arXiv:2212.06246*, 2022.
- [90] T. Kudo, "Subword regularization: Improving neural network translation models with multiple subword candidates," *arXiv preprint arXiv:1804.10959*, 2018.
- [91] A. v. d. Oord, Y. Li, and O. Vinyals, "Representation learning with contrastive predictive coding," *arXiv preprint arXiv:1807.03748*, 2018.



**Jeong Hun Yeo** received the B.S. and M.S. degrees in electrical & electronic engineering from Korea Advanced Institute of Science and Technology (KAIST), Daejeon, South Korea, in 2020 and 2022, respectively. He is currently pursuing the Ph.D. degree in electrical engineering at KAIST, Daejeon, South Korea. His research interests include deep learning, image/video analysis, visual speech recognition, and multi-modal learning.



**Minsu Kim** received the B.S. degree in electrical & electronic engineering from Yonsei University, Seoul, South Korea, in 2019. He is currently pursuing the Ph.D. degree in electrical engineering at Korea Advanced Institute of Science and Technology (KAIST), Daejeon, South Korea. His research interests include deep learning, image/video analysis, audio-visual speech recognition, and multi-modal analysis.



**Jeongsoo Choi** received the B.S. degree in electrical & electronic engineering from Korea Advanced Institute of Science and Technology (KAIST), Daejeon, South Korea, in 2020. He is currently pursuing the Ph.D. degree in electrical engineering at KAIST, Daejeon, South Korea. His research interests include deep learning, image/video analysis, speech synthesis, and multi-modal analysis.



visual recognition, video action localization, machine learning, and computer vision.

**Dae Hoe Kim** received the B.S. degree from Hanyang University, Seoul, Korea, in 2010 and the M.S. and Ph.D. degrees from the Korea Advanced Institute of Science and Technology (KAIST), Daejeon, South Korea, in 2012 and 2017, respectively. He was a Senior Researcher with Agency for Defense Development, Daejeon, South Korea. Since 2019, he has been working as a Senior Researcher at the Superintelligence Creative Research Laboratory of the Electronics and Telecommunications Research Institute, Daejeon. His research interests include



University of Toronto, Canada. He is currently a Professor with the Department of Electrical Engineering and the Director of the Center for Applied Research in Artificial Intelligence (CARAI), KAIST. Among the years, he has been conducting research in a wide spectrum of image and video systems research topics. Among those topics, his interests include image processing, computer vision, visual recognition, multimodal learning, video representation/compression, and object detection. He received the Young Investigator Finalist Award of ISMRM, in 1992, and the Year's Scientist Award (Korea), in 2003. He served as an Associate Editor for IEEE Signal Processing Letters. He currently serves as an Associate Editor for IEEE Transactions on Circuits and Systems for Video Technology. He served as a TPC in many international conferences, including the program chair and organized special sessions.

**Yong Man Ro** (Senior Member, IEEE) received the B.S. degree from Yonsei University, Seoul, South Korea, and the M.S. and Ph.D. degrees from the Korea Advanced Institute of Science and Technology (KAIST), Daejeon, South Korea. He was a Researcher at Columbia University, a Visiting Researcher at the University of California at Irvine, Irvine, CA, USA, and a Research Fellow of the University of California at Berkeley, Berkeley, CA, USA. He was a Visiting Professor with the Department of Electrical and Computer Engineering,