

v1 - A Multi-view Auto Lip-reading The Kip-reader's (Progress Report)

Michael Shell
School of Electrical and
Computer Engineering
Georgia Institute of Technology
Atlanta, Georgia 30332-0250

Email: <http://www.michaelshell.org/contact.html>

Homer Simpson
Twentieth Century Fox
Springfield, USA
Email: homer@thesimpsons.com

James Kirk
and Montgomery Scott
Starfleet Academy
San Francisco, California 96678-2391
Telephone: (800) 555-1212
Fax: (888) 555-1212

Abstract—Auto lip-reading is a promising method for enhancing speech recognition, by combining the audio input with a visual. In this work we evaluate the usage of a combined method of the classical approach to ALR where handcrafted feature extraction is used, with the more modern approach where end-to-end neural network is used. In particular we investigate the combination of optical flow and key feature extraction in combination with neural networks, such as convolutional networks and Long short-term memory. The performance of each of these systems is compared to the current state-of-the-art architecture.

I. INTRODUCTION

Lip-reading is a technique to understand speech by visually interpreting the movement of the lips, face and thought¹. This technique is not limited to deaf or hear-of-hearing people but is also used by people which have a normal hearing process. A phenomena known as the McGurk effect[1] show this relation, where the interpretation of speech for the same sound is changed with the image. Just as people use lip-reading for speech recognition it is also seen to have its application in artificial, where a higher accuracy can be obtained by combining the acoustic and visual information [2]. In artificial intelligent the combination of acoustic and visual information is know as audio-visual speech recognition (AVSR) and system where only the visual information is used in commonly known as automatic lip-readig (ALR) or visual speech recognition (VSR). ALR also have other promising applications, beside the combination with acoustic information, such as visual password, silent speech interface and forensic video analysis.

The main challenge in ALR is duo to a large variation in visual factors both from the recording such as changes in illumination and camera angle[3], but also from factors that is person specific such as mouth shape an visual pronunciation.

In order to address each of these challenges, different experimental seups is proposed for ALR such as: Speaker dependent (SD) or speaker independent (SI), single-view or multi-view. SD is the simplest setting where the personal variation from the speaker is removed since data from one speaker is used both for the training and evaluation. In SI the variation from the speaker such as mouth shape and visual

pronunciation is included where unseen speakers is used for the evaluation. Single-view eliminate the change in camera angle where a multi-view setting include this dependency.

Previous work in ALR can in generally be grouped in to two, one with a classical approach and a more recent approach where deep neural networks is used.

In the classical approach the visual feature extraction is based on methods as component analysis, discrete wavelet transform, discrete cosine transform, active appearance model [5], local binary pattern [6], optical flow [7], Eigenlips [8], histograms of oriented gradients [9], internal motion histograms, motion boundary histograms and their mixed models [8, 10]. For multi-viewpoint lip-reading [11] adopt a minimum cross-pose variance analysis technique.²

Better performing systems is later obtained by the use of neural networks, where neural networks both have its application for both the feature extraction and the temporal correlation. In [2] a deep autoencoder for the feature extraction (Are there any temporal model in this architecture?). A long short-term memory (LSTM) is used in [4] to make a temporal time correlation of the features. A state-of-the art performance is then obtained in [5], where a convolutional neural network (CNN) is used for the feature extraction together with a LSTM for the temporal correlation.

In this work we propose a combined solution between the classical approach and the neural network approach, by using techniques from the classical approach for the feature extraction in combination with a CNN and LSTM. We limit our scope to focus on speaker independent and multi-view setting, where we use the OuluVS2 database[6] to evaluate our design, in relation to the error rate of word or phrase classification. This task is related to the challenges given at the ACCV 2016 workshop, multi-view lip-reading/audio-visual challenge³ (MLAC 2016).

¹Wiki

²Taken from Multi-View Automatic

³<http://ouluvs2.cse.oulu.fi/ACCVE.html>

II. RELATED WORK

III. BACKGROUND

A. Convolutional Neural Network

Convolutional neural networks (CNN) is usually used in many fields of computer vision. CNN is classically applied to classification and detection task. CNN is inspired by multi-layer perceptron containing small sub-regions of a visual field called receptive field [16]. Unlike fully connected layered network, CNN has sparse connectivity and shared weights for the purpose of increasing computational efficiency and global representation power. CNN is now the most popular and effective selection for learning visual features in computer vision and machine learning fields. We obtain a feature map at layer h with input x pixel at coordinates (i, j) as the following equation:

$$h_{ij} = a((W \cdot x)_{ij} + b) \quad (1)$$

, where weight matrix W and bias vector b is the filter of this feature map, a is activation function for non-linearities.

B. Long Short-Term Memory

Among numerous methodologies, recurrent neural network (RNN) and its variants are now common in handling sequential data with their promise of performance and ease of use. The fundamental neural network does not consider the dependency of all inputs and outputs. However, in various tasks, there exist dependency in inputs and outputs, such as sentence analysis. RNN recurrently use the previous computation result to compute the current output. Long-short term memory (LSTM) [?], one of the most popular RNN variants that is able to capture long-range dependencies, is commonly adopted. LSTM is RNN with gates, which is proposed to prevent the vanishing gradient problem, becoming more effective in dealing with long sequences.

The basic structure of LSTM unit consists of a cell state with three essential gates: input gate, forget gate and output gate. The cell controls the information storing for a long period via gates. Given an input vector \mathbf{x}_t at time step t , the formal equation for updating gates, output and cell state are defined as follows:

$$\mathbf{i}_t = \sigma(\mathbf{x}_t \mathbf{U}^i + \mathbf{h}_{t-1} \mathbf{W}^i)$$

$$\mathbf{f}_t = \sigma(\mathbf{x}_t \mathbf{U}^f + \mathbf{h}_{t-1} \mathbf{W}^f)$$

$$\mathbf{o}_t = \sigma(\mathbf{x}_t \mathbf{U}^o + \mathbf{h}_{t-1} \mathbf{W}^o)$$

$$\mathbf{c}_t = \mathbf{c}_{t-1} \circ \mathbf{f}_t + \mathbf{i}_t \circ \tanh(\mathbf{x}_t \mathbf{U}^c + \mathbf{h}_{t-1} \mathbf{W}^c)$$

$$\mathbf{h}_t = \tanh(\mathbf{c}_t) \circ \mathbf{o}_t$$

where $\mathbf{W}^i, \mathbf{W}^f, \mathbf{W}^o, \mathbf{W}^c \in \mathbb{R}^{N \times N}$, $\mathbf{U}^i, \mathbf{U}^f, \mathbf{U}^o, \mathbf{U}^c \in \mathbb{R}^{N \times N}$ are weight matrices, \mathbf{h}_t is output vector and i, f and o represent input (i), forget (f) and output (o) gates.

We are planning to apply the bi-directional LSTM, which considers the both directional(forward and backward) sequence of data. There are some research about outperforming result of bi-directional LSTM([?]).

C. Optical Flow

Optical flow estimation is one of the key problems in the computer vision. In the previous approaches, Horn and Schunk suggest the original optical flow([?]). After the original optical flow, many improvements are suggested as variational model optical flow. In our approach, we are planning to apply the optical flow to feature extraction. Optical flow represents the characteristic and it can be used as a feature of the specific part in the image.

D. Key features

IV. DATASET

To train and evaluate our approach we use the public available OuluVS2 database. The recording environment used is illustrated in figure 1, where the multi-view setting can be seen. The setup make use of five cameras located at different angels in relation to the subject. For each recorded subject three different scenarios is performed, namely the pronunciation of: (1) a sequence of ten fixed digit sequences, (2) ten daily-use short English phrases and (3) five randomly selected TIMIT sentences. Examples of each of these can be seen in table I A total of 52 different test subject is used, where each subject is pronouncing each phrase three times.

A preprocessed version is available where the recording from different angles is synchronized and the region of interest (ROI) is both segmented, with the recording of interest, and cropped to only contain the mouth region. An example of the preprocessed data can be seen in figure 2.

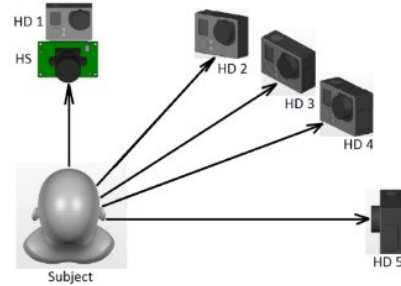


Fig. 1. Illustration of the multi-view setup used in the OuluVS2 recording, reference to figure???



Fig. 2. Example of the preprocessed data with with the region of interest.

(1) Digits	"1 7 3 5 1 6 2 6 6 7"
	"4 0 2 9 1 8 5 9 0 4"
(2) Phrases	"Thank you"
	"Have a good time"
(3) TIMIT	"Chocolate and roses never fail as a romantic gift"

TABLE I

EXAMPLES OF THE THREE DIFFERENT SCENARIOS USED

V. METHOD

In this section our approach to solve the Automatic Lip-reading is presented. The high-level architecture is first introduced, where the main components and their functions is described. The particular architectures which we like to evaluate is then presented.

A. High-level Architecture

In figure 3 an illustration of the high-level architecture can be seen with the three main components; *Visual model*, *temporal model* and classifier. In the figure the flow of information can also be seen, from the sequence of images inputted to the architecture and to the outputted class probabilities. The *visual model* is used to extract features from the input image. These features is then passed on to the *temporal model* where they are correlated in time. The time correlated features is then inputted to the classifier which is outputting the different class probabilities. For each of the three main components, different implementations is of interest where each of these is listed in figure 3. Each of the different methods is presented in the following sections.

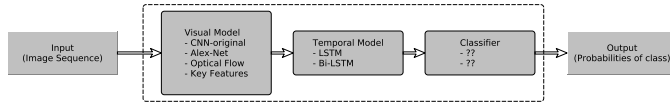


Fig. 3. High-level architecture used

B. Visual Model

For the feature extraction 4 different methods is of interest, where two of these make use of neural network's and the other two is computer vision related.

a) *CNN-original*: This visual model is corresponding to the one used in the state-of-the-art model, presented in [5]. An illustration of the model can be seen in figure 4. It make use of two convolutional layers with 16 to 256 filters in the shape of (3,3). Each convolutional layer has a successive max-pooling layer in the shape of (2,2). Lastly a fully connected layer with a dimension of 8 to 64 outputs is used.

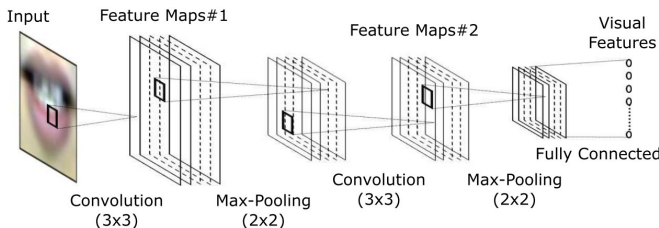


Fig. 4. CNN-original

b) *Alex-Net*: A popular neural network architecture for image recognition is an Alex-net[7]. This architecture is a deep CNN as depicted in figure 5.

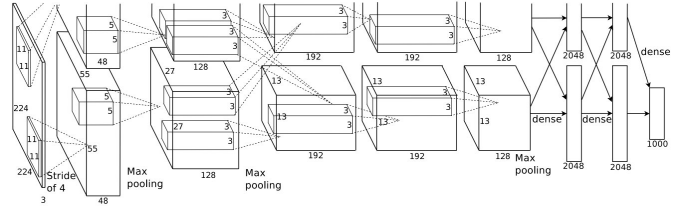


Fig. 5. Alex-net

c) *Optical flow*: For the calculation of optical flow we expect the differential method commonly denoted as the Lucas-Kanade method[8]. The method have the assumption of the flow begin essentially constant in a local neighbourhood.

d) *Key features*: For the key feature extraction a method similar to the one presented in [9] is used. The mouth is divided in to four parts, lower and upper part of the upper lip and lower and upper part of the lower lip. A vertical grid is then made, where the coordinates of the intersection points is used. An illustration of this can be seen in figure 6.

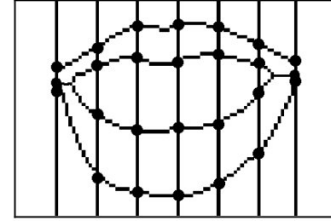


Fig. 6. Key Feature

C. Temporal model

Two different temporal models is used, which is *LSTM* and *bidirectional-LSTM*. For both of these models a two layered design is used, where the size of each layer is here varying with the number of output features from the visual model.

D. Classifier

What classifiers are we expecting to use?????

VI. RESULT

VII. CONCLUSION

The conclusion goes here.

REFERENCES

- [1] H. McGurk and J. MacDonald, "Hearing lips and seeing voices," vol. 264, pp. 746–748, Dec. 1976.
- [2] J. Ngiam, A. Khosla, M. Kim, J. Nam, H. Lee, and A. Y. Ng, "Multimodal Deep Learning," *Proceedings of The 28th International Conference on Machine Learning (ICML)*, pp. 689–696, 2011.
- [3] G. Potamianos and C. Neti, "Audio-visual speech recognition in challenging environments," *Interspeech*, pp. 1293–1296, 2003. [Online]. Available: http://www.tsi.enst.fr/~chollet/Biblio/Articles/Domaines/BIOMET/AudioVisual/Old/EURO03_CHALLENGING.pdf
- [4] M. Wand, J. Koutnik, and J. Schmidhuber, "Lipreading with long short-term memory," in *2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, mar 2016, pp. 6115–6119. [Online]. Available: <http://ieeexplore.ieee.org/document/7472852/>

- [5] D. Lee, J. Lee, and K.-e. Kim, "Multi-View Automatic Lip-Reading using Neural Network," pp. 1–14.
- [6] I. Anina, Z. Zhou, G. Zhao, and M. Pietikainen, "OuluVS2: A multi-view audiovisual database for non-rigid mouth motion analysis," *2015 11th IEEE International Conference and Workshops on Automatic Face and Gesture Recognition, FG 2015*, 2015.
- [7] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet Classification with Deep Convolutional Neural Networks," *Advances In Neural Information Processing Systems*, pp. 1–9, 2012.
- [8] B. D. Lucas and T. Kanade, "An Iterative Image Registration Technique with an Application to Stereo Vision," *Imaging*, vol. 130, no. x, pp. 674–679, 1981. [Online]. Available: <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.49.2019&rep=rep1&type=pdf>
- [9] M. Li and Y. M. Cheung, "A novel motion based lip feature extraction for lip-reading," *Proceedings - 2008 International Conference on Computational Intelligence and Security, CIS 2008*, vol. 1, pp. 361–365, 2008.