

CAN DNNs LEARN TO LIPREAD FULL SENTENCES?

George Sterpu, Christian Saam, Naomi Harte

SigmaMedia, ADAPT Centre, School of Engineering, Trinity College Dublin, Ireland

ABSTRACT

Finding visual features and suitable models for lipreading tasks that are more complex than a well-constrained vocabulary has proven challenging. This paper explores state-of-the-art Deep Neural Network architectures for lipreading based on a Sequence to Sequence Recurrent Neural Network. We report results for both hand-crafted and 2D/3D Convolutional Neural Network visual front-ends, online monotonic attention, and a joint Connectionist Temporal Classification-Sequence-to-Sequence loss. The system is evaluated on the publicly available TCD-TIMIT dataset, with 59 speakers and a vocabulary of over 6000 words. Results show a major improvement on a Hidden Markov Model framework. A fuller analysis of performance across visemes demonstrates that the network is not only learning the language model, but actually learning to lipread.

Index Terms— Lipreading, Sequence to Sequence Recurrent Neural Networks, TCD-TIMIT

1. INTRODUCTION

Automatic lipreading of continuous and large vocabulary speech is a promising technology with many applications, recovering the information in speech from a different modality than the acoustic one. The traditional approaches have largely followed early approaches in speech recognition, using hand-crafted features and Hidden Markov Models (HMM). These have been so far unsuccessful at modelling the complex patterns of visual speech [1, 2, 3, 4, 5], and several research problems, such as finding good representations, remain open in the lipreading community.

The Sequence to Sequence Recurrent Neural Network (Seq2seq RNN) architecture has seen a surge in popularity since it was first introduced in [6] for machine translation. It has an elegant formulation, makes minimal assumptions about the modelled sequences, requires less domain knowledge and has obtained competitive results on many benchmarks. Together with the Connectionist Temporal Classification (CTC) method [7], these represent the main end-to-end trainable approaches for transcribing temporal patterns. In

this work we prefer Seq2seq for its additional property of implicitly learning a language model, as CTC performance is limited by the conditional independence of its predictions [7].

Several recent advancements in machine learning have not been explored by the lipreading community. These include the monotonic attention [8] and the joint CTC-Sequence loss [9]. In addition, several successful applications in Automatic Speech Recognition (ASR) [10, 11, 12, 13], using both medium-sized (TIMIT, WSJ) and large (Google Speech Commands) datasets, give us some useful insights from a different, though correlated modality. Yet experience has shown that techniques successful for audio-only speech recognition don't automatically translate well to a lip-reading task [1, 2, 3, 4, 5]. Thus our contribution is an exploration of state of the art Seq2seq techniques within the domain of lipreading, to determine what approaches hold the greatest potential in this domain and identify where further challenges remain.

There are, to our best knowledge, only two papers in the literature to date that address the problem of lipreading at the sub-word level using DNNs. The first one [14] uses a spatio-temporal Convolutional Neural Network (CNN) and the CTC loss in order to produce a phonetic transcription of the input sentence. However, the model was applied on a low perplexity dataset, GRID [15] where it can be argued that the model can heavily rely on the predictable structure of the sentences. In addition, the CTC loss has its own shortcomings due to the independence assumption for the predicted labels. We address this by testing the algorithms on TCD-TIMIT [1], a dataset of phonetically-balanced sentences and a vocabulary of approx. 6000 words. The second paper [16] makes use of a spatial only CNN, and applies the Seq2seq network to produce sentence transcriptions at the character level. The dataset used for evaluation, LRS, is larger than TCD-TIMIT but not public, and has recently been superseded by a more challenging, public version, MV-LRS. Our work differs from these two by making predictions at the viseme level, which is a unit choice that avoids ambiguities. In this way, the language model does not have to be well trained in advance, as in [16]. In addition, we explore a wider range of architectures, such as both hand-crafted and 2D/3D CNN visual front-ends, online monotonic attention and a joint CTC-Seq2seq loss. The paper is organised as follows: In Section 2 we describe the general network architecture. Section 3 presents our experiments, and we discuss our findings in Section 4.

Supported by the ADAPT Centre for Digital Content Technology, which is funded under the SFI Research Centres Programme (Grant 13/RC/2106) and is co-funded under the European Regional Development Fund.

2. MODEL ARCHITECTURE

Our lipreading pipeline has a video processing front-end and a Seq2seq RNN, learning from variable-length videos and producing variable-length transcriptions at the viseme level. The system is trained end-to-end.

2.1. Visual front-end

The visual front-end involves segmenting the lip region from a visual stream and computing a feature vector for each frame. We consider both handcrafted features and learnt CNN-based visual representations. At this stage, we can also take advantage of the temporal dimension by appending derivatives or by using 3D convolution kernels.

2.2. Sequence modelling

Next, the extracted visual features are fed to a Seq2Seq model, which consists of two RNNs termed as the *encoder* and the *decoder*. With each input timestep, the encoder updates its internal state and produces one output. We collect all the outputs in a *memory* and retain only the final state, known as a *thought vector* that summarises the input sentence. The decoder is initialised from the thought vector and starts producing output symbols from a designated start-of-sentence token until it finally produces an end-of-sentence symbol. As the temporal dimension is warped onto the one-dimensional thought vector, the decoder is allowed to peek into the memory and soft-select the vectors that are correlated with its current internal state. This mechanism is known as *attention*, and the soft-selection temporal pattern is called *alignment*.

With speech signals, enforcing this alignment to be monotonic with respect to the encoded inputs may alleviate the problem of attending to the repetitions of a word in the same sentence. The impact may be more significant for visemes, where the number of classes is typically much lower than for phonemes or characters. In addition, scanning only a past history of the memory enables the on-line application of the lipreading system, further reducing the time complexity. We consider the implementation of [8], which was shown to outperform related strategies with a minimal loss in accuracy over the softmax attention baseline.

2.3. Training and decoding

In the training stage, the entire transcription is available to the decoder. The embedding of the ground-truth symbol gets fed at every time step, but from time to time we replace it with the previously decoded symbol in order to increase the robustness of the network to recover from mistakes. This training process implies that the predicted output transcription has an identical length with the ground-truth transcription, thus a cross-entropy loss function can be applied. In the evaluation stage, the ground-truth transcription cannot be used, and

the decoder is likely to produce a transcription of a different length. Hence, we evaluate the quality of the prediction by computing the Levenshtein edit distance with respect to the ground truth.

Combining the Seq2Seq cross entropy loss with the CTC loss could lead to several benefits. First, the CTC loss forces the encoder to better focus on the input signal, as it tends to become "lazy" due to the power of the implicitly learnt language model on the decoding side. In addition, the encoder should now learn representations that are more closely related to the class labels, as the CTC first predicts a class for each frame, and only later it merges the repeated symbols.

3. EVALUATION

3.1. Dataset

We performed our experiments on TCD-TIMIT [1], a publicly available audio-visual dataset with 59 subjects, each reciting 98 phonetically balanced sentences from a vocabulary of 6000 words, totalling around 8 hours of recordings. Sentences vary from 10 to 65 visemes in length. Evaluation was done on the speaker-dependent protocol of [1], choosing 67 and 31 sentences from each speaker for train and test respectively. The dialect-dependent sentences (name begins with SA) were removed. As in the original TIMIT database, these two sentences were common across all speakers. Early results demonstrated that the models quickly learned the structure of these sentences, giving misleading high performance. Our labels are the same viseme level transcriptions as in [1], which were obtained from a phonetic transcription by mapping phonemes into 12 viseme clusters.

3.2. Setup

Visual features. As the lip region coordinates were already provided in [1, 2], we used them to crop this region from the video frames, downsampled to 36x36 pixels and converted it to grey scale as a preprocessing step. We first considered handcrafted features and kept 44 low frequency coefficients of the lip region 2D DCT transform, plus their first two derivatives, as in [1, 2]. To check the impact of the implicitly learnt language model alone, we also present the results in the absence of a visual stream by replacing the features with zeros.

Next, we tested multiple CNN architectures on the previously cropped region, additionally using a 36x36 RGB version and a 64x64 grey one to check the benefits of color and a larger window size. Our 2D CNNs have 4 layers with 16, 32, 64 and 128 feature detectors respectively, a small 3x3 convolution kernel and rectified linear activations. After the first layer, our convolutions use a stride of 2 to reduce the dimensionality. The activations of the last layer are flattened and fully connected to a new layer of 128 units, producing our frame-wise feature vectors. The 3D CNN is of the same structure, differing only in the use of a 3x3x3 convolution kernel.

Encoder-decoder RNN. For our Seq2Seq model we start with two unidirectional recurrent layers of 128 Long Short-term Memory (LSTM) cells each, for both the encoder and the decoder. The one layer version was not performing well and we do not report these results. However, we test a one layer bidirectional LSTM (BiLSTM) version, processing the sentence both in the forward and backward directions, while maintaining the same number of parameters. Decoding was performed using a beam search strategy of width equal to 4.

Attention. Our default attention mechanism was the Luong [17] version with the energy term scaled, and we obtained significantly worse results with the more popular Bahdanau attention style [18]. We also tested the online monotonic attention strategy of [8]. To make it work, we found it was essential to turn off the pre-sigmoid noise and set the scalar bias to a negative value.

Joint CTC-Seq2seq loss. As the Seq2seq language model exhibits a strong early influence in training, we try to add a CTC loss over the encoder’s outputs, inserting a softmax layer over the vocabulary size plus 1, and training jointly with the cross-entropy loss on the decoder side. Since [9] obtained the best results for a mixing coefficient of 0.2 for the CTC loss, we only consider this case here.

3.3. Practical aspects

Input pipeline. We noticed a consistent improvement when randomly shuffling the train files with each dataset iteration. Grouping sentences of similar lengths together, a concept known as bucketing, leads to a smaller zero padding of batches, noticeably reducing the RNN processing time. Our bucket width was 15 frames, or approximately 0.5 seconds.

Regularisation. We generally obtained good results with dropout applied to the recurrent cells [19], keeping the inputs, the states and the outputs with a probability of 0.9. For the best results with the CNN architectures, we interleaved dropout layers with a rate of 50% between convolutions. We also applied L2-norm regularisation on the recurrent and the convolutional weights, scaled by 0.0001 and 0.01 respectively. We enable gradient clipping to a maximum norm of 10.0 and we also clip the LSTM cells between -10.0 and 10.0.

4. DISCUSSION AND CONCLUSION

The results of our study are shown in Table 1. We first observe a massive improvement over the HMM baseline. However, a large part is owed to the implicitly learnt RNN-based language model, as hypothesised in [12] and revealed by system **D**. In comparison, a bi-gram language model only increased the accuracy of the HMM system **A** up to 35% [2] on the same dataset, using the same DCT features. Looking at the predictions, we note that the model quickly learns to output only two visemes in an interleaved pattern, surrounded by the silence visemes delimiting the start and the end of each sen-

Table 1. Lipreading accuracy on TCD-TIMIT. The right column shows the number of iterations needed to reach convergence (or *nc* for *no convergence*).

Feature	Accuracy	Iters
A. DCT + HMM baseline [2]	31.59 %	-
B. AAM + HMM baseline [2]	25.28 %	-
C. Eigenlips + DNN-HMM [4]	46.61 %	-
D. zeros + LSTMs	45.87 %	160
E. DCT + LSTMs	61.52 %	250
F. DCT + BiLSTMs	60.72 %	180
G. E w/o attention	48.29 %	270
H. E w/ monotonic attention	61.58 %	170
I. DCT + joint CTC-Seq2seq	61.18 %	180
J. 2D CNN + LSTMs		<i>nc</i>
K. 2D CNN + BiLSTMs	66.27 %	400
L. J on RGB + joint CTC-Seq2seq	66.20%	150
M. J on 64x64 + joint CTC-Seq2seq		<i>nc</i>
N. Gray 3D CNN + LSTMs		<i>nc</i>
O. 2D CNN + joint CTC-Seq2seq	64.61%	260

tence. These correspond to the *Lips relaxed*, *narrow opening* and *Tongue up or down* classes, and together they account for 52.56% of the occurrences in TCD-TIMIT scripts. Since the scripts were phonetically balanced, the viseme distribution only reflects a natural speech pattern.

We identified this matter in all our experiments, typically taking at least 100 iterations before the predictions start to look diverse. This suggests that the language model might slow down training convergence, as the system will learn the patterns from the input signal more slowly.

The use of DCT features with a Seq2seq model led to a substantial improvement over the state of the art on the TCD-TIMIT dataset [4]. There is a noticeable boost in convergence speed from unidirectional to bidirectional LSTMs, yet it does not always translate into higher accuracy, as demonstrated by **E** and **F**. This could be explained by the fact that two single-layer networks are less powerful than a single two-layer variant. We tried another variant of two-layered bidirectional LSTM which did not improve the performance.

As noted by [16], the attention-less system **G** could not learn meaningful patterns from the input, predicting a similar transcription for most sentences. This could imply that either the temporal information vanishes during encoding, or the decoding process relies heavily on the language model. The attention-based system **E** alleviates these aspects, obtaining an absolute 13.23% improvement over this variant.

Replacing the Luong-style softmax attention with the monotonic attention of [8] maintains the performance at the same level. This is also demonstrated by the alignments in Figure 1, where the softmax attention learns to align monotonically, producing a sharp peak in the weight distribution.

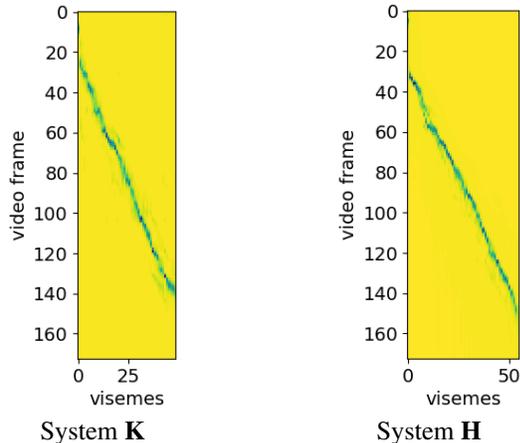


Fig. 1. Typical alignments learnt by our systems

Consequently, the enforced monotonic attention would represent a suitable choice for lipreading, further reducing the time complexity and enabling online decoding. Cited as a possible extension in [16], our benchmark shows the first successful application of online and monotonic attention to lipreading.

The use of 2D-CNN features led to an additional $\approx 5\%$ absolute improvement over the best performing DCT-based system, as is the case with system **K**. In this case, using BiLSTMs was crucial to prevent the system from getting stuck in a local minimum, as in **J**. However, our experiments on images of increased resolution (64x64) and with 3D convolutions did not reach convergence, showing the limits of a shallow CNN architecture.

The use of the joint CTC-Seq2seq loss function significantly accelerates the training process. However, in our case, the test set accuracy was lower than for the cross-entropy loss function alone. The impact of the CTC loss may be twofold. It enforces a frame-wise classification on the encoder’s outputs, which could lead to better gradients for the CNN layers. This is demonstrated by the performance achieved with systems **L** and **O**, which could not converge without the additional CTC loss. On the other hand, the two loss functions could have competing requirements for the state representation, and a proper weighting may be vital for optimal performance, as shown in [9].

On the alignments produced by the decoder we could observe that they tend to get fuzzy towards the end of the sentence, sometimes resembling to a river delta. This suggests that the thought vector is quite good at summarising the recent past, and the attention is only needed to boost the decoding of early events. We hypothesise that a different assignment of the thought vector and attention duties, where the first encodes a rather short history and the latter attends to more distant key frames, could enhance the overall performance.

We have compared the viseme confusion matrices of systems **A**, the DCT + HMM baseline, and **K**, our top perform-

Table 2. Viseme accuracy of the best DNN system **K** and relative change from HMM baseline (**A**). Visemes sorted by decreasing visibility.

Viseme	TIMIT Phonemes	Accuracy K [%]	Δ Accuracy K - A [%]
Lips to teeth	/f/ /v/	85.6	21.25
Lips puckered	/er/ /ow/ /r/ /q/ /w/ /uh/ /uw/ /axr/ /ux/	83.4	50.81
Lips together	/b/ /p/ /m/ /em/	94.8	30.40
Lips relaxed moderate opening to lips narrow-puckered	/aw/	45.7	25.90
Tongue between teeth	/dh/ /th/	58.4	27.79
Lips forward	/ch/ /jh/ /sh/ /zh/	65.4	18.26
Lips rounded	/oy/ /ao/	31.6	-8.41
Teeth Approximated	/s/ /z/	81.6	52.24
Lips relaxed narrow opening	/aa/ /ae/ /ah/ /ay/ /ey/ /ih/ /iy/ /y/ /eh/ /ax-h/ /ax/ /ix/	95.6	73.50
Tongue up or down	/d/ /l/ /n/ /t/ /el/ /nx/ /en/ /dx/	84.8	56.17
Tongue back	/g/ /k/ /ng/ /eng/	63.2	24.41
Silence	/sil/ /pcl/ /tcl/ /kcl/ /bcl/ /dcl/ /gcl/ /h#/ /#h/ /pau/ /epi/	93.6	0.21

ing DNN-based lipreading system. Table 2 shows the relative performance increase across the viseme classes for these two systems. The table also shows the TIMIT phonemes mapped to each viseme class and their visibility, or ease of observation for a human. The improvement from **A** to **K** is ubiquitous with the exception of a single viseme corresponding to the *Lips rounded* shape. This viseme is most frequently confused with the *Lips relaxed narrow opening* viseme, suggesting that it is difficult even for the CNN to learn features that disambiguate them. Lower improvements are seen for *Lips forward* and *Tongue back*. The frontal view used as input does not capture any depth information, however the database includes a second view at 30° which could be useful for such visemes.

Overall, the Seq2seq model greatly outperforms HMM and hybrid DNN-HMM systems even without CNN-based feature extraction. The fully neural architectures achieved the highest accuracies in our experiments. Additionally, the use of the joint loss function boosted the training convergence and enabled learning visual features on higher dimensional inputs. Lastly, we demonstrate the efficiency of online monotonic attention on this task, a necessary step towards online decoding.

5. ACKNOWLEDGEMENTS

We are grateful to Eugene Brevdo, Marco Forte, Oriol Vinyals and Li Deng for their helpful comments and suggestions.

6. REFERENCES

- [1] Naomi Harte and Eoin Gillen, “TCD-TIMIT: An audio-visual corpus of continuous speech,” *IEEE Transactions on Multimedia*, vol. 17, no. 5, pp. 603–615, May 2015.
- [2] George Sterpu and Naomi Harte, “Towards lipreading sentences using active appearance models,” in *AVSP*, Stockholm, Sweden, August 2017.
- [3] Kwanchiva Thangthai and Richard Harvey, “Improving computer lipreading via dnn sequence discriminative training techniques,” in *The 18th Annual Conference of the International Speech Communication Association Interspeech 2017*. May 2017.
- [4] Kwanchiva Thangthai, Helen L Bear, and Richard Harvey, “Comparing phonemes and visemes with dnn-based lipreading,” in *Workshop on Lip-Reading using deep learning methods*, 2017, BMVC 2017.
- [5] G. Potamianos, C. Neti, G. Gravier, A. Garg, and A. W. Senior, “Recent advances in the automatic recognition of audiovisual speech,” *Proceedings of the IEEE*, vol. 91, no. 9, pp. 1306–1326, Sept 2003.
- [6] Ilya Sutskever, Oriol Vinyals, and Quoc V Le, “Sequence to sequence learning with neural networks,” in *Advances in Neural Information Processing Systems 27*, Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, Eds., pp. 3104–3112. Curran Associates, Inc., 2014.
- [7] Alex Graves, Santiago Fernández, Faustino Gomez, and Jürgen Schmidhuber, “Connectionist temporal classification: Labelling unsegmented sequence data with recurrent neural networks,” in *Proceedings of the 23rd International Conference on Machine Learning*, New York, NY, USA, 2006, ICML ’06, pp. 369–376, ACM.
- [8] Colin Raffel, Minh-Thang Luong, Peter J. Liu, Ron J. Weiss, and Douglas Eck, “Online and linear-time attention by enforcing monotonic alignments,” in *Proceedings of the 34th International Conference on Machine Learning*, Doina Precup and Yee Whye Teh, Eds., International Convention Centre, Sydney, Australia, 06–11 Aug 2017, vol. 70 of *Proceedings of Machine Learning Research*, pp. 2837–2846, PMLR.
- [9] S. Kim, T. Hori, and S. Watanabe, “Joint ctc-attention based end-to-end speech recognition using multi-task learning,” in *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, March 2017, pp. 4835–4839.
- [10] Jan Chorowski, Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio, *End-to-end continuous speech recognition using attention-based recurrent NN: First results*, 2014.
- [11] Jan K Chorowski, Dzmitry Bahdanau, Dzmitry Serdyuk, Kyunghyun Cho, and Yoshua Bengio, “Attention-based models for speech recognition,” in *Advances in Neural Information Processing Systems 28*, C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, Eds., pp. 577–585. Curran Associates, Inc., 2015.
- [12] D. Bahdanau, J. Chorowski, D. Serdyuk, P. Brakel, and Y. Bengio, “End-to-end attention-based large vocabulary speech recognition,” in *2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, March 2016, pp. 4945–4949.
- [13] W. Chan, N. Jaitly, Q. Le, and O. Vinyals, “Listen, attend and spell: A neural network for large vocabulary conversational speech recognition,” in *2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, March 2016, pp. 4960–4964.
- [14] Yannis M. Assael, Brendan Shillingford, Shimon Whiteson, and Nando de Freitas, “Lipnet: Sentence-level lipreading,” vol. abs/1611.01599, 2016.
- [15] Martin Cooke, Jon Barker, Stuart Cunningham, and Xu 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.
- [16] Joon Son Chung, Andrew Senior, Oriol Vinyals, and Andrew Zisserman, “Lip reading sentences in the wild,” in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [17] Minh-Thang Luong, Hieu Pham, and Christopher D. Manning, “Effective approaches to attention-based neural machine translation,” *CoRR*, vol. abs/1508.04025, 2015.
- [18] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio, “Neural machine translation by jointly learning to align and translate,” *CoRR*, vol. abs/1409.0473, 2014.
- [19] Wojciech Zaremba, Ilya Sutskever, and Oriol Vinyals, “Recurrent neural network regularization,” *CoRR*, vol. abs/1409.2329, 2014.