

Towards pathological speech synthesis from articulation

Bence Halpern¹², Rob J. J. H. van Son¹², Michiel W. M. van den Brekel¹²

¹NKI-AVL, Amsterdam

²ACLC, University of Amsterdam, The Netherlands

b.halpern@nki.nl, r.v.son@nki.nl

Abstract

This paper presents a technique to synthesise speech that is pathological on the articulatory level. This technique combines a vocoder with a speaker-independent articulatory to acoustic neural network using electromagnetic articulography recordings from the three largest articulatory datasets. A visualisation technique is described to shed light on what these neural networks learn. It is shown that speech with manipulated articulation can readily be synthesised. However, the baseline quality of the vocoder used turned out to be low, accounting for 88% of the output’s variance. The baseline vocoder quality is currently too low to evaluate the pathological aspects of the manipulated speech.

Index Terms: computational paralinguistics, articulatory-to-acoustic speech synthesis, deep learning, pathological speech

1. Introduction

Synthesising pathological speech could enable many potential applications. One of them is the creation of synthetic data which could be used as training data for pathological speech detection. It could be also used as a clinical tool for counselling patients about post-treatment speech outcomes. Additionally, a deeper understanding of the articulatory to speech relationship could offer improvement in speech therapy tools.

Understanding how articulation affects speech is a central question in speech research. The source-filter model was one of the first models to tackle this problem by discovering that speech production could be described by the geometry of the vocal tract and the glottal wave [1] [2]. A significant drawback of this method however is that it does not model the movement of articulators.

Recently, deep learning methods became popular to model articulation. These methods use a measurement tool, called electromagnetic articulography (EMA) to obtain articulation data [3] [4] [5] along with recurrent neural networks, which are function approximators that are able to deal with the sequential nature of data [6]. Data-driven methods became of interest also in real-time speech synthesis, using a technique called permanent magnetic articulography by [7], resulting in intelligible speech. The conclusion of these endeavours were that while it is possible to predict some of the pitch from articulation, the quality suffers. However, it is possible to obtain satisfactory values for the cepstrum.

This indicates, that this technique could be a good candidate for synthesising pathological speech where the pitch of the voice is natural. For example, in the case of oral tumours, the laryngeal function remains intact, meaning the pitch remains natural. Thus, it is proposed that the F_0 could be simply obtained through vocoder analysis and only predict the cepstral values through articulation to model pathologies.

In this paper, a technique is described which combines healthy speech from the three largest articulatory datasets, MNGU0 [8], MOCHA-TIMIT [9] and TORGO [10], in order to create a general speaker-independent articulatory to acoustic model and introduce a framework for pathological speech synthesis.

The main contributions of this paper are,

- a description of a method for speaker-independent MFCC prediction in Section 2
- a technique to incorporate articulation domain-knowledge into pathological speech synthesis in Section 2.4
- a discussion of the current limitations of this framework in Section 3.4
- an attempt to shed light on what these neural networks might learn in Section 3.3

The code of the experiments is available as a Github repository online. [11].

2. Method

2.1. Dataset preprocessing

2.1.1. Electrode preprocessing

Articulators recorded are slightly different in each dataset, meaning particular attention has to be paid to align these. An example of EMA recording locations are shown on Figure 1. Seven electrodes were used for this experiment out of the total eight, Table 1 includes the alignment of the channels that were used. This ensures that each input channel records reasonably similar information, meaning that the channels should have similar variance. These are then standardised on a per speaker basis. These steps alleviate some of the speaker-wise deviations, but does not alleviate problems if an electrode falls off during the experiment or if an electrode needs to be changed.

In the case of the TORGO dataset, some of the channels contained artifacts, these have been excluded. The signal to noise ratio in these spiky regions was low enough to affect training.

Previously [12], the effect of delay on the output signal were investigated. It has been found that delay is beneficial for the case of causal models. In Section 2.3, a bidirectional recurrent model will be introduced which is not causal, meaning there is no need for delays.

2.1.2. Speech data processing

The total dataset contains speech from six British male and three British female speakers, with a total of 6117 utterances, approximately 10 hours of recorded speech with a sampling frequency of 16kHz. Only the healthy speech has been included from the TORGO dataset. There are 1263 utterances from the

Table 1: Articulatory information recorded in datasets

MNGU0	MOCHA-TIMIT	TORG0
Tongue dorsum (T3)	Tongue dorsum (T3)	Tongue back
Tongue blades (T2)	Tongue blades (T2)	Tongue middle
Tongue tip (T1)	Tongue tip (T1)	Tongue tip
Lower incisor (T3)	Jaw	Lower incisor
Upper incisor	Nose	Upper incisor
Upper lip	Upper lip	Upper lip
Lower lip	Lower lip	Lower lip

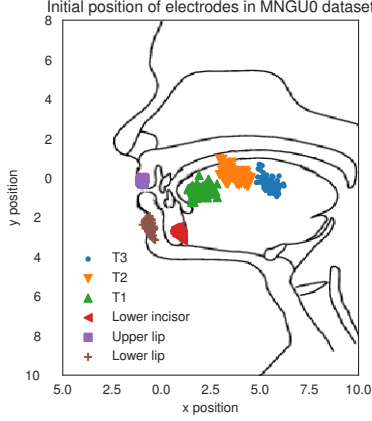


Figure 1: The visualisation of electrode locations for 300 samples from the MNGU0 dataset at their initial position. The drawing is only indicative of real positions.

MNGU0, 920 from the MOCHA-TIMIT and 3934 from the TORG0 dataset.

Vocoder features were extracted with the PyWORLD vocoder [13] and compressed with the PySPTK toolkit [14]. The period between consecutive frames was 5 milliseconds. The resulting 40 MFCC and 1 power parameters were used to generate static and delta parameters, resulting in 82 parameters for the training. As the first step of the MFCC extraction $\alpha = 0.42$ were used as a pre-emphasis coefficient. The PyWORLD vocoder also provides the F_0 and BAP values, which were not used for training.

Preprocessing techniques of previous publications are summarised in Table 2.

2.1.3. Sampling

The sampling frequency of the original EMA signals was 500 Hz, however the MNGU0 was provided to us downsampled to 200 Hz. To match this frequency, the sampling frequency of the other datasets was also downsampled to 200 Hz.

For the MNGU0 dataset, NaN (not a number) values occurred when the measurement precision was low. These values were interpolated linearly.

To ease training, the input signals were either truncated or padded so there were a total of $T = 1000$ samples for each training example. For input signals which are shorter, it is assumed that the last part is silence, so it is padded with the last element. These are not propagated back during training, to avoid the neural network making inference based on the length of the last element.

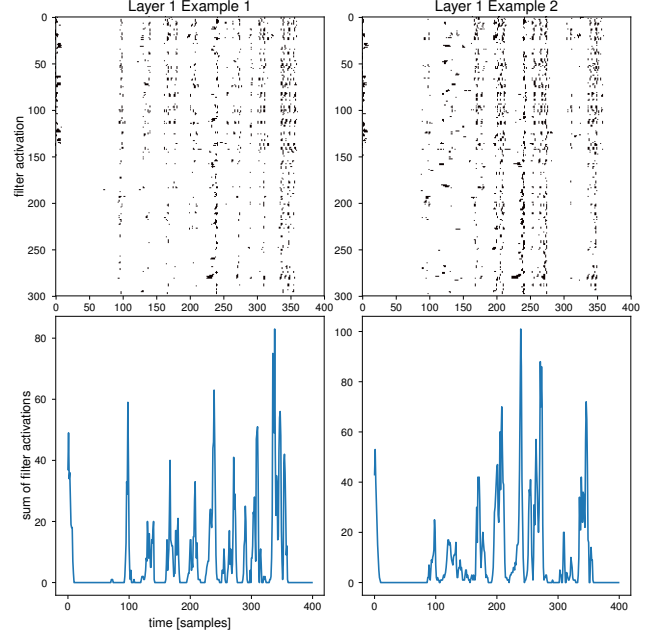


Figure 2: Thresholded Sobel mask of activations indicates that a boundary phenomena is learned by the neural network.

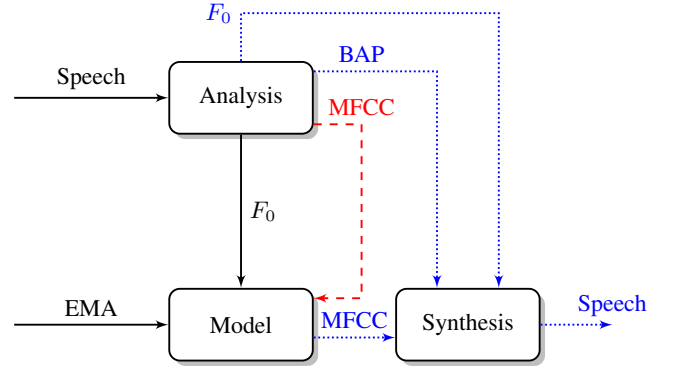


Figure 3: Red dashed line indicates training-only setup, and blue dotted lines indicate inference for speech synthesis. Best viewed in colour.

2.1.4. Fundamental frequency interpolation

In this framework, the F_0 is also used for prediction, in order to mitigate the error due to the residual pitch information in the MFCCs. Thus, it needs to be processed to be used by the neural network. Previously, it has been found beneficial to take the logarithm of the pitch to obtain a continuous F_0 curve in the prediction setting. When the logarithm is not defined, linear interpolation has been done. [7] An alternative method with exponential interpolation is described in [15].

2.2. Synthesis setup

The setup for inference and training can be seen in Figure 3. In the training setup, only the MFCCs are given. The pitch and band aperiodicities (BAP) are directly fed to the vocoder during synthesis time, as these don't contain information about articulation.

Table 2: Comparison of preprocessing techniques of some previous studies. In the case of EMA data, there is a clear consensus of 40 MFCC channels.

Author	Liu	Taguchi	Gonzalez
EMA/PMA	EMA	EMA	PMA
MFCC	40 + 1	40 + 1	24 + 1
Delta	No	Yes	Yes
EMA sampling	200 Hz	200 Hz	100 Hz*
Standardisation	Yes	Yes	Yes
Smoothing	No	Yes	No
Vocoder	STRAIGHT [16]	WORLD	STRAIGHT

*upsampled to 200 Hz to match analysis rate

2.3. Neural network design

In this paper, a recurrent neural network will be used in order to approximate the articulatory to acoustic mapping. To construct this speaker-independent network, previous speaker-dependent architectures have been studied, to conclude on an appropriate design.

It has been concluded that the optimisation schedules of the publications studied were very different, most likely due to the problem of vanishing and exploding gradients in recurrent neural networks. This is the reason why [4] used incremental training along with gradient clipping, and probably the reason why [5] used a learning rate scheduler. However, [7] used Adam optimiser [17] which is known to manage both of these problems with the minor disadvantage of the absence of good convergence guarantees. The fact that Adam was able to obtain similar results without careful parameter-tuning indicated that it will be an appropriate candidate as an optimiser for our model.

Previous publications on speaker-dependent models reported best performance on bidirectional architectures, however it was unclear whether BLSTM or BGru architectures are better. Also, [4] resorted to a combination of fully connected and recurrent layers.

In order to determine the best architecture, a pilot study has been performed on all three neural networks which are summarised in Table 4, however all of them were trained with an Adam optimiser for the reasons mentioned above, and a learning rate of 0.003, and a batch size of 100 without noise on MNGU0 dataset. The best performing neural network was then trained on the entire dataset.

Table 3: Performance of speaker-independent articulatory to acoustic neural network for 10-fold cross-validation with 95 % confidence intervals. In the TORGO dataset, different recording sessions were kept in different datasets.

Dataset	Multi-speaker MCD	Single-speaker MCD
Combined result	5.31 ± 0.09 dB	N/A
MNGU0	5.93 ± 0.31 dB	4.77 dB
Female MOCHA-TIMIT	5.02 ± 0.06 dB	5.23 dB
Male MOCHA-TIMIT	4.06 ± 0.06 dB	5.83 dB
TORGO Female 1	4.48 ± 0.03 dB	N/A
TORGO Female 2A	4.23 ± 0.06 dB	N/A
TORGO Female 2B	4.81 ± 0.14 dB	N/A
TORGO Female 3	4.94 ± 0.09 dB	N/A
TORGO Male 1A	4.64 ± 0.04 dB	N/A
TORGO Male 1B	4.70 ± 0.05 dB	N/A
TORGO Male 2A	4.62 ± 0.04 dB	N/A
TORGO Male 2B	15 ± 0.86 dB	N/A
TORGO Male 3	4.63 ± 0.11 dB	N/A
TORGO Male 4	4.85 ± 0.12 dB	N/A

Table 4: Comparison of different training methods used in previous publications with the results of the pilot study using held-out validation. The method described in the paper of Gonzalez performed best.

Author	Liu	Taguchi	Gonzalez
BLSTM layers	4 (128)	2 (256)	4 (150) GRU
Dense layers	1	3+1	1
Regularisation	No	LayerNorm	Noise 0.05
Dropout	No	Yes (50 %)**	No
Optimiser	SGD	RMSProp	Adam
Learning rate	0.01*	0.01	0.003
Gradient clipping	No	5	No
Early stopping	Yes	Yes	Yes
MLPG [18]	No	Yes	Yes
Maximum epochs	32	N/A	100
Batch size	N/A	8	100
Incremental training	No	Yes	Yes
MCD***	4.84 dB	7.28 dB	4.77 dB

* with decay after Epoch 11 ** from author communication

*** results of our training with Adam optimiser

For training the mean squared error loss function was used, and for evaluation the Mel cepstral distortion (MCD) have been employed. [19]

For the speaker-independent experiments, ten fold cross-validation was performed to estimate the out-of-sample generalisation capability of the neural networks.

2.4. Articulatory space modification

Using this framework, the problem of making pathological speech can be traded for the problem of making pathological articulation and feeding pathological articulation through the neural network.

Consider the EMA signal $\mathbf{x} \in \mathbb{R}^{t \times m}$, where t is the number of samples and m are the number of electrode channels recorded. The problem of articulatory space modification is about finding $\hat{\mathbf{x}} := f(\mathbf{x})$. The pathological framework presented effectively allows the interested people to design their own functions.

The aim here is to give some ideas to the reader. In some cases it happens that a certain articulator cannot reach a certain target or the articulator cannot move at all.

In that case, it is possible to model this by taking,

$$\hat{x}_{i,k} = \begin{cases} c & \text{if } x_{i,k} > c \\ -c & \text{if } x_{i,k} < -c \\ x_{i,k} & \text{otherwise,} \end{cases}$$

where $k \in [1, m]$ is the index of an articulator channel. The first method we used simply fixes $c = 0$.

It is often the case that pathological domain knowledge is known about the velocity, or the acceleration of the tongue. In that case, it is possible to model it as difference equation. For example, if it is assumed that the speed cannot exceed a certain $c > 0$,

$$\hat{\dot{x}}_{i,k} = \begin{cases} c & \text{if } x_{i+1,k} - x_{i,k} > c \\ -c & \text{if } x_{i+1,k} - x_{i,k} < -c \\ x_{i+1,k} - x_{i,k} & \text{otherwise,} \end{cases}$$

could be used to model the articulation. Finally, the pathological signal could be obtained using,

$$\hat{x}_{i,m} = \sum_{j=1}^i x_{j,m}.$$

3. Results and discussion

3.1. Pilot study

The pilot study results are summarised in Table 4.

Based on our training, it seems clear that the GRU architecture was superior to an LSTM architecture in our case, when used with an Adam optimiser. There is no general consensus whether GRU or LSTM is better for particular datasets. [20]

3.2. Prediction of MFCC values

The prediction results for the MFCC values are summarised in Table 3. Results were all in similar range as previously reported values for speaker-dependent datasets, and in our framework the speaker-independent architectures clearly performed better than the speaker-dependent architectures on the MOCHA-TIMIT datasets.

3.3. What do these neural networks learn?

Recently, there have been many advancements in understanding what neural networks learn. Convolutional neural networks can be analysed via conventional methods in filter analysis [21], Classification neural networks can propagate back gradients to find the most important inputs for the prediction [22]. These techniques are not applicable for recurrent neural networks in a regression context, so we resort to exploring the temporal activations of the layers.

To make these intelligible, a Sobel mask is thresholded to find peaks in the activations. On Figure 2. We observe that line-like boundaries are learned, and their duration indicate these might approximate phone to word syllable representations.

3.4. The current limitations of the synthesis

The quality of the synthesised speech depends on the quality of vocoding and the quality of prediction.

According to our observations, the quality is bounded more by the quality of the vocoder, than the synthesis itself. It has been found that mean squared error (MSE) between the vocoder resynthesised speech and the predicted speech is 11 on the MNGU0 dataset. The MSE between the analysis-resynthesis and the vocoder is 80. That means that 88% of the variance is due to the vocoder analysis-resynthesis. This indicates that future improvements should focus on better vocoding rather than better acoustic mapping.

Despite this performance loss, we wanted to investigate the sensitivity of the model performance to more training data. The neural network was retrained using incrementally more training data in twenty percent batches. The MSE was calculated at all epochs of training for the validation set, which can be seen on Figure 4.

A paired t-test confirmed that there is a statistically significant ($p < 10^{-10}$) improvement with each addition of the training data, meaning more data would significantly help generalisation performance.

3.5. Pathological speech examples

Synthesised pathological speech examples can be found on the webpage of the author, see [23] and the media files. Informal discussions with speech language pathologist indeed confirmed that some of these synthesised samples resemble dysarthric or disordered speech, but these simple heuristics usually don't incorporate enough knowledge about a particular pathology to show it consistently.

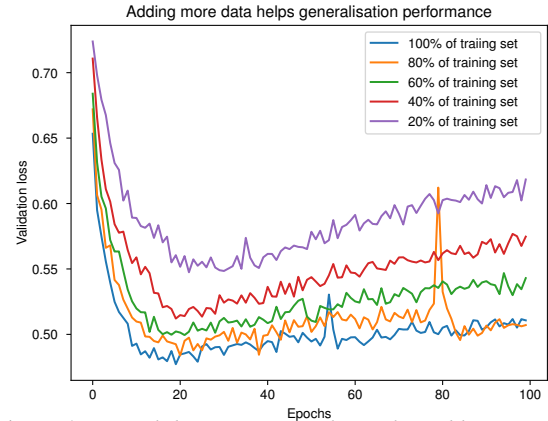


Figure 4: Partial data retraining shows that adding more data would decrease loss

This confirms that the framework can be used as a platform to implement a pathological articulation by reflecting the physiological changes in the articulatory space. This is most easily done in a data-driven fashion, by recording healthy and pathological articulation and creating a pathology-dependent mapping in the articulatory space. Because it is a low dimensional space, this is a much easier problem than learning the same mapping in the high dimensional cepstral space.

4. Conclusions

This paper is a proof of concept that it is possible to make pathological speech by incorporating changes in articulatory domain. Benchmarks have been also established and an open source repository is also available in order to reproduce these results.

Further work needs to be done on improving vocoder speech quality and creating models which consistently show a certain pathology.

5. Acknowledgements

This project has received funding from the European Union's Horizon 2020 research and innovation programme under Marie Skłodowska-Curie grant agreement No 766287. This work was carried out on the Dutch national e-infrastructure with the support of SURF Cooperative. The Department of Head and Neck Oncology and surgery of the Netherlands Cancer Institute receives a research grant from Atos Medical (Malm, Sweden), which contributes to the existing infrastructure for quality of life research.

6. References

- [1] J. Benesty, M. M. Sondhi, Y. Huang, and S. Greenberg, *Springer Handbook of Speech Processing*, 2009, vol. 126, no. 4. [Online]. Available: <http://scitation.aip.org/content/asa/journal/jasa/126/4/10.1121/1.3203918>
- [2] G. Fant, "The source filter concept in voice production," *Quarterly Progress and Status Report*, 1981. [Online]. Available: <http://www.speech.kth.se/gpsr>
- [3] S. Aryal and R. Gutierrez-Osuna, "Data driven articulatory synthesis with deep neural networks," *Computer Speech and Language*, vol. 36, pp. 260–273, 2016. [Online]. Available: <http://dx.doi.org/10.1016/j.csl.2015.02.003>
- [4] F. Taguchi and T. Kaburagi, "Articulatory-to-speech conversion using bi-directional long short-term memory." [Online]. Avail-

able: https://www.isca-speech.org/archive/Interspeech{_}2018/pdfs/0999.pdf

- [5] Z. C. Liu, Z. H. Ling, and L. R. Dai, "Articulatory-to-acoustic conversion using BLSTM-RNNs with augmented input representation," *Speech Communication*, vol. 99, no. February, pp. 161–172, 2018. [Online]. Available: <https://doi.org/10.1016/j.specom.2018.02.008>
- [6] S. Hochreiter and J. Schmidhuber, "Long Short-Term Memory," *Neural Computation*, 1997.
- [7] J. A. Gonzalez, L. A. Cheah, P. D. Green, J. M. Gilbert, S. R. Ell, R. K. Moore, and E. Holdsworth, "Evaluation of a silent speech interface based on magnetic sensing and deep learning for a phonetically rich vocabulary," *Proceedings of the Annual Conference of the International Speech Communication Association, INTERSPEECH*, vol. 2017-Augus, pp. 3986–3990, 2017.
- [8] K. Richmond, P. Hoole, and S. King, "Announcing the electromagnetic articulography (day 1) subset of the mngu0 articulatory corpus," in *Proceedings of the Annual Conference of the International Speech Communication Association, INTERSPEECH*, 2011.
- [9] A. Wrench, "The MOCHA-TIMIT articulatory database." 1999. [Online]. Available: <http://www.speech.kth.se/gpsr>
- [10] F. Rudzicz, A. K. Namasivayam, and T. Wolff, "The TORGO database of acoustic and articulatory speech from speakers with dysarthria," *Language Resources and Evaluation*, 2012.
- [11] B. Halpern, "Articulatory vocoder." [Online]. Available: <http://github.com/karkirole/vocoder-clean>
- [12] J. A. Gonzalez, L. A. Cheah, J. M. Gilbert, J. Bai, S. R. Ell, P. D. Green, and R. K. Moore, "A silent speech system based on permanent magnet articulography and direct synthesis," *Computer Speech and Language*, 2016.
- [13] M. Morise, F. Yokomori, and K. Ozawa, "WORLD: A vocoder-based high-quality speech synthesis system for real-time applications," *IEICE Transactions on Information and Systems*, 2016.
- [14] "PySPTK toolkit." [Online]. Available: <http://github.com/r9y9/pysptk>
- [15] C. J. Chen, R. A. Gopinath, M. D. Monkowski, M. A. Picheny, and K. Shen, "New methods in continuous Mandarin speech recognition." *5th European Conference on Speech Communication and Technology (EUROSPEECH)*, 1997.
- [16] H. Kawahara, "STRAIGHT, exploitation of the other aspect of VOCODER: Perceptually isomorphic decomposition of speech sounds," *Acoustical Science and Technology*, 2006.
- [17] D. P. Kingma and J. L. Ba, "Adam: A method for stochastic gradient descent," *ICLR: International Conference on Learning Representations*, 2015.
- [18] Z. Wu and S. King, "Improving Trajectory Modelling for DNN-Based Speech Synthesis by Using Stacked Bottleneck Features and Minimum Generation Error Training," *IEEE/ACM Transactions on Audio Speech and Language Processing*, 2016.
- [19] R. F. Kubichek, "Mel-cepstral Distance Measure for Objective Speech Quality Assessment," *IEEE Pacific Rim Conference on Communications, Computers and Signal Processing*, pp. 125–128, 1993.
- [20] R. Jozefowicz, W. Zaremba, and I. Sutskever, "An Empirical Exploration of Recurrent Network Architectures," *JMLR*, 2015.
- [21] D. Palaz, M. Magimai-Doss, and R. Collobert, "Analysis of CNN-based speech recognition system using raw speech as input," in *Proceedings of the Annual Conference of the International Speech Communication Association, INTERSPEECH*, 2015.
- [22] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, "Grad-CAM: Visual Explanations from Deep Networks via Gradient-Based Localization," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017.
- [23] B. Halpern, "Speech Examples." [Online]. Available: <http://karkirole.github.io/paper1>