

The Impact of Deep Learning on Speech Synthesis with Embedded or Mobile Devices

Hannes Bohnengel

Technical University of Munich

hannes.bohnengel@tum.de

ABSTRACT

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultricies augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetur adipiscing elit.

KEYWORDS

NECESSARY??? Deep Learning, Deep Neural Network, Embedded System, Mobile Device, Speech Synthesis, Text-to-Speech

1 INTRODUCTION

Virtual personal assistants (VPA) like Siri, Cortana or Google Now start having a huge impact on the way of interacting with electronic devices like smartphones or notebooks. Up to now the VPAs help only with rather simple tasks like search queries, starting phone calls or setting a clock, but according to a recent survey from the IT research firm Gartner [3], this will change in the near future. With the Facebook Messenger it is already possible to make purchases or to order an Uber car and new use cases are expected soon. The survey also states, that through the vast increase of devices in the scope of the Internet of Things (IoT) the way of interacting with machines will go towards minimal or zero touch. Instead of interacting through common touch-displays or buttons, the user simply speaks to the device, like to another person. To enable this, both Automatic Speech Recognition (ASR) and speech synthesis are essential technologies.

In this paper I will only focus on the speech synthesis part. A widely spread technique to synthesize human speech from a given text or from linguistic descriptions is Statistical Parametric Speech Synthesis (SPSS); also referred to as Statistical Parametric Speech Generation (SPSG) [10]. This technique is based on the usage of Hidden Markov Models (HMMs). Black *et al.* [4] show that it has several advantages over its predecessor, the concatenative speech synthesis, for example the flexibility in changing voice characteristics and a smaller memory footprint. However the quality of the generated speech still has potential for improvement. Due to over-smoothing the voice sounds muffled in comparison to natural speech.

This is where recent achievements in deep learning come in. Deep learning is usually referred to as a class of machine learning techniques that achieve tasks like feature extraction or pattern analysis by using many connected layers of non-linear information processing [7, 10]. Since 2006 advances in the training algorithms of Deep Neural Networks (DNNs) have enabled the field of deep learning applications to emerge [5]. Most machine learning models until then had used shallow structures, like for example HMMs, Gaussian Mixture Models (GMMs), Conditional Random

Fields (CRFs) or Support Vector Machines (SVMs). In these structures only one layer is responsible for generating features out of the raw input signals. While achieving quite good results with rather simple problems, they reach their limit when it comes to more complex tasks like processing human language or natural images [7]. In the tutorial survey [7] the author also states four different approaches to improve speech synthesis through deep learning models, whereof three are dealing with SPSS. One of those three approaches is described in [14], where the authors implemented a part of the speech synthesis system by using a DNN and observed an improved performance in predicting output features. In [8] a more general approach is conducted by investigating what effects the deployment of a DNN on different parts of the SPSS system has. An improvement of the naturalness of the generated speech was one of the main results.

For implementing speech synthesis on resource-constrained devices like smartphones or tablets, SPSS is considered the best solution due to the tradeoff between voice quality and acceptable footprint size [12]. Since the computational costs of SPSS are often high, some optimization steps like applying fewer conditional calls are conducted in [12] to make the HMM-based speech synthesis technique SPSS more suitable for mobile devices. Going one step further, in [5] an approach to adapt SPSS for embedded devices by using a deep learning model, an Auto Encoder (AE), is employed. Four tasks (syllabification, phonetic transcription, part-of-speech tagging and lexical stress prediction) are examined and tested with the use of this deep learning model. As results the authors highlight highly reduced model sizes, higher training times, very close performance and a similar run time in comparison to the state-of-the-art models. This shows that the usage of deep learning models for speech synthesis on embedded systems is a reasonable step, not only to improve performance and voice quality, but also towards the independability on online databases for speech synthesis.

The remaining paper is structured as follows: Section 2 first states the motivation, why speech synthesis is a useful technology. Then it describes the conventional approach without deep learning models for speech synthesis and gives an overview of advantages and drawbacks of the used models and techniques. This is followed by a brief explanation of the probably most common used technique SPSS, where the paper [4] has been chosen as commonly cited reference. Thereafter two possibilities how SPSS can be improved by deploying deep learning models are characterized, whereof the papers [8, 14] are reviewed. In Section 4 the motivation, why speech synthesis is important on embedded or mobile devices is given, followed by two examples on how speech synthesis can be implemented on an embedded system, once without [12] and once with deep learning models [5]. Finally Section 5 summarizes the essential points of this paper and gives some future directions.

2 CONVENTIONAL SPEECH SYNTHESIS

2.1 Motivation & Approaches

Speech synthesis has emerged over the last 10 years due to a vast contribution by the global community of researchers and the increasing computational power for data processing. Its quality and naturalness has increased steadily and different approaches have been developed so far [11]. The typical applications like navigation systems in cars or telephone-based dialogue systems are nowadays widely established. But also as reading aid for visually impaired people [1] or as in the case of the famous scientist Stephen Hawking, who has been using a synthesized voice to communicate since 1997 [2], speech synthesis has proven to be very useful. Another very interesting application of speech synthesis is shown in [XX]. The author proposed to introduce synthetic speech as means of communication between pilots, since there have been many accidents due to misunderstandings at radio-based communication.

According to [9] speech synthesis can be divided into three types: Canned speech, Context-to-Speech (CTS) and Text-to-Speech (TTS). Canned speech more or less is the replay of prerecorded spoken sentences or words with none or very little adjustments. A typical example are the announcements on train stations. Because of the high effort of recording everything (almost) exactly as it is replayed, this approach is limited to only a few simple applications. With CTS the waveform is generated out of a linguistic description without any information of the respective text. In this way no natural language processing is required, but nevertheless CTS nowadays has not made any important impact. The last and most promising type is TTS.

A TTS system consists of a Natural Language Processing (NLP) part, where the text is analysed and the word and sentence structure and accents are extracted. In the next step these accents are used to generate the prosody of the given text like duration, intensity and pitch. Then the created phonetic representations with prosody information are stringed together to a continuous stream of signal parameters. The last task, the speech generation, uses this stream to generate the respective waveform. This function block can be implemented in different ways. In [9] three general approaches are named as follows: Parametric Speech Synthesis (formant-based synthesis), Concatenative Speech Synthesis (unit-selection synthesis) and Statistical Parametric Speech Synthesis (HMM-based synthesis). The methods in brackets are the respective implementations, which are most commonly used.

The formant-based synthesis is the oldest approach. To generate a voice waveform, an excitation signal is fed into multiple formant filters which describe the characteristics of the human vocal tract. The output of the filters then forms the voice waveform. This technique is the only one which does not need any recorded speech, but instead generates the synthesized voice only by modeling the human vocal tract. The quality of the generated voice is the lowest in comparison to the other techniques, but therefore formant-based synthesizers have the smallest footprint and the voice characteristics can easily be modified by just changing the filter parameters [9].

With the development of concatenative speech synthesizers the quality of the generated speech improved tremendously. Very similar to CTS prerecorded speech is used as reference. Very basically said, the recorded speech is divided into units and these units are then stringed together to form the new speech signal according to a given text. Hereby the chosen size of the units determines both the footprint size and the voice quality. With larger units a higher voice quality can be achieved, but this also results in a much

larger database. The challenge in unit-selection is to ensure, that the transitions between the units are as natural as possible [9]. In Section 2.2 some concepts on how to achieve this as well as the third implementation, the HMM-based synthesis, a specific instance of SPSS are described in detail.

2.2 HMM-based Speech Synthesis

In this section the HMM-based, the most recent approach for speech synthesis will be described further and both advantages and drawbacks compared to unit-selection synthesis will be highlighted. Therefore the work of Black *et al.* [4] will be taken as reference.

The quality of unit-selection synthesis directly depends from the quality of the prerecorded speech. But even with a database of excellent quality sporadic errors can still not be avoided totally. If a specific phonetic or prosodic part of a generated sentence is not well represented in the database the output quality of this sentence suffers immensely. To try to avoid this a huge effort in specifically designing the database for the required application can be performed, but still there is no guarantee that such bad joins happen. In addition the fact, that in unit-selection no or only very little adaptations of the voice characteristics are possible without an enormous increase of the database size, the ambition towards seamless speech synthesis leads to the HMM-based approach. There a statistic representation of some sets of speech segments is used to generate arbitrary synthetic speech.

Details about unit-selection synthesis ???

In Figure 1 the structure of a typical HMM-based synthesizer is shown. The whole system can be divided into two parts, the training and the synthesis part. The connection of these two parts is a set of context-dependent (what is that?) HMMs. In the training part spectrum and excitation parameters of the recorded speech are used to generate acoustic models represented by the HMMs. Thereby phonetic, linguistic and prosodic parameters are considered. Details about spectrum and excitation parameters? In comparison to a unit-selection system the database only is needed in the training part.

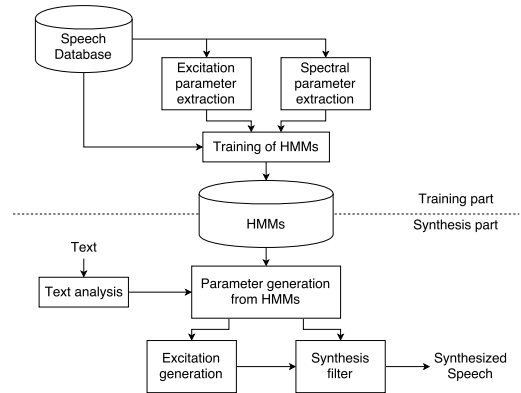


Figure 1: Function blocks of HMM-based synthesis [4]

In the synthesis part first the text which is to be synthesized is transformed to a sequence of parameters containing information about the context (what context?). According to this sequence the respective HMMs (what is that?) are concatenated in order to form an utterance HMM. Then after determining the state durations of the HMMs a sequence of coefficients is created, which finally

is used to construct the speech waveform using a special filter (details?).

The main disadvantage of this approach compared to unit-selection synthesis is the quality of the synthesized speech. Three factors are accountable for the lack of quality: the vocoder, the modeling accuracy, and an effect called over-smoothing (details?).

However there are some essential advantage, which makes the HMM-based approach a competitive alternative to unit-selection synthesis. First, the voice characteristics can be modified without much effort. Thus the implementation of different languages and the realization of different speaking styles with emotional emphasize is possible. Second, these just named aspects require a much smaller database than in unit-selection synthesis, since only a statistic representation of speech segments rather than raw speech data is stored.

In Table 1 the up to here mentioned techniques are compared regarding the most prominent advantages and drawbacks.

Table 1: Comparison of speech generation methods [4, 9]

Technique	Advantages	Drawbacks
Formant-based	No prerecorded speech required	Very artificial and metallic voice
Unit-selection	Very high voice quality possible	Large database required
HMM-based	Adjustable voice and small footprint	Voice sounds muffled

- (Relationship between the two approaches)
- (Hybrid approaches)
- Why there is need to further improve this technology?
- What is a HMM?
- Bring up MLPG (Maximum Likelihood Parameter Generation)
- Bring up Decision Tree
- Robustness as additional advantage (source [10] in [14])

3 SPSS WITH DEEP LEARNING MODELS

In this section first different approaches of deploying deep learning models are investigated, as shown in [8] and then one specific approach is explained further, wherefore [14] is used as reference.

3.1 One specific approach for improvement

In the previous section we learned that the quality issue of HMM-based speech synthesis is caused by three aspects, the vocoder, the accuracy of acoustic models and the over-smoothing effect. Zen *et al.* [14] suggest a specific approach to eliminate one of these causes, the accuracy of acoustic models, by allocating this task to a DNN. In conventional HMM-based systems the mapping between context features (phonetic and linguistic properties) and speech parameters is done by decision tree based context clustering. Thereby the context-dependent HMMs are assigned to different clusters depending on the combination of contexts using binary decision trees. Each cluster is characterized by a specific set of speech parameters. In this way it is possible to estimate all HMMs in a robust way, with a typically sized training database.

However decision trees soon reach their limits when handling complex contexts. Only by increasing the size and in this way

decreasing the efficiency of the decision tree more complicated contexts (e.g. XOR) can be dealt with.

In addition to that decision require partitioned input data with each partition having a different set of parameters. In this way with less data per region overfitting (weak generalization) is likely to happen, which then results in a lack of quality. These downsides can be avoided by using a DNN instead of multiple decision trees.

But the use of a DNN also introduces two disadvantages: The first one arises in terms of computational power. Both in the training and in the prediction stage decision trees require much less operations (total amount and level of complexity) than DNNs. The second one has to do with the decision process in its basic form. With decision trees a binary question has to be answered, whilst a DNN consists of weighted neurons, which use non-linear activation functions (e.g. sigmoid, tanh, ReLU [6]) to determine their state. As consequence of that interpretable rules are far easier to produce with decision trees, than with DNNs.

In Figure 2 the structure of a DNN-based speech synthesis system is shown. First a sequence of input features is generated after analysing the input text. These parameters contain numeric values like the number of words in a sentence or the duration of a phoneme as well as binary answers to questions like "is the current phoneme aa?". Then this parameter sequence is fed into the DNN where a mapping to output features is deployed by using forward propagation. The DNN has to be trained sometime before with pairs of input and output features from a database. In the following steps the speech parameters are extracted from the statistics of the output features and the voice waveform in turn generated from the speech parameters. this is done in the same way as in the HMM-based system. For this system the function blocks of text analysis, parameter generation, and waveform synthesis can be reused from a HMM-based system. Only the mapping from input features (e.g. linguistic contexts) to output features (spectral and excitation parameters) is implemented in a different way.

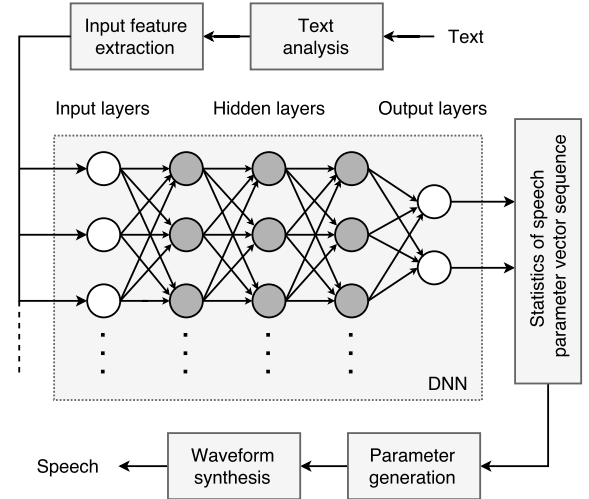


Figure 2: Speech synthesis based on DNN [14]

To compare the output of the above describe framework with that of a HMM-based, Zen *et al.* conducted some experiments with each a HMM-based and a DNN-based speech synthesis system in [14]. Therefore they used the same speech data in US English which includes about 33 000 utterances for both systems.

The HMM-based system used 2 554 questions for the decision tree-based context clustering. To influence the size of the decision trees the scaling factor α is used, where a high value results in a small decision tree and $\alpha = 1$ denotes a typical HMM-based system.

342 binary features (e.g. phonemes identities) and 25 numerical features (e.g. number of syllables in a word) represent the input features of the DNN-based system. As activation function the sigmoid function is used, since the authors experienced superior performance of this type in previous tests. In total one network with different number of layers (1, 2, 3, 4, or 5) and units per layer (256, 1 024, or 2 048) are used.

Objective¹

As subjective evaluation the generated speech is played back to a number of listeners, who then choose which synthesis is preferred. If no difference is perceived the option "neutral" can be selected. In this test the same number of parameters for both systems are used. For the structure of the DNN-based approach 4 layers with a different number of units per layer were applied. In Table 2 the outcome of the subjective evaluation is shown. Unmistakably the speech generated with the DNN-based systems were preferred, regardless of the amount of units per layer. The listeners described them as less muffled.

Table 2: Subjective scores (in %) of speech samples in [14]

HMM (α)	DNN (layers \times units)	Neutral
15.8 (16)	38.5 (4 \times 256)	45.7
16.1 (4)	27.2 (4 \times 512)	56.8
12.7 (1)	36.6 (4 \times 1 024)	50.7

Conclusions: Improved performance for predicting spectral and excitation parameters. More natural sounding voice. BUT: Higher computational costs both at training and prediction stage.

3.2 Other ways for improvement

Although, the approach of SPSS has brought many advantages over unit-selection synthesis, as shown in the previous section, the generated voice is still not as natural as desired. Therefore deep learning models recently have been used to further improve SPSS. Since DNNs have proven to be very effective in speech recognition, they have found their way into speech synthesis. With the use of DNNs it is possible to represent higher dimensional and correlated features in an efficient way as well as compactly modeling complex mapping functions. Exactly these properties can be used in SPSS, where different features representing the context (linguistic, prosody, etc.) have to be considered for the acoustic modeling [8].

In [8] the effects of deep learning methods on SPSS are investigated. Therefore the different parts of a HMM-based speech synthesis system are modeled with a DNN and then the output is compared to the conventional approach.

¹blablabla

-> other approach than DNN as acoustic model as in [14]

4 SPEECH SYNTHESIS ON EMBEDDED DEVICES

4.1 Motivation

Why is it important to implement speech synthesis on embedded platform?

What needs to be thought about when dealing with embedded or mobile devices?

4.2 HMM-based Approach

4.3 Deep Learning-based Approach

5 CONCLUSIONS

Here the core points will be repeated and concluded.
Some future aspects will be highlighted.
What should be done in the future?

See [13]

- Voice cloning
- Voice reconstruction
- Personalised speech-to-speech translation
- Articulatory-controllable speech synthesis

REFERENCES

- [1] 2017. The Benefits of Text to Speech. (02 June 2017). <http://www.readspeaker.com/benefits-of-text-to-speech/>
- [2] 2017. Stephen Hawking - The Computer. (02 June 2017). <http://www.hawking.org.uk/the-computer.html>
- [3] 2017. Survey on usage of virtual personal assistants by Gartner. (22 May 2017). <http://www.gartner.com/newsroom/id/3551217>
- [4] Alan W. Black, Heiga Zen, and Keiichi Tokuda. 2007. Statistical Parametric Speech Synthesis. In *2007 IEEE International Conference on Acoustics, Speech and Signal Processing - ICASSP '07*, Vol. 4. IV-1229-IV-1232. <https://doi.org/10.1109/ICASSP.2007.367298>
- [5] Tiberiu Boros and Stefan Daniel Dumitrescu. 2015. Robust Deep-learning Models for Text-to-speech Synthesis Support on Embedded Devices. In *Proceedings of the 7th International Conference on Management of Computational and Collective Intelligence in Digital EcoSystems (MEDES '15)*. ACM, New York, NY, USA, 98-102. <https://doi.org/10.1145/2857218.2857234>
- [6] Hoon Chung, Sung Joo Lee, and Jeon Gue Park. 2016. Deep neural network using trainable activation functions. In *2016 International Joint Conference on Neural Networks (IJCNN)*. 348-352. <https://doi.org/10.1109/IJCNN.2016.7727219>
- [7] Li Deng. 2014. A tutorial survey of architectures, algorithms, and applications for deep learning. *APSIPA Transactions on Signal and Information Processing* 3 (January 2014). <https://doi.org/10.1017/atsip.2013.9>
- [8] K. Hashimoto, K. Oura, Y. Nankaku, and K. Tokuda. 2015. The effect of neural networks in statistical parametric speech synthesis. In *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. 4455-4459. <https://doi.org/10.1109/ICASSP.2015.7178813>
- [9] Florian Hinterleitner. 2017. *Quality of Synthetic Speech: Perceptual Dimensions, Influencing Factors, and Instrumental Assessment (T-Labs Series in Telecommunication Services)* (1st ed. 2017 ed.). Springer, Singapore. <https://doi.org/10.1007/978-981-10-3734-4>
- [10] Z. H. Ling, S. Y. Kang, H. Zen, A. Senior, M. Schuster, X. J. Qian, H. M. Meng, and L. Deng. 2015. Deep Learning for Acoustic Modeling in Parametric Speech Generation: A systematic review of existing techniques and future trends. *IEEE Signal Processing Magazine* 32, 3 (May 2015), 35-52. <https://doi.org/10.1109/MSP.2014.2359987>
- [11] David Suendermann, Harald Höge, and Alan Black. 2010. *Challenges in Speech Synthesis*. Springer US, Boston, MA, 19-32. https://doi.org/10.1007/978-0-387-73819-2_2
- [12] Bálint Tóth and Géza Németh. 2012. Optimizing HMM Speech Synthesis for Low-Resource Devices. In *Journal of Advanced Computational Intelligence and Intelligent Informatics*. 327-334.
- [13] Junichi Yamagishi. 2011. New and emerging applications of speech synthesis. (February 2011). www.cstr.ed.ac.uk
- [14] Heiga Zen, Andrew Senior, and Mike Schuster. 2013. Statistical parametric speech synthesis using deep neural networks. In *2013 IEEE International Conference on Acoustics, Speech and Signal Processing*. 7962-7966. <https://doi.org/10.1109/ICASSP.2013.6639215>
- [15] Heiga Zen, Keiichi Tokuda, and Alan W. Black. 2009. Statistical parametric speech synthesis. (April 2009), 1039 - 1064 pages. <https://doi.org/10.1016/j.specom.2009.04.004>

Remarks starting here:

Page count estimation:

• Section 1	1	(incl. Abstract)
• Section 2.1	0.75	
• Section 2.2	0.5	
• Section 3.2	0.5	
• Section 3.1	0.5	
• Section 4.1	0.25	
• Section 4.2	0.75	
• Section 4.3	0.75	
• Section 5	0.25	
• References	0.5	

Total: 6 pages ???

To do:

- cite both long and short version of [15]
- check main title
- check (sub)section titles
- check references (bibtex warnings ???)
- check correct citing
- check for typos
- insert footnotes if further explanations are necessary
- keywords necessary ???
- check page count (6 pages)
- disable screen mode on last draft
- explain acronyms first time they appear in abstract AND first time they appear in body

Questions:

- (1) Why speech synthesis is important? What are its applications?
- (2) What are the conventional techniques of speech synthesis? What are the drawbacks of such techniques?
- (3) What is deep learning? What improvements do deep learning algorithms bring?
- (4) How some algorithms are modified to suit speech synthesis?
- (5) Why is it important to implement speech synthesis on embedded platform?
- (6) An example of how speech synthesis can be implemented on embedded platform without deep learning.
- (7) How the 3 can be combined?
- (8) Future works.