

# NEURAL NET WORD REPRESENTATIONS FOR PHRASE-BREAK PREDICTION WITHOUT A PART OF SPEECH TAGGER

*O. Watts, S. Gangireddy,  
J. Yamagishi, S. King, S. Renals*

The Centre for Speech Technology Research  
University of Edinburgh, UK

*A. Stan, M. Giurgiu*

Communications Department  
Technical University of Cluj-Napoca, Romania

## ABSTRACT

The use of shared projection neural nets of the sort used in language modelling is proposed as a way of sharing parameters between multiple text-to-speech system components. We experiment with pretraining the weights of such a shared projection on an auxiliary language modelling task and then apply the resulting word representations to the task of phrase-break prediction. Doing so allows us to build phrase-break predictors that rival conventional systems without any reliance on conventional knowledge-based resources such as part of speech taggers.

*Index Terms*— Speech synthesis, TTS, unsupervised learning, neural net language modelling, multitask learning.

## 1. INTRODUCTION

Neural networks (NNs) have re-emerged as a popular paradigm for the construction of text-to-speech (TTS) systems in recent years [1, 2, 3], much of their popularity being due to their successful use in learning ‘deep’ representations of data. In TTS (as in related work in speech recognition [4]) the emphasis has been on learning representations of acoustic data. For example, [3] uses Restricted Boltzmann Machines to operate directly on spectral envelopes, in effect replacing conventional representations of speech derived with expert knowledge (mel cepstral coefficients, line spectral pairs, etc.) with ones learned by the deep model. In this paper we focus instead on the text part of TTS, which has so far received less attention, and treat NNs as a way of obtaining representations of text which are optimised directly for the prediction of speech features.

Our previous work [5, 6] on obtaining features in an unsupervised way for use in TTS systems adopted a vector space model (VSM) approach, using a two-stage approach to semi-supervised learning. In the first – unsupervised – stage, we extract high-dimensional continuous-valued features to characterise textual or linguistic units of interest (letters, phonemes, words, utterances, etc.). This is done by compiling matrices of cooccurrence counts and then applying low rank matrix factorisation to obtain lower-dimensional distributional representations of the items in question. Then in a second

supervised step, these features are used in place of conventional features (part of speech (POS) tags, phonetic categories etc.) as predictors for various supervised tasks such as phrase-break prediction and acoustic state clustering. The advantage of this approach is that the first stage allows the system to take advantage of large quantities of unannotated text. However, although we have shown empirically that the representations obtained in the first step in many cases improve the predictors trained in the second, the weakness of this approach is that the two steps are nonetheless unlinked, and there is no guarantee that the features from the unsupervised phase will be useful in the supervised one.

In the current work, we investigate the use of *shared projection feed-forward NNs* of the sort that have become popular for language modelling [7] to learn word representations. This architecture has the potential to overcome the problems of our VSM approach as it allows word representations to be learned jointly with a classifier on some supervised task of interest. However, we want to retain the benefit of being able to exploit large text resources which has been a key feature of our previous work, and so consider an approach similar to that of [8, 9]. There, word representations obtained within a neural network language model are subsequently used for a variety of natural language processing tasks in a *multitask learning* (MTL) framework. The idea of MTL [10] is to train predictors for several related tasks in such a way that some parameters of the various task-specific models are shared, and can therefore be estimated on more data. Classically, some classifiers are devoted to *auxiliary tasks*: the predictions of these classifiers are of no direct interest, and they are learned purely to improve the estimation of the shared parameters. In [8, 9] the shared parameters are word representations from a NN language model (NNLM): the language modelling task is an auxiliary task whose only role is to initialise word representations which are subsequently refined on the other tasks. This is an attractive paradigm for TTS where multiple system components are trained on small amounts of expensive manually-labelled (usually disjoint) data and optimised in isolation. One goal of our on-going work is to develop a TTS system with shared representations for all units of interest (letters, phonemes, words, phrases, etc.) which are trained

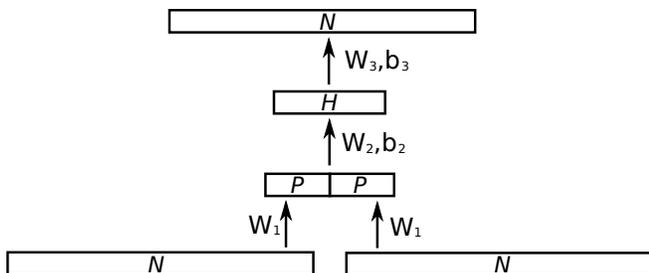


Fig. 1. Topology of neural net language model used.

jointly with all system components (relating both to text processing and acoustic modelling) in the manner of [11].

The scope of the current work is limited in comparison to the longer-term plans outlined in the previous paragraph, but is a first step in that direction. In [5] we took the task of phrase-break (PB) prediction for rapid testing of VSM features before applying them to other tasks in later work. Here we take the same task: its credible objective measure and the small size of its training data enable a faster turnaround of experiments than would be allowed, for example, by the training and subjective evaluation of acoustic models. We take advantage of the relatively small scale of the PB task to pose some questions which it seems desirable to answer before progressing to more complex systems. Firstly, are features taken from a standard feed-forward NNLM equally as useful for our task as VSM features were previously shown to be? Secondly, how big an increase in performance is afforded by updating word representations on the supervised task? Conversely, what happens if we learn representations from scratch on the supervised task?

## 2. SHARED PROJECTION FEED-FORWARD NNS

What we term a *shared projection feed-forward NN* was used for letter-to-sound conversion in [12], and has become popular for language modelling [7, 13, 14]. In all cases, shared weights in an initial hidden *projection layer* (in [12] termed a *self-organising code layer*) map from orthogonal inputs encoding the identities of textual units (letters, words) in some history or context to continuous representations of them. The fact that weights are shared between all units in the context or history means that the representations of those units are invariant to their positions in the context or history. Furthermore, the representations are learned along with the other parameters of the network so that the letter/word representations which are obtained are optimised on the task of interest (predicting the correct phoneme in [12], predicting the following word in [7]). In [12] a direct comparison is made with the knowledge-based coding of letters used in the NET-Speak system: the learned representations outperformed the knowledge-based coding.

Figure 1 shows the topology of a feed-forward trigram NNLM.  $N$  denotes the size of the vocabulary: input consists of a  $2N$  sparse vector encoding the 2 word trigram history. Shared weights  $W_1$  transform the inputs for individual tokens into their  $P$ -dimensional representations: the concatena-

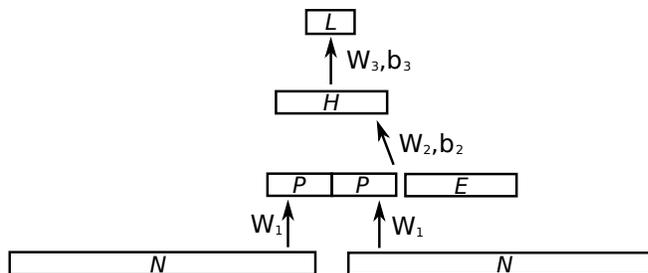


Fig. 2. Topology of phrase-break predictor used.

tion of the representations of the words in the  $n$ -gram history constitutes the output of the projection layer. A hidden layer with  $H$  units and a non-linear activation function transforms the outputs of the projection layer. Finally, an  $N$ -dimensional output layer produces  $n$ -gram probabilities for each word in the vocabulary given the input history. Use of the *softmax* normalisation function to constrain the network's outputs to be in the range  $[0, 1]$  and to sum to 1 ensures a valid probability distribution.

There is nothing about this model topology which restricts it to language modelling: as already mentioned, in [12] it is used for LTS conversion. We here use the slightly modified version of it shown in Figure 2 for phrase-break prediction. Here, two words of context are also used, although in the current experiments these are the words preceding and following a possible break, in contrast to the language model history. Indeed, the use of the shared projection means that words' representations are invariant to position in context, so that arbitrarily different contexts (in terms of the words' relative position to the target and their number) can be used for different tasks within the same experiment. The reason we restrict the context to preceding and following words is to ensure results which are comparable to those of previous work, where the same context was used. Another difference between the language model and PB predictor is that the size  $L$  of the output layer is not the size of the vocabulary, and is determined by the number of possible values of the variable to be predicted. In the experiments presented here,  $L$  is 2 to encode *break* and *no break* (which we note is equivalent to using a single sigmoidal unit). A final difference is that an extra  $E$  inputs are fed directly into the second hidden layer: these extra inputs encode additional contextual information beyond the 2 words' context such as distance to punctuation (see the following section for details).

## 3. EXPERIMENTS

### 3.1. Data

The data used for the supervised part (PB prediction) of present experiments is identical to that described in [5]: a subset of 39 stories from the SEC/MARSEC corpus [15, 16, 17] which consists of annotated radio broadcast transcriptions. Punctuation marks and two levels of PB (break/no break) were associated with the token that precedes them as a feature of that token. The data consist of c.35,000 non-punctuation tokens; 11% are from the 'overlap' section (annotated prosod-

**Table 1.** Summary of benchmark systems built

Neural net systems	Decision tree systems	Feature description
B	B'	Basic
G	G'	Guessed POS
T	T'	Topline POS

ically by both the corpus’s transcribers) which we use as in [5] as a test set. No single fixed development set was used in the current work as in [5]; instead, partition into training and validation sets was carried out randomly per model built, as explained in Section 3.2.

For the unsupervised pretraining part of the present experiments (NNLM training) we used the same 1.2 million tokens of Wall Street Journal text from which VSM features were obtained in [5]. The tokenised text was lowercased, but no sophisticated text normalisation was applied. A small set of tokens was found using the procedure described in [5] and rewritten with the  $\langle unk \rangle$  token used to handle unseen words at validation and run time.

### 3.2. Systems built

**Benchmark systems** Table 1 summarises the benchmark systems used in this work. Systems B, G and T are NN benchmark systems which in terms of the features they use are directly comparable with decision tree-based systems of the same names from previous work [6], which are here denoted B', G' and T' to disambiguate.<sup>1</sup> None of systems B, G or T in effect makes use of the network’s projection layer – for those systems, an input identical for each example is fed into the two words’ context input of the NN, in effect resulting in a standard feed-forward NN with a single hidden layer.

All systems make use of the following basic punctuation and positional features:

- The identity of the word’s punctuation symbol;
- The number of words {since, until} a word with a *strong punctuation* mark (i.e. excluding quotes);
- The number of words {since, until} the beginning/end of the utterance;

For training NNs, the punctuation feature was represented using 1-of- $k$  coding, and the positional features were normalised to have zero mean and unit variance.

**System B** makes use of these features only. All other systems additionally make use of some representation of the chosen *context words*: this was limited to the words preceding and following a possible phrase-break for consistency with our previous work.

Despite the great variety of machine learning methods that have been used for PB prediction in the past, all of them have used POS tags as independent variables. A system using POS tag representations of the context words is therefore an appropriate topline system, and **system T** supplements the basic

<sup>1</sup>Results for the slightly refined systems presented in [6] are used, rather from the nearly identical systems presented in [5]; note that the system here called T' corresponds to system Tt (using the TnT tagger) in [6].

**Table 2.** Summary of systems built with induced word representations; ● and ○ indicate *true* and *false* respectively.

NN system	DT system	Feature description	Distributional pretraining	Representations tuned for PBs	Vocabulary from PB data
U	U'	Unsupervised	●	○	○
R	-	Randomly init'd	○	●	●
F	-	Fine-tuned	●	●	○
S	-	Subset vocabulary	●	●	●

features with the output of a high-quality POS tagger (*TnT*: [18]). The tagger’s output for the context words was collapsed to the 23-tag set used in [19, 20, 5] and was presented (coded as two 1-of- $k$  vectors) to the network along with the basic features of system B as the  $E$  extra inputs shown in Figure 2.

A full POS tagger can be approximated with a few simple rules; **system G** makes use of such rules, consulting 9 lists of different types of function words at run time. Words not found in any of the lists are tagged as *content* words; others are replaced with the tag associated with the list in which they appear (e.g. *pronoun*, *modal verb*, etc.). This is a suitable benchmark against which to compare unsupervised methods as it is a system built with only minimal expert knowledge.<sup>2</sup>

**Experimental systems** Experimental systems U, R, F and S are summarised in Table 2. They all make use of the projection layer as well the  $E$  extra inputs shown in Figure 2, which are used to input the basic positional and punctuation features already described.

Several different configurations using induced word representations were explored. Most of them used the training of a trigram NNLM as an auxiliary task: for this the 1.2M tokens of text described in Section 3.1 were used. **System U** made use of weights  $W_1$  which were pretrained on the full vocabulary of the unsupervised (language modelling) task, but kept them frozen during training on the PB task. Unseen words are handled with the  $\langle unk \rangle$  token of the NNLM’s vocabulary. This system is the nearest equivalent to the decision tree system (here called U') from [6] using a VSM (learned on the same 1.2 million tokens of text) and where the classifier was not able to transform individual words’ representations on the PB task.

The full vocabulary and projection weights  $W_1$  were similarly taken from the trained NNLM and used to initialise the projection layer weights  $W_1$  of **system F**. Instead of then being kept frozen, however, these weights were then updated with backpropagation along with the NN’s other weights when the NN is trained on the PB data. Comparison of the

<sup>2</sup>External benchmarks: [21, 19, 20] report F scores of 74.4%, 78.3%, and 81.6% on MARSEC data, although possible differences of data preparation mean this comparison needs to be made with caution.

results of systems U and F therefore will allow us to test the usefulness of fine-tuning existing word representations on the task of interest.

**Systems R** is designed to determine the usefulness of pre-training weights  $\mathbf{W}_1$  on the language modelling task. Its projection weights  $\mathbf{W}_1$  were randomly initialised and then tuned only on the PB task. The vocabulary is therefore limited to that of the PB training set: to handle unseen words at testing time, words with a single occurrence in the training data were found, and  $n\%$  of them were randomly selected and replaced with the  $\langle unk \rangle$  token. Three versions of this configuration were tried, with the value of  $n$  set to 100, 50 and 10.

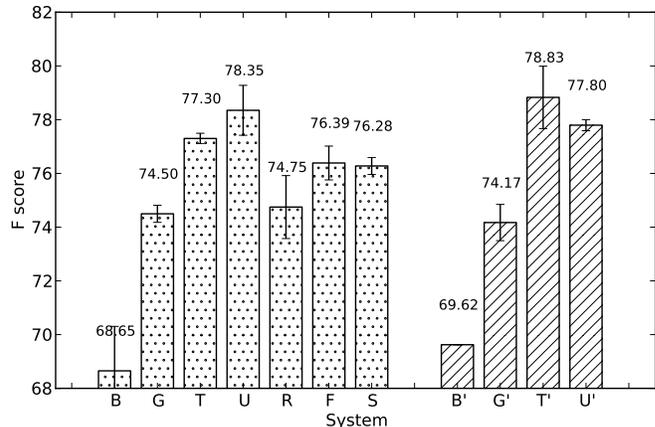
To investigate the impact of the mismatch caused by fine-tuning only the representations of tokens seen in the PB training data (and not the tokens seen in the LM data but absent in the PB training data), **system S** was built. This is identical to F, except that it uses only a subset of the language model vocabulary, limited to the words seen in the PB training data. The representation of the  $\langle unk \rangle$  token – initialised in the NNLM and updated on the PB task – is therefore used to handle all words absent from the PB training data at validation and test time.

### 3.3. Network training

Weights and biases of all models were initialised uniformly at random according to the *normalised initialisation* suggested by [22] (except weights  $\mathbf{W}_1$  of systems U, F and S where pre-trained weights are used). 10% of the training examples were chosen at random for each model trained for use as a validation set. The remaining 90% were resampled to balance class probabilities – resampling was done with a new random seed for each model trained. Stochastic gradient descent with minibatches was used to train NNs using negative log likelihood of the training examples as the cost function. A learning rate was used which decayed exponentially when improvement in validation set negative log likelihood between successive epochs fell beneath a prespecified threshold. Training finished when the negative log likelihood of this validation set stopped decreasing, or when 15 epochs had been completed. A small  $L_2$  regularisation term was used for training the NNLM, but no explicit regularisation was used in training the PB predictors. In all NNs built, projection dimension  $P$  was fixed at 50 (consistent with the dimension of the VSM features used in previous work). For NNLM training, hidden layer size  $H$  was set to 100. For PB predictors, 5 different values of  $H$  (10, 50, 100, 150, 200) were tried. Five PB predictors were trained with different weight initialisations for each combination of configuration and  $H$  value.

### 3.4. Results

Results are reported as F scores on phrase-breaks. Mean F scores on validation sets were used to determine 50 as the optimal setting of parameter  $n$  in system R. The optimal setting for  $H$  for all configurations was also found in this way. A single set of 5 systems in each of these 7 configurations was then evaluated on the test set; means and standard deviations of F scores are plotted in Figure 3 alongside corresponding results



**Fig. 3.** Mean and st. deviation of F scores (test set) of final systems (left) & decision tree systems for reference (right) from earlier work (where the statistics summarise the scores of 10 systems per configuration).

## 4. DISCUSSION

The overall trend among baseline systems B, G and T is the same as that of decision tree benchmarks (B', G', T'), except that T makes less good use of the POS features than the corresponding decision tree system T'. Adding the word representations pretrained on the LM task and keeping them fixed (U) gives the best performance, outperforming T and equalling that of T'. Thus using NNLM features allows us to close the gap in performance between topline and unsupervised word feature systems while using only 1.2M tokens of unannotated text data. This is a pleasing outcome: in previous work [6] 20M tokens of text were needed to close this gap when using VSMs. Pretraining representations in the NNLM is clearly beneficial to training them from scratch on the PB task, although doing only that allows system R to rival the minimally supervised system G. The slightly inferior performance of system F is surprising, as a major motivation for this work was the expectation that updating the word representations on the PB task directly would yield an improvement. We hypothesise that this is due to the mismatch between the fine-tuned vocabulary and the items whose NNLM representations were 'left behind' during fine-tuning. Using these 'left behind' representations gives similar results to the case where they are simply replaced with the  $\langle unk \rangle$  token (system S). If true, this has important implications for MTL as it suggests that a large and consistent inventory of representations is better than one which has been partially tuned on a task of interest. Several ways to enforce representations' consistency during fine-tuning could be proposed, such as tuning a shared-weight adaptation layer instead of the word representations themselves [14, 23], or by using a different training curriculum, such as interleaving batches of PB examples and LM examples.

## 5. ACKNOWLEDGEMENTS

This work is supported by the European Community's Seventh Framework Programme (FP7/2007-2013) under grant agreement N° 287678, and by the JST CREST project uDialogue.

## 6. REFERENCES

- [1] Shiyin Kang, Xiaojun Qian, and Helen Meng, “Multi-distribution deep belief network for speech synthesis,” in *Acoustics, Speech and Signal Processing (ICASSP), 2013 IEEE International Conference on*, 2013, pp. 8012–8016.
- [2] Heiga Zen, Andrew Senior, and Mike Schuster, “Statistical parametric speech synthesis using deep neural networks,” in *Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, 2013, pp. 7962–7966.
- [3] Zhen-Hua Ling, Li Deng, and Dong Yu, “Modeling spectral envelopes using Restricted Boltzmann Machines and Deep Belief Networks for statistical parametric speech synthesis,” *Audio, Speech, and Language Processing, IEEE Transactions on*, vol. 21, no. 10, pp. 2129–2139, 2013.
- [4] G. Hinton, Li Deng, Dong Yu, G.E. Dahl, A. Mohamed, N. Jaitly, A. Senior, V. Vanhoucke, P. Nguyen, T.N. Sainath, and B. Kingsbury, “Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups,” *Signal Processing Magazine, IEEE*, vol. 29, no. 6, pp. 82–97, 2012.
- [5] Oliver Watts, Junichi Yamagishi, and Simon King, “Unsupervised continuous-valued word features for phrase-break prediction without a part-of-speech tagger,” in *Proc. Interspeech*, Florence, Italy, Aug. 2011.
- [6] Oliver Watts, *Unsupervised Learning for Text-to-Speech Synthesis*, Ph.D. thesis, University of Edinburgh, 2012.
- [7] Yoshua Bengio, Réjean Ducharme, Pascal Vincent, and Christian Jauvin, “A neural probabilistic language model,” *Journal of Machine Learning Research*, vol. 3, pp. 1137–1155, 2003.
- [8] R. Collobert and J. Weston, “A unified architecture for natural language processing: Deep neural networks with multitask learning,” in *International Conference on Machine Learning, ICML*, 2008.
- [9] Ronan Collobert, Jason Weston, Léon Bottou, Michael Karlen, Koray Kavukcuoglu, and Pavel P. Kuksa, “Natural language processing (almost) from scratch,” *CoRR*, vol. abs/1103.0398, 2011.
- [10] Rich Caruana, “Multitask learning,” *Machine Learning*, vol. 28, pp. 41–75, July 1997.
- [11] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, “Gradient-based learning applied to document recognition,” *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [12] Kåre Jean Jensen and Søren Riis, “Self-organizing letter code-book for text-to-phoneme neural network model,” in *INTERSPEECH*, 2000, pp. 318–321.
- [13] Holger Schwenk and Jean-Luc Gauvain, “Connectionist language modeling for large vocabulary continuous speech recognition,” in *International Conference on Acoustics, Speech and Signal Processing*, 2002, pp. 765–768.
- [14] Holger Schwenk, “Continuous space language models,” *Computer Speech & Language*, vol. 21, no. 3, pp. 492 – 518, 2007.
- [15] Gerry Knowles, Briony Williams, and L. Taylor, *A Corpus of Formal British English Speech: The Lancaster/IBM Spoken English Corpus*, Longman, 1996.
- [16] Gerry Knowles, Anne Wichmann, and Peter Alderson, *Working with Speech: Perspectives on Research into the Lancaster/IBM Spoken English Corpus*, Longman, 1996.
- [17] Peter Roach, Gerry Knowles, Tamas Varadi, and Simon Arnfield, “Marsec: A machine-readable spoken english corpus,” *Journal of the International Phonetic Association*, vol. 23, no. 2, pp. 47–54, 1993.
- [18] Thorsten Brants, “TnT: a statistical part-of-speech tagger,” in *Proc. 6th Conf. Applied Natural Language Processing*, 2000, pp. 224–231.
- [19] Paul Taylor and Alan W. Black, “Assigning phrase breaks from part-of-speech sequences,” *Computer Speech & Language*, vol. 12, no. 2, pp. 99 – 117, 1998.
- [20] Ian Read and Stephen Cox, “Stochastic and syntactic techniques for predicting phrase breaks,” *Computer Speech and Language*, vol. 21, no. 3, pp. 519 – 542, 2007.
- [21] Bertjan Bussler, Walter Daelemans, and Antal van den Bosch, “Predicting phrase breaks with memory-based learning,” in *Proc. 4th ISCA Speech Synthesis Workshop*, 2001.
- [22] Xavier Glorot and Yoshua Bengio, “Understanding the difficulty of training deep feedforward neural networks,” *Journal of Machine Learning Research - Proceedings Track*, vol. 9, pp. 249–256, 2010.
- [23] Junho Park, Xunying Liu, Mark J. F. Gales, and Philip C. Woodland, “Improved neural network based language modelling and adaptation,” in *INTERSPEECH*, 2010, pp. 1041–1044.