



# Context-Dependent Acoustic Modeling without Explicit Phone Clustering

Tina Raissi, Eugen Beck, Ralf Schlüter, Hermann Ney

Human Language Technology and Pattern Recognition Group  
RWTH Aachen University

tina.raissi@rwth-aachen.de, {beck, schluter, ney}@cs.rwth-aachen.de

## Abstract

Phoneme-based acoustic modeling of large vocabulary automatic speech recognition takes advantage of phoneme context. The large number of context-dependent (CD) phonemes and their highly varying statistics require tying or smoothing to enable robust training. Usually, Classification and Regression Trees are used for phonetic clustering, which is standard in Hidden Markov Model (HMM)-based systems. However, this solution introduces a secondary training objective and does not allow for end-to-end training. In this work, we address a direct phonetic context modeling for the hybrid Deep Neural Network (DNN)/HMM, that does not build on any phone clustering algorithm for the determination of the HMM state inventory. By performing different decompositions of the joint probability of the center phoneme state and its left and right contexts, we obtain a factorized network consisting of different components, trained jointly. Moreover, the representation of the phonetic context for the network relies on phoneme embeddings. The recognition accuracy of our proposed models on the Switchboard task is comparable and outperforms slightly the hybrid model using the standard state-tying decision trees.

**Index Terms:** automatic speech recognition, context-dependent acoustic modeling, hybrid DNN/HMM system

## 1. Introduction

The realization of the phonetic co-articulation effect in large vocabulary continuous speech recognition (LVCSR) systems with standard Gaussian Mixture Model/Hidden Markov Model (GMM/HMM) takes into account a context-dependent (CD) representation of phones, usually triphones [1]. The extension of each phoneme with its left and right contexts leads to a considerable growth of the number of possible states. Finding the right trade-off between the model complexity and the available data can become complicated, on the grounds that during training many triphones are unevenly distributed or never observed. In order to overcome sparsity issues, for long, Classification and Regression Trees (CART) marked the state-of-the-art in ASR for tying CD phone states into generalized triphone states [2]. The successful advent of neural-based models in LVCSR paved the way for the hybrid Deep Neural Network (DNN)/HMM architecture [3], where the Gaussian mixture based emission probabilities are replaced by normalized scaled generalized triphone state posteriors, predicted by a discriminative model.

The introduction of CART labels as output targets of the NN model has given an important contribution to the improvement of the performance of the recognition systems, maintaining at the same time a two-fold dependency to the GMM system. The frame-level state alignment for training CART derives normally from a GMM. Furthermore, there is a mismatch between the features used for the estimation of Gaussian mixture

parameters and the one used for learning the posterior probabilities of the tied-states in the neural network component.

The majority of the research works on the CD acoustic modeling in connection with the hybrid approach aims to either integrate the context directly into the neural network [4, 5], or to eliminate the dependency to the GMM system. It is shown that the initial alignments to the context-independent (CI) states for the standard tree-based clustering approach can be provided by a flat-started DNN [6]. Similar approaches design the set of CD targets by clustering the activations of a CI DNN [7, 8]. There are also different possible training criteria for the state-tying algorithm, such as Kullback-Leibler divergence [9, 10], entropy [11], and based on DNN and classification error [12]. The elimination of the state-tying decision trees is the topic of research also for end-to-end models such as CTC where a CD embedding network is applied instead [13].

The common trait between most of the mentioned works is a phone clustering principle and the necessity of having one more training and optimization step. It is important to underline that in addition to the supplementary time and resource effort, another crucial concern regarding this further modeling approach is how the set of clustered states can affect the decision boundaries in the final neural network, which learns the probability distribution over their posteriors. This is especially true when the classic phonetic decision trees are involved. The relative heuristics regarding the choice of the questions or maximum number of leaves can affect directly the definition of the set of class labels, which, if not well-defined, can lead to over-fitting problems [14].

In this work, we propose a CD acoustic modeling for the hybrid approach, which disposes of the necessity of an additional phone clustering step for the determination of the HMM state inventory. The resulting model is partitioned into separate components, trained conjointly, and corresponding to one of the factorized elements of the joint probability of the center phoneme state with its left and right phonetic contexts. Depending on the type of decomposition, each component learns a posterior probability distribution over phonemes and phoneme states in mono-, di- and triphone context. We show that for the Switchboard task the recognition system built upon our direct context integration approach with no state-tying clustering can obtain a similar performance to a hybrid model using standard tying based on CART.

## 2. Formulation of the Problem

The statistical formulation of automatic speech recognition task maximizes the a-posteriori probability of a word sequence  $w_1^N$  of length  $N$  given the acoustic feature sequence  $x_1^T$  of length  $T$ , with  $T \gg N$ , based on Bayes decision rule [15]:

$$x_1^T \rightarrow \tilde{w}_1^N(x_1^T) = \operatorname{argmax}_{w_1^N} \left\{ p(x_1^T | w_1^N) \cdot p(w_1^N) \right\}. \quad (1)$$

The acoustic-phonetic component  $p(x_1^T | w_1^N)$  of Eq. (1), in the standard HMM with generative approach and involving a sequence of triphone states  $s_1^T$  is formulated as:

$$\begin{aligned} p(x_1^T | w_1^N) &= \sum_{s_1^T} \prod_{t=1}^T p(x_t | s_t, w_1^N) \cdot p(s_t | s_{t-1}, w_1^N) \\ &= \sum_{s_1^T} \prod_{t=1}^T p(x_t | s_t, \phi_1^M, w_1^N) \cdot p(s_t | s_{t-1}, \phi_1^M, w_1^N), \end{aligned}$$

where  $\phi_1^M$  represents a suitable triphone sequence of length  $M$  corresponding to the word sequence.

### 3. Integration of the Context

Denote by  $\{\phi_\ell, \phi_c, \phi_r\}_t$  the set of left, center and right phonemes of the aligned triphone at time frame  $t$ . Each phoneme consists of three HMM states, and each state can be associated with a state class  $c(s_t, w_1^N) = \{\phi_\ell, \phi_c, \phi_r, i\}_t$ , where  $i$  enumerates the HMM state of the corresponding triphone. The likelihood of generating a feature vector  $x$  at time frame  $t$  given a triphone, can be written as:

$$p(x_t | s_t, \phi_1^M, w_1^N) = p(x_t | c(s_t, w_1^N)) = p_t(x | \phi_\ell, \phi_c, \phi_r, i).$$

For simplicity, we use the parametrized probability distribution  $p$  and its further denotation  $p_t$  at time frame  $t$ , interchangeably. By applying Bayes identity we have:

$$p(x | \phi_\ell, \phi_c, \phi_r, i) = \frac{p(\phi_\ell, \phi_c, \phi_r, i | x) \cdot p(x)}{p(\phi_\ell, \phi_c, \phi_r, i)}. \quad (2)$$

Let  $\sigma_c$  be the current HMM state within the center phoneme, the CD neural network should ideally model the joint probability of  $\sigma_c$  with the left and right phonetic contexts appearing in the nominator of Eq. (2), which can be written through the following mapping as:

$$p(\phi_\ell, \phi_c, \phi_r, i | x) \rightarrow p(\phi_\ell, \sigma_c, \phi_r | x). \quad (3)$$

### 4. Different Decompositions

The joint posterior probability distribution of Eq. (3) would demand a high number of parameters and an infeasible memory requirement, if conceived as the output of a neural network. One possible solution is to obtain a factorization into CD probabilities, by applying the classic Markov chain rule [16].

#### 4.1. Diphone

The emission probability defined for the diphone model, as shown in Fig. 1a, is obtained by conditioning only on the left phonetic context. Starting with the modified version of Eq. (2), which takes also into consideration the mapping Eq. (3) we have:

$$\begin{aligned} p(x | \sigma_c, \phi_\ell) &= \frac{p(\sigma_c, \phi_\ell | x) \cdot p(x)}{p(\sigma_c, \phi_\ell)} \\ &= \frac{p(\sigma_c | \phi_\ell, x) \cdot p(\phi_\ell | x) \cdot p(x)}{p(\sigma_c | \phi_\ell) \cdot p(\phi_\ell)}. \end{aligned} \quad (4)$$

#### 4.2. Triphone Forward

In case of all triphone models, depicted in Figs. 1b to 1d, it is possible to achieve different decompositions by having as the start point all three entities, namely right and left contexts along

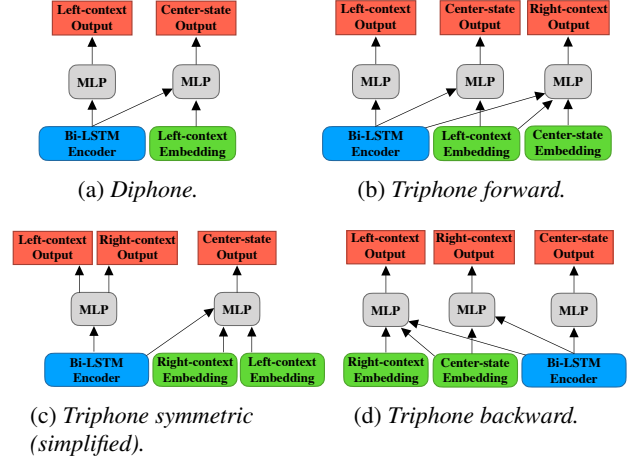


Figure 1: The architecture of different models defined in Eqs. (4) to (7). The left and right output layers have the respective phoneme identities  $\phi_\ell$  and  $\phi_r$  as targets. The target for the center phoneme output is the CI state referred to the phoneme inventory of the vocabulary.

with the center phoneme state. beginning with the right context, the chain rule will produce a left-to-right trigram, as below:

$$\begin{aligned} p(x | \phi_\ell, \sigma_c, \phi_r) &= \frac{p(\phi_\ell, \sigma_c, \phi_r | x) \cdot p(x)}{p(\phi_\ell, \sigma_c, \phi_r)} \\ &= \frac{p(\phi_r | \phi_\ell, \sigma_c, x) \cdot p(\sigma_c | \phi_\ell, x) \cdot p(\phi_\ell | x) \cdot p(x)}{p(\phi_r | \phi_\ell, \sigma_c) \cdot p(\sigma_c | \phi_\ell) \cdot p(\phi_\ell)}. \end{aligned} \quad (5)$$

#### 4.3. Triphone Symmetric

Another possible decomposition starts with the center phoneme state given the left and right contexts. The context-dependency in the remaining factors is not taking into account the center phoneme. By assuming that there is no interdependency between the left and right contexts, we drop the dependency to the right context  $\phi_r$  in the probability value  $p(\phi_\ell | \phi_r, x)$  of Eq. (6a), ending up with Eq. (6b). This independence assumption is valid also for the respective prior  $p(\phi_\ell | \phi_r)$ .

$$\begin{aligned} p(x | \phi_\ell, \sigma_c, \phi_r) &= \frac{p(\phi_\ell, \sigma_c, \phi_r | x) \cdot p(x)}{p(\phi_\ell, \sigma_c, \phi_r)} \\ &= \frac{p(\sigma_c | \phi_\ell, \phi_r, x) \cdot p(\phi_\ell | \phi_r, x) \cdot p(\phi_r | x) \cdot p(x)}{p(\sigma_c | \phi_\ell, \phi_r) \cdot p(\phi_\ell | \phi_r) \cdot p(\phi_r)} \quad (6a) \\ &= \frac{p(\sigma_c | \phi_\ell, \phi_r, x) \cdot p(\phi_\ell | x) \cdot p(\phi_r | x) \cdot p(x)}{p(\sigma_c | \phi_\ell, \phi_r) \cdot p(\phi_\ell) \cdot p(\phi_r)}. \end{aligned} \quad (6b)$$

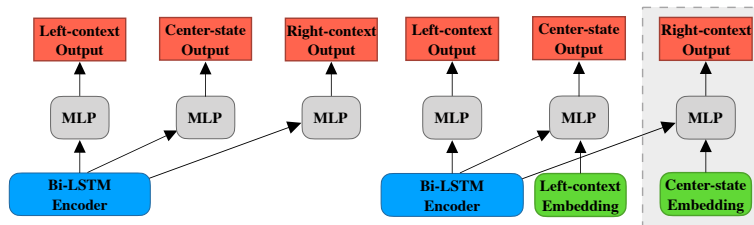
#### 4.4. Triphone Backward

A different possible factorization starts with the left context, leaving the center phoneme state as the single entity [4].

$$\begin{aligned} p(x | \phi_\ell, \sigma_c, \phi_r) &= \frac{p(\phi_\ell, \sigma_c, \phi_r | x) \cdot p(x)}{p(\phi_\ell, \sigma_c, \phi_r)} \\ &= \frac{p(\phi_\ell | \sigma_c, \phi_r, x) \cdot p(\phi_r | \sigma_c, x) \cdot p(\sigma_c | x) \cdot p(x)}{p(\phi_\ell | \sigma_c, \phi_r) \cdot p(\phi_r | \sigma_c) \cdot p(\sigma_c)} \end{aligned} \quad (7)$$

Table 1: Different pre-training procedures for the proposed triphone model with forward decomposition (Fwd) of Fig. 1b. The outputs of the trained architectures at monophone and diphone stages, depicted in Figs. 2a and 2b, are marked in the respective network columns. Decoding for diphone models and Fwd follows Eqs. (4) and (5), respectively. The recognition results in terms of Word Error Rate (WER) are over 300h Switchboard using 4-gram language model. The experiments of each row can be described as follows: (1) - Baseline triphone model with standard state-tying, (2) - Fwd with no pre-training, (3) - Diphone model of Fig. 1a trained with no pre-training and used for the initialization of Fwd, (4) - Pre-trained Fwd with only monophone stage, (5) Fwd with monophone and diphone pre-training stages, (6): Similar to experiment (5) with optional inclusion of the network branch having output distribution  $p(\phi_r|\sigma_c, x)$ .

#	Model	Monophone stage			Diphone stage			[%WER]	Triphone stage
		Network			Network				
		$\phi_\ell$	$\phi_c$	$\phi_r$	$\phi_\ell$	$\phi_c$	$\phi_r$		
1	Base	not applicable (n/a)							13.9
2	Fwd	-	-	-	-	-	-	n/a	13.9
3		-	-	-	✓	✓	-	14.8	13.9
4		✓	✓	✓	-	-	-	n/a	13.9
5		✓	✓	✓	✓	✓	-	14.4	13.8
6		✓	✓	✓	✓	✓	✓	14.2	13.6



(a) Monophone stage of rows 4 to 6. (b) Diphone stage rows 3, 5 and 6.

Figure 2: The network architectures used in monophone and diphone stages for pre-training of proposed triphone model with forward decomposition. The highlighted branch of the diphone network in Fig. 2b with output distribution  $p(\phi_r|\sigma_c, x)$  is included only for Experiment 6 of Table 1.

## 5. Modeling and Training Details

### 5.1. Model Architecture

The architecture of each model is divided into two separate constituting parts: (1) a bidirectional Long Short Time Memory (Bi-LSTM) network which obtains an encoding of the input features following the relatively well-established acoustic modeling background [17, 18], (2) a factorized neural network which integrates the context into the whole model. Regarding the CD component, there are three aspects to be underlined. The left and right phonemes of each triphone along with the center phoneme state are represented by using an embedding layers. Each output layer is preceded by a Multi-Layer Perceptron (MLP). It is possible to organize the MLP layers with different settings. They could be used as a shared combination component or be located in independent flows. Experimental results over different architectures show that in case of backward and forward model, these arbitrary choices do not implicate significant changes in the model performance.

### 5.2. Multi-Stage Phonetic Training

The final models are improved by using pre-training. The whole procedure can be defined as a multi-stage training which builds on incremental phonetic dependencies. We start with a monophone network and at each stage augment the context-dependency relations to adhere to a higher acoustic-phonetic n-gram scheme. From one stage to the following, the outputs needed for the final factorized model are kept without further modification. For this set of pre-training experiments on the

Switchboard task, and reported in Table 1, we took advantage of the fact that the diphone model is actually a complete sub-architecture of the triphone model with the forward decomposition. The experiments outcome show that with no pre-training it is possible to obtain the same performance of the baseline model using CART. The comparison between Experiments (3) with only diphone and (4) with only monophone pre-training stages, against Experiments (5) and (6) shows that the model benefits from the three-stage training. For the proposed diphone and triphone models the WER is consequently decreased from 14.8% to 14.2% and from 13.9% to 13.6%, respectively. Furthermore, the optional inclusion of the additional factor  $p(\phi_r|\sigma_c, x)$  during the diphone stage of Experiment (6), boosts both diphone and triphone models' performance. For more details about the experimental setup, see Sec. 6.

## 6. Experimental Setting and Results

We compare the CD acoustic models described in Sec. 4 with a baseline hybrid model using the standard state-tying with CART. All models are trained and evaluated over 300h Switchboard-1 Release 2 (LDC97S62) [19] and Hub500 data (LDC2002S09), respectively, with the aid of RETURNN and RASR toolkits [20, 21].

### 6.1. Setting

The frame-wise state alignment for training derives from a GMM/HMM system relying on CART. Our proposed approach makes use of a state inventory consisting of 136 state labels

Table 2: Different real time factor (RTF) values for comparable average number of states per frame during time synchronous prefix tree search for the triphone models with CART-based state-tying (Base) and forward decomposition (Fwd).

LM	Model	#States	#Trees	RTF
4-gram	Base	17345	113	0.5
	Fwd	15617	75	12.18
LSTM	Base	58504	179	5.59
	Fwd	61753	199	13.0

corresponding to 45 phonemes with three states and the single-state silence entity. For the baseline model, a set of 9001 CART labels are considered.

Both baseline and proposed CD models use a Bi-LSTM encoder comprising 6 forward and backward layers of size 500 with 10% dropout probability [22]. The input speech signal to the encoder is represented by 40-dimensional Gammatone Filterbank based features [23], optionally concatenated with i-vectors of dimension 200 for a subset of the conducted experiments [24]. All models share the same set of training hyper-parameters and are trained with frame-wise cross-entropy criteria and Adam optimizer with Nesterov momentum [25]. By means of Newbob scheduling, the initial learning rate of  $5 \times 10^{-4}$  with a decay factor of  $\sqrt{0.8}$  is controlled by using the cross-validation frame error rate and decreased to a minimum value of  $5 \times 10^{-6}$ . For the regularization, an  $L_2$  weight decay with constant 0.01, gradient noise with a variance of 0.3 and the focal loss factor of 2.0 are adopted [26, 27]. Each CD model is trained for 80 hours. The pick performance for the baseline model is obtained after 24 epochs, requiring 8.5h less than the best proposed CD model. Concerning the proposed approach, the one-hot encoding of the left and/or right phonemes and the center phoneme states are projected by using linear layers of dimension 10 and 30, respectively. Furthermore, The prior quantities appearing in the denominator of Eqs. (4) to (7) are estimated by averaging over the output activations of the network using a subset of the training set.

On the recognition side, we considered both 4-gram and LSTM language models [28, 29, 30]. The prior scales for each factor and the LM scale are tuned by using a grid search. The real time factors of two baseline and forward models are compared in Table 2. Forwarding all possible context combinations in batch mode gives an important contribution to the optimization of our approach. However, we aim to proceed with other optimization methods as a future work.

## 6.2. Results

The experimental results for the CD models of the proposed approach show that different decompositions obtain similar performance. The triphone model with forward decomposition outperforms slightly the hybrid baseline model. The improvement is maintained also when a different LM or i-vector adaptation are applied. We believe that the performance drop in case of symmetric model derives from the simplifying assumption regarding no interdependency between the two contexts, as discussed in Sec. 4.

## 7. Discussion

For the proposed CD models, the identity of each HMM state is uniquely defined by the identity of the center phoneme and the

Table 3: Comparison of WERs between the baseline system using standard state-tying with CART (Std. tying) and the proposed CD models with forward (Fwd), backward (Bwd) and symmetric (sym) decompositions, using 4-gram and LSTM LMs. A subset of experiments are carried out using i-vectors (I-vec).

Context-Dependency	LM	I-vec	[% WER]			
			Std. Tying	Proposed Approach		
Triphone	4-gram	no	13.9	Fwd	Bwd	Sym
		yes	13.3	12.9	—	—
	LSTM	no	12.7	12.6	12.8	13.8
		yes	11.9	11.7	—	—
Diphone	4-gram	no	15.0	14.2		
Monophone			17.3			

position within it, along with its right and left phonemes. The state labels in this case correspond to the set of CI states. The consideration of a subset of factors and not the full factorized model during the decoding leads to a considerable performance degradation. As an example, for a symmetric model, if we use only the normalized posterior  $\frac{p(\sigma_c|\phi_\ell, \phi_r, x)}{p(\sigma_c|\phi_\ell, \phi_r)}$  from Eq. (6), we observe up to 48% relative WER deterioration. Furthermore, by including an additional target belonging to a larger context span during the training and choosing a subset of the factors during decoding it is possible to obtain improvement. This is for example the case of the pre-trained diphone model of Experiment 6 of Table 1 having also the  $p(\phi_r|\sigma_c, x)$  factor during training, against the pre-trained diphone of Experiment 5. These observations suggest two aspects about the CD models: (1) the model learns the context-dependencies during joint training of the factors, (2) the decision rule carried out with respect to the defined theoretical framework is consistent and sound.

## 8. Conclusions

We showed that in acoustic modeling for the hybrid approach it is possible to discard the phone clustering step. Our results indicate that direct modeling of context provides sufficient smoothing ability with respect to the variability in context-dependent phoneme statistics and performs as well as the former clustering-based approach. However, at this stage of the work, the training of the models is still based on the frame-wise alignment derived from a separate GMM/HMM system. Future work concentrates on training this direct modeling approach from scratch, in order to also eliminate this secondary dependence on phonetic clustering.

## 9. Acknowledgements

This work has received funding from the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation programme (grant agreement No 694537, project ”SEQCLAS”) and from a Google Focused Award. The work reflects only the authors’ views and none of the funding parties is responsible for any use that may be made of the information it contains. We thank Markus Kitza for providing us with the i-vectors for model training.

## 10. References

- [1] L. Bahl, R. Bakis, J. Bellegarda, P. Brown, D. Burshtein, S. Das, P. De Souza, P. Gopalakrishnan, F. Jelinek, D. Kanevsky, R. Mercer, A. Nadns, D. Nahamnn, and M. Pichrny, "Large vocabulary natural language continuous speech recognition," in *Proc. IEEE Intern. Conf. on Acoustics, Speech and Signal Process. (ICASSP)*, pp. 465–467, 1989.
- [2] S. J. Young, J. J. Odell, and P. C. Woodland, "Tree-based state tying for high accuracy acoustic modelling," in *Proc. Workshop on Human Language Technology*, pp. 307–312. Association for Computational Linguistics, 1994.
- [3] H. A. Bourlard and N. Morgan, *Connectionist speech recognition: a hybrid approach*. Springer Science & Business Media, 2012.
- [4] H. Bourlard, N. Morgan, C. Wooters, and S. Renals, "Cdn: A context dependent neural network for continuous speech recognition," in *Proc. IEEE Intern. Conf. on Acoustics, Speech and Signal Process. (ICASSP)*, 1992.
- [5] G. E. Dahl, D. Yu, L. Deng, and A. Acero, "Context-dependent pre-trained deep neural networks for large-vocabulary speech recognition," in *IEEE Transactions on Audio, Speech, and Language Process.*, Vol. 20, No. 1, pp. 30–42, 2011.
- [6] A. Senior, G. Heigold, M. Bacchiani, and H. Liao, "GMM-free DNN acoustic model training," in *Proc. IEEE Intern. Conf. on Acoustics, Speech and Signal Process. (ICASSP)*, pp. 5602–5606, 2014.
- [7] C. Zhang and P. C. Woodland, "Standalone training of context-dependent deep neural network acoustic models," in *Proc. IEEE Intern. Conf. on Acoustics, Speech and Signal Process. (ICASSP)*, pp. 5597–5601, 2014.
- [8] M. Bacchiani and D. Rybach, "Context dependent state tying for speech recognition using deep neural network acoustic models," in *Proc. IEEE Intern. Conf. on Acoustics, Speech and Signal Process. (ICASSP)*, pp. 230–234, 2014.
- [9] M. Razavi, R. Rasipuram, and M. Magimai-Doss, "On modeling context-dependent clustered states: Comparing HMM/GMM, hybrid HMM/ANN and KL-HMM approaches," in *Proc. IEEE Intern. Conf. on Acoustics, Speech and Signal Process. (ICASSP)*, pp. 7659–7663, 2014.
- [10] G. Gosztolya, T. Grósz, L. Tóth, and D. Imseng, "Building context-dependent DNN acoustic models using Kullback-Leibler divergence-based state tying," in *Proc. IEEE Intern. Conf. on Acoustics, Speech and Signal Process. (ICASSP)*, pp. 4570–4574, 2015.
- [11] L. Zhu, K. Kilgour, S. Stüker, and A. Waibel, "Gaussian free cluster tree construction using Deep Neural Network," in *Proc. Interspeech*, 2015.
- [12] S. Wiesler, G. Heigold, M. Nußbaum-Thom, R. Schlüter, and H. Ney, "A discriminative splitting criterion for phonetic decision trees," in *Proc. Interspeech*, pp. 54–57, 2010.
- [13] J. Chorowski, A. Lancucki, B. Kostka, and M. Zpotoczny, "Towards using context-dependent symbols in CTC without state-tying decision trees," *arXiv:1901.04379*, 2019.
- [14] P. Bell, P. Swietojanski, and S. Renals, "Multitask learning of context-dependent targets in deep neural network acoustic models," *IEEE/ACM Transactions on Audio, Speech, and Language Process.*, Vol. 25, No. 2, pp. 238–247, 2016.
- [15] T. Bayes, "An essay towards solving a problem in the doctrine of chances. by the late rev. mr. bayes, frs communicated by mr. price, in a letter to john canton, AMFR S," *Philosophical Transactions of The Royal Society of London*, No. 53, pp. 370–418, 1763.
- [16] N. Morgan and H. Bourlard, *Factoring networks by a statistical method*. MIT Press, 1992.
- [17] A. Graves, N. Jaitly, and A. R. Mohamed, "Hybrid speech recognition with deep bidirectional LSTM," in *Proc. IEEE Workshop on Automatic Speech Recognition and Understanding (ASRU)*, pp. 273–278, 2013.
- [18] A. Zeyer, P. Doetsch, P. Voigtlaender, R. Schlüter, and H. Ney, "A comprehensive study of deep bidirectional LSTM RNNs for acoustic modeling in speech recognition," in *Proc. IEEE Intern. Conf. on Acoustics, Speech and Signal Process. (ICASSP)*, pp. 2462–2466, 2017.
- [19] J. J. Godfrey, E. C. Holliman, and J. McDaniel, "Switchboard: Telephone speech corpus for research and development," in *Proc. IEEE Intern. Conf. on Acoustics, Speech and Signal Process. (ICASSP)*, Vol. 1, pp. 517–520, 1992.
- [20] A. Zeyer, T. Alkhouli, and H. Ney, "RETURNN as a generic flexible neural toolkit with application to translation and speech recognition," *arXiv:1805.05225*, 2018.
- [21] S. Wiesler, A. Richard, P. Golik, R. Schlüter, and H. Ney, "RASR/NN: The RWTH neural network toolkit for speech recognition," in *Proc. IEEE Intern. Conf. on Acoustics, Speech and Signal Process. (ICASSP)*, pp. 3281–3285, 2014.
- [22] N. Srivastava, G. E. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: a simple way to prevent neural networks from overfitting," *Journal of machine learning research*, Vol. 15, No. 1, pp. 1929–1958, 2014.
- [23] R. Schlüter, I. Bezrukov, H. Wagner, and H. Ney, "Gamma-tone features and feature combination for large vocabulary speech recognition," in *Proc. IEEE Intern. Conf. on Acoustics, Speech and Signal Process. (ICASSP)*, Vol. 4, pp. IV–649, 2007.
- [24] M. Kitza, P. Golik, R. Schlüter, and H. Ney, "Cumulative adaptation for BLSTM acoustic models," *arXiv:1906.06207*, 2019.
- [25] T. Dozat, "Incorporating nesterov momentum into adam," 2016.
- [26] A. Neelakantan, L. Vilnis, Q. V. Le, I. Sutskever, L. Kaiser, K. Kurach, and J. Martens, "Adding gradient noise improves learning for very deep networks," *arXiv:1511.06807*, 2015.
- [27] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *Proc. IEEE Intern. Conf. on Computer Vision*, pp. 2980–2988, 2017.
- [28] R. Kneser and H. Ney, "Improved backing-off for m-gram language modeling," in *Proc. IEEE Intern. Conf. on Acoustics, Speech and Signal Process. (ICASSP)*, Vol. 1, pp. 181–184, 1995.
- [29] M. Sundermeyer, R. Schlüter, and H. Ney, "Lstm neural networks for language modeling," in *Proc. Interspeech*, 2012.
- [30] E. Beck, W. Zhou, R. Schlüter, and H. Ney, "Lstm language models for lvcsr in first-pass decoding and lattice-rescoring," *arXiv:1907.01030*, 2019.