



Sequential recurrent neural networks for language modeling

Youssef Oualil^{1,2}, Clayton Greenberg^{1,2,3}, Mittul Singh^{1,3}, Dietrich Klakow^{1,2,3}

¹Spoken Language Systems (LSV)

²Collaborative Research Center on Information Density and Linguistic Encoding

³Graduate School of Computer Science

Saarland University, Saarbrücken, Germany

{firstname.lastname}@lsv.uni-saarland.de

Abstract

Feedforward Neural Network (FNN)-based language models estimate the probability of the next word based on the history of the last N words, whereas Recurrent Neural Networks (RNN) perform the same task based only on the last word and some context information that cycles in the network. This paper presents a novel approach, which bridges the gap between these two categories of networks. In particular, we propose an architecture which takes advantage of the explicit, sequential enumeration of the word history in FNN structure while enhancing each word representation at the projection layer through recurrent context information that evolves in the network. The context integration is performed using an additional word-dependent weight matrix that is also learned during the training. Extensive experiments conducted on the Penn Treebank (PTB) and the Large Text Compression Benchmark (LTCB) corpus showed a significant reduction of the perplexity when compared to state-of-the-art feedforward as well as recurrent neural network architectures.

Index Terms: Recurrent neural networks, language modeling

1. Introduction

A high quality Language Model (LM) is considered to be an integral component of many systems for language technology applications, such as speech recognition [1], machine translation [2], etc. The goal of an LM is to identify probable sequences of predefined linguistic units, which are typically words. Semantic and syntactic properties of the language, encoded by the LM, guide these predictions.

Intrinsically, the performance of an LM can be evaluated based upon its ability to predict the next word given its context. The most common approach to build such models is the word count-based method, which is commonly known as N -gram language modeling [3, 4]. By simply enumerating all possibilities over a short span of words and assigning probabilities to them directly, N -grams were difficult to outperform for a very long time.

The introduction of neural networks for language modeling led to a significant improvement over these standard models. This was mainly due to the continuous word representations they provide, which typically overcome the exponential growth of parameters that N -gram models require to enumerate possibilities. Bengio et al. [5] proposed a Feedforward Neural Network (FNN) for language modeling, as an alternative to

N -grams, to estimate the probability of a given word sequence while considering a fixed context (word history) size. This approach was very successful and has been shown to outperform a mixture of different other models [6], and to significantly improve speech recognition performance [7].

In order to overcome the fixed context size constraint and to capture long range dependencies known to be present in language, Mikolov et al. [8, 9] proposed a Recurrent Neural Network (RNN) which allows context information to cycle in the network. Another recurrence-based network architecture, Long-Short Term Memory (LSTM) [10], addresses some learning issues from the original RNN and explicitly controls the longevity of context information in the network.

Contrary to FNN, recurrent models such as RNN and LSTM predict the next word based only on the current word and the context representation. Therefore, they lose information about word position rather quickly and cannot model short range dependencies as well as FNN and N -grams. For example, English has position-dependent patterns such as “he * he” (“he said he”, “he mentioned he”, ...). The position of “he” is essential for making the right prediction in this case, and the recurrent models are not designed to encode that. Rather, they are better for smooth incremental updates and hence for longer range dependencies.

This paper proposes a novel approach that models short range dependencies like FNN and long range dependencies like RNN. In particular, the hidden layers combine explicit encoding of the local context and a recurrent architecture, which allows the context information to sequentially evolve in the network at the projection layer. In the first step, the word representation are enhanced using the context information. This step maps the word representations from a universal embedding space into a context-based space. Then, the system performs the next word prediction as it is typically done in FNN. The learning of the network weights uses the Back-Propagation Through Time (BPTT) algorithm similarly to RNN. The main difference here is the additional network error resulting from the additional sequential connections. This paper also shows that learning of word-dependent sequential connections can substantially improve the performance of the proposed network.

We proceed as follows. Section 2 presents a brief overview of FNN and RNN models. Section 3 introduces the proposed architecture which combines these two models. Then, Section 4 evaluates the proposed network in comparison to different state-of-the-art language models for perplexity on the PTB and the LTCB corpus. Finally, we conclude in Section 5.

This research was funded by the German Research Foundation (DFG) as part of SFB 1102.

2. Neural Network Language Models

The goal of a language model is to estimate the probability distribution $p(w_1^T)$ of word sequences $w_1^T = w_1, \dots, w_T$. Using the chain rule, this distribution can be expressed as

$$p(w_1^T) = \prod_{t=1}^T p(w_t | w_1^{t-1}) \quad (1)$$

The rest of this section shows how FNN and RNN are used to approximate this probability distribution.

2.1. Feedforward Neural Networks

Similarly to N -gram models, FNN uses the Markov assumption of order $N-1$ to approximate (1) according to

$$p(w_1^T) \approx \prod_{t=1}^T p(w_t | w_{t-N+1}^{t-1}) \quad (2)$$

Subsequently, each of the terms involved in this product, i.e., $p(w_t | w_{t-N+1}^{t-1})$, is estimated, separately, in a single bottom-up evaluation of the network according to

$$P_{t-i} = X_{t-i} \cdot U, \quad i = N-1, \dots, 1 \quad (3)$$

$$H_t = f \left(\sum_{i=1}^{N-1} P_{t-i} \cdot V_i \right) \quad (4)$$

$$O_t = g(H_t \cdot W) \quad (5)$$

X_{t-i} is a one-hot encoding of the word w_{t-i} , whereas the rows of U encode the continuous word representations (i.e., embeddings). Thus, P_{t-i} is the continuous representation of the word w_{t-i} . W and $V = [V_1, \dots, V_{N-1}]$ are the network connection weights, which are learned during training in addition to U . Moreover, $f(\cdot)$ is an activation function, whereas $g(\cdot)$ is the softmax function. Figure 1(a) shows an example of an FNN with a fixed context size $N-1=3$ with a single hidden layer.

2.2. Recurrent Neural Networks

An RNN attempts to capture the complete history in a context vector h_t , which represents the state of the network and evolves in time. Therefore, it approximates (1) according to

$$p(w_1^T) \approx \prod_{t=1}^T p(w_t | w_{t-1}, h_{t-1}) = \prod_{t=1}^T p(w_t | h_t) \quad (6)$$

RNN evaluates this distribution similarly to FNN. The main difference occurs in Equations (3) and (4) which are combined into

$$H_t = f(X_{t-1} \cdot U + H_{t-1} \cdot V) \quad (7)$$

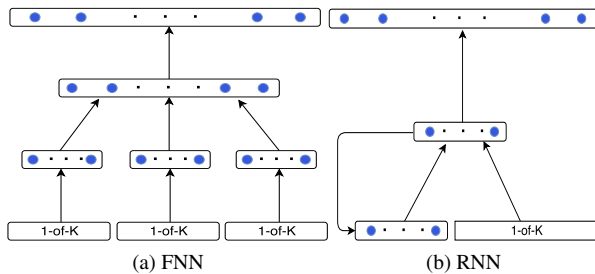


Figure 1: FNN vs RNN Architecture.

Figure 1(b) shows an example of a standard RNN. The next Section will show how an RNN can be extended to explicitly model short range dependencies through additional sequential connections.

3. Sequential Recurrent Neural Network

The main difference between an RNN and an FNN is the context representation. More precisely, The context layer H_t of an FNN is estimated based on a fixed context size i.e., the last $N-1$ words, whereas in an RNN, H_t is constantly updated (at each time iteration) using only the last word and context at time $t-1$.

3.1. The proposed Neural Architecture

We propose in this paper an architecture which captures short range dependencies over the last $N-1$ word positions as it is done in FNN, and the long range context through recurrence, similarly to RNN. The design of this structure is motivated by the inefficiency of RNN to model position dependent patterns, which are particularly frequent in conversational speech. RNN loses information about word position quickly and therefore cannot efficiently model short range dependencies. FNN and N -gram models, however, are designed as position-dependent models, which deal only with short-term context. Extending RNN structure to explicitly represent the short term history as it is done in FNN will 1) help improve the modeling of short range context, as it will 2) allow the network to capture any residual/additional context information that may be present in the past $i = t-N+1, \dots, t-2$ time iterations but which may have been lost during the last context update, which is based only on the last word at $t-1$ (See illustration in Figure 2). In the worst case scenario, the context information will be simply redundant and is expected not to harm the performance. The rest of this Section introduces the mathematical formulation of this approach.

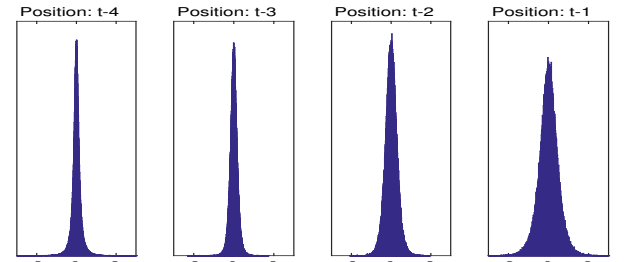


Figure 2: Histograms of the projection-to-hidden weights V_1, V_2, V_3 and V_4 (see Figure 3) for each of the 4 word positions of an SRNN ($N=5$) trained on LTCB. These histograms show that the magnitude of the weights decays with the word position (from $t-1$ to $t-4$) but does not nullify. Thus, the farther word positions still capture some residual/additional context.

The proposed Sequential Recurrent Neural Network (SRNN) approximates (1) according to

$$p(w_1^T) \approx \prod_{t=1}^T p(w_t | w_{t-N+1}^{t-1}, h_{t-N+1}) = \prod_{t=1}^T p(w_t | h_{t-N+2}^t) \quad (8)$$

The proposed architecture to estimate (8) explicitly represents the history over the last $N-1$ word positions as it is done in FNN to approximate (2) while it enhances the actual word representations using the recurrent context information, which propagates sequentially within the network. Furthermore, restricting the context to a 1-word history window ($N=2$) in (8) leads to the RNN approximation in (6). Therefore, the proposed approach can be seen as an extension of the standard RNN to explicitly model and capture short range context.

The additional sequential connections allow the context information to propagate from the past to the future within the network. These connections can be defined as a Word-Independent (WI) recurrence vector, which fixes the amount of context information allowed to propagate in the network, as they can be designed as Word-Dependent (WD) vectors. In this case, each word will have its own context weight vector, which will typically learn which context “neurons” are relevant for that particular word and therefore scales each context unit accordingly.

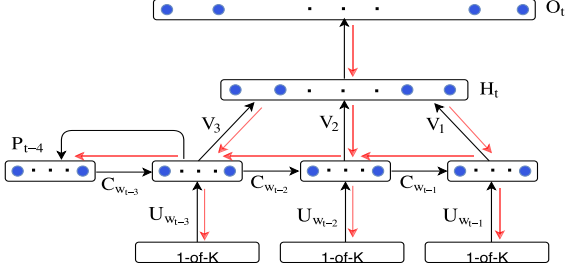


Figure 3: Sequential Recurrent Neural Network architecture. The backward path (red arrows) shows the error propagation during training (this figure does not include BPTT).

The network evaluation is performed similarly to FNN, the main difference occurs in Equation (3), which becomes in the case of the word-independent model

$$P_{t-i} = f_s(X_{t-i} \cdot U + C \odot P_{t-i-1}), \quad i = N-1, \dots, 1 \quad (9)$$

as it becomes in the case of the word-dependent model

$$P_{t-i} = f_s(X_{t-i} \cdot U + C_{w_{t-i}} \odot P_{t-i-1}), \quad i = N-1, \dots, 1 \quad (10)$$

where $f_s(\cdot)$ is an activation function and \odot is the element-wise product operator. C is the word-independent recurrence weight vector, whereas $C_{w_{t-i}}$ is the word-dependent context weight corresponding to the word w_{t-i} . Figure (3) shows an example of an SRNN with three additional sequential connections ($N-1=3$) and a single hidden layer.

The proposed SRNN model is a general architecture that includes different networks. In particular, setting $C = [0, \dots, 0]$ and $f_s(x) = x$ results in the classical FNN architecture, whereas setting $N = 2$ leads to a standard RNN with a diagonal recurrence matrix and an additional non-recurrent layer. Moreover, setting C to a fixed value in $[0, 1]$ and $f_s(x) = x$ leads to the Fixed-size Ordinally-Forgetting Encoding (FOFE) [11] architecture, which was proposed to uniquely encode word sequences.

The proposed model replaces the universal word embeddings at the projection layer of an FNN by context-dependent word embeddings. More particularly, both Equations (9) and (10) show that each word representation is enhanced using the context information before proceeding to the next word prediction. Therefore, we can see this particular step as a transformation from the universal embedding space into a context-dependent space with a better discrimination of words.

3.2. SRNN Training

The parameters to train for an SRNN are the word embeddings U , the project-to-hidden connection weights $V = [V_1, \dots, V_{N-1}]$, the hidden-to-output connection weights W and the context weight vector C for the WI model, or $C = [C_1^T, \dots, C_K^T]^T$ (K is the vocabulary size) for the WD model.

In this case, each word w in the vocabulary will be characterized by two learnable vectors, namely, the continuous representation (embedding) U_w and the context weight C_w .

Similarly to RNN, the parameter learning of an SRNN architecture follows the standard Back-Propagation Through Time (BPTT) algorithm. The main difference occurs at the projection layer, where the additional error vectors resulting from the sequential connections should be taken into account (See example of error propagation in Figure 3) before unfolding the network in time.

4. Experiments and Results

4.1. Experimental Setup

We evaluated the proposed architecture on two different benchmark tasks. The first set of experiments was conducted on the Penn Treebank (PTB) corpus using the standard division, e.g. [9, 11]: sections 0-20 are used for training while sections 21-22 and 23-24 are used for validation and testing. The vocabulary was limited to the most 10k frequent words while the remaining words were all mapped to the token $\langle \text{unk} \rangle$. In order to evaluate how the proposed approach scales to large corpora, we run a set of experiments on the Large Text Compression Benchmark (LTCB) [12]. This corpus is based on the enwik9 dataset which contains the first 10^9 bytes of enwiki-20060303-pages-articles.xml. We adopted the same training-test-validation data split and preprocessing from [11]. All but the 80k most frequent words were replaced by $\langle \text{unk} \rangle$. Details about the sizes of these two corpora and the percentage of Out-Of-Vocabulary (OOV) words that were mapped to $\langle \text{unk} \rangle$ can be found in Table 1.

Table 1: Corpus size in number of words and $\langle \text{unk} \rangle$ rate.

Corpus	Train		Dev		Test	
	#W	$\langle \text{unk} \rangle$	#W	$\langle \text{unk} \rangle$	#W	$\langle \text{unk} \rangle$
PTB	930K	6.52%	82K	6.47%	74K	7.45%
LTCB	133M	1.43%	7.8M	2.15%	7.9M	2.30%

The proposed approach (SRNN) is compared to different systems including the N -gram Kneser-Ney (KN) model and different feedforward and recurrent neural architectures. For feedforward networks, the baseline systems include 1) the FNN-based LM [5] as well as the 2) Fixed-size Ordinally Forgetting Encoding (FOFE) approach, which was implemented as a feedforward sentence-based model [11]. The FOFE results were obtained using the FOFE toolkit [11]. The results are reported for different context sizes ($N=1, 2$ and 4) and different numbers of hidden layers (1 or 2). Regarding recurrent models, we compare the proposed approach to 3) the full RNN (without classes) [9], 4) to a deep RNN [13], which investigates different ways of adding hidden layers to RNN, and finally 5) to the LSTM architecture [10], which explicitly regulates the amount of information that propagates in the network.

4.2. PTB Experiments

For the PTB experiments, the FNN, FOFE and SRNN architectures have similar configurations. That is, the hidden layer(s) size is 400 with all hidden units using the Rectified Linear Unit (ReLU) i.e., $f(x) = \max(0, x)$, as an activation function, whereas the word representation (embedding) size was set to 200 for FNN, FOFE and LSTM and 100 for SRNN. The latter uses $f_s = \tanh(\cdot)$ as sequential activation function. The hidden layer size of RNN and LSTM were set to 400 and follow the original configuration proposed in [9] and [10], respectively.

We also use the same learning setup adopted in [11]. Namely, we use the stochastic gradient descent algorithm with a mini-batch size of 200, the learning rate is initialized to 0.4, the momentum is set to 0.9, the weight decay is fixed to 4.10^{-5} and the training is done in epochs. The weights initialization follows the normalized initialization proposed in [14]. Similarly to [8], the learning rate is halved when no significant improvement in the log-likelihood of the validation data is observed. Then, we continue with seven more epochs while halving the learning rate after each epoch. The BPTT was set to 5 time steps. In the tables below, WI-SRNN refers to the word-independent SRNN model proposed in (9), whereas WD-SRNN refers to the word-dependent model in (10). For both models, the context connection weights, C , were randomly initialized in $[0, 1]$. In order to compare to the FOFE approach, we also report results where C is reduced to a scalar forgetting factor that is fixed at 0.7. This is denoted as WI-SRNN* in the tables below. We report the results in terms of perplexity (PPL), Number of model Parameters (NoP) and the training speed, which is defined as the number of words processed per second (w/s) on a GTX TITAN X GPU.

Table 2: *LMs performance on the PTB test set.*

N-1=	model			model+KN5			NoP	w/s
	1	2	4	1	2	4	4	4
1 Hidden Layer								
FNN	176	131	119	132	116	107	6.32M	24.3K
FOFE	123	111	112	108	100	101	6.32M	17.2K
WI-SRNN*	117	110	109	105	100	99	5.16M	12.9K
WI-SRNN	112	107	107	102	98	97	5.16M	11.2K
WD-SRNN	109	106	106	99	96	95	6.16M	10.4K
2 Hidden Layers								
FNN	176	129	114	132	114	102	6.48M	21.8K
FOFE	116	108	109	104	98	97	6.48M	16.6K
WI-SRNN*	114	108	107	102	98	96	5.32M	10.8K
WI-SRNN	109	105	104	99	96	94	5.32M	9.6K
WD-SRNN	108	103	104	97	94	94	6.32M	9.2K
Recurrent Models								
RNN	123			107			8.16M	20.6K
Deep RNN	107.5			—			6.96M	—
LSTM	114			99			6.96M	7.6K

Table 2 shows the LMs evaluation on the PTB test set. We can clearly see that the proposed approach outperforms all other models using the lowest Number of model Parameters (NoP) among all configurations. This also includes other models that were reported in the literature, such as RNN with maximum entropy [15], random forest LM [16], structured LM [17] and syntactic neural network LM [18]. More particularly, SRNN with two hidden layers achieves a comparable performance to a mixture of RNNs [19]. We can also conclude that the explicit modeling of short range dependencies through sequential connections improves the performance. More precisely, the results show that increasing the history window (1, 2 and 4) improves the performance for all SRNN models. Table 2 also shows that using a fixed scalar forgetting factor (WI-SRNN*) leads to a slight improvement over the FOFE approach, which is mainly due to the additional non-linear activation function f_s . Furthermore, the word-dependent (WD-SRNN) model slightly outperforms the word-independent model (WI-SRNN) but with a non-negligible increase in the number of parameters. Regarding the training speed, we can conclude that training an SRNN model requires approximately twice the time needed for FFN and RNN, whereas it needs less time compared to LSTM.

4.3. LTCB Experiments

The LTCB experiments use the same PTB setup with minor changes. The results shown in Table 3 follow the same experimental setup used in [11]. More precisely, these results were obtained without usage of momentum or weight decay whereas the mini-batch size was set to 400. The FNN and FOFE architectures contain 2 hidden layers of size 600 (or 400) whereas RNN and SRNN have a single hidden layer of size 600. In order to compare to [11], the forgetting factor C of WI-SRNN* is fixed at 0.6.

Table 3: *LMs Perplexity on the LTCB test set.*

Context Size M=N-1	model			NoP
	1	2	4	4
KN	239	156	132	—
FNN [M*200]-600-600-80k	235	150	114	64.84M
FOFE [M*200]-400-400-80k	120	115	108	48.48M
FOFE [M*200]-600-600-80k	112	107	100	64.84M
WI-SRNN* [M*200]-600-80k	110	102	94	64.48M
WI-SRNN [M*200]-600-80k	85	80	77	64.48M
WD-SRNN [M*200]-600-80k	77	74	72	80.48M
RNN [600]-600-80k	85			96.36M

The LTCB results shown in Table 3 generally confirm the PTB conclusions. In particular, we can see that SRNN models outperform all other models while requiring comparable or fewer model parameters. Moreover, the WI-SRNN* model with a single hidden layer slightly outperforms FOFE (2 hidden layers). These results, however, show a more significant improvement for the WD-SRNN model and for the increased window size (from 1 to 4) compared to the improvement obtained on the PTB. This is mainly due to the large amount of LTCB training data, which allows us to train richer WD context vectors.

Table 4: *Examples of top 5 similar words.*

in		strictly		germany	
U_w	C_w	U_w	C_w	U_w	C_w
into	at	solely	purely	italy	japan
throughout	on	rigidly	totally	france	russia
through	for	broadly	physically	britain	italy
during	their	purely	solely	switzerland	france
within	to	ostensibly	technically	england	spain

Table 4 shows some word examples with their top 5 cosine similarities for word embeddings U_w and Euclidean distance for context weights C_w . These examples show a general trend, not valid for every example, that the embeddings capture semantic (conceptual) similarities and the context weights model syntactic (functional) similarities.

5. Conclusion and Future Work

We have presented a sequential recurrent neural network which captures short range dependencies using short history windows, and models long range context through recurrent connections. Experiments on PTB and LTCB corpora have shown that this architecture substantially outperforms many state-of-the-art neural systems, due to its successful combination of the motivating features of its feedforward and recurrent predecessors. Further gains could be made by more optimally controlling the amount of information evolving in the network, as it is done in LSTM, and by more thoroughly addressing long range dependencies. These will be investigated in future work.

6. References

- [1] S. Katz, "Estimation of probabilities from sparse data for the language model component of a speech recognizer," *IEEE Transactions on Acoustics, Speech, and Signal Processing*, vol. 35, no. 3, pp. 400–401, Mar. 1987.
- [2] P. F. Brown, J. Cocke, S. A. D. Pietra, V. J. D. Pietra, F. Jelinek, J. D. Lafferty, R. L. Mercer, and P. S. Roossin, "A statistical approach to machine translation," *Comput. Linguist.*, vol. 16, no. 2, pp. 79–85, Jun. 1990.
- [3] R. Rosenfeld, "Two decades of statistical language modeling: Where do we go from here?" in *Proceedings of the IEEE*, vol. 88, 2000, pp. 1270–1278.
- [4] R. Kneser and H. Ney, "Improved backing-off for m-gram language modeling," in *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, Detroit, Michigan, USA, May 1995, pp. 181–184.
- [5] Y. Bengio, R. Ducharme, P. Vincent, and C. Jauvin, "A neural probabilistic language model," *J. Mach. Learn. Res.*, vol. 3, pp. 1137–1155, Mar. 2003.
- [6] J. Goodman, "A bit of progress in language modeling, extended version," Microsoft Research, Tech. Rep. MSR-TR-2001-72, 2001.
- [7] H. Schwenk and J. Gauvain, "Training neural network language models on very large corpora," in *Human Language Technology Conference and Conference on Empirical Methods in Natural Language Processing (EMNLP)*, Oct. 2005, pp. 201–208.
- [8] T. Mikolov, M. Karafiát, L. Burget, J. Cernocký, and S. Khudanpur, "Recurrent neural network based language model," in *11th Annual Conference of the International Speech Communication Association (INTERSPEECH)*, Makuhari, Chiba, Japan, Sep. 2010, pp. 1045–1048.
- [9] T. Mikolov, S. Kombrink, L. Burget, J. Cernocký, and S. Khudanpur, "Extensions of recurrent neural network language model," in *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, May 2011, pp. 5528–5531.
- [10] M. Sundermeyer, R. Schlüter, and H. Ney, "LSTM neural networks for language modeling," in *13th Annual Conference of the International Speech Communication Association (INTERSPEECH)*, Portland, OR, USA, Sep. 2012, pp. 194–197.
- [11] S. Zhang, H. Jiang, M. Xu, J. Hou, and L. Dai, "The fixed-size ordinally-forgetting encoding method for neural network language models," in *53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing of the Asian Federation of Natural Language Processing ACL*, vol. 2, July 2015, pp. 495–500.
- [12] M. Mahoney, "Large text compression benchmark," 2011. [Online]. Available: <http://mattmahoney.net/dc/textdata.html>
- [13] R. Pascanu, Ç. Gülçehre, K. Cho, and Y. Bengio, "How to construct deep recurrent neural networks," *CoRR*, vol. abs/1312.6026, 2013.
- [14] X. Glorot and Y. Bengio, "Understanding the difficulty of training deep feedforward neural networks," in *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics (AISTATS)*, Chia Laguna Resort, Sardinia, Italy, May 2010, pp. 249–256.
- [15] T. Mikolov, A. Deoras, D. Povey, L. Burget, and J. Cernocký, "Strategies for training large scale neural network language models," in *IEEE Workshop on Automatic Speech Recognition & Understanding (ASRU)*, Waikoloa, HI, USA, Dec. 11–15, 2011, pp. 196–201.
- [16] P. Xu and F. Jelinek, "Random forests and the data sparseness problem in language modeling," *Computer Speech & Language*, vol. 21, no. 1, pp. 105–152, 2007.
- [17] D. Filimonov and M. P. Harper, "A joint language model with fine-grain syntactic tags," in *Conference on Empirical Methods in Natural Language Processing (EMNLP), A meeting of SIGDAT, a Special Interest Group of the ACL*, Singapore, Aug. 2009, pp. 1114–1123.
- [18] A. Emami and F. Jelinek, "Exact training of a neural syntactic language model," in *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, Montreal, Quebec, Canada, May 2004, pp. 245–248.
- [19] T. Mikolov, A. Deoras, S. Kombrink, L. Burget, and J. Cernocký, "Empirical evaluation and combination of advanced language modeling techniques," in *12th Annual Conference of the International Speech Communication Association (INTERSPEECH)*, Florence, Italy, Aug. 27–31, 2011, pp. 605–608.