# RETURNN: THE RWTH EXTENSIBLE TRAINING FRAMEWORK FOR UNIVERSAL RECURRENT NEURAL NETWORKS

Patrick Doetsch, Albert Zeyer, Paul Voigtlaender, Ilia Kulikov, Ralf Schlüter, Hermann Ney

Human Language Technology and Pattern Recognition, Computer Science Department,
RWTH Aachen University, 52062 Aachen, Germany
{doetsch,zeyer,voigtlaender,kulikov,schlueter,ney}@cs.rwth-aachen.de

## ABSTRACT

In this work we release our extensible and easily configurable neural network training software. It provides a rich set of functional layers with a particular focus on efficient training of recurrent neural network topologies on multiple GPUs. The source of the software package is public and freely available for academic research purposes and can be used as a framework or as a standalone tool which supports a flexible configuration. The software allows to train state-of-the-art deep bidirectional long short-term memory (LSTM) models on both one dimensional data like speech or two dimensional data like handwritten text and was used to develop successful submission systems in several evaluation campaigns.

**Index Terms**: recurrent neural networks, lstm, rnn, speech recognition, software package, multi-gpu

## 1. INTRODUCTION

Recurrent neural networks (RNNs) and in particular LSTMs [1] now dominate most sequential learning tasks including automatic speech recognition (ASR) [2, 3], statistical machine translation (SMT) [4], and image caption generation [5]. The training of deep recurrent neural networks is considerably harder compared to pure feed-forward structures due to the accumulation of gradients over time. For a long time there were only very few implementations of the methods and topologies that are required for RNN training. This changed rapidly when solutions for automatic differentiation and symbolic representations were combined into powerful computing libraries [6]. In the machine learning community the most prominent example during that time was Theano [7, 8], which provides an extensive Python package to compute derivatives using symbolic mathematical expressions. While most of these packages allow to comfortably design neural network architectures on a low level, they do not serve as ready-to-use solutions for large scale tasks. Instead, their primary focus is generality in order to allow for various system designs without introducing constraints due to performance or usability issues. There is also naturally not much interest in getting the best performance for a particular hardware setup, but instead to keep compatibility at a maximum.

RETURNN draws on Theano as an additional layer on top of the Theano library which aims to fill in the gap between research oriented software packages and application driven machine learning software like Caffe [9]. Our software provides highly optimized LSTM kernels written in CUDA, as well as efficient training in a multi-GPU setup. We simplify the construction of new topologies using a JSON based network configuration file while also providing a way to extend the software by functional layers. The software comes with few dependencies and it is furthermore tightly integrated into the RASR software package that is also developed at our institute [10, 11]. The aim of this paper is to provide an overview of the most important aspects of the software.

The paper is organized as follows: In Section 2 we give an overview over the software RETURNN is based on, as well as competing implementations that were used for tasks in ASR. Section 3 describes how to design a neural network training setup within RETURNN. Section 4 gives an overview of the components of the tool. Section 5 then provides further information on how to extend RETURNN through additional functional layers. Finally we demonstrate the efficiency of RETURNN empirically by comparing it to TensorFlow and Torch.

## 2. RELATED WORK

Theano [7] is a Python based framework for symbolic mathematical tensor expressions with support for automatic differentiation. Expressions are modeled in a computational dependency graph which can further be augmented through an automatic optimization procedure. The implementation of each graph node is abstract and can be defined for various types of hardware like CPUs or GPUs. These properties make Theano particularly useful for neural network training tasks. By providing the required building blocks Theano allows to define complex connectionist structures that are fully differentiable. Keras [12] is a high-level Theano based framework for data-driven machine learning. It is maybe the most similar software package to RETURNN. Keras started as a pure Theano based framework but now it also supports TensorFlow as back-end with minimal restrictions. Similar projects that are built on top of the Theano library include Lasagne [13] and Blocks [14].

TensorFlow is the most recent open source machine learning package by Google [15]. It is actively developed and comes already with many predefined solutions such as LSTMs, end-to-end systems and others. TensorFlow is similar to Theano as it also works with symbolic computation graphs and automatic differentiation.

Torch [16] uses the Lua programming language and consists of many flexible and modular components that were developed by the community. In contrast to Theano, Torch does not use symbolic expressions and all calculations are done explicitly.

Other notable frameworks are C++ -based Caffe [17], Python-based Neon [18] and Brainstorm [19]. In [20] a comparison between Caffe, Neon, Theano, and Torch was done.

```
"fw_0" : { "class":"rec", "n_out":300," direction ":1 },
"bw_0" : { "class":"rec", "n_out":300," direction ":−1 },
"fw_1" : { "class":"rec", "n_out":300," direction ":1,
          "from" : ["fw_0", "bw_0"]},
"bw_1" : { "class":"rec", "n_out":300," direction ":−1,
          "from" : ["fw_0", "bw_0"]},
"output" : { "class":"softmax","from" : ["fw_1", "bw_1"]}
```

**Fig. 1**: An example network specification JSON file that realizes a bidirectional LSTM-RNN with two layers containing 300 nodes in forward and backward direction correspondingly.

Task specific software packages like RASR [11] or Kaldi [21] which are both for developing speech recognition systems, contain modules to train and decode ASR systems, including neural networks. While the EESEN package [22] extends Kaldi by adding rudimentary support for LSTMs, RETURNN extends RASR to support various recurrent neural networks architectures in ASR systems.

## 3. GENERAL USAGE

RETURNN provides a fully functional training software, which includes user interaction, a multi-batch trainer and the possibility to extract the network activations for further processing. No other dependencies besides Theano are required and network topologies will always run on CPU or GPU. During execution, RETURNN writes useful information with configurable verbosity to the standard output and a log file. Network activations can be forwarded into a HDF5 [23] file or directly be passed to the RASR decoder as described in Section 4. It is further possible to execute RETURNN in a daemon mode which allows to access model evaluation using web services.

### 3.1. Configuration

Network architectures are described using a JSON format. Each network hereby is a map from layer identification names to layer descriptions. A layer description is simply a dictionary containing a *class* parameter which specifies the layer class, an optional list of incoming layers, and layer specific parameters. When constructing the network, RETURNN looks for a layer with a specified loss and then recursively instantiates all layers that are directly or indirectly connected to it. See Figure 1 for an example of a bidirectional LSTM network.

The remaining configuration parameters are provided as simple parameters. A typical configuration file contains the task, the path or descriptor of the input data, a learning rate together with suitable adjustment methods, and information on how batches should be crafted. The configuration parameters can also be merged into the network JSON description file or even provided fully in Python format, such that a single configuration file can be used. The Python format further allows to define custom layer types and other functions in the configuration file and to use them in the network. We provide some demo setups and configuration files together with the release of the software.

### 3.2. Layers

Layers are the fundamental building blocks of RETURNN. Each layer is a named class which is callable in the JSON description file by specifying its constructor parameters. We already provide a rich set of feed-forward layers including convolutional operations [24], and support the most common activation functions. Convolutional layers hereby make use
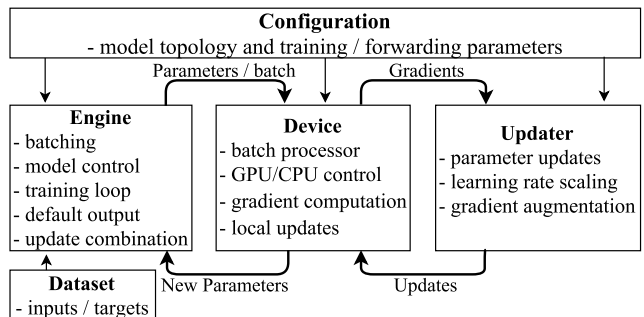


**Fig. 2**: The RETURNN processing pipeline. Sequences are generated and passed over to the main engine. The engine combines sequences into batches and performs an epoch-wise training of the network parameters using one or more devices. Each device computes the error signal for its batch and generates updates according to the optimizer.

of the efficient kernel implementations within CuDNN. Network behavior can further be augmented by using functional components like sampling and windowing. Output layers with various loss functions are available, including the cross-entropy, the mean-squared-error and Connectionist Temporal Classification (CTC).

The main focus however lies in the recurrent layers. Different cell implementations including (one- and two-dimensional) LSTM, gated recurrent units (GRU) [25], associative LSTM [26] and many more variants are available. Recurrent layers can further be connected by passing over the final state from one RNN to another one, allowing for encoder / decoder topologies with attention as described in [4]. A configurable attention mechanism is available to calculate expected inputs from encoder networks. Using a similar method, these recurrent layers further allow for basic language modeling. An example can be seen in Figure 1. Several layers can further be composed into sub-networks and then used as regular layers, which allows to model high order and circular dependencies between layers.

## 4. ENGINE

Large vocabulary speech recognition tasks have memory requirements that are significantly larger than the memory limitations of the operating hardware. This is particularly true for GPUs. We therefore implemented a data caching technique that minimizes hard disk usage while keeping the amount of allocated memory below a configurable threshold. Also in many tasks the lengths of the sequences deviate by a large amount and combining sequences into batches requires to process many additional frames that were added by zero-padding. The software therefore provides an option to chunk sequences into (possibly overlapping) segments of constant length. By sacrificing contextual information provided in the end states of the chunks, chunking allows to make a much more efficient use of the GPU memory [27]. RETURNN also supports a generic pre-training scheme where simpler network topologies are automatically generated based on a given network topology. A currently experimental Torch-Theano bridge which will be released with this software further allows to run Torch code within RETURNN.

In the recent release 0.8 of RASR [10], we added several generic Python interfaces, which allow to pass data in between of RASR and RETURNN efficiently. These interfaces can
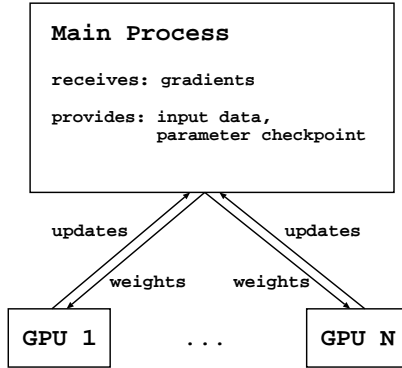
```
┌─────────────────────────────────────┐
│ Main Process                        │
│                                     │
│ receives: gradients                 │
│                                     │
│ provides: input data,               │
│           parameter checkpoint      │
│                                     │
│                                     │
└─────────────────────────────────────┘
        ↗           ↖
   updates            updates

      weights   weights

┌──────────┐              ┌──────────┐
│  GPU 1   │    ...       │  GPU N   │
└──────────┘              └──────────┘
```

**Fig. 3**: The processing pipeline in multi-GPU training. An initial parameter set $\theta$ is passed to all workers. The workers make consecutive updates without synchronizing the parameters after every update. After processing three batches, the workers send their current parameter estimate back to the CPU process where they get combined into a single set of parameters.

be used to perform the feature extraction within RASR while passing the resulting inputs to the network in real-time. They also provide a method to send the output activations of a network to RASR in order to perform decoding or to retrieve an error signal which was calculated based on discriminative training criteria available in RASR.

### 4.1. Multi-GPU training

Modern machines consist of several GPUs where each of these cards defines an isolated computation system. These computation systems can be used as independent sub-batch processors. Unfortunately, the library internally only allows to handle a single device context. We therefore chose to implement the multi-GPU functionality as an interaction of several independent system processes similar to [28]. Each GPU is attached to its own sub-process. Only the main process, which is scheduled on the CPU, has access to the real network parameters. When data batches are processed, a user specified number of batches is assigned to each device and the corresponding data is copied to the GPU memory. The main process then provides an image of the current network parameters to each of the GPU workers, which will apply their updates asynchronously batch by batch. After processing a specific amount of batches the GPU workers send their modified network parameters back to the main process where they are combined into a single set of parameters by averaging. The overall process is depicted in Figure 3. The processes communicate via sockets using a simple self-defined protocol. Weight matrices are transferred as serialized arrays, which significantly slows down training if the workers are synchronized too often. However, in our experiments we observe very stable convergence even if we only synchronize once per epoch. In fact, we often observe a regularizing effect and measure a better generalization error when keeping the GPUs asynchronous for several hundred batches.

### 4.2. CUDA Kernels for 1D and 2D LSTM Layers

We noticed that a straightforward LSTM implementation in Theano using scan (as used in Keras and Lasagne) is not very efficient in terms of both speed and memory. We therefore chose to implement the LSTM kernels directly using CUDA and cuBLAS [24]. The non-recurrent part of the LSTM for-

ward computations are performed in a single matrix multiplication for the whole mini-batch of sequences. The same applies to the back propagation step with respect to the weights and the inputs, after the recurrent part is back propagated through time. Furthermore, we reuse memory wherever possible and use custom CUDA kernels for the LSTM gating mechanism.

To the best of our knowledge, we provide the first publicly available GPU-based implementation of multidimensional long short-term memory (MDLSTM) [29]. In an MDLSTM layer, the hidden state $h(u, v)$ for position $(u, v)$ is calculated based on the predecessor hidden activations $h(u - 1, v)$ and $h(u, v - 1)$, which means that it can only be computed after both predecessor states are known. As a consequence, previous (CPU-based) implementations of MDLSTM [30] only process one pixel at a time by traversing the image column-wise in an outer loop and row-wise in an inner loop. We noticed that the activations for all positions on a common diagonal can be computed at the same time, which allows us to exploit the massive parallelism offered by modern GPUs. Additionally, we process multiple images and also the four directions of a multi-directional MDLSTM layer simultaneously by using batched cuBLAS operations and custom CUDA kernels. Optionally, the stable cell described in [31] can be used to improve convergence.

### 4.3. Optimization

When optimizing deep RNNs, regular stochastic gradient descent may not allow the network to converge to a fixed point in weight space and more sophisticated methods are needed. In those cases, learning rate scaling schedules aim to estimate a better parameter dependent update step. In RETURNN many well known learning rate schedules are implemented, including Adagrad, Adadelta and Adam [32, 33, 34]. Furthermore RETURNN allows for both the classical momentum term and also the simplified Nesterov accelerated gradient [35]. Decreasing the learning rate during training can be done based on the validation error. In particular noise addition, norm constraints and outlier detection mechanisms allow for a better convergence and avoid numerical instabilities. Note that batches gradients are not scaled in RETURNN and the dimension of the batch has a direct influence on the norm of the gradients. Regularization is possible using dropout on the layer inputs of any layer or by penalizing large L2 norms of the weight matrices.

## 5. EXTENSIBILITY

RETURNN is mostly written in Python with some parts extended by modules using the C++ CUDA API (see Section 4.2) and follows an object-oriented design. Any layer described in Section 3.2 can be used as base class to extend the package by new functional elements. Each layer is hereby considered as a black box that reads a batch of sequences and writes a batch of sequences, possibly of different shape. In order to avoid influences of zero-padding when multiple sequences of different lengths are processed together, we use an index tensor which indicates for each time step and batch, whether the frame should be considered as part of the sequence or not. The layer definition itself can be any kind of Theano expression. Each layer class is provided with a list of incoming layers and the index tensor. The layer is expected to create a 3D tensor, with time (or sequence progress) as first, the batch index as second and the layer output size as third dimension. Likewise each layer in the list of incoming layers will provide a member

**Table 1**: Comparison of runtime and memory requirement for different software packages. The numbers were averaged over 10 training epochs. Note that precise memory usage estimates for Torch and TensorFlow can not be obtained due to their internal memory management.

| Toolkit | Runtime [sec] | Memory [GB] | FER [%] |
|---|---|---|---|
| RETURNN | 198 | 2.4 | 42.51 |
| Theano LSTM | 366 | 3.2 | 42.63 |
| Keras | 619 | 5.4 | 44.36 |
| TensorFlow | 693 | $\sim 7.2$ | 47.41 |
| Torch (CuDNN) | 164 | $\sim 2.6$ | 43.02 |

called *output* with above defined shape.

A newly written layer class can directly be executed using the JSON description file, where the variables of the corresponding JSON object are passed on as arguments to the constructor of the layer.

### 5.1. Data Handling

The dataset is abstracted as a generic interface. Any dataset can provide multiple inputs and output targets of variable dimensionality and shape, where inputs and outputs can be encoded sparsely. We have a wide range of dataset implementations. Most prominently we support the HDF5 hierarchical data format, which is also used as format for models produced in RETURNN. Moreover, features from RASR can directly be used within RETURNN as described in Section 4. The release of this software contains several examples of dataset usages.

## 6. EXPERIMENTS

We demonstrate the performance of RETURNN on framewise labeled speech data from the CHiME dataset [36]. Our aim is to show that RETURNN successively converges during training and compare it to other implementations. Each frame consists of 17 consecutive speaker-adapted 16-dimensional MFCC vectors reduced to 45 dimensions by LDA. The vectors were labeled with 1501 tied allophone states using a Viterbi alignment obtained from a previously trained hidden Markov model. The segments have an average length of 738 frames with a variance of 291 frames. To ease the processing in recurrent neural networks we divided all observation sequences into constant chunks of 250 frames, padding zero-frames if required. We compare RETURNN to Keras, Torch and TensorFlow. In Torch we use a recently published CuDNN based LSTM implementation [1], which is also planned to be migrated into Theano based frameworks. For each package we measure the average runtime, average memory consumption and relative number of misclassified frames on 10% of the training data. Training is performed on 81 chunks in parallel on a NVIDIA GTX 1080. Three bidirectional LSTM layers were used in the experiments containing 512 units in each direction. A similar configuration achieved a word error rate of 6.49% and 8.43% on the development set and evaluation set of the corpus.

It can be seen in Table 1 that the internal LSTM kernel of RETURNN outperforms all competitors except the CuDNN implementation w. r. t. runtime and memory usage. In order to provide a more direct comparison of the LSTM implementations, we also present the runtime of RETURNN with an LSTM version that does not make use of our optimized LSTM kernels ("Theano LSTM" in Table 1). The internal memory
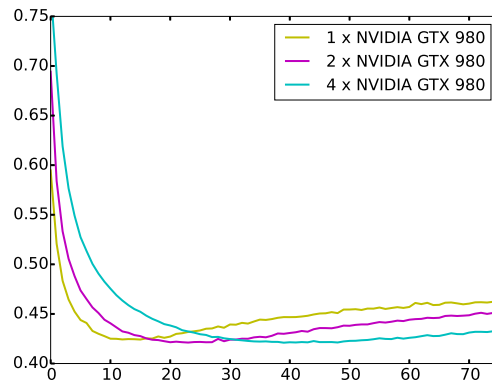


**Fig. 4**: Frame-error rate on the CHiME dataset over 75 epochs when training with multiple devices. The arrows indicate the minimal frame error and the total training time for one, two or four GPUs.

management of TensorFlow and Torch make it difficult to obtain exact measurements of their memory usage, but we can see that 25% less memory is required in our LSTM kernel compared to the Theano based kernels, including Keras.

We also conducted experiments to evaluate the runtime and classification performance of RETURNN on multiple GPUs. Here, the training time per epoch was 140, 80 and 41 seconds for one, two or four NVIDIA GTX 980 GPUs respectively. The evolution of the frame error rate (FER) and the corresponding minima are shown in Figure 4. We can see that convergence time can be significantly reduced by using multiple devices. We further observe a smoothing effect from the model averaging, such that the system trained on four GPUs achieved the lowest frame error in this experiment.

## 7. CONCLUSIONS

We presented RETURNN, a highly configurable training framework for neural networks. The software is only based on Theano and CUDA and provides very fast training procedures for recurrent neural networks upon others. It further includes a rich set of functional layers which can be applied in new network designs using a convenient JSON syntax. RETURNN was successfully used in several recent evaluation campaigns, including the READ handwriting competition on ICFHR 2016, the IWSLT 2016 German speech transcription task, and the CHiME speech separation and recognition challenge in 2016 where we ranked first, first, and second respectively.

By providing an RNN training framework, which allows to train neural networks with minimal configuration effort, we hope to increase interest in this research area and to allow more people to access these methods. RETURNN can be downloaded on our institute's website[2] and is freely available for academic research purposes.

## 8. ACKNOWLEDGEMENTS

---

[1] https://github.com/soumith/cudnn.torch

[2] https://www-i6.informatik.rwth-aachen.de/web/Software/index.html

## 9. REFERENCES

[1] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, 1997.

[2] Haşim Sak, Andrew Senior, and Françoise Beaufays, "Long short-term memory based recurrent neural network architectures for large vocabulary speech recognition," *arXiv preprint arXiv:1402.1128*, 2014, `http://arxiv.org/pdf/1402.1128`.

[3] Albert Zeyer, Patrick Doetsch, Paul Voigtlaender, Ralf Schlüter, and Hermann Ney, "A comprehensive study of deep bidirectional LSTM RNNs for acoustic modeling in speech recognition," in *International Conference on Acoustics, Speech, and Signal Processing*, 2017.

[4] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio, "Neural machine translation by jointly learning to align and translate," *arXiv preprint arXiv:1409.0473*, 2014.

[5] Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan, "Show and tell: A neural image caption generator," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 3156–3164.

[6] Atilim Gunes Baydin, Barak A. Pearlmutter, and Alexey Andreyevich Radul, "Automatic differentiation in machine learning: a survey," *CoRR*, vol. abs/1502.05767, 2015.

[7] Frédéric Bastien, Pascal Lamblin, Razvan Pascanu, James Bergstra, Ian J. Goodfellow, Arnaud Bergeron, Nicolas Bouchard, and Yoshua Bengio, "Theano: new features and speed improvements," Deep Learning and Unsupervised Feature Learning NIPS 2012 Workshop, 2012.

[8] James Bergstra, Olivier Breuleux, Frédéric Bastien, Pascal Lamblin, Razvan Pascanu, Guillaume Desjardins, Joseph Turian, David Warde-Farley, and Yoshua Bengio, "Theano: a CPU and GPU math expression compiler," in *Proceedings of the Python for Scientific Computing Conference (SciPy)*, June 2010.

[9] Yangqing Jia, Evan Shelhamer, Jeff Donahue, Sergey Karayev, Jonathan Long, Ross Girshick, Sergio Guadarrama, and Trevor Darrell, "Caffe: Convolutional architecture for fast feature embedding," *arXiv preprint arXiv:1408.5093*, 2014.

[10] David Rybach, Stefan Hahn, Patrick Lehnen, David Nolden, Martin Sundermeyer, Zoltan Tüske, Simon Wiesler, Ralf Schlüter, and Hermann Ney, "RASR - the RWTH Aachen university open source speech recognition toolkit," in *IEEE ASRU*, Waikoloa, HI, Dec. 2011.

[11] Simon Wiesler, Alexander Richard, Pavel Golik, Ralf Schlüter, and Hermann Ney, "RASR/NN: The RWTH neural network toolkit for speech recognition," in *IEEE International Conference on Acoustics, Speech, and Signal Processing*, Florence, Italy, May 2014, pp. 3313–3317.

[12] François Chollet, "Keras," `https://github.com/fchollet/keras`, 2015.

[13] "Lasagne, lightweight library to build and train neural networks in Theano," `http://lasagne.readthedocs.org/`, Accessed: 2016-03-17.

[14] Bart van Merriënboer, Dzmitry Bahdanau, Vincent Dumoulin, Dmitriy Serdyuk, David Warde-Farley, Jan Chorowski, and Yoshua Bengio, "Blocks and Fuel: Frameworks for deep learning," *arXiv preprint arXiv:1506.00619*, 2015.

[15] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dan Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng, "TensorFlow: Large-scale machine learning on heterogeneous systems," 2015, Software available from tensorflow.org.

[16] R. Collobert, S. Bengio, and J. Mariéthoz, "Torch: a modular machine learning software library," Tech. Rep., IDIAP, 2002.

[17] Yangqing Jia, Evan Shelhamer, Jeff Donahue, Sergey Karayev, Jonathan Long, Ross Girshick, Sergio Guadarrama, and Trevor Darrell, "Caffe: Convolutional architecture for fast feature embedding," in *Proceedings of the ACM International Conference on Multimedia*. ACM, 2014, pp. 675–678.

[18] "Neon: Fast, scalable, easy-to-use Python based deep learning framework by Nervana™," `https://github.com/nervanasystems/neon`, Accessed: 2016-03-17.

[19] "Brainstorm," `https://github.com/IDSIA/brainstorm`, 2015.

[20] Soheil Bahrampour, Naveen Ramakrishnan, Lukas Schott, and Mohak Shah, "Comparative study of Caffe, Neon, Theano, and Torch for deep learning," *arXiv preprint arXiv:1511.06435*, 2015.

[21] Daniel Povey, Arnab Ghoshal, Gilles Boulianne, Lukas Burget, Ondrej Glembek, Nagendra Goel, Mirko Hannemann, Petr Motlicek, Yanmin Qian, Petr Schwarz, Jan Silovsky, Georg Stemmer, and Karel Vesely, "The Kaldi speech recognition toolkit," in *IEEE 2011 Workshop on Automatic Speech Recognition and Understanding*, Rue Marconi 19, Martigny, Dec. 2011, number Idiap-RR-04-2012, IEEE Signal Processing Society, IEEE Catalog No.: CFP11SRW-USB.

[22] Yajie Miao, Mohammad Gowayyed, and Florian Metze, "EESEN: end-to-end speech recognition using deep RNN models and WFST-based decoding," *CoRR*, vol. abs/1507.08240, 2015.

[23] The HDF Group, "Hierarchical data format version 5," 2000-2010.

[24] Paul Voigtlaender, Patrick Doetsch, and Hermann Ney, "Handwriting recognition with large multidimensional long short-term memory recurrent neural networks," in *International Conference on Frontiers in Handwriting Recognition*, Shenzhen, China, Oct. 2016, pp. 228–233.

[25] Junyoung Chung, Çaglar Gülçehre, KyungHyun Cho, and Yoshua Bengio, "Gated feedback recurrent neural networks," *CoRR*, vol. abs/1502.02367, 2015.

[26] Ivo Danihelka, Greg Wayne, Benigno Uria, Nal Kalchbrenner, and Alex Graves, "Associative long short-term memory," *CoRR*, vol. abs/1602.03032, 2016.

[27] Patrick Doetsch, Michal Kozielski, and Hermann Ney, "Fast and robust training of recurrent neural networks for offline handwriting recognition," in *International Conference on Frontiers in Handwriting Recognition*, Crete, Greece, Sept. 2014, pp. 279–284.

[28] Jeffrey Dean, Greg Corrado, Rajat Monga, Kai Chen, Matthieu Devin, Mark Mao, Marc'aurelio Ranzato, Andrew Senior, Paul Tucker, Ke Yang, Quoc V. Le, and Andrew Y. Ng, "Large scale distributed deep networks," in *Advances in Neural Information Processing Systems 25*, F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, Eds., pp. 1223–1231. 2012.

[29] Alex Graves and Jürgen Schmidhuber, "Offline handwriting recognition with multidimensional recurrent neural networks," in *Neural Information Processing Systems Foundation*, 2008, pp. 545–552.

[30] Alex Graves, "RNNLIB: A recurrent neural network library for sequence learning problems," http://sourceforge.net/projects/rnnl/.

[31] Gundram Leifert, Tobias Strauß, Tobias Grüning, and Roger Labahn, "Cells in multidimensional recurrent neural networks," *arXiv preprint arXiv:1412.2620*, 2014.

[32] John Duchi, Elad Hazan, and Yoram Singer, "Adaptive subgradient methods for online learning and stochastic optimization," *J. Mach. Learn. Res.*, vol. 12, pp. 2121–2159, July 2011.

[33] Matthew D. Zeiler, "ADADELTA: an adaptive learning rate method," *CoRR*, vol. abs/1212.5701, 2012.

[34] Diederik P. Kingma and Jimmy Ba, "Adam: A method for stochastic optimization," *CoRR*, vol. abs/1412.6980, 2014.

[35] Ilya Sutskever, James Martens, George E. Dahl, and Geoffrey E. Hinton, "On the importance of initialization and momentum in deep learning," in *Proceedings of the 30th International Conference on Machine Learning, ICML*, June 2013, pp. 1139–1147.

[36] Jon Barker, Ricard Marxer, Emmanuel Vincent, and Shinji Watanabe, "The third 'CHiME' speech separation and recognition challenge: Dataset, task and baselines," in *IEEE Workshop on Automatic Speech Recognition and Understanding (ASRU), Scottsdale, AZ*, 2015, pp. 504–511.