

End-to-End Training of Acoustic Models for Large Vocabulary Continuous Speech Recognition with TensorFlow

Ehsan Variani, Tom Bagby, Erik McDermott, Michiel Bacchiani

Google Inc, Mountain View, CA, USA

{variiani, tombagby, erikmcd, michiel}@google.com

Abstract

This article discusses strategies for end-to-end training of state-of-the-art acoustic models for Large Vocabulary Continuous Speech Recognition (LVCSR), with the goal of leveraging TensorFlow components so as to make efficient use of large-scale training sets, large model sizes, and high-speed computation units such as Graphical Processing Units (GPUs). Benchmarks are presented that evaluate the efficiency of different approaches to batching of training data, unrolling of recurrent acoustic models, and device placement of TensorFlow variables and operations. An overall training architecture developed in light of those findings is then described. The approach makes it possible to take advantage of both data parallelism and high speed computation on GPU for state-of-the-art sequence training of acoustic models. The effectiveness of the design is evaluated for different training schemes and model sizes, on a 20,000 hour Voice Search task.

Index Terms: speech recognition, tensorflow

1. Introduction

In recent years there has been an explosion of research into new acoustic models for LVCSR. Large performance gains have resulted from the shift from Gaussian Mixture Models (GMMs) to feed-forward Deep Neural Networks (DNNs) trained at the frame level [1], and further gains have been obtained from utterance-level sequence training of DNNs [2, 3]. Significant additional gains have resulted from switching from feed-forward DNNs to Recurrent Neural Networks (RNNs), in particular Long Short Term Memory Models (LSTMs), again trained both at the frame level and sequence level [4]. Further architecture exploration has included evaluation of Convolutional Neural Networks [5] and much deeper networks such as Residual Networks [6].

A common thread in the rapid evolution of modeling architectures is increasing model size, and use of larger training sets, creating new challenges for the efficient use of computational resources. There is a clear need for well engineered, carefully designed acoustic model trainers which can leverage both data and model parallelism and high performance computational resources such as Graphical Processing Units (GPUs), while still allowing the use of the complex, non-cohesive decoder or lattice operations necessary for state-of-the-art sequence training, which easily scales up to much larger models and larger training data sets, and that produces state-of-the-art performance.

A number of general purpose neural network training packages such as TensorFlow [7] are now available to the research community, as are off-the-shelf fast computation hardware resources such as GPUs. However, speech recognition has its own specific needs. For example, the current state-of-the-art model in speech is the LSTM, which presents unique challenges for efficient computation. The specific choices for unrolling, e.g. full

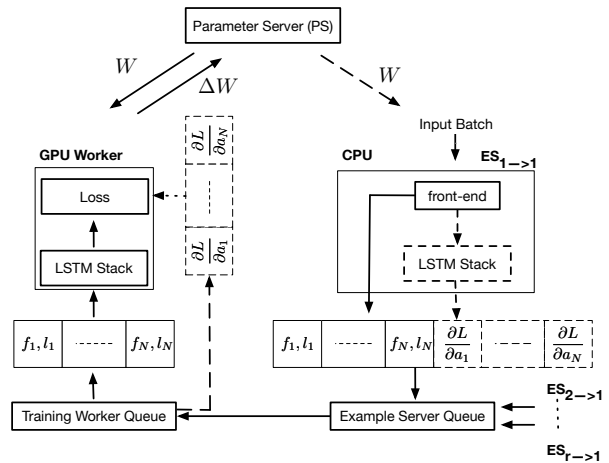


Figure 1: The end-to-end training system: Parameter Server (PS), Example Server (ES) and workers. The Example Servers pass training examples to server queues which then can be parsed and batched by worker queues. ES's have access to PS's, necessary for sequence training.

unrolling up to the utterance length [8] vs truncated unrolling [9], optionally with state-saving, and the choice of batching, have significant consequences for training efficiency. In addition to the modeling complexities, there are ASR-specific optimization challenges. In particular, sequence training of acoustic models is a critical step, but the decoder and lattice operations it involves are not easy to implement on GPUs efficiently, in a manner that scales to larger models. This has led to hybrid use of CPU and GPU (placing decoding/lattice operations on the CPU, and computation of local acoustic posterior scores on the GPU) [10], “lattice free” use of a simplified training LM with forward-backward based sequence training on GPUs [11], or biting the bullet on porting complex lattice operations to GPUs [12], among many studies.

In this article we present a comprehensive, “end-to-end” view of these critical issues in acoustic model trainer design, in light of the functions provided by the TensorFlow platform. We benchmark different strategies for recurrent model unrolling, data batching, and device placement of TensorFlow operations. Building on a previous approach to asynchronous sequence training of DNNs and LSTMs [4, 13], we describe an architecture in which the computation of outer and inner derivatives is decoupled between CPU and GPU, while still leveraging data parallelism using e.g. 100s of workers. We then present experimental results evaluating this TensorFlow-based training architecture on a 20,000 hour Voice Search training set, for different model sizes, achieving high final word accuracy in much reduced training times.

2. Acoustic Modeling Training with TensorFlow

TensorFlow distinguishes itself from other machine learning platforms such as DistBelief [14] in the way it expresses computation as stateful dataflow graphs. The dataflow computation paradigm and ability to distribute computation across machines gives us a powerful tool for acoustic model training. With careful placement of different parts of this graph onto different machines and devices, we can train complete state-of-the-art LVCSR systems using a mix of GPU and CPU devices. Figure 1 presents one way of constructing an acoustic model trainer which uses three major components, Example Servers (ES), Training Workers and Parameter Servers (PS). The level of parallelization of this system is specified by the number of training workers, n , the number of ES’s per worker, r and also efficient parametrization of TF queues used to communicate between example servers and training workers.

2.1. Distributed Feature Extraction

A typical ASR front-end consists of a sequence of computations such as Framing, Windowing, Short Term Fourier Transform, and finally log Mel filtering. In addition, we add noise and distortion to training examples on-the-fly, which is expensive to compute and requires I/O from additional data sources. We distribute front-end feature extraction and noisification onto a set of machines, as shown in the right side of Figure 1. Each machine, denoted as an Example Server (ES), is a multi-core CPU machine with access to a shard of training data, and optionally for sequence training, access to the parameter servers in order to compute sequence loss outer derivatives. For each training worker i , there are r example servers, $ES_{1 \rightarrow i}, \dots, ES_{r \rightarrow i}$ extracting features using a shared computation graph to serve worker i . Each example server converts extracted information into `tf.SequenceExample` format and streams the result back to the corresponding training worker using TF native queues. On the client side, standard parsing and queuing functions in TF are used to produce batches of training data. For truncated unrolling, `batch_sequences_with_states` creates batches of segments of sequences and propagates states between segments. For full unrolling, `bucket_by_sequence_length` groups training examples by their length and batching examples of similar length into buckets. Correct tuning of the input queue parameters, such as queue capacities, is very important to the performance of the training setup.

2.2. Parallelization methods

Data parallelism uses copies of the model, processing different batches from the training data in parallel. In a standard TensorFlow data parallel training setup, each training replica computes gradients for a mini-batch and then asynchronously applies gradient updates to weights stored in a parameter server.

For model parallelism, each layer of the stacked LSTM architecture is assigned to its own GPU. As step n of a given layer only depends on step $n - 1$ of that layer and step n of the previous layer, it is possible to pipeline computation and make effective use of multiple GPUs with close to linear scaling. Parallelizing data movement is handled transparently by TensorFlow and latency of data transfers overlaps well with compute. The limitation of this approach is that GPU use is dictated by architecture, each layer must map to a single compute device.

2.3. Training Recurrent architecture

TensorFlow provides two functions for unrolling RNNs: `static_rnn` and `dynamic_rnn`. The weights and architecture of the RNN are separated from the mechanism of unrolling them in time. `static_rnn` creates an unrolled graph for a fixed RNN length; a complete subgraph of the RNN operations is repeated for each time step. The limitations of this are excess memory use, slow initial creation of a large graph, and sequence length cannot be longer than the fixed unrolling. However, this simple unrolling can better pipeline timesteps across multiple GPUs more effectively. `dynamic_rnn` solves memory and sequence length limitations by using a `tf.WhileLoop` to dynamically unroll the graph when it is executed. That means graph creation is faster and batches can be of variable sequence length. Performance for both is similar, and both can be used with state saving for truncated unrolling.

Figure 2 shows the number of frames per second for different batch and unrolling sizes. It can be seen that increasing batch size increases throughput much more than longer unrolling. The fixed unrolling setup can also be made faster using model parallelization placing each layer of the model on one GPU.

3. End-to-end training scheme

Cross Entropy (CE) training is performed using the tandem of example servers described earlier, running on CPU, and training workers, running on either CPU or GPU, that compute model updates using the standard CE loss function, calculated with the features and targets provided by the example servers. The `tf.nn.softmax_cross_entropy_with_logits` can be used for CE training. Alternatively, Connectionist Temporal Classification (CTC) [15] can be used; this is provided in TF as the `tf.nn.ctc_loss` op.

Sequence training uses the same tandem of example servers and training workers, illustrated in Figure 1, but with outer derivatives instead of CE targets, and a simple auxiliary loss function that implements an asynchronous chaining of outer and inner derivatives [13] constituting the total sequence-level loss gradient,

$$\frac{\partial \mathcal{L}(\mathbf{X}, \theta)}{\partial \theta_i} = \sum_{t=1}^T \sum_k^N \frac{\partial \mathcal{L}(\mathbf{X}, \theta)}{\partial a(\mathbf{x}_t, k, \theta)} \frac{\partial a(\mathbf{x}_t, k, \theta)}{\partial \theta_i}, \quad (1)$$

for a sequence \mathbf{X} of T feature vectors \mathbf{x}_t , unfolded RNN logits $a(\mathbf{x}_t, k, \theta)$, for all N network output classes, and all network parameters θ_i [3]. The example servers compute the outer derivatives,

$$w(\mathbf{x}_t, k, \theta') = \left. \frac{\partial \mathcal{L}(\mathbf{X}, \theta)}{\partial a(\mathbf{x}_t, k, \theta)} \right|_{\theta=\theta'}, \quad (2)$$

for a snapshot of the parameters θ' obtained from the parameter server; and the training workers use the auxiliary function

$$\mathcal{L}_{AUX}(\mathbf{X}, \theta) = \sum_{t=1}^T \sum_{k=1}^N w(\mathbf{x}_t, k, \theta') \log p_{\theta}(k|\mathbf{x}_t), \quad (3)$$

where $p_{\theta}(k|\mathbf{x}_t)$ is the network output, i.e. the usual softmax over the logits. It is easy to show that the gradient of this auxiliary function with respect to the live parameters θ then approximates the chaining of outer and inner derivatives in Eq. (1) [13]. The computation of sequence loss outer derivatives on the

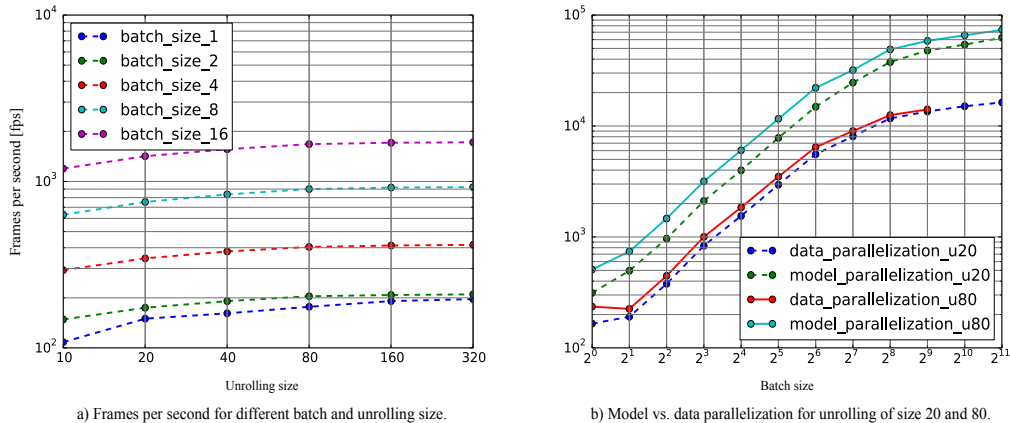


Figure 2: Benchmark of different training schemes for training 5-layer LSTMs with 768 cells per layer using 5 GPUs.

example servers uses on-the-fly lattice generation and forward-backward over the error expectation semi-ring.

Using difference machines (example servers vs. training workers) to compute outer and inner derivatives allows for faster computational throughput, but results in a significantly higher degree of asynchrony than in [3, 13]. Tuning the number of workers, queue size, and filtering out outer derivative sequences computed for a θ' deemed overly stale, were found to be important to obtaining good convergence.

The snapshot parameters θ' must be periodically refreshed; this can be implemented in TensorFlow by adding to the inference graph used in the example server sequence training pipeline a set of operations copying live variables on the parameter server to the snapshot variables.

This asynchronous approach to sequence training offers flexibility in the choice of batching and unrolling schemes for RNN training discussed in Section 2.3, as it decouples computation of outer and inner derivatives. As outer derivative computation for the sequence-level loss requires the entire utterance, `dynamic_rnn` can be used. The issues with full unrolling previously discussed are mitigated by the fact that the outer derivatives are computed with just a forward pass. On the inner derivative side, large batches with limited unrolling can leverage the efficiency of GPU-based computation.

Sequence training is performed starting from a CE-trained acoustic model. The sequence training criterion in this study is state-level Minimum Bayes Risk (sMBR) [2, 3]. Lattices are generated on-the-fly using a bigram LM estimated over the training utterance transcripts.

4. Experiments

Training a recurrent neural network involves the choice of many parameters, including unrolling and batch size. Hardware can impose constraints on the optimal choice of these parameters. The experiments are designed to investigate the effect of these choices on Word Error Rate (WER).

Data: The training data consists of about 20,000 hours of spontaneous speech from anonymized, human-transcribed Voice Search queries. For noise robustness training, each utterance in the training set is artificially noisified using 100 different styles of noises with varying degrees of noise and reverbera-

tion. The test sets are separate anonymized, human-transcribed Voice Search datasets of about 25 hours each. Evaluation is presented on two sets, *clean*, from real Google voice search traffic and *noise*, an artificially noisified version of the clean test set, with the same configuration used for training set noisification.

Front-end: The training examples are 480-dimensional frames output by the front end every 30 ms, generated by stacking four frames of 128 log mel filterbank energies extracted from a 32 ms window shifted over the audio samples in 10 ms increments [16]. During CE training, to emulate the effect of the right context, we delay the output state label by 5 frames [17]. Each frame is labeled with one out of 8192 context-dependent output phoneme states; each feature frame and target label are bundled by the example server in `tf.SequenceExample` format. For sequence training, the example server bundles outer derivatives instead of target labels, along with the feature frame and sMBR-related statistics, e.g. the utterance level loss value itself, for summarization in TensorBoard. The outer derivatives are computed over the entire utterance using the model parameters loaded in the example server at the time of processing the utterance, as described in Section 3.

Baseline: For the baseline, an integrated architecture performing feature extraction and training on the same machine, and parallelized across CPUs, was used to produce the state-of-the-art system. This system was trained with a batch size of 1 and full unrolling. The training used a total of 1600 CPU machines for all training workers and parameter servers. The training parameters for this setup have been highly tuned for the same task, with performances reported in first row of Table 1. The architecture shown in Figure 1 was used for the rest of the experiments presented in this paper.

Training: ASGD was used for all experiments. For CE training, a ratio of 10 to 1 is used for the number of example servers per worker. For sequence training, a ratio of 30 to 1 is used. All trainer workers are Nvidia K40 GPUs and use data parallelism. In total, 32 GPU workers were used for CE training; 16 GPU workers were used for sequence training. The parameter servers are multi-core CPU machines; For GPU experiments, 13 parameter servers were used.

Choice of unrolling scheme: The choice of full unrolling vs. truncated unrolling can be evaluated from different perspectives. In Table 1, these are compared in terms of WER performance, convergence time, maximum number of frames per mini-batch and *Average Padding Ratio* (APR), the average number of padded frames per step of training. The model used for comparison was a 5 layer LSTM with 600 cells per layer, denoted as 5xLSTM(600). The fully unrolled model was trained with batch 1, the setting for our best full unrolling system. For truncated unrolling, the model was trained with batch 256 and unrolling of 20. The fully unrolled model was trained using 500 CPU workers and 97 parameter servers. The default GPU setup described above was used for training the truncated model. With respect to WER, both models converge to similar WERs for both clean and noise test sets. The convergence time of the truncated model is significantly shorter than the fully unrolled model with batch size 1, due to the effectiveness of large batches in the GPU setup. Though full unrolling can be sped up using batching (with bucketing of sequences by length), the potential speed up depends on the distribution of sequence length over the training data. When bucketing, the number of frames in each step of SGD varies with the bucket sequence length, affecting step time proportionally. As step time changes, the number of asynchronous steps computed by other replicas changes, introducing more variable parameter staleness. Attempting to control for this, with e.g. variable batch size, introduces more training parameters and complexity.

With truncated unrolling, the maximum number of frames per step is bounded by batch size times unrolling size. The main drawback to batching sequence data is the need for padding each batch to the maximum sequence length in the batch. In our case the average number of padded frames per mini-batch (Table 1) is about one sixth of the mini-batch size. The number of wasted padded frames increases as the unrolling is increased. For full unrolling, padding depends on the bucket size used and on the distribution of sequence length over the training set. For the training set used for this paper, truncated unrolling was preferred for its simplicity and is used for the rest of experiments presented here.

	WER [%]		Conv. [days]	Max frames	APR [%]
	clean	noise			
full	12.07	19.52	28 d	45k	0
truncated	12.10	19.76	6 d	5120	15.6

Table 1: *WER, convergence time and mini-batch metrics (maximum # of frames & Average Padding Ratio (APR)) for full vs truncated unrolling.*

Choice of batching: Truncated unrolling is desirable, as it allows the use of very large batch sizes, which significantly improves the throughput of GPU parallel computation. However, the choice of batch size is not completely independent of the choice of unrolling size.

To examine this behavior, three pairs of batch and unrolling sizes were chosen such that the total number of frames in each mini-batch, $b \times u$, is constant. This allows us to avoid learning rate tuning for each model. Table 2 compares the WER of a 5xLSTM(768) model trained with three batch sizes, 512, 256 and 128 with corresponding unrolling sizes of 10, 20, and 40. The performances are presented for CE. The model trained with unrolling of 20 and 40 performs similarly, while the model trained with larger batch size of 512 and unrolling of 10 shows some performance degradation. This might be due to the fact that the unrolling size is not sufficient for learning

longer temporal dependencies, or to optimization issues such as the padding ratio introduced by batching for truncated unrolling. Table 2 presents the mean and standard deviation of the total number of padded frames over training steps. The mean and stddev for an unrolling of 10 and a batch size of 512 are the largest over the setups examined. This means that the number of frames used for learning varies significantly across steps, which might explain the performance degradation.

(b, u)	(512, 10)	(256, 20)	(128, 40)
clean	11.70	11.40	11.45
noise	19.21	18.51	18.56
avg padding ratio	20.6 ± 19.2	15.6 ± 10.0	17.6 ± 7.4

Table 2: *WERs for different batching and unrolling schemes.*

Choice of model parameters: Table 3 summarizes WERs after sequence training for three models with different number of LSTM cells per layer. These models were trained with batch of 256 and unrolling of 20, with learning rate tuned for all models. As discussed in Section 3, outer derivative staleness is an issue. To address that, for each model we zeroed out the 32 most stale examples in each mini-batch.

The model parameters can be chosen to make the best use of the available hardware. In our examples, increasing the layer size together with the batch size allows more efficient use of GPU hardware. In Table 3, two LSTM topologies, one with 600 units per layer and one with 768 units per layer, are trained with same resources; the convergence time for both models is similar. However, the larger model is significantly better in terms of WER. Furthermore, the wider model with 1024 cells per layer shows extra gains but of course this leads to extra training time.

	cross_entropy		sequence_training		Num params	Conv. [days]
	clean	noise	clean	noise		
5xLSTM(600)	12.10	19.76	10.43	15.56	22 M	6 + 3
5xLSTM(768)	11.40	18.51	10.10	15.01	28 M	6.5 + 3
5xLSTM(1024)	11.04	17.35	9.88	14.13	38 M	10 + 4

Table 3: *Comparison of WER after sequence training.*

5. Conclusion

This article discussed different approaches to efficient distributed training of RNNs in TensorFlow, with a focus on different strategies for data batching, recurrent model unrolling, and device placement of TensorFlow variables and operations. A training architecture was proposed that allows for flexible design of those strategies, enabling the use of a hybrid distributed CPU/GPU training scheme that extends previous work on asynchronous sequence training. Experimental results on a 20,000 hour Voice Search task show that this training architecture suffers no loss compared to a previous LSTM baseline of the same model size that uses conventional settings (full unrolling, no batching, synchronous sequence training), but in much reduced time, thanks to the GPU-based implementation enabled by the proposed design. Additional gains were obtained with significantly wider models, again in much improved training time.

6. Acknowledgements

The authors would like to thank all members of the Google speech team.

7. References

- [1] F. Seide, G. Li, X. Chen, and D. Yu, "Feature engineering in context-dependent deep neural networks for conversational speech transcription," in *Proc. ASRU*, 2011, pp. 24–28.
- [2] B. Kingsbury, "Lattice-based optimization of sequence classification criteria for neural-network acoustic modeling," in *Acoustics, Speech and Signal Processing, 2009. ICASSP 2009. IEEE International Conference on*. IEEE, 2009, pp. 3761–3764.
- [3] E. McDermott, G. Heigold, P. J. Moreno, A. W. Senior, and M. Bacchiani, "Asynchronous stochastic optimization for sequence training of deep neural networks: towards big data." in *Interspeech*, 2014, pp. 1224–1228.
- [4] H. Sak, O. Vinyals, G. Heigold, A. Senior, E. McDermott, R. Monga, and M. Mao, "Sequence discriminative distributed training of long short-term memory recurrent neural networks," *entropy*, vol. 15, no. 16, pp. 17–18, 2014.
- [5] O. Abdel-Hamid, A.-r. Mohamed, H. Jiang, L. Deng, G. Penn, and D. Yu, "Convolutional neural networks for speech recognition," *IEEE/ACM Transactions on audio, speech, and language processing*, vol. 22, no. 10, pp. 1533–1545, 2014.
- [6] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.
- [7] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin *et al.*, "Tensorflow: Large-scale machine learning on heterogeneous distributed systems," *arXiv preprint arXiv:1603.04467*, 2016.
- [8] P. J. Werbos, "Backpropagation through time: what it does and how to do it," *Proceedings of the IEEE*, vol. 78, no. 10, pp. 1550–1560, 1990.
- [9] R. Williams and J. Peng, "An efficient gradient-based algorithm for on-line training of recurrent network trajectories," *Neural Comput.*, vol. 2, no. 4, pp. 490–501, 1990.
- [10] K. Vesely, A. Ghoshal, L. Burget, and D. Povey, "Sequence-discriminative training of deep neural networks." in *Interspeech*, 2013, pp. 2345–2349.
- [11] D. Povey, V. Peddinti, D. Galvez *et al.*, "Purely sequence-trained neural networks for asr based on lattice-free mmi." *Interspeech*, 2016.
- [12] H. Su, G. Li, D. Yu, and F. Seide, "Error back propagation for sequence training of context-dependent deep networks for conversational speech transcription," in *Acoustics, Speech and Signal Processing (ICASSP), 2013 IEEE International Conference on*. IEEE, 2013, pp. 6664–6668.
- [13] G. Heigold, E. McDermott, V. Vanhoucke, A. Senior, and M. Bacchiani, "Asynchronous stochastic optimization for sequence training of deep neural networks," in *Acoustics, Speech and Signal Processing (ICASSP), 2014 IEEE International Conference on*. IEEE, 2014, pp. 5587–5591.
- [14] J. Dean, G. Corrado, R. Monga, K. Chen, M. Devin, M. Mao, A. Senior, P. Tucker, K. Yang, Q. V. Le *et al.*, "Large scale distributed deep networks," in *Advances in neural information processing systems*, 2012, pp. 1223–1231.
- [15] A. Graves, S. Fernández, F. Gomez, and J. Schmidhuber, "Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks," in *Proceedings of the 23rd international conference on Machine learning*. ACM, 2006, pp. 369–376.
- [16] G. Pundak and T. N. Sainath, "Lower frame rate neural network acoustic models," *Interspeech 2016*, pp. 22–26, 2016.
- [17] H. Sak, A. W. Senior, and F. Beaufays, "Long short-term memory recurrent neural network architectures for large scale acoustic modeling." in *Interspeech*, 2014, pp. 338–342.