

---

# Recurrent Memory Array Structures

## Technical Report

Kamil M Rocki\*

*IBM Research, San Jose, 95120, USA*

(July 11, 2016)

### Abstract

The following report introduces ideas augmenting standard Long Short Term Memory (LSTM) architecture with multiple memory cells per hidden unit in order to improve its generalization capabilities. It considers both deterministic and stochastic variants of memory operation. It is shown that the nondeterministic Array-LSTM approach improves state-of-the-art performance on character level text prediction achieving 1.402 BPC<sup>†</sup> on enwik8 dataset. Furthermore, this report establishes baseline neural-based results of 1.12 BPC and 1.19 BPC for enwik9 and enwik10 datasets respectively.

## 1 Background

It has been argued that the ability to compress arbitrary redundant patterns into short, compact representations may require an understanding that is equivalent to general artificial intelligence (MacKay, 2003; Hutter, 2005). One example of such a process is demonstrated by learning to predict text a letter at a time. A strong connection between compression and prediction was shown (Shannon, 1951). Therefore, this report considers experiments on natural wikipedia text corpora, however the algorithms can be in principle applied to any sequences of patterns.

## 2 Simple RNN

A standard recurrent neural network (so called simple RNN, cite) is composed of a matrix of connections between its inputs and hidden states  $W$ , and a matrix  $U$ , connecting hidden states in consecutive time steps. In such a simple RNN architecture the entire history of observations is aggregated in hidden states of neurons (fig 2.1). States ( $h^t$ ) are determined by previous states ( $h^{t-1}$ ) and feedforward immediate inputs ( $x^t$ ). This architecture is deterministic, for

---

\*kmrocki@us.ibm.com

<sup>†</sup>bits per character

2 identical  $h^{t-1}$ ,  $x^t$  there will be 2 identical outputs  $h^t$ . A single time step update can be expressed with an equation (1) or equivalently using the following graphical representation as shown in fig. 2.1.

$$h^t = \tanh(Wx^t + Uh^{t-1} + b) \quad (2.1)$$

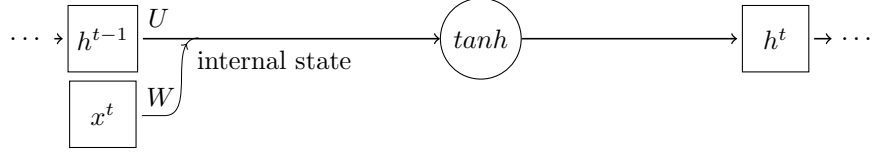


Figure 2.1: Simple RNN unit (omitted implementation specifics);  $h^t$  - internal (hidden) state at time step  $t$ ;  $x$  are inputs,  $y$  are optional outputs to be emitted; All connections are learnable parameters. No explicit asynchronous memory, implicit history aggregation only through hidden states  $h$ . Omitted bias terms for brevity.

### 3 Memory structures

A simple RNN architecture does not handle long-range interactions and multiple simultaneous context well (Bengio et al., 1994). However, it can be modified in order to make learning long-term dependencies easier. One solution to the problem is to change the way of interactions between hidden units, i.e. add multiplicative connections (Sutskever et al., 2011). Another is by adding explicit memory cells. Networks involving register-like functionality allowing hidden states to store and load its contents in an asynchronous way have been very successful recently. Examples of such architectures include Long-Short Term Memory (Hochreiter and Schmidhuber, 1997) and Gated Recurrent Unit (GRU) (Chung et al., 2014) networks.

#### 3.1 LSTM

This section describes a standard LSTM network used in our experiments and serving as a foundation for array extensions. Equations 3.1-3.6 define a single LSTM time step update.

$$f^t = \sigma(W_f x^t + U_f h^{t-1} + b_f) \quad (3.1)$$

$$i^t = \sigma(W_i x^t + U_i h^{t-1} + b_i) \quad (3.2)$$

$$o^t = \sigma(W_o x^t + U_o h^{t-1} + b_o) \quad (3.3)$$

$$\tilde{c}^t = \tanh(W_c x^t + U_c h^{t-1} + b_c) \quad (3.4)$$

$$c^t = f_t \odot c^{t-1} + i_t \odot \tilde{c}^t \quad (3.5)$$

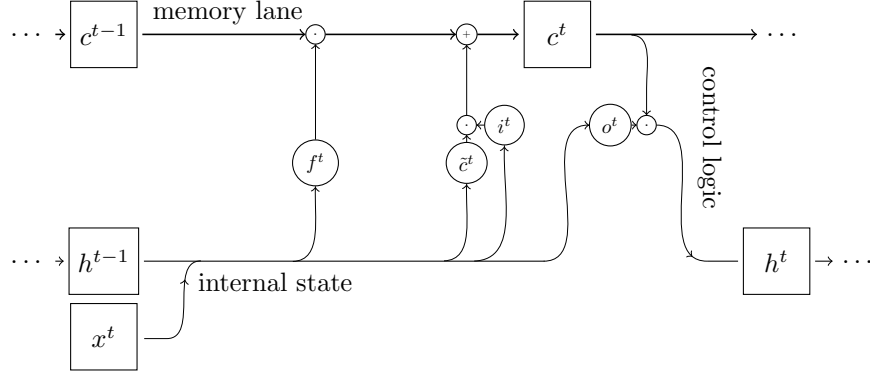


Figure 3.1: LSTM, biases and nonlinearities omitted for brevity; The memory content  $c^t$  is set according to the gates' activations which are in turn driven by the bottom-up input  $x^t$  and previous internal state  $h^{t-1}$ .

$$h^t = o_t \odot \tanh(c^t) \quad (3.6)$$

### 3.2 State-sharing Memory: Array-LSTM

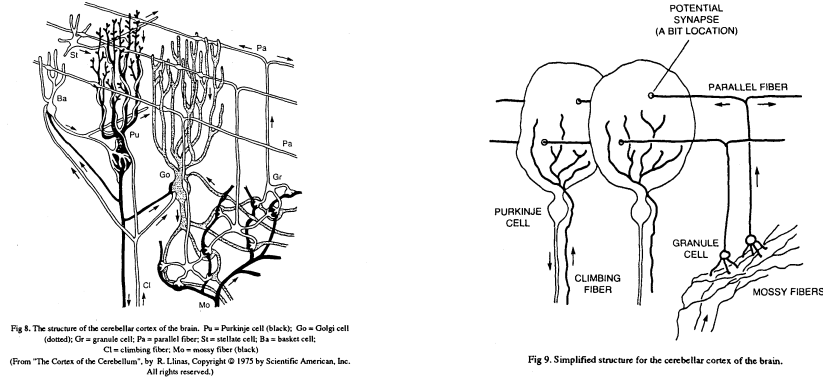


Figure 3.2: The structure of the cerebellar cortex (Kanerva, 1988)

Cerebellar cortex structure as described in (Kanerva, 1988) among others: Fig. 3.2 shows array-like structure in what seems to be Random-Access Memory in cerebellum. It has served as an inspiration for the Array-LSTM architecture. The main idea is that instead of building hierarchies of layers (as in stacked LSTM (Graves, 2013), Gated-Feedback RNN (Chung et al., 2015)) keep a single layer, but build more complex memory structures inside a RNN unit

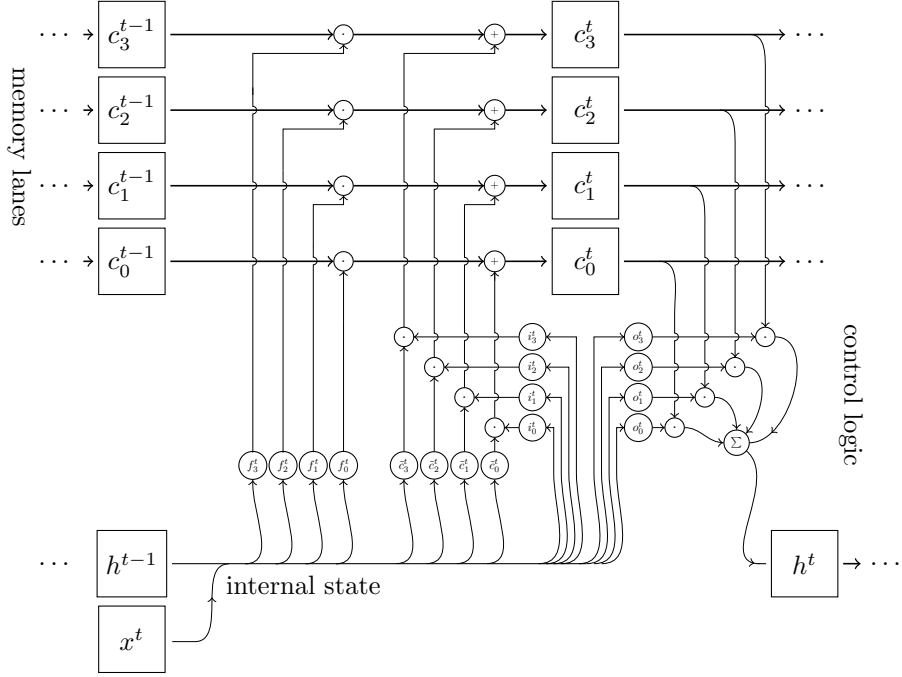


Figure 3.3: Array-LSTM with 4 memory cells per hidden control unit; modulated vs modulating connections, omitted nonlinearities and biases for brevity; It becomes a standard LSTM when only 1 memory cell per hidden unit is present

(a similar line of thinking has been explored by constructing a more complex transition function inside a layer (Pascanu et al., 2014)). We want to create a *bottleneck* by sharing internal states, forcing the learning procedure to *pool* similar or interchangeable content using memory cells belonging to one hidden unit (analogy would be that a word related to car would activate a particular hidden state and each memory cell could be interchangeably used if it represents a particular car type, therefore externally the choice of a particular memory cell would not be relevant. Similar concepts exist already in convnets, i.e. spatial pooling, the hidden state should work as a complex cell pooling multiple possible substates. Figure 3.3 and equations 3.1-3.6 describe a single Array-LSTM architecture. Note the similarity between figures 3.2 and 3.3. Furthermore, We consider modifications of the presented simple Array-LSTM architecture which are meant to improve capacity, memory efficiency, learning time or generalization. This report shows two types of changes, one pertains to deterministic family of architectures and the other one to stochastic operations (working in a dropout-like fashion).

$$f_k^t = \sigma(W_{fk}x^t + U_{fk}h^{t-1} + b_{fk}) \quad (3.7)$$

$$i_k^t = \sigma(W_{ik}x^t + U_{ik}h^{t-1} + b_{ik}) \quad (3.8)$$

$$o_k^t = \sigma(W_{ok}x^t + U_{ok}h^{t-1} + b_{ok}) \quad (3.9)$$

$$\tilde{c}_k^t = \tanh(W_{ck}x^t + U_{ck}h^{t-1} + b_{ck}) \quad (3.10)$$

$$c_k^t = f_k^t \odot c_k^{t-1} + i_k^t \odot \tilde{c}_k^t \quad (3.11)$$

$$h^t = \sum_k o_k^t \odot \tanh(c_k^t) \quad (3.12)$$

## 4 Deterministic Array-LSTM extensions

### 4.1 Lane selection: Soft attention

We allow hidden state to *choose* a memory cell that it *wants* to *read* from and *write* to. This choice will be invisible to other hidden states, so that multiple memory cell choices can be mapped to the same internal state, giving some space to learning invariant temporal patterns. Compared to the original LSTM and simple Array-LSTM, one additional gate activation per memory lane is computed. We call it *selection* gate. See Fig. 4.1 and equations (4.1-4.8) for details. The intuition behind this approach is that the  $s$  gates control how much a particular memory lane should be used during current time step. If  $s^t$  is 1 then all other activations remain unchanged, we fully *process* this memory cell. If  $s^t$  is 0, then the contents are carried over from the previous time step and the memory cell is not affected during that time step (no read/no write). The idea is that such a mechanism should allow less leaky memory cells, effectively not relying entirely on forget gate in order to preserve its contents. Each memory lane has a selection gate  $s^t$  associated with it, controlling information flow through or bypassing current time step (carrying over memory content from previous time step). The idea is that  $s$  gates should be responsible for controlling the timescale (as proposed in the Zoneout paper (Krueger et al., 2016)).

Attention signals  $k$ .

$$a_k^t = \sigma(W_{ak}x^t + U_{ak}h^{t-1} + b_{ak}) \quad (4.1)$$

Softmax normalization.

$$s_k^t = \frac{e^{a_k^t}}{\sum_k e^{a_k^t}} \quad (4.2)$$

$$f_k^t = s_k^t \odot \sigma(W_{fk}x^t + U_{fk}h^{t-1} + b_{fk}) \quad (4.3)$$

$$i_k^t = s_k^t \odot \sigma(W_{ik}x^t + U_{ik}h^{t-1} + b_{ik}) \quad (4.4)$$

$$o_k^t = s_k^t \odot \sigma(W_{ok}x^t + U_{ok}h^{t-1} + b_{ok}) \quad (4.5)$$

$\tilde{c}_k^t$  is not affected by  $s_k^t$ .

$$\tilde{c}_k^t = \tanh(W_{ck}x^t + U_{ck}h^{t-1} + b_{ck}) \quad (4.6)$$

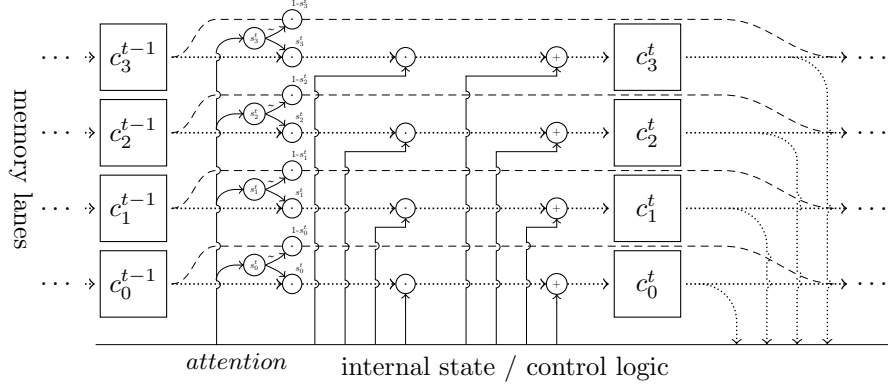


Figure 4.1: Lane selection through soft attention, solid lines are control logic signals (from/to gates), dotted lines are memory cell lanes used upon selection, dashed lines represent the carry lanes used otherwise; dotted and dashed lanes are mutually exclusive

Note the inverted  $f$  gate effect for clearer notation (cell  $k$  is reset for  $f_k = 1$ ). If  $s_k^t$  is 0, then  $f_k^t$  is 0 and effectively  $k$  memory cell's contents are entirely transferred for the previous time step.

$$c_k^t = (1 - f_k^t) \odot c_k^{t-1} + i_k^t \odot \tilde{c}_k^t \quad (4.7)$$

$$h^t = \sum_k o_k^t \odot \tanh(c_k^t) \quad (4.8)$$

**Max pooling version** In this version the algorithm used a hard, deterministic selection, by choosing the lane with highest  $s^t$  value. It ignores other lanes and backpropagates errors only through that lane, It is analogous to max pooling in CNNs.

## 5 Non-deterministic Array-LSTM extensions

### 5.1 Stochastic Output Pooling

This is the simplest of the considered stochastic architectures (Fig. 5.1). It works by treating initial  $o$  gate activations as inputs to a softmax output distribution and sampling from this distribution. Therefore it enforces normalization of output response and sparse binary outputs ( $\leq 1/2$ ). During backpropagation, the algorithm uses only the selected gate (as in the max attention algorithm in 4.1). All other steps are exactly the same as the standard Array-LSTM approach in 3.2.

Probability that memory cell  $i$  will be used during  $h^t$  update (other cells' outputs are not used):

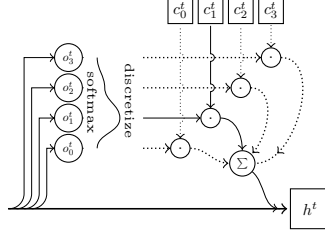


Figure 5.1: Stochastic Output Pooling: Sample open output gate index  $i$  from this distribution calculated using *softmax* normalization. Dotted lines represent inactive connections (set to zero due to the sampling procedure).

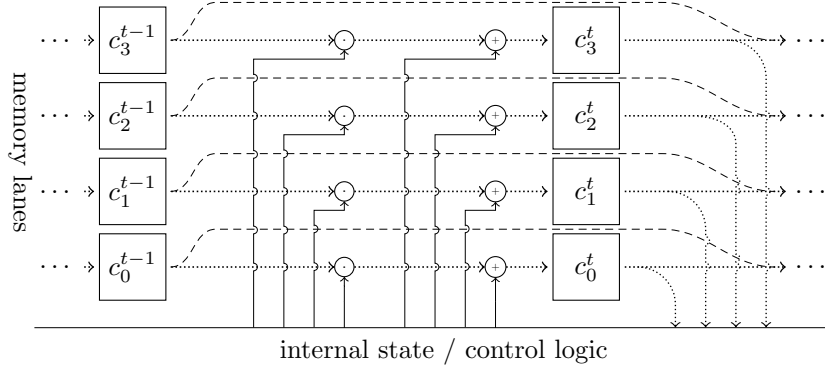


Figure 5.2: Stochastic Memory Array: solid lines are control logic signals (from/to gates), dotted lines are memory cell lanes used upon selection, dashed lines represent the carry lanes used otherwise; dotted and dashed lanes are mutually exclusive

$$p(i = k) = \frac{e^{o_k^t}}{\sum_k e^{o_k^t}} \quad (5.1)$$

The hidden state  $h^t$  is computed using the active output cell index.

$$h^t = o_i^t \odot \tanh(c_i^t) \quad (5.2)$$

## 5.2 Stochastic Memory Array

This is a more complex architecture, an *active* cell is being chosen for the entire cycle of computation. Instead of selecting only the output memory cell, the algorithm chooses one memory lane to be used during the entire forward and backward loop. If a lane is not selected, carry over cell contents. This idea is very similar to the algorithm described as Zoneout (Krueger et al., 2016). The main difference between two implementations is the fact that here, the hidden state is not affected directly by the noise injecting procedure. Instead, by choosing

exactly one memory lane, it ensures that the output of the memory is non-zero and the content of  $h^t$  is going to be non-zero (I found that the main instability of all stochastic approaches, including Recurrent Dropout (Zaremba et al., 2014) is caused by the fact that all or almost all hidden states’ activations can be zeros. Figure 5.2 depicts the changes to the original Array-LSTM architecture from 3.2. Each memory lane has an additional bypass connection which is used when memory cell is not selected for computation (as in the soft attention algorithm in 4.1). We consider 2 variants of this approach: (a) 1 out of  $K$  cells is active (b)  $1/2$  cells are active, i.e. 4 out of 8. In the first implementation, the algorithm randomly chooses index of the memory cell to be used. In the latter one, the same procedure applies, only to groups of 2 cells (odd or even).

### 5.2.1 Stochastic hard attention

Here this report describes different implementations of combining the soft attention mechanism with stochastic lane selection. We would like to select memory cell in a semi-random way, according to some selection distribution as described in 4.1.

**Semi-Hard** One memory lane is selected as being active (as in Stochastic Memory Array from 5.2). However, unlike the previous fully random selection mechanism, in this approach the lane is selected according to softmax distribution where  $s$  are inputs controlling this distribution. The backward step is the same as in the soft attention version, ignoring the fact that in the forward pass a sample was used and computes derivatives using differentiable softmax outputs.

**Hard** Same as v1, but the backward pass uses only the selected lane (as in 5.2)

## 5.3 Implementation considerations

Array-LSTM brings many advantages from the computation cost perspective.

1. **More memory cells** Given the same amount of memory for allocation, Array-LSTM has effectively use more memory cells due to smaller matrix between hidden units and gates.
2. **More parallelism** More cells give raise to more independent elementwise computation which is very cheap on SIMD hardware like GPUs.
3. **More data locality** Cells belonging to the same hidden unit will be somehow correlated which might improve cache hit ratio. In addition to that, sparse hidden to hidden connectivity might be possible for larger number of cells per unit.
4. **Approximate** Stochastic Array Memory is naturally resilient to noisy input.



## 6 Related Works

There are many other approaches to improving baseline LSTM architecture including deterministic approaches such as Depth gated LSTM (Yao et al., 2015), Grid LSTM (Kalchbrenner et al., 2015) or Adaptive Computation Time (Graves, 2016). Examples of stochastic variants include recurrent dropout (Zaremba et al., 2014; Semeniuta et al., 2016) Batch normalization (Cooijmans et al., 2016) and Zoneout (Krueger et al., 2016).

## 7 Experiments

### 7.1 Methodology

The learning algorithm used was backpropagation through time (BPTT), it proceeded by selecting random sequences of length 10000 randomly from a given corpus. The learning algorithm used was Adagrad<sup>‡</sup> with a learning rate of 0.001. Weights were initialize using so-called Xavier initialization Glorot and Bengio (2010). Sequence length for BPTT was 75 and batch size 128, states were carried over for the entire sequence of 10000 simulating full BPTT. Forget bias was set initially to 1. These values were not tuned, so it is possible that better results can be easily obtained with different settings. The algorithm was written in C++ and CUDA 8 and ran on GTX Titan GPU for up to 20 days. Link to the code is at the end.

### 7.2 Data

#### 7.2.1 enwik8

It constitutes first  $10^8$  bytes of English Wikipedia dump (with all extra symbols present in xml), also known as Hutter Prize challenge dataset\*.

#### 7.2.2 enwik9

This dataset is used in Large Text Compression Benchmark - the first  $10^9$  bytes of the English Wikipedia dump\*.

#### 7.2.3 enwik10

This dataset is an extension of enwik9 dataset and was created by taking first  $10^{10}$  bytes of english wikipedia dump from June 1, 2016<sup>†</sup>, entire dump<sup>‡</sup> (57G).

---

<sup>‡</sup>with a modification taking into consideration only recent window of gradient updates

<sup>†</sup><https://www.dropbox.com/s/kzb5a0bih99ltui/enwik10.txt>

<sup>‡</sup><https://www.dropbox.com/s/1wigbsjpwxtwh2k/enwiki-20160601-pages-articles.xml>

### 7.2.4 Test data

First 90% of each corpus was used for training, the next 5% for validation and the last 5% for reporting test accuracy.

## 7.3 Results

It is not known what the limit on the compression ratio is beforehand. It is not known if a pattern is compressible, before actually being able to compress it or proving that it cannot be compressed. It was shown that for humans the perceived prediction quality of natural english text depends largely on the amount of context given and the ability to incorporate previous knowledge into making predictions and ranges from 0.6 to 1.3 depending on the case (Shannon, 1951). Another estimate was obtained by bzip and cmix9 algorithms used for text compression which prove the upper limit on the compressibility of these datasets (Table 7.1).

Table 7.1: Bits per character (test data). The results were collected using networks of approximately equal size (66M parameters)

	enwik8	enwik9	enwik10
bzip2	2.32	2.03	-
cmix9	1.25	0.99	-
mRNN*(Sutskever et al., 2011)	1.60	1.55	-
GF-RNN (Chung et al., 2015)	1.58	-	-
Grid LSTM (Kalchbrenner et al., 2015)	1.47	-	-
LSTM (my impl)	1.45	1.13	1.24
Stacked 2-LSTM (my impl)	1.468	<b>1.12</b>	<b>1.19</b>
Array-4 LSTM	1.445	-	-
SoftAtt-Array-2 LSTM	1.43	-	-
HardAtt-Array-2 LSTM	1.43	-	-
Pooling-Array-2	1.422	-	-
Stochastic-Array-2	<b>1.407</b>	-	-

## 7.4 Observations

### 7.4.1 Worked

1. **Array-LSTM performs as well as regular LSTM**; it seems that sharing control states does not affect negatively training as long as there is the same number of memory cells), despite having limited number of hidden nodes (but same number of parameters). However, like LSTM it exhibits strong overfitting on enwik8 datasets (after about 24h of training). In

\*preprocessed data, with 86-symbol alphabet, 2G dataset instead of enwik9

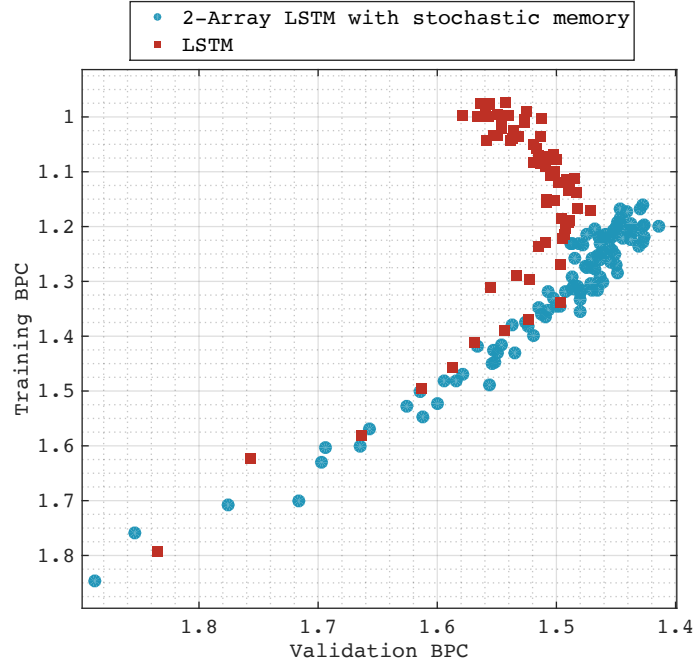


Figure 7.1: Training progress on enwik8 corpus (validation vs training), bits/character) - overfitting

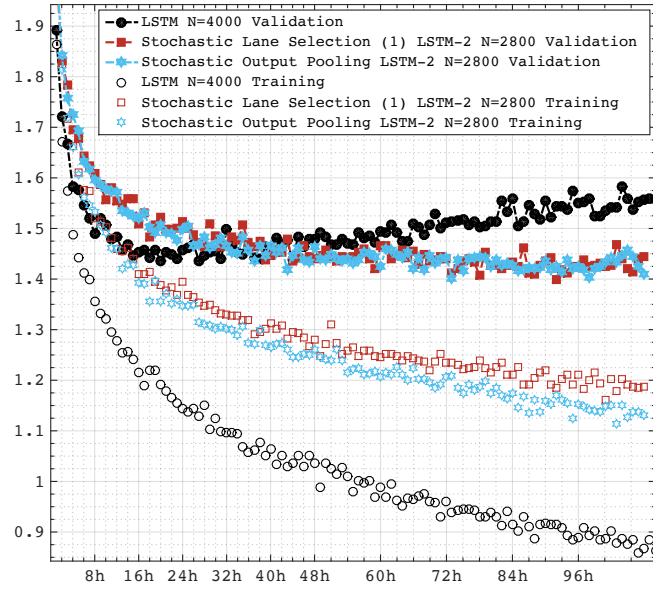


Figure 7.2: Training progress on enwik8 corpus, bits/character)

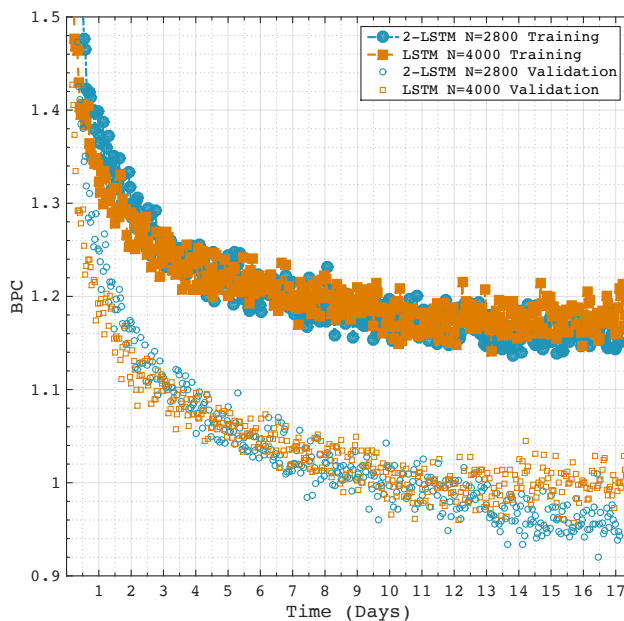


Figure 7.3: Training progress (x-axis is time) on enwik9 corpus. Problem seems to be capacity limited for this dataset, after 2 weeks of learning both the training loss and validation loss are still decreasing. Learning proceeds in a very similar way on enwik10 dataset (both curves are slightly higher, approximately 0.05 bits). Array-LSTM architectures did not provide any boost in learning convergence, comparable performance.

some cases convergence speed is better when using state sharing, however it does not seem to generalize better than standard LSTM contrary to expectations.

2. **Stochastic Memory Array** - need to make sure that at least 1 cell is on, otherwise it blows up.

#### 7.4.2 Inconclusive

1. **Hard attention** - couldn't make it work well, converged slowly, hard to debug. Max variant works slightly better, but not results are not convincing.
2. **Attention modulated LSTM** comparable to Array-LSTM, not much of an improvement.
3. **Stochastic Output Pooling** works if  $n = 2$ , for  $n > 2$  there is too much randomness (for  $n = 4$ , effectively the dropout rate is 75 percent, training is slow), Stochastic Memory Array works better with 1/2 cells active.

### 7.4.3 Did not work

1. Unconstrained Dropout (allowing all cells or states to be zeros - causes unstable behavior).
2. Multiplicative interactions between memory cells within one column/array or stack (example - 7.4)

### 7.4.4 Other observations

1. 2-LSTM or multilayer LSTM do not really improve things, no gain in generalization, slightly better capacity observed on enwik10; a better way of injecting compositionality is required. In short, one large layer is equally capable, which may be explained by inherent deep structure of recurrent nets.
2. No visible overfitting with standard LSTM on enwik9 and enwik10 after 2 weeks of training - 2 problems.
3. On enwik8, the main problem is overfitting, even small nets (less than 500 units) can basically memorize small corpora, observed that during generation of sequences it can *recite* fragments of text, need to improve generalization - incorporate temporal invariance into architecture.
4. results from compression challenge (cmix9) show that there is some redundancy left, so it should be possible to get better results
5. Batch size affects convergence speed, sequence length (possibly because we carry last state over emulating full BPTT) and epoch length have low impact
6. Performance analysis supports persistent RNN paper (Diamos et al., 2016); Most time is spent moving data, pointwise OPs are very cheap
7. Qualitative evaluation Generated sequences give some hints, these are only hypotheses which need to be confirmed. cells may be sensitive to: ...

(a) position, dates, city names, languages, topics, numbers, quantities

An analysis as in (Karpathy et al., 2015) may be required. See the Appendix for generated samples, we also provide nearly 5000 samples of length 5000 each as a by-product of the experiments for analysis.

8. Even BPC of 1.0 seems to be insufficient in order to generate human-level text, but getting close (heavily structured text is OK)

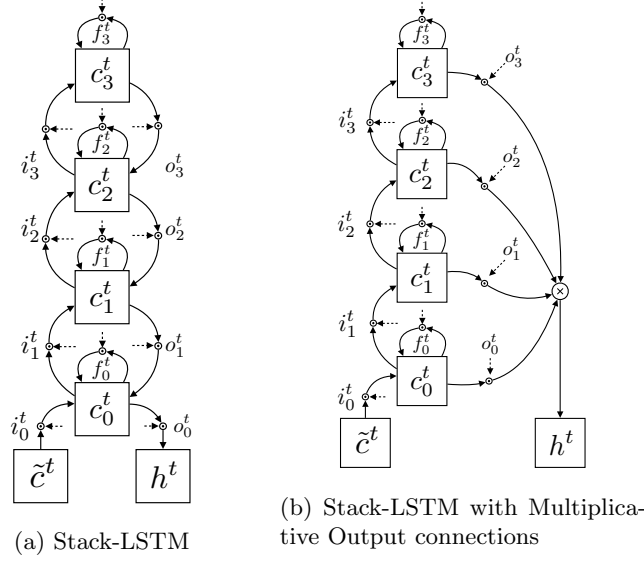


Figure 7.4: The main idea is that a stack-like structure enforces ordering in the data-flow between low and high frequency patterns; here we assume that bottom memory cells deal with high frequency and top cell with least frequent changes

#### 7.4.5 Further work

1. How to incorporate compositionality in a better way to reflect recursive structure - feedback: how to implement it efficiently?
2. How to do learning in a more efficient way which reflects the asynchronous nature of event. BPTT works, but the time horizon is fixed, so how to make parts of network *see* different past sequences, incorporate only relevant symbols, i.e. discarding exact timing and preserving ordering.
  - (a) **Stack vs Array, Tape** The data structure should enforce representation compositionality, we are experimenting with the following structures injecting feedback signal, but no significant improvement. We got it to work in practice, and there is some experimental evidence that it may increase capacity, but it took too much time on enwik9 and enwik10 datasets using current implementation to draw any definitive conclusions.
  - (b) The relationship between gating-like mechanism in RNNs and thalamocortical circuits in the brain - resonant columns; there seem to be some clues that cortex incorporates some form of gating mechanism though thalamus, but we don't know exactly how it works and how feedback is implemented - some research is needed - see Appendix B

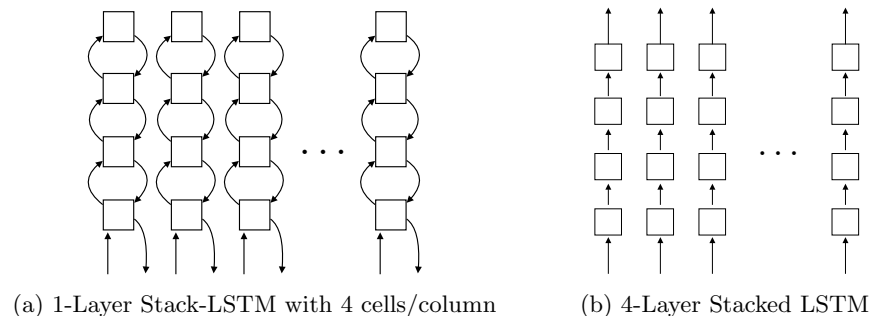


Figure 7.5: Stack-LSTM with the memory hierarchy built-in into the layer (1 layer with columnar cell stack) feedback architecture vs popular stacked-LSTM which is a feedward multilayer architecture without feedback from higher layers

3. What is the hard limit on BPC - will improving BPC lead to better samples? See Appendix A for low entropy samples from enwik9 and enwik10 datatests

## 8 Conclusions

This report presented Array-LSTM approaches. Stochastic memory operation seems to be necessary in order to mitigate overfitting effects. We tried many deterministic variants, but they all overfit as easily as baseline LSTM. Based on results from enwik9 and enwik10 datasets, the best regularization strategy is simply more data. Furthermore the Stochastic Memory Array approach extended state-of-the-art on enwik8 and set reference results for neural based algorithms on enwik9 and enwik10 datasets.

## Acknowledgements

This work has been supported in part by the Defense Advanced Research Projects Agency (DARPA). The source code is available<sup>§</sup> (contains implementations of the described ideas and was used to obtain the reported results).

## References

- Y. Bengio, P. Simard, and P. Frasconi. Learning long-term dependencies with gradient descent is difficult. *IEEE Transactions on Neural Networks*, 5(2):157–166, 1994. URL <http://www.iro.umontreal.ca/~lisa/pointeurs/ieeetrnn94.pdf>.

<sup>§</sup><https://github.com/krocki/ArrayLSTM>

- J. Chung, Ç. Gülçehre, K. Cho, and Y. Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. Technical Report Arxiv report 1412.3555, Université de Montréal, 2014. Presented at the Deep Learning workshop at NIPS2014.
- J. Chung, Ç. Gülçehre, K. Cho, and Y. Bengio. Gated feedback recurrent neural networks. *CoRR*, abs/1502.02367, 2015. URL <http://arxiv.org/abs/1502.02367>.
- T. Cooijmans, N. Ballas, C. Laurent, and A. Courville. Recurrent batch normalization. *CoRR*, abs/1603.09025, 2016. URL <http://arxiv.org/abs/1603.09025>.
- G. Diamos, S. Sengupta, B. Catanzaro, M. Chrzanowski, A. Coates, E. Elsen, J. Engel, A. Hannun, and S. Satheesh. Persistent rnns: Stashing recurrent weights on-chip. In *Proceedings of The 33rd International Conference on Machine Learning*, page 20242033, 2016.
- X. Glorot and Y. Bengio. Understanding the difficulty of training deep feed-forward neural networks. In *In Proceedings of the International Conference on Artificial Intelligence and Statistics (AISTATS10)*. Society for Artificial Intelligence and Statistics, 2010.
- A. Graves. Generating sequences with recurrent neural networks. *CoRR*, abs/1308.0850, 2013. URL <http://dblp.uni-trier.de/db/journals/corr/corr1308.html#Graves13>.
- A. Graves. Adaptive computation time for recurrent neural networks. *CoRR*, abs/1603.08983, 2016. URL <http://arxiv.org/abs/1603.08983>.
- S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural Comput.*, 9(8):1735–1780, Nov. 1997. ISSN 0899-7667. doi: 10.1162/neco.1997.9.8.1735. URL <http://dx.doi.org/10.1162/neco.1997.9.8.1735>.
- M. Hutter. *Universal Artificial Intelligence: Sequential Decisions based on Algorithmic Probability*. Springer, Berlin, 2005. ISBN 3-540-22139-5. doi: 10.1007/b138233. URL <http://www.hutter1.net/ai/uaibook.htm>. 300 pages, <http://www.hutter1.net/ai/uaibook.htm>.
- N. Kalchbrenner, I. Danihelka, and A. Graves. Grid long short-term memory. *CoRR*, abs/1507.01526, 2015. URL <http://arxiv.org/abs/1507.01526>.
- P. Kanerva. *Sparse Distributed Memory*. MIT Press, Cambridge, MA, USA, 1988. ISBN 0262111322.
- A. Karpathy, J. Johnson, and F. Li. Visualizing and understanding recurrent networks. *CoRR*, abs/1506.02078, 2015. URL <http://arxiv.org/abs/1506.02078>.



- D. Krueger, T. Maharaj, J. Kramár, M. Pezeshki, N. Ballas, N. R. Ke, A. Goyal, Y. Bengio, H. Larochelle, A. C. Courville, and C. Pal. Zoneout: Regularizing rnns by randomly preserving hidden activations. *CoRR*, abs/1606.01305, 2016. URL <http://arxiv.org/abs/1606.01305>.
- D. J. C. MacKay. *Information Theory, Inference, and Learning Algorithms*. Cambridge University Press, 2003. URL <http://www.cambridge.org/0521642981>. Available from <http://www.inference.phy.cam.ac.uk/mackay/itila/>.
- R. Pascanu, Ç. Gülçehre, K. Cho, and Y. Bengio. How to construct deep recurrent neural networks. In *International Conference on Learning Representations 2014 (Conference Track)*, Apr. 2014. URL <http://arxiv.org/abs/1312.6026>.
- S. Semeniuta, A. Severyn, and E. Barth. Recurrent dropout without memory loss. *CoRR*, abs/1603.05118, 2016. URL <http://arxiv.org/abs/1603.05118>.
- C. E. Shannon. Prediction and entropy of printed english. *Bell System Technical Journal*, 30:50–64, Jan. 1951. URL <http://languagelog.ldc.upenn.edu/myl/Shannon1950.pdf>.
- I. Sutskever, J. Martens, and G. Hinton. Generating text with recurrent neural networks. In L. Getoor and T. Scheffer, editors, *Proceedings of the 28th International Conference on Machine Learning (ICML-11)*, ICML ’11, pages 1017–1024, New York, NY, USA, June 2011. ACM. ISBN 978-1-4503-0619-5.
- K. Yao, T. Cohn, K. Vylomova, K. Duh, and C. Dyer. Depth-gated LSTM. *CoRR*, abs/1508.03790, 2015. URL <http://arxiv.org/abs/1508.03790>.
- W. Zaremba, I. Sutskever, and O. Vinyals. Recurrent neural network regularization. *CoRR*, abs/1409.2329, 2014. URL <http://arxiv.org/abs/1409.2329>.

# Appendices

## A Generated samples

A gallery of samples generated by networks, standard procedure was used: initialize all  $h$  and  $c$  memory units with 0, generate one symbol at a time according to the output probability distribution conditioned on  $h^t$ , treat that input as the next input, iterate a few thousand times, see the link to the source code for more samples.

## Listing 1: enwik8.txt Stochastic Array-2 LSTM N-4000 Validation BPC-1.395

```

1 well s lie to capture political activities created in and 2000 would be made up- if the Scripturals reached between
  ngation and decent from three to five, and when the painting followed, about 68,000 members of societies had
  not yet deployed its monumental calliets as I continued to come back towards them, illiced to familiarity.
2 Most hinders therefore, concerning the behavioral science of the latter was done in the 1980s.
  *Third participation did not imply that close relationships were annoyal because of the distinctive emotions, both
  intrinsic and realism.
3
4 (see [[Dark mass disaster]]). In particular, Darwinisms that arise out by meipsycological treatment because many is
  reeated in he response limit of inn. In some further cases, particularly apogetative sounds do not prefix
  conditions, such as results of non-conditional analogy in cles of gaining chemical in different ways will
  almost completely lead to death, or require some extent by emphasizing that in general times or death, types
  or etching.
5
6 This is not associated with [[Negative testiny]] or [[Gaut in the ulture]].
7
8 ==Usage==
9 Within the [[Chinese experiment of sexual abuse]], '''evolution''' is also called &quot;[[Generalissimple kines
  object]]&quot;;. It however is at least essentially understood and alleged by the philosophers of [[societal
  planning]] as a result of analogous working-nature expressions in overrun variables based on the expectations
  of which aid is inevitable.
10
11 ==History of the name==
12 Criticism of fundamental differentialism by modern day liberties now refers to individual organizations such as [[
  Black particular]]s. In [[Germany]], which examples as either a consonant or non-GFT, a correspondence against
  the people of the [[West End of Britain]]. Prior to the case was beginning in 1983 in relation to the
  subdivision of the cultural landforms of many years of academic life (at least once, according to 16th
  century chronology law neither language 'ar' Esperanto had less so in [[Old Norse literature]] than on
  other balancing words.
13
14 [[Old Scottish Gaelic language|Old English variations]] are be&quot;liked&quot; in [[Taggart Stanley]]'s &quot;New
  Yorkers.&quot;
15
16 *Romantic Cross-and-Sip Modern-Literature generators would retain him in [[Christmas]] notice as the composer of a
  child's chick in the Charles W. Burgh6 story in books or history; two L5 centuries.
17 The Dish Studios describe a giant as Ann She given in [[Weether University]]. This lings the linement in the semi-
  final dayline: in his ack of main characters, Jeffrey was indeed his lyrical heart-town. However, the general
  rule would find confusing well, the film, the odds being exceeded. Another popular transgressionist [[Evan
  Tobl]], openly Wet artist, inspired Moore for the time. He also made many successes, including the sequel,
  '''[[Rolling Stone]]''', and a number of pseudoblats. '''[[John Gapling Madden]]''' also starred [[Elizabeth
  Jones]] as '''Back and East'''.
18
19 Born in [[Scagblaw]], Buckingham Carolini were among the only serious visual concerns over one such theeto ad
  continued to exchange dance bolbloonfriends their involvement in devesong her being. A young Bruce Goldwater
  wife promoted William Randolly to I.L. in this year. In the end was Bell running back in Parliament with a
  bug in the [[odorn slug administration]].
20
21 During which three former top-ever, filed and bigenting factories in the upstate New York, factories family were
  typically severely defered; she sometimes refused.
22
23 Coupland also appeared in '''The Party''' and '''[[RuneQuest]]''' ('''[[300 Begs For Studies]]''') in [[1968]]. It took
  the characters found at Berkeley in theater in [[Rochester, New York|Rochester]] on [[June 3]], [[1923]] at
  Anny Bryant London as an staple, after the [[1904 in film|1945]] interview that the [[film|movie]] played
  the play.
24
25 After this peregression, Beatle came from a group of [[Rock and Roll Hall of Fame]]rs. Much of his life is
  considered cross-time and less extreme and successful in popular drink. Her novel '''Places in Amoryo music''',
  penned for four plays selling the opening band [[Kazuma Ahmet]] ('''[[Amalghend]]''') rolled out the [[United
  Kingdom|British]] movie of the 20th century and is most name out of the role.
26
27 ==External links==
28 {{wikibooks}}
29 * {{indb name|id=0416003|name=Charles Lindberg}}
30 * [http://www.farleyfnet.com/ Far Easter.com] (unofficial) Quote
31 * [http://www.fantagnalsilot.co.uk Frankie Mobina, Andy Warhol]
32 * [http://www.friendsfootballsholdon.com/shaganshish.html Fall for Anarchy Fancifully Charted from The Golden Dawn
  of the 1913 General Theater]
33 * [http://www.cwvtc.us/archives/hackman/050802.pdf '''And there Were Today's ' Portrait of Andrew Jackson'''] from a
  musical account of written work on Camus ]
34 * {{cite book | first=Francis | last= Catted | author=Astounding Spooks | other= | others= | publisher=Freedman |
  year=1986 | id=ISBN 1-854889-04-7}}
35
36 '''G&amp;SCollingd:
37 * St. Claire Clarke, &quot;Casualties of Eden Young Check in Closets&quot; in '''W

```

## Listing 2: enwik9.txt LSTM N-3200 BPC-1.116

```

1 spread allith the age of 18 year. The [[population density]] is 218.1/km&sup2; (560.5/mi&sup2;). There are
  3,298 housing units at an average density of 153.5/km&sup2; (390.7/mi&sup2;). The racial makeup of
  the town is 98.62% [[White (U.S. Census)|White]], 0.07% [[African American (U.S. Census)|African American]],
  0.13% [[Native American (U.S. Census)|Native American]], 0.11% [[Asian (U.S. Census)|Asian]], 0.02% [[Pacific
  Islander (U.S. Census)|Pacific Islander]], 0.09% from [[Race (U.S. Census)|other races]], and 0.48% from two
  or more races. 0.16% of the population are [[Hispanic (U.S. Census)|Hispanic]] or [[Latino (U.S. Census)|
  Latino]] of any race.
2
3 There are 2,000 households out of which 45.5% have children under the age of 18 living with them, 85.6% are [[
  Marriage|married couples]] living together, 3.2% have a female householder with no husband present, and 14.9%

```

```

4         are non-families. 12.9% of all households are made up of individuals and 3.8% have someone living alone who
5         is 65 years of age or older. The average household size is 3.02 and the average family size is 3.38.
6
7     In the town the population is spread out with 26.2% under the age of 18, 8.0% from 18 to 24, 29.8% from 25 to 44,
8         22.4% from 45 to 64, and 12.7% who are 65 years of age or older. The median age is 36 years. For every 100
9         females there are 100.8 males. For every 100 females age 18 and over, there are 101.5 males.
10
11     The median income for a household in the township is $35,621, and the median income for a family is $35,729. Males
12         have a median income of $29,463 versus $16,071 for females. The [[per capita income]] for the township is $12
13         ,675. 3.4% of the population and 2.4% of families are below the [[poverty line]]. Out of the total
14         population, 1.1% of those under the age of 18 and 3.9% of those 65 and older are living below the poverty
15         line.
16
17     == External links ==
18     {{Mapit-US-cityscale|41.830224|-92.713274}}
19
20     [[Category:Villages in Illinois]]
21     [[Category:Fayette County, Illinois|"Grand Isle,"]]</text>
22     </revision>
23     </page>
24     <page>
25         <title>Rethshawn, Illinois</title>
26         <id>111452</id>
27         <vision>
28             <id>16001292</id>
29             <timestamp>2004-12-21T22:31:38Z</timestamp>
30             <contributor>
31                 <username>Rambot</username>
32                 <id>6120</id>
33             </contributor>
34             <comment>Updated internal [[Geographic References|references]]. Added [[Template:Mapit-US-cityscale]] external
35                 links. Reorganized and cleaned up the page.</comment>
36             <text xml:space="preserve">'''Otho''' is a town located in [[Lyon County, North Carolina]]. As of the
37                 [[2000]] census, the town had a total population of 380.
38
39             == Geography ==
40             [[Image:NCMap-doton-Saxonville.PNG|right|Location of Saxonville, North Carolina]]
41             Stockyard is located at 35&deg;20'20&quot; North, 81&deg;36'1&quot; West (35.543920, -81.635405){{GR|1}}.
42
43             According to the [[United States Census Bureau]], the town has a total area of 7.3 [[square kilometer|km&sup2;]]
44                 (2.8 [[square mile|mi&sup2;]]). 7.3 km&sup2; (2.8 mi&sup2;) of it is land and none of it is
45                 covered by water.
46
47             == Demographics ==
48             As of the [[census]]{{GR|2}} of [[2000]], there are 9,851 people, 4,251 households, and 2,736 families residing in
49                 the town. The [[population density]] is 233.0/km&sup2; (607.3/mi&sup2;). There are 3,102 housing
50                 units at an average density of 226.4/km&sup2; (583.2/mi&sup2;). The racial makeup of the city is
51                 86.82% [[White (U.S. Census)|White]], 5.93% [[African American (U.S. Census)|African American]], 0.04% [[
52                 Native American (U.S. Census)|Native American]], 15.67% [[Asian (U.S. Census)|Asian]], 0.07% [[Pacific
53                 Islander (U.S. Census)|Pacific Islander]], 1.04% from [[Race (U.S. Census)|other races]], and 1.64% from two
54                 or more races. 7.17% of the population are [[Hispanic (U.S. Census)|Hispanic]] or [[Latino (U.S. Census)|
55                 Latino]] of any race.
56
57             There are 408 households out of which 26.8% have children under the age of 18 living with them, 52.2% are [[Marriage
58                 |married couples]] living together, 8.5% have a female householder with no husband present, and 35.2% are non
59                 -families. 31.0% of all households are made up of individuals and 14.0% have someone living alone who is 65
60                 years of age or older. The average household size is 2.43 and the average family size is 3.03.
61
62             In the town the population is spread out with 25.6% under the age of 18, 6.0% from 18 to 24, 26.8% from 25 to 44,
63                 24.9% from 45 to 64, and 14.9% who are 65 years of age or older. The median age is 39 years. For every 100
64                 females there are 95.3 males. For every 100 females age 18 and over, there are 91.4 males.
65
66             The median income for a household in the city is $25,708, and the median income for a family is $31,398. Males have
67                 a median income of $29,589 versus $16,750 for females. The [[per capita income]] for the town is $13,053.
68                 22.3% of the population and 18.0% of families are below the [[poverty line]]. Out of the total population,
69                 37.3% of those unde

```

### Listing 3: enwik10.txt 2-LSTM N-2304 BPC-1.27

```

1     =
2     |image/ParkinmapsfromlIMWebTourism.UH-CIPTIFICAIN.jpg
3     |image_size = 240
4     |image_caption = View of the Republic of Pitza Park in Nottingham.
5     |image_flag =
6     |image_seal =
7     |image_map = NSW Map of Indiana cities locations in Indiana.svg
8     |map_caption = Location in Indiana
9     |image_map1 =
10    |mapsize1 =
11    |map_caption1 =
12    |subdivision_type = [[List of countries|Country]]
13    |subdivision_name = [[United States]]
14    |subdivision_type1 = [[Political divisions of the United States|State]]
15    |subdivision_name1 = [[Indiana]]
16    |subdivision_type2 = [[List of counties in Indiana|County]]
17    |subdivision_name2 = [[Mississippi County, Indiana|Mississippi]]
18    |subdivision_type3 = [[List of townships in Indiana|Township]]
19    |subdivision_name3 = [[Holden Township, Coshocton County, Indiana|Holden]]
20    &lt;!-- Politics -----&gt;

```

```

21 |government_footnotes =
22 |government_type =
23 |leader_title =
24 |leader_name =
25 |leader_title1 =
26 |leader_name1 =
27 |established_title =
28 |established_date =
29
30 <!--gme1901|photomft.<!--ref name=Mw=102<!--[[Hibernian|''Michigan Historical Markers'']]. [http://www.ks.gov/
    facts/hebstr/geog/grossmap.htm Former Government History database]<!--/ref>
31 |established_title3 = Town
32 |established_date3 = 1863
33 |area_magnitude =
34 |area_total_km2 = 1.48
35 |area_total_sq_mi = 0.47
36 |area_land_sq_mi = 0.17
37 |area_water_sq_mi = 0.02
38 |area_water_percent = 0
39 |area_urban_sq_mi =
40 |area_metro_sq_mi =
41 |area_blank1_title =
42 |area_blank1_title =
43 |area_blank1_km2 =
44 |area_blank1_sq_mi =
45 <!-- Population ----->
46 |population_as_of = [[2010 United States Census|2010]]
47 |population_est = 9211
48 |pop_est_as_of = 2012<!--ref name="2012 Pop Estimate">>{{cite web|title=Population Estimates|
    url=http://www.census.gov/popest/data/cities/totals/2012/SUB-EST2012.html|publisher=[[United States Census
    Bureau]]|accessdate=2013-05-29}}<!--/ref>
49 |population_note =
50 |population_total = 118
51 |population_density_km2 = 662.1
52 |population_density_sq_mi = 1832.3
53 |population_metro =
54 |population_density_metro_km2 =
55 |population_density_metro_sq_mi =
56 |population_urban =
57 |population_density_urban_km2 =
58 |population_density_urban_sq_mi =
59 |population_blank1_title =
60 |population_blank1 =
61 |population_density_blank1_km2 =
62 |population_density_blank1_sq_mi =
63 <!-- General information ----->
64 |timezone = [[North American Central Time Zone|Central (CST)]]
65 |utc_offset = -6
66 |timezone_DST = CDT
67 |utc_offset_DST = -5
68 |area_land_km2 = 1.38
69 |area_water_km2 = 0
70 |area_footnotes = <!--ref name="census-g001">>
71 |area_total_km2 = 1.23
72 |area_total_sq_mi = 0.60
73 |area_land_sq_mi = 0.59
74 |area_water_sq_mi = 0
75
76 <!-- Population -->
77 |population_as_of = [[2010 United States Census|2010]]
78 |population_est = 655
79 |pop_est_as_of = 2012<!--ref name="2012 Pop Estimate">>{{cite web|title=Population Estimates|
    url=http://www.census.gov/popest/data/cities/totals/2012/SUB-EST2012.html|publisher=[[United States Census
    Bureau]]|accessdate=2013-05-28}}<!--/ref>
80 |population_footnotes = <!--ref name="FactFinder">>
81 |population_total = 364
82 |population_density_km2 = 650.2
83 |population_density_sq_mi = 1840.0
84
85 <!-- General information -->
86 |timezone = [[North American Central Time Zone|Central (CST)]]
87 |utc_offset = -6
88 |timezone_DST = CDT
89 |utc_offset_DST = -5
90 |elevation_footnotes =
91 |elevation_m = 472
92 |elevation_ft = 1546
93 |coordinates_display = inline,title
94 |coordinates_type = region:US_type:city
95 |latd = 49 |latm = 18 |lats = 45 |latNS = N
96 |longd = 101 |longm = 21 |longs = 33 |longEW = W
97
98 <!-- Area/postal codes & others -->
99 |postal_code_type = [[ZIP code]]
100 |postal_code = 58023
101 |area_code = [[Area code 508|508]]
102 |blank_name = [[Federal Information Processing Standard|FIPS code]]
103 |blank_info = 48-15224<!--ref name="GR2">>{{cite web|url=http://factfinder2.census.gov|
    publisher=[[United States Census Bureau]]|accessdate=2008-01-31|title=American FactFinder}}<!--/ref>
104 |blank1_name = [[Geographic Names Information System|GNIS]] feature ID

```

105 |blank1\_info = 0475524 &lt;ref name=&quot;GR3&quot;&gt;{{cite web|url=http://geonames.usgs.gov|  
accessdate=2008-01-31|title=US Board on Geographic Names|p