# Progress and Tradeoffs in Neural Language Models

**Raphael Tang** and **Jimmy Lin**

David R. Cheriton School of Computer Science
University of Waterloo
{r33tang, jimmylin}@uwaterloo.ca

## Abstract

In recent years, we have witnessed a dramatic shift towards techniques driven by neural networks for a variety of NLP tasks. Undoubtedly, neural language models (NLMs) have reduced perplexity by impressive amounts. This progress, however, comes at a substantial cost in performance, in terms of inference latency and energy consumption, which is particularly of concern in deployments on mobile devices. This paper, which examines the quality–performance tradeoff of various language modeling techniques, represents to our knowledge the first to make this observation. We compare state-of-the-art NLMs with "classic" Kneser-Ney (KN) LMs in terms of energy usage, latency, perplexity, and prediction accuracy using two standard benchmarks. On a Raspberry Pi, we find that orders of increase in latency and energy usage correspond to less change in perplexity, while the difference is much less pronounced on a desktop.

## 1 Introduction

Deep learning has unquestionably advanced the state of the art in many natural language processing tasks, from syntactic dependency parsing (Chen and Manning, 2014) to named-entity recognition (Lample et al., 2016) to machine translation (Luong et al., 2015). The same certainly applies to language modeling, where recent advances in neural language models (NLMs) have led to dramatically better approaches as measured using standard metrics such as perplexity (Melis et al., 2018; Merity et al., 2018b).

Specifically focused on language modeling, this paper examines an issue that to our knowledge has not been explored: advances in neural language models have come at a significant cost in terms of increased computational complexity. Computing the probability of a token sequence using non-neural techniques requires a number of phrase lookups and perhaps a few arithmetic operations, whereas model inference with NLMs require large matrix multiplications consuming perhaps millions of floating point operations (FLOPs). These performance tradeoffs are worth discussing.

In truth, language models exist in a quality–performance tradeoff space. As model quality increases (e.g., lower perplexity), performance as measured in terms of energy consumption, query latency, etc. tends to decrease. For applications primarily running in the cloud—say, machine translation—practitioners often solely optimize for the lowest perplexity. This is because such applications are embarrassingly parallel and hence trivial to scale in a data center environment.

There are, however, applications of NLMs that require less one-sided optimizations. On mobile devices such as smartphones and tablets, for example, NLMs may be integrated into software keyboards for next-word prediction, allowing much faster text entry. Popular Android apps that enthusiastically tout this technology include SwiftKey and Swype. The greater computational costs of NLMs lead to higher energy usage in model inference, translating into shorter battery life.

In this paper, we examine the quality–performance tradeoff in the shift from non-neural to neural language models. In particular, we compare Kneser–Ney smoothing, widely accepted as the state of the art prior to NLMs, to the best NLMs today. The decrease in perplexity on standard datasets has been well documented (Melis et al., 2018), but to our knowledge no one has examined the performances tradeoffs. With deployment on a mobile device in mind, we evaluate energy usage and inference latency on a Raspberry Pi (which shares the same ARM architecture as nearly all smartphones today). We find that a 2.5× reduction in perplexity on PTB comes at a staggering cost in terms of performance: inference with

NLMs takes $49\times$ longer and requires $32\times$ more energy. Furthermore, we find that impressive reductions in perplexity translate into at best modest improvements in next-word prediction, which is arguable a better metric for evaluating software keyboards on a smartphone. The contribution of this paper is the first known elucidation of this quality–performance tradeoff. Note that we refrain from prescriptive recommendations: whether or not a tradeoff is worthwhile depends on the application. Nevertheless, NLP engineers should arguably keep these tradeoffs in mind when selecting a particular operating point.

## 2 Background and Related Work

Melis et al. (2018) evaluate recent neural language models; however, their focus is not on the computational footprint of each model, but rather the perplexity. To further reduce perplexity, many neural language model extensions exist, such as continuous cache pointer (Grave et al., 2017) and mixture of softmaxes (Yang et al., 2018). Since our focus is on comparing "core" neural and non-neural approaches, we disregard these extra optimizations techniques in all of our models.

Other work focus on designing lightweight models for resource-efficient inference on mobile devices. Liu et al. (2018) explore LSTMs (Hochreiter and Schmidhuber, 1997) with binary weights for language modeling; Botha et al. (2017) examine shallow feedforward neural networks for natural language processing.

**AWD-LSTM.** Merity et al. (2018b) show that a simple three-layer LSTM, with proper regularization and optimization techniques, can achieve state of the art on various language modeling datasets, surpassing more complex models. Specifically, Merity et al. (2018b) apply randomized backpropagation through time, variational dropout, activation regularization, embedding dropout, and temporal activation regularization. A novel scheduler for optimization, non-monotonically triggered ASGD (NT-ASGD) is also introduced. Merity et al. (2018b) name their three-layer LSTM model trained with such tricks, "AWD-LSTM."

**Quasi-Recurrent Neural Networks.** Quasi-recurrent neural networks (QRNNs; Bradbury et al., 2017) achieve current state of the art in word-level language modeling (Merity et al., 2018a). A quasi-recurrent layer comprises two separate parts: a convolution layer with three

weights, and a recurrent pooling layer. Given an input $\mathbf{X} \in \mathbb{R}^{k \times n}$, the convolution layer is

$$\mathbf{Z} = \tanh(\mathbf{W}_z \cdot \mathbf{X})$$
$$\mathbf{F} = \sigma(\mathbf{W}_f \cdot \mathbf{X})$$
$$\mathbf{O} = \sigma(\mathbf{W}_o \cdot \mathbf{X})$$

where $\sigma$ denotes the sigmoid function, $\cdot$ represents masked convolution across time, and $\mathbf{W}_{\{z,f,o\}} \in \mathbb{R}^{m \times k \times r}$ are convolution weights with $k$ input channels, $m$ output channels, and a window size of $r$. In the recurrent pooling layer, the convolution outputs are combined sequentially:

$$\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + (1 - \mathbf{f}_t) \odot \mathbf{z}_t$$
$$\mathbf{h}_t = \mathbf{o}_t \odot \mathbf{c}_t$$

Multiple QRNN layers can be stacked for deeper hierarchical representation, with the output $\mathbf{h}_{1:t}$ being fed as the input into the subsequent layer: In language modeling, a four-layer QRNN is a standard architecture (Merity et al., 2018a).

**Perplexity–Recall Scale.** Word-level perplexity does not have a strictly monotonic relationship with recall-at-$k$, the fraction of top $k$ predictions that contain the correct word. A given R@$k$ imposes a weak minimum perplexity constraint—there are many free parameters that allow for large variability in the perplexity given a certain R@$k$. Consider the corpus, "choo choo train," with an associated unigram model $P(\text{"choo"}) = 0.1$, $P(\text{"train"}) = 0.9$, resulting in an R@1 of $1/3$ and perplexity of $4.8$. Clearly, R@1 $= 1/3$ for all $P(\text{"choo"}) \leq 0.5$; thus, perplexity can drop as low as 2 without affecting recall.

## 3 Experimental Setup

We conducted our experiments on Penn Treebank (PTB; Marcus et al., 1993) and WikiText-103 (WT103; Merity et al., 2017). Preprocessed by Mikolov et al. (2010), PTB contains 887K tokens for training, 70K for validation, and 78K for test, with a vocabulary size of 10,000. On the other hand, WT103 comprises 103 million tokens for training, 217K for validation, and 245K for test, spanning a vocabulary of 267K unique tokens.

For the neural language model, we used a four-layer QRNN (Bradbury et al., 2017), which achieves state-of-the-art results on a variety of datasets, such as WT103 (Merity et al., 2018a) and PTB. To compare against more common

LSTM architectures, we also evaluated AWD-LSTM (Merity et al., 2018b) on PTB. For the non-neural approach, we used a standard five-gram model with modified Kneser-Ney smoothing (Chen and Goodman, 1996), as explored in Mikolov and Zweig (2012) on PTB. We denote the QRNN models for PTB and WT103 as `ptb-qrnn` and `wt103-qrnn`, respectively.

For each model, we examined word-level perplexity, R@3 in next-word prediction, latency (ms/q), and energy usage (mJ/q). To explore the perplexity–recall relationship, we collected individual perplexity and recall statistics for each sentence in the test set.

## 3.1 Hyperparameters and Training

The QRNN models followed the exact training procedure and architecture delineated in the official codebase from Merity et al. (2018a). For `ptb-qrnn`, we trained the model for 550 epochs using NT-ASGD (Merity et al., 2018b), then fine-tuned for 300 epochs using ASGD (Polyak and Juditsky, 1992), all with a learning rate of 30 throughout. For `wt103-qrnn`, we followed Merity et al. (2018a) and trained the QRNN for 14 epochs, using the Adam optimizer with a learning rate of $10^{-3}$. We also applied regularization techniques from Merity et al. (2018b); all the specific hyperparameters are the same as those in the repository. Our model architecture consists of 400-dimensional tied embedding weights (Inan et al., 2017) and four QRNN layers, with 1550 hidden units per layer on PTB and 2500 per layer on WT103. Both QRNN models have window sizes of $r = 2$ for the first layer and $r = 1$ for the rest.

For the KN-5 model, we trained an off-the-shelf five-gram model using the popular SRILM toolkit (Stolcke, 2002). We did not specify any special hyperparameters.

## 3.2 Infrastructure

We trained the QRNNs with PyTorch (0.4.0; commit `1807bac`) on a Titan V GPU. To evaluate the models under a resource-constrained environment, we deployed them on a Raspberry Pi 3 (Model B) running Raspbian Stretch (4.9.41-v7+). The Raspberry Pi (RPi) is not only a standard platform, but also a close surrogate to mobile phones, using the same Cortex-A7 in many phones. We then transferred the trained models to the RPi, using the same frameworks for evaluation. We plugged the RPi into a Watts Up Pro meter, a power me-

| Model | Val. | Test |
|---|---|---|
| **Penn Treebank** | | |
| Skip LSTM (Melis et al., 2018) | 60.9 | 58.3 |
| AWD-LSTM (Merity et al., 2018b) | 60.0 | 57.3 |
| QRNN | 59.1 | 56.8 |
| **WikiText-103** | | |
| Rae-LSTM (Rae et al., 2018) | 36.0 | 36.4 |
| QRNN | 31.9 | 32.8 |

Table 1: Comparison of neural language models on Penn Treebank and WikiText-103.

ter that can be read programatically over USB at a frequency of 1 Hz. For the QRNNs, we used the first 350 words of the test set, and averaged the ms/query and mJ/query. For KN-5, we used the entire test set for evaluation, since the latency was much lower. To adjust for the base power load, we subtracted idle power draw from energy usage.

For a different perspective, we further evaluated all the models under a desktop environment, using an i7-4790k CPU and Titan V GPU. Because the base power load for powering a desktop is much higher than running neural language models, we collected only latency statistics. We used the entire test set, since the QRNN runs quickly.

In addition to energy and latency, another consideration for the NLP developer selecting an operating point is the cost of underlying hardware. For our setup, the RPi costs $35 USD, the CPU costs $350 USD, and the GPU costs $3000 USD.

## 4 Results and Discussion

To demonstrate the effectiveness of the QRNN models, we present the results of past and current state-of-the-art neural language models in Table 1; we report the Skip- and AWD-LSTM results as seen in the original papers, while we report our QRNN results. Skip LSTM denotes the four-layer Skip LSTM in Melis et al. (2018). Rae et al. (2018) focus on Hebbian softmax, a model extension technique—Rae-LSTM refers to their base LSTM model without any extensions. In our results, KN-5 refers to the traditional five-gram model with modified Kneser-Ney smoothing, and AWD is shorthand for AWD-LSTM.

**Perplexity–recall scale.** In Figure 1, using KN-5 as the model, we plot the log perplexity (cross entropy) and R@3 error (1 − R@3) for every sentence in PTB and WT103. The horizontal clusters arise from multiple perplexity points representing
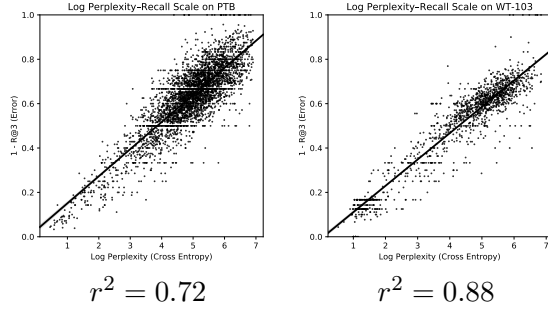
Figure 1: Log perplexity–recall error with KN-5.

$$r^2 = 0.72 \qquad r^2 = 0.88$$



Figure 2: Log perplexity–recall error with QRNN.

$$r^2 = 0.77 \qquad r^2 = 0.86$$

| # Method | Model Quality | | | RPi | | CPU \| GPU | |
|---|---|---|---|---|---|---|---|
| | Val. | Test | R@3 | ms/q | mJ/q | ms/q | ms/q |
| **Penn Treebank** | | | | | | | |
| 1 KN-5 | 148.4 | 141.5 | 36.7% | 7 | 6 | 0.8 | – |
| 2 AWD | 59.2 | 56.8 | 44.9% | 223 | 295 | 7.9 | 1.7 |
| 3 QRNN | 59.1 | 56.8 | 44.7% | 224 | 296 | 7.5 | 1.6 |
| **WikiText-103** | | | | | | | |
| 4 KN-5 | 145.2 | 152.7 | 39.8% | 264 | 229 | 37 | – |
| 5 QRNN | 31.9 | 32.8 | 53.5% | 1240 | 1480 | 59 | 3.5 |

Table 2: Language modeling results on performance and model quality.

the same R@3 value, as explained in Section 2. We also observe that the perplexity–recall scale is non-linear—instead, log perplexity appears to have a moderate linear relationship with R@3 error on PTB ($r = 0.85$), and an even stronger relationship on WT103 ($r = 0.94$). This is partially explained by WT103 having much longer sentences, and thus less noisy statistics.

From Figure 2, we find that QRNN models yield strongly linear log perplexity–recall plots as well, where $r = 0.88$ and $r = 0.93$ for PTB and WT103, respectively. Note that, due to the improved model quality over KN-5, the point clouds are shifted downward compared to Figure 1. We conclude that log perplexity, or cross entropy, provides a more human-understandable indicator of R@3 than perplexity does. Overall, these findings agree with those from Chen et al. (1998), which explores the log perplexity–word error rate scale in language modeling for speech recognition.

**Quality–performance tradeoff.** In Table 2, from left to right, we report perplexity results on the validation and test sets, R@3 on test, and finally per-query latency and energy usage. On the RPi, KN-5 is both fast and power-efficient to run, using only about 7 ms/query and 6 mJ/query for PTB (Table 2, row 1), and 264 ms/q and 229 mJ/q on WT103 (row 5). Taking 220 ms/query and consuming 300
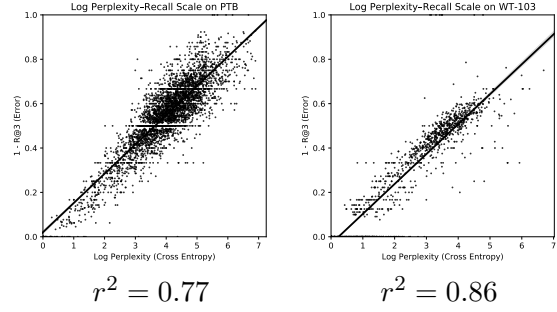
mJ/query, AWD-LSTM and ptb-qrnn are still viable for mobile phones: The modern smartphone holds upwards of 10,000 joules (Carroll et al., 2010), and the latency is within usability standards (Miller, 1968). Nevertheless, the models are still 49× slower and 32× more power-hungry than KN-5. The wt103-qrnn model is completely unusable on phones, taking over 1.2 seconds per next-word prediction. Neural models achieve perplexity drops of 60–80% and R@3 increases of 22–34%, but these improvements come at a much higher cost in latency and energy usage.

In Table 2 (last two columns), the desktop yields very different results: the neural models on PTB (rows 2–3) are 9× slower than KN-5, but the absolute latency is only 8 ms/q, which is still much faster than what humans perceive as instantaneous (Miller, 1968). If a high-end commodity GPU is available, then the models are only twice as slow as KN-5 is. From row 5, even better results are noted with wt103-qrnn: On the CPU, the QRNN is only 60% slower than KN-5 is, while the model is faster by 11× on a GPU. These results suggest that, if only latency is considered under a commodity desktop environment, the QRNN model is humanly indistinguishable from the KN-5 model, even without using GPU acceleration.

## 5  Conclusion

In the present work, we describe and examine the tradeoff space between quality and performance for the task of language modeling. Specifically, we explore the quality–performance tradeoffs between KN-5, a non-neural approach, and AWD-LSTM and QRNN, two neural language models. We find that with decreased perplexity comes vastly increased computational requirements: In one of the NLMs, a perplexity reduction by 2.5× results in a 49× rise in latency and 32× increase in energy usage, when compared to KN-5.

# References

Jan A Botha, Emily Pitler, Ji Ma, Anton Bakalov, Alex Salcianu, David Weiss, Ryan McDonald, and Slav Petrov. 2017. Natural language processing with small feed-forward networks. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2879–2885.

James Bradbury, Stephen Merity, Caiming Xiong, and Richard Socher. 2017. Quasi-recurrent neural networks. In *International Conference on Learning Representations*.

Aaron Carroll, Gernot Heiser, et al. 2010. An analysis of power consumption in a smartphone. In *USENIX Annual Technical Conference*, volume 14, pages 21–21.

Danqi Chen and Christopher Manning. 2014. A fast and accurate dependency parser using neural networks. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing*, pages 740–750.

Stanley F Chen, Douglas Beeferman, and Ronald Rosenfeld. 1998. Evaluation metrics for language models. In *DARPA Broadcast News Transcription and Understanding Workshop*, pages 275–280.

Stanley F Chen and Joshua Goodman. 1996. An empirical study of smoothing techniques for language modeling. In *Proceedings of the 34th annual meeting on Association for Computational Linguistics*, pages 310–318.

Edouard Grave, Armand Joulin, and Nicolas Usunier. 2017. Improving neural language models with a continuous cache. In *International Conference on Learning Representations*.

Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural Computation*, 9(8):1735–1780.

Hakan Inan, Khashayar Khosravi, and Richard Socher. 2017. Tying word vectors and word classifiers: A loss framework for language modeling. In *International Conference on Learning Representations*.

Guillaume Lample, Miguel Ballesteros, Sandeep Subramanian, Kazuya Kawakami, and Chris Dyer. 2016. Neural architectures for named entity recognition. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 260–270.

Xuan Liu, Di Cao, and Kai Yu. 2018. Binarized LSTM language model. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, volume 1, pages 2113–2121.

Thang Luong, Hieu Pham, and Christopher D. Manning. 2015. Effective approaches to attention-based neural machine translation. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 1412–1421.

Mitchell P Marcus, Mary Ann Marcinkiewicz, and Beatrice Santorini. 1993. Building a large annotated corpus of English: The Penn Treebank. *Computational linguistics*, 19(2):313–330.

Gábor Melis, Chris Dyer, and Phil Blunsom. 2018. On the state of the art of evaluation in neural language models. In *International Conference on Learning Representations*.

Stephen Merity, Nitish Shirish Keskar, and Richard Socher. 2018a. An analysis of neural language modeling at multiple scales. *arXiv:1803.08240*.

Stephen Merity, Nitish Shirish Keskar, and Richard Socher. 2018b. Regularizing and optimizing LSTM language models. In *International Conference on Learning Representations*.

Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. 2017. Pointer sentinel mixture models. In *International Conference on Learning Representations*.

Tomáš Mikolov, Martin Karafiát, Lukáš Burget, Jan Černockỳ, and Sanjeev Khudanpur. 2010. Recurrent neural network based language model. In *Eleventh Annual Conference of the International Speech Communication Association*.

Tomas Mikolov and Geoffrey Zweig. 2012. Context dependent recurrent neural network language model. *SLT*, 12:234–239.

Robert B. Miller. 1968. Response time in man-computer conversational transactions. In *Proceedings of the December 9-11, 1968, Fall Joint Computer Conference, Part I*, AFIPS '68 (Fall, part I), pages 267–277.

B. T. Polyak and A. B. Juditsky. 1992. Acceleration of stochastic approximation by averaging. *SIAM J. Control Optim.*, 30(4):838–855.

Jack W Rae, Chris Dyer, Peter Dayan, and Timothy P Lillicrap. 2018. Fast parametric learning with activation memorization. *arXiv:1803.10049*.

Andreas Stolcke. 2002. SRILM — an extensible language modeling toolkit. In *Seventh International Conference on Spoken Language Processing*.

Zhilin Yang, Zihang Dai, Ruslan Salakhutdinov, and William W. Cohen. 2018. Breaking the softmax bottleneck: A high-rank RNN language model. In *International Conference on Learning Representations*.