

# Paper review of "Exploring the Limits of Language Modeling"

Alexandra Ioana Neagu

TU Delft

neagu.alexandra0206@gmail.com

## Abstract

This paper review critically evaluates the research paper titled *Exploring the Limits of Language Modeling* (Jozefowicz et al., 2016). The review assesses the strengths and weaknesses of the paper, considering factors such as originality, clarity, and potential impact. The review highlights the research questions addressed in the experimental evaluation and suggests potential areas for improvement in the presented work. Additionally, it evaluates the reproducibility of the findings and the confidence level in the assessment. Overall, the review aims to provide a comprehensive analysis of the paper's contributions to the field of natural language processing (NLP).

## 1 What is this paper about, and what contributions to NLP does it make?

This paper focuses on language modelling (LM) in the context of NLP. Language models that accurately predict distributions over sentences not only capture grammatical structures but also encapsulate knowledge present in the training data. The paper highlights the importance of LMs in various NLP tasks such as speech recognition, machine translation, and text summarization. Despite the prevalence of simpler models like N-grams, recent advancements in deep learning and recurrent neural networks (RNNs) have driven LM research, enabling the exploration of tasks where strong conditional independence assumptions are unrealistic.

The authors emphasize the significance of training better LMs, as it often improves the performance of downstream tasks and enables the extraction of knowledge encoded in vast amounts of training data. They advocate for research on larger-scale LM tasks, like the One Billion Word Benchmark dataset, which offers a more realistic representation of language compared to smaller datasets like the Penn Tree Bank. By releasing their models and

training recipes, the authors aim to encourage further research in LM beyond existing benchmarks.

In my opinion, the main contributions of the paper are outlined as follows:

1. Exploration, extension, and unification of current research on large-scale LM.
2. Design of a Softmax loss based on character-level convolutional neural networks (CNNs), which is efficient to train and as precise as a full Softmax but with significantly fewer parameters.
3. Significant improvements to the state-of-the-art on the One Billion Word Benchmark dataset, achieving a perplexity reduction from 51.3 to 30.0 for single models while reducing the number of parameters by a factor of 20.
4. Demonstration that an ensemble of different models further reduces perplexity to 23.7, surpassing current state-of-the-art results.
5. Lastly, the sharing of models and training recipes greatly contributes to future research in this domain.

### 1.1 Paper summary

In summary, the paper addresses the challenges and advancements in large-scale language modelling, proposing novel techniques and achieving notable improvements in LM performance on benchmark datasets.

*The related work section* provides a comprehensive overview of related background in the field of LM and NN-based approaches. It discusses the significance of LM in NLP and outlines various approaches, including parametric and non-parametric methods. Count-based approaches, such as Kneser-Ney smoothed 5-gram models, have been successful in capturing language statistics and challenging NN-based approaches, particularly in scenarios with large training data. The paper primarily

focuses on RNN models, specifically Long Short-Term Memory (LSTM) networks, which excel in capturing long-term dependencies in language sequences. The authors also highlight the importance of large-scale datasets, like the **One Billion Word Benchmark**, in training effective LMs. Furthermore, the section introduces convolutional embedding models, which incorporate character-level inputs to build word embeddings, and discusses the computational challenges associated with assigning probability distributions over large vocabularies, proposing solutions such as importance sampling (IS) and Noise Contrastive Estimation (NCE).

In *Section 3*, the paper delves into LM improvements, focusing on RNN-based models and various strategies to address computational challenges associated with large-scale softmax layers. Firstly, the relationship between NCE and IS is discussed in the context of approximating the softmax layer for efficient training. While both methods propose surrogate tasks to approximate the softmax, IS optimizes a multiclass classification task and may offer advantages over NCE. The paper then introduces the concept of a **CNN Softmax layer**, which leverages character-level features to reduce the number of parameters in the softmax layer while maintaining efficiency. Additionally, the paper explores **Char LSTM Predictions**, which combines word-level and character-level models to improve computational efficiency by predicting words character by character using an LSTM network. Though these models are computationally attractive, they may not yet outperform traditional softmax-based models, but further research is anticipated to enhance their effectiveness.

**Section 4** details the experimental setup and methodology employed in the study. Using the TensorFlow system, experiments are conducted on the 1B Word Benchmark dataset, which serves as a widely accepted benchmark for statistical LM progress evaluation. The dataset comprises approximately 0.8 billion words with a vocabulary size of 793,471 words, where out-of-vocabulary words are denoted with a special token. Model architectures and training procedures are meticulously described, with a focus on evaluating various RNN LM architectures and exploring techniques to optimize performance. The experiments involve testing different LSTM configurations, including dimensions of embedding layers, state sizes, projection sizes, and the number of LSTM layers. The training

procedure involves using an AdaGrad optimizer with a learning rate of 0.2, unrolling RNNs for 20 steps without resetting LSTM states, and employing gradient clipping to ensure stable training. Additionally, details about training large LSTM models and approximating the softmax using importance sampling are elucidated, highlighting the computational considerations and strategies employed to handle the large target vocabulary effectively.

**Section 5** presents the results and analysis of the experiments conducted. The authors express hope that their results, including achieving a new record in perplexity with ensembles, will facilitate rapid progress in LM, and they pledge to release model weights and recipes upon publication. Subsections delve into various aspects of the findings:

1. **Size Matters:** Larger LSTM models tend to perform better, with the best results achieved by the largest models that could fit into GPU memory. Increasing the LSTM layer size and embedding dimensions also contributes to improved performance.
2. **Regularization Importance:** Dropout regularization proves effective in improving results, even for relatively small models, helping to mitigate overfitting issues.
3. **Importance Sampling is Data Efficient:** IS significantly enhances speed and overall model performance compared to NCE.
4. **Word Embeddings vs. Character CNN:** Replacing the embedding layer with a character-level NN allows the model to handle arbitrary words and is beneficial for datasets with conversational or informal text. Furthermore, using character-level embeddings does not degrade performance and reduces the number of input layer parameters.
5. **Smaller Models with CNN Softmax:** Substituting the Softmax layer with a CNN Softmax sub-network reduces the model size considerably while maintaining performance. Moreover, introducing a correction term further narrows the gap between regular and CNN Softmax.
6. **Training Speed:** Training speed varies with model size, with smaller models reaching competitive perplexity levels in a matter of

hours, while the best models require several days to weeks to achieve optimal performance.

7. **Ensembles:** Averaging multiple models yields significant improvements in perplexity, with the ensemble achieving a perplexity of 23.7, a notable advancement over previous work.
8. **LSTMs are Best on Tail Words:** LSTM models outperform N-gram models, especially on rare words, suggesting their superiority for languages or datasets with a large number of rare words.

Lastly, qualitative evaluation through sampled sentences demonstrates the high quality of the generated text, reflecting the impressive perplexities attained by the models. These findings collectively underscore the efficacy of LSTM-based language models and highlight avenues for further research and improvement in LM. The paper concludes by suggesting that future research should focus on reasonably sized datasets, drawing inspiration from successes in computer vision research, such as the ImageNet dataset.

## 2 What strengths does this paper have?

This paper exhibits several strengths. Firstly, it explores novel techniques in LM, such as employing large, regularized LSTM LMs trained with IS to approximate the true Softmax. This approach leads to a significant reduction in perplexity compared to traditional N-gram models. The authors provide detailed descriptions of the experimental setup, including data sets used, model architectures, training procedures, and results. Additionally, they indicate their intention to release the accompanying code and models upon publication, enhancing the replicability of their work. The findings of this paper could have a substantial impact on the field of LM. By demonstrating the superiority of LSTM models over N-grams and providing insights into effective training techniques, the paper sets a precedent for future research in large-scale language understanding.

In my opinion, the paper is well-written and effectively communicates the methodologies and findings. Each section is logically structured, making it easy for readers to follow the research process and understand the implications of the results.

## 3 What weaknesses does this paper have?

While the paper has many strengths, it also has some weaknesses. It does not extensively discuss the limitations of the proposed approach or potential areas for improvement. Acknowledging and addressing these limitations could strengthen the paper's credibility and provide more explicit avenues for future research. While the paper mentions comparisons with existing N-gram models, it could benefit from a more comprehensive comparison with other SOTA approaches in LM.

Furthermore, the paper briefly mentions occasional mistakes in the generated sentences but does not delve into a detailed analysis of these errors. Understanding the limitations of the model in specific scenarios could provide valuable insights for future improvements. Lastly, while the paper outlines the experimental setup and intends to release code and models, replicating the results may still be challenging due to the complexity of training large LSTM LMs and the computational resources required. Providing more guidance or resources for replication could address this issue.

## 4 What are the research questions that are answered in the experimental evaluation?

In my opinion, the experimental evaluation in the paper aims to answer several research questions, out of which three of the most important ones are, in no particular order:

1. **RQ1: Impact of Model Size:** One research question revolves around the influence of model size on LM performance. The key result in this regard is the demonstration that larger LSTM models, such as the 2-layer LSTM with 8192+1024 dimensional recurrent states, outperform smaller models, indicating that size indeed matters in LM.
2. **RQ2: Effectiveness of Regularization Techniques:** Another research question focuses on the effectiveness of regularization techniques, particularly dropout, in mitigating overfitting. The experimental results show that dropout improves the model's generalization performance, even for relatively small models. However, it also highlights that overfitting can still occur, suggesting the need for careful regularization strategies.

3. **Comparison of Sampling Methods:** The experimental evaluation also compares the effectiveness of NCE and IS in approximating the partition function during training. The results demonstrate that IS significantly improves speed and overall model performance compared to NCE, providing insights into efficient training strategies for large-scale language models.

## 5 Can you think of how you would improve the work presented in this paper?

One potential improvement that I believe to the work presented in this paper could be to **investigate the incorporation of external knowledge sources or context into the LM framework**. This could be approached through the following research questions:

1. **Incorporating Semantic Embeddings: How does integrating semantic embeddings or knowledge graphs into the LM impact its performance in capturing semantic relationships and improving context understanding?** By incorporating external semantic information, such as *WordNet* or embeddings trained on domain-specific corpora, the LMI may be able to better understand and generate text that reflects nuanced semantic meanings.
2. **Adaptive Contextual Learning: Can the LM dynamically adapt its context window based on the specific linguistic properties or discourse structures of the input text?** By incorporating mechanisms for adaptive context learning, the model could adjust its attention or memory mechanisms to focus on relevant contextual information, leading to more accurate predictions and improved coherence in the generated text.

These potential improvements could enhance the robustness, flexibility, and performance of the LM framework, leading to more advanced natural language understanding and generation capabilities.

## 6 How reproducible is the work?

I would rate the reproducibility of the work described in the paper as a **3 out of 4**. Researchers and practitioners could mostly reproduce the results by utilizing the provided methodologies, datasets,

and code, which were made available upon publication. However, there might be some challenges in reproducing certain aspects of the experiments due to the complexity of the model architectures and the specific choices of hyperparameters. Additionally, while the paper provides detailed descriptions of the methodologies employed, there may still be some subjective elements in parameter tuning or model design that could introduce variability in reproduction efforts. Overall, with access to the provided resources and clear documentation, reproducing the core findings of the paper should be feasible for most researchers and practitioners (Raff, 2019).

## 7 How confident are you in your assessment of this paper?

I would rate my confidence in the assessment of this paper as a **4 out of 5**. I carefully reviewed the important points and tried to ensure I didn't miss any crucial details. However, there's always a chance that I might have overlooked something that could affect my ratings. While I have a good understanding of the general area, I didn't delve deeply into some of the paper's details, such as the intricacies of the experimental design. Therefore, I'm fairly confident in my evaluation, but there's still a possibility of missing some nuances.

## References

- Rafal Jozefowicz, Oriol Vinyals, Mike Schuster, Noam Shazeer, and Yonghui Wu. 2016. Exploring the limits of language modeling. *arXiv preprint arXiv:1602.02410*.
- Edward Raff. 2019. A step toward quantifying independently reproducible machine learning research. *Advances in Neural Information Processing Systems*, 32.