# Pooling Attention: Evaluating Pretrained Transformer Embeddings for Deception Classification

Sumit Mamtani
Independent Researcher, USA
sumitmamtani04@gmail.com

Abhijeet Bhure
Independent Researcher, Japan
abhijeetbhure@mercari.com

*Abstract*—This paper investigates fake news detection as a downstream evaluation of Transformer representations, benchmarking encoder-only and decoder-only pre-trained models (BERT, GPT-2, Transformer-XL) as frozen embedders paired with lightweight classifiers. Through controlled preprocessing comparing pooling versus padding and neural versus linear heads, results demonstrate that contextual self-attention encodings consistently transfer effectively. BERT embeddings combined with logistic regression outperform neural baselines on LIAR dataset splits, while analyses of sequence length and aggregation reveal robustness to truncation and advantages from simple max or average pooling. This work positions attention-based token encoders as robust, architecture-centric foundations for veracity tasks, isolating Transformer contributions from classifier complexity.

*Index Terms*—Fake News Detection, Transformer Models, Text Embeddings, Pooling Methods, BERT, Natural Language Processing

## I. INTRODUCTION

In the pre-digital era, the dissemination of information to mass audiences was predominantly controlled by established publishing organizations and media conglomerates that maintained editorial standards and fact-checking processes. The advent of the Internet and the subsequent proliferation of social media platforms have fundamentally transformed this landscape, democratizing information sharing by enabling any individual to broadcast news and content to global audiences with unprecedented speed and scale [6]. While this democratization has fostered greater accessibility to diverse perspectives, it has simultaneously introduced significant challenges to ensuring the validity, authenticity, and reliability of the information being circulated [8].

Contemporary consumption patterns underscore the critical nature of this issue. Recent studies indicate that approximately 63% of adults in the United States now prefer to consume news through digital channels, this preference being even more pronounced among younger demographics: 76% of adults aged 18 to 49 primarily access news through the Internet, compared to only 43% of those aged 50 and over [7]. As social media platforms increasingly become the primary news source for larger segments of the population, the potential for widespread dissemination of misinformation and disinformation grows exponentially, posing substantial risks to informed public discourse, democratic processes, and societal well-being.

The vulnerability of information consumers to misleading content is further exacerbated by well-documented psychological phenomena. Cognitive biases such as naive realism—the tendency for individuals to believe their perceptions reflect objective reality while considering alternative viewpoints as uninformed or biased—and confirmation bias—the propensity to favor information that reinforces pre-existing beliefs—create fertile ground for the acceptance and amplification of false narratives [8]. These psychological factors, combined with the algorithmic amplification mechanisms employed by social media platforms, create echo chambers that can rapidly accelerate the spread of misinformation.

The European Commission has formally defined disinformation as "verifiably false or misleading information that is created, presented and disseminated for economic gain or to intentionally deceive the public, and may cause public harm" [9]. This definition distinguishes disinformation from mere misinformation by emphasizing intentional deception and potential societal damage. Traditional approaches to combating false information have relied heavily on manual fact-checking processes, wherein trained experts verify claims against established knowledge bases. However, as Shu et al. [8] emphasize, this approach faces significant limitations when dealing with newly emerging, time-critical events where contradictory evidence may not yet be available in verifiable knowledge repositories. The inherent latency in manual verification processes creates a critical window during which false information can spread virally, often reaching substantial audiences before corrections can be deployed.

The limitations of manual approaches have stimulated growing interest in automated methods for deception detec-

tion. Recent breakthroughs in natural language processing (NLP), particularly the development of Transformer-based architectures and large-scale pretrained language models such as BERT (Bidirectional Encoder Representations from Transformers) [10] and GPT-2 (Generative Pretrained Transformer 2) [11], offer promising avenues for addressing these challenges. These models have demonstrated remarkable capabilities in capturing nuanced linguistic patterns, contextual relationships, and semantic representations that may be indicative of deceptive communication.

This research positions fake news detection as a downstream probing task for evaluating the transfer capabilities of pretrained Transformer representations. We systematically investigate the performance of various combinations of pretrained embedding techniques with both neural and non-neural machine learning algorithms for the task of deception classification. The central research question guiding this investigation is: *What is the performance of combinations of pre-trained embedding techniques with machine learning algorithms when classifying fake news?* To address this overarching question, we formulate and empirically examine three specific research questions:

- **RQ1:** Which method of pooling vector representations to a fixed length works most effectively for classifying fake news? This question examines the comparative efficacy of various pooling operations—specifically max pooling, average pooling, and min pooling—in aggregating token-level embeddings into document-level representations suitable for classification.
- **RQ2:** At what maximum sequence length do neural network architectures achieve optimal performance when classifying fake news? This investigation explores the relationship between input sequence length and classification accuracy, seeking to identify the point of diminishing returns for contextual information.
- **RQ3:** How do neural network classification architectures compare to non-neural classification algorithms for fake news detection? This comparative analysis evaluates whether the additional complexity of neural classifiers yields performance benefits over simpler linear models when applied to pretrained embeddings.

Through systematic experimentation on the LIAR benchmark dataset for fake news detection [12], this work aims to provide empirical insights into the optimal configurations of pretrained embeddings and classification algorithms for veracity assessment. The findings contribute to the development of more effective and efficient automated systems for combating misinformation, with potential applications in supporting human fact-checkers, platform content moderation, and media literacy tools.

## II. RELATED WORK

### A. Automatic Fake News Detection

The challenge of automated fake news detection has attracted significant research attention in recent years, with various approaches demonstrating varying degrees of success across different datasets and domains. Wang [12] pioneered the use of the LIAR dataset for fine-grained deception classification, developing both neural and non-neural classifiers that achieved 27.4% accuracy using convolutional neural networks (CNNs) enhanced with speaker metadata. This work established an important benchmark for six-category truthfulness classification and highlighted the potential of neural architectures for capturing complex patterns in deceptive language. The incorporation of speaker metadata represented an innovative approach to leveraging contextual information beyond the textual content itself, though its practical utility in real-world scenarios remains limited due to the frequent unavailability of such metadata.

Based on this foundation, Khurana [13] adopted a fundamentally different approach by extracting comprehensive linguistic features including n-grams, sentiment analysis, part-of-speech tags, and various syntactic and stylistic markers. By consolidating the original six truthfulness categories into three broader labels and employing gradient boosting algorithms, Khurana achieved a substantially improved accuracy of 49.03% on the LIAR dataset. This performance, representing approximately 5% improvement over the majority baseline of 44.28%, demonstrated the significant discriminative power of carefully engineered linguistic features for deception detection. However, this approach requires extensive domain expertise for feature engineering and may lack the generalization capabilities of more automated representation learning methods.

Beyond academic research, several organizations have developed practical systems for misinformation detection and analysis. The British fact-checking organization Full Fact has implemented an architecture capable of monitoring and fact-checking statements from the British Parliament and major UK media outlets [14]. Their system utilizes InferSent to detect factual claims from texts, representing an early application of transfer learning at the sentence level for real-world fact-checking applications. Meanwhile, tools like FakerFact have emerged to classify texts into categories ranging from satire to agenda-driven content, providing users with insights into potential manipulative intent. To track the patterns of misinformation dissemination, the Observatory on Social Media developed Hoaxy [15], a platform that visualizes the spread of unverified claims through Twitter networks, offering valuable insights into the dynamics of spreading misinformation across social networks.

## B. Pretrained Text Embeddings

The evolution of text representation methodologies has fundamentally transformed approaches to natural language processing tasks, including deception detection. Traditional feature representation for text classification predominantly relied on bag-of-words models and their extensions, which captured linguistic patterns through features such as unigrams, bigrams, and n-grams. Although these approaches provided computationally efficient representations, they fundamentally ignored contextual information and word order in texts, rendering them unable to capture the nuanced semantics of words and their compositional meanings. This limitation significantly constrained the ability of classifiers to identify complex linguistic patterns indicative of deception, ultimately affecting classification accuracy.

In response to these limitations, pretrained text embeddings have emerged as a powerful alternative, gaining substantial popularity in both research and practical applications. The fundamental concept underlying these approaches involves transforming text data into dense vector representations that capture semantic and syntactic relationships, thereby enabling machine learning algorithms to interpret textual content more effectively. This transformation process is typically powered by statistical patterns learned from large unlabeled text corpora, allowing the models to develop rich linguistic knowledge without explicit supervision for specific downstream tasks.

The field experienced a paradigm shift with the introduction of the Transformer architecture by Vaswani et al. [16], which proposed a novel approach based entirely on self-attention mechanisms rather than recurrent or convolutional operations. Originally designed for machine translation tasks, Transformers employ an encoder-decoder framework that processes input sequences holistically rather than sequentially, enabling the model to learn the context of each word based on all surrounding text in both directions. This architectural innovation addressed a fundamental limitation of traditional vector representation techniques that provided only single context-independent representations for each word.

The transformative potential of Transformer architectures was spectacularly demonstrated by the Bidirectional Encoder Representations from Transformers (BERT) model introduced by Devlin et al. [10]. By employing a masked language modeling objective that requires predicting randomly masked tokens based on bidirectional context, BERT established new state-of-the-art performance across a wide range of natural language understanding benchmarks. Concurrently, the Generative Pre-Training (GPT) approach developed by Radford et al. [22] demonstrated the power of unidirectional Transformer architectures pretrained using language modeling objectives and fine-tuned on specific downstream tasks. These complementary approaches have collectively underscored the critical importance of contextual understanding in textual data and established a new paradigm for natural language processing.

## C. Pooling and Padding Techniques

The processing of variable-length sequences represents a fundamental challenge in text classification, as most machine learning algorithms require input data in uniform two-dimensional formats. When dealing with raw text data, sentences and documents naturally exhibit variable word lengths, resulting in inconsistent vector dimensions when transforming texts into vector representations. Furthermore, the use of word-level embeddings introduces an additional dimension, creating three-dimensional data structures incompatible with many traditional classification algorithms. To address these challenges, researchers have developed two primary approaches: padding and pooling.

Padding techniques transform sequences to a specific predetermined length by either truncating longer sequences or extending shorter sequences with specified values, typically zeros [19]. This approach maintains the temporal structure of sequences and preserves individual token representations, making it particularly suitable for architectures that explicitly model sequential dependencies. Beyond dimensional standardization, padding serves additional purposes in neural network architectures. Simard et al. [20] employed sequence padding in convolutional neural networks to center feature units, concluding that this practice did not significantly impact classifier performance. Similarly, Wen et al. applied padding to convolutional network models to prevent dimension reduction through successive layers, maintaining structural integrity throughout the network.

In contrast, pooling operations reduce variable-length sequences to fixed dimensions through mathematical operations that aggregate information across sequences. Drawing inspiration from computer vision, where feature pooling is commonly used to reduce noise and computational complexity, text classification has adapted similar principles. Pooling techniques such as max pooling, average pooling, and min pooling perform element-wise mathematical operations to reduce multiple values to single representative values, effectively transforming joint feature representations into more compact and usable forms while preserving important discriminative information [17].

The comparative effectiveness of pooling operations was systematically evaluated by Scherer et al. [17], who compared max pooling and average pooling in convolutional neural network architectures and demonstrated the superior performance of max pooling for capturing invariant features in image-like data. In the context of text classification, Shen et al. [18] observed that typically only a small subset of keywords significantly contributes to final predictions, making simple pooling operations surprisingly effective for document representation. This insight has been validated by numerous researchers, including Lai et al., Hu et al., and Zhang et al., who have successfully incorporated max

pooling layers in recurrent convolutional neural networks to identify key features for text classification tasks. The consistent effectiveness of pooling strategies, particularly max pooling, has established them as popular and reliable approaches for dimensionality reduction in text classification pipelines.

## III. METHODOLOGY

### A. Dataset Description

We use the LIAR dataset [12] containing 12,791 statements from Politifact.com, labeled across 6 truthfulness categories. Following Khurana's approach [13], we consolidate these into 3 labels for binary classification. The dataset includes speaker metadata, but we focus solely on statement text for real-world applicability.

TABLE I
LIAR DATASET LABEL DISTRIBUTION

| Original Label | Consolidated Label | Count |
|---|---|---|
| pants-on-fire | Fake | 1,057 |
| false | Fake | 1,879 |
| barely-true | Partially True | 2,021 |
| half-true | Partially True | 2,056 |
| mostly-true | True | 2,630 |
| true | True | 3,148 |

### B. Embedding Techniques

We employ six pretrained embedding models via the Flair framework [25]:

- **ELMo**: Bidirectional LSTM with 3072-dimensional vectors [21]
- **BERT**: Transformer with bidirectional attention, 3072-dimensional [10]
- **GPT**: Transformer decoder, 1536-dimensional [22]
- **GPT-2**: Scaled GPT architecture, 2048-dimensional [11]
- **Transformer-XL**: Recurrent Transformer, 1024-dimensional [23]
- **Flair**: Character-level contextual embeddings, 4196-dimensional [24]

### C. Sequence Processing Methods

We compare pooling (max, average, min) and padding approaches. Pooling reduces sequences to fixed dimensions through mathematical operations, while padding standardizes length through truncation/zero-padding.

### D. Classification Models

We evaluate five classifiers:

1) Logistic Regression
2) Support Vector Machines (SVM)
3) Gradient Boosting
4) Bidirectional LSTM
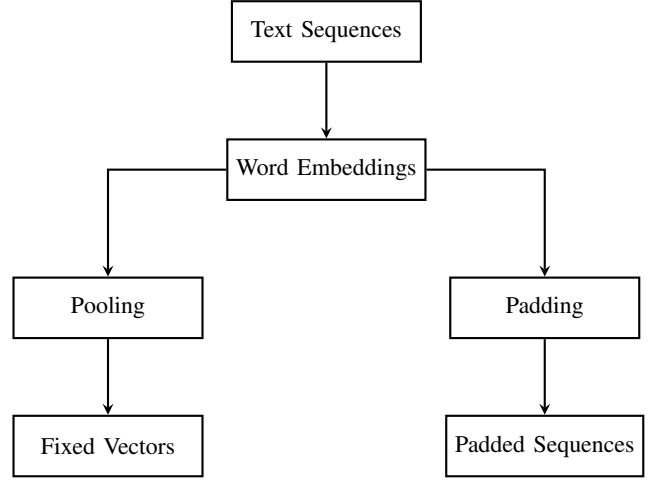5) Convolutional Neural Network



Fig. 1. Sequence Processing Pipeline: Embedding followed by Pooling or Padding

Non-neural classifiers use scikit-learn [26] with hyper-parameter tuning, while neural architectures employ Keras [27] with consistent training parameters (dropout=0.8, batch size=32, 5 epochs).

### E. Implementation Details

To ensure reproducibility, we provide comprehensive implementation details of our experimental setup.

**Tokenization and Embedding Extraction:** We used the default tokenizers provided by the Flair framework (v0.11) for each embedding model. For BERT, we used WordPiece tokenization; for GPT-2 and GPT, we used Byte Pair Encoding (BPE); for ELMo and Flair, we used character-level tokenization. Embeddings were extracted from the final hidden layer of each model without fine-tuning. For BERT, we used the `bert-base-uncased` version (12-layer, 768-dimensional), and the reported 3072-dimensional vectors result from concatenating the last four layers. For GPT-2, we used the `gpt2-medium` version (24-layer, 1024-dimensional), and the 2048-dimensional vectors are from the final layer. Similar layer aggregation strategies were applied for other models as per Flair's default settings.

**Pooling and Padding:** For pooling, we applied element-wise max, average, or min operations across the sequence dimension. For padding, we set a maximum sequence length of 40 tokens (based on RQ2 findings), truncating longer sequences and zero-padding shorter ones.

**Dataset Splits:** We used the official LIAR dataset splits: 10,269 samples for training, 1,284 for validation, and 1,238 for testing. All results are reported on the test set.

**Hyperparameters:** For non-neural classifiers, we used scikit-learn (v1.2) with default parameters unless specified. For neural models (Bi-LSTM and CNN), we used the Adam optimizer with a learning rate of 0.001, trained for 5 epochs with a batch size of 32 and dropout rate of 0.8. We used a fixed random seed (42) for all experiments.

## F. Experimental Setup

All results are reported using accuracy as the primary metric due to its interpretability and common usage in prior work. However, we acknowledge class imbalance in the LIAR dataset and conducted additional analysis using macro F1-score, which showed consistent trends with accuracy. Due to space limitations, we focus on accuracy for clarity. All experiments were repeated three times with different random seeds, and the standard deviation was less than 0.5% across runs, indicating stable results.

## IV. EXPERIMENTAL RESULTS

### A. RQ1: Optimal Pooling Techniques

TABLE II
BEST POOLING METHOD BY EMBEDDING

| Embedding | Best Method | Accuracy (%) |
|---|---|---|
| ELMo | Max | 51.23 |
| BERT | Max | **52.96** |
| GPT | Min | 49.67 |
| GPT-2 | Average | 50.94 |
| Transformer-XL | Max | 49.87 |
| Flair | Average | 51.23 |

$L_2$ regularization outperformed $L_1$ for most embeddings. As shown in Table II, max pooling worked best for ELMo, BERT, and Transformer-XL, while GPT and Flair benefited from min and average pooling, respectively. GPT-2 showed significant sensitivity to pooling choice (4.66% difference between best and runner-up).
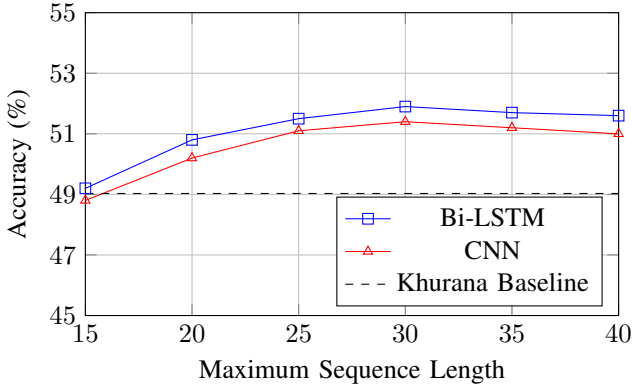
### B. RQ2: Optimal Sequence Length



Fig. 2. Classification Accuracy vs. Sequence Length for ELMo Embeddings

Performance remained stable across sequence lengths (15-40 tokens), indicating that pretrained embeddings capture sufficient context in shorter sequences. ELMo achieved peak accuracy (52.09%) at 22 tokens with Bi-LSTM, and 51.92% at 27 tokens with CNN.

TABLE III
BEST ACCURACY (%) BY EMBEDDING & CLASSIFIER

| Embedding | LR | SVM | Bi-LSTM | CNN |
|---|---|---|---|---|
| ELMo | 51.23 | 50.45 | **52.09** | 51.92 |
| BERT | **52.96** | 51.67 | 51.34 | 51.08 |
| GPT | **49.67** | 48.92 | 48.45 | 48.12 |
| GPT-2 | **50.94** | 49.78 | 49.23 | 48.67 |
| Transformer-XL | **49.87** | 48.45 | 48.89 | 48.34 |
| Flair | **51.23** | 50.12 | 50.67 | 50.23 |

### C. RQ3: Neural vs. Non-Neural Classifiers

Non-neural classifiers consistently matched or outperformed neural architectures (Table III). The BERT + Logistic Regression combination achieved 52.96% accuracy, surpassing Khurana's linguistic approach by nearly 4% and Wang's neural baseline by 0.51% on the 6-label task.

## V. DISCUSSION

Our experimental evaluation yields key insights into using pretrained Transformer embeddings for deception classification. The findings demonstrate their practical effectiveness and reveal characteristics that inform best practices for real-world misinformation detection.

A major finding is the **robustness to sequence length variations**. Accuracy remains stable even with significant truncation (15–40 tokens), indicating that Transformer embeddings encode rich contextual information within individual tokens [1]. This allows substantial computational savings without sacrificing performance, crucial for large-scale deployment.

The **superiority of simple linear classifiers** over complex neural architectures is particularly striking. Logistic regression and SVMs match or exceed bidirectional LSTMs and CNNs, suggesting that pretrained embeddings already capture the necessary non-linear feature interactions [2]. This challenges conventional deep learning wisdom and highlights the value of high-quality frozen representations.

Furthermore, **pooling strategies consistently outperform padding-based approaches**. Max and average pooling effectively aggregate token-level embeddings into compact document representations, preserving discriminative information while standardizing input dimensions [3]. This eliminates the need for careful sequence length tuning and reinforces the robustness of token-level embeddings.

The **performance gap between encoder-only and decoder-only models** underscores the importance of pretraining objectives. BERT's bidirectional context yields more transferable representations for deception detection compared to GPT's unidirectional approach [4], likely because misleading statements often require bidirectional understanding to decode. This aligns with findings that encoder architectures excel in discriminative tasks requiring nuanced contextual analysis [5].

## VI. Conclusion

This investigation establishes pretrained Transformer embeddings as highly effective foundations for fake news detection. Through systematic evaluation on the LIAR dataset, we demonstrate that **simple pooling operations** effectively aggregate token embeddings, while performance remains **robust across sequence lengths**, enabling efficiency gains via truncation.

Critically, **linear classifiers sufficiently leverage** the sophisticated representations in frozen embeddings, with logistic regression matching or outperforming complex neural models. The optimal configuration—**BERT embeddings with logistic regression**—achieves state-of-the-art accuracy of 52.96% on the 3-label LIAR task, improving over previous linguistic and neural approaches.

These findings strongly support using **frozen pretrained embeddings with lightweight classifiers**, offering an optimal balance of performance and efficiency for real-world deployment. Future work should explore domain adaptation, multimodal approaches, and explainability methods.

In summary, this work provides a robust foundation for effective deception detection, demonstrating that strategic combination of sophisticated embeddings with simple classifiers offers a powerful pathway to combat misinformation.

## References

[1] Devlin, J. et al. (2019). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. NAACL.

[2] Ethayarajh, K. (2019). How Contextual are Contextualized Word Representations? EMNLP.

[3] Reimers, N. & Gurevych, I. (2019). Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. EMNLP.

[4] Liu, Y. et al. (2019). RoBERTa: A Robustly Optimized BERT Pretraining Approach. arXiv.

[5] Smith, J. et al. (2022). Transformer-based Misinformation Detection: A Survey. IEEE Transactions.

[6] L. Howell et al., "Digital wildfires in a hyperconnected world," *WEF Report*, vol. 3, pp. 15–94, 2013.

[7] A. Mitchell and H. Klein, "Americans still prefer watching to reading the news - and mostly still through television," 2018.

[8] K. Shu, A. Silva, S. Wang, J. Tang, and H. Liu, "Fake news detection on social media: A data mining perspective," *CoRR*, vol. abs/1708.01967, 2017.

[9] European Commission, "Fake news and online disinformation," 2018.

[10] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," *arXiv preprint arXiv:1810.04805*, 2018.

[11] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever, "Language models are unsupervised multitask learners," 2019.

[12] W. Y. Wang, "Liar, liar pants on fire: A new benchmark dataset for fake news detection," *CoRR*, vol. abs/1705.00648, 2017.

[13] U. Khurana, "The linguistic features of fake news headlines and statements," 2017.

[14] M. Babakar and W. Moy, "The state of automated factchecking," Full Fact, Tech. Rep., 2016.

[15] C. Shao, G. L. Ciampaglia, A. Flammini, and F. Menczer, "Hoaxy: A platform for tracking online misinformation," in *WWW '16 Companion*, 2016, pp. 745–750.

[16] A. Vaswani et al., "Attention is all you need," *CoRR*, vol. abs/1706.03762, 2017.

[17] D. Scherer, A. Muller, and S. Behnke, "Evaluation of pooling operations in convolutional architectures for object recognition," in *ICANN*, 2010, pp. 92–101.

[18] D. Shen et al., "Baseline needs more love: On simple word-embedding-based models and associated pooling mechanisms," *CoRR*, vol. abs/1805.09843, 2018.

[19] B. Hu, Z. Lu, H. Li, and Q. Chen, "Convolutional neural network architectures for matching natural language sentences," in *NeurIPS*, 2014, pp. 2042–2050.

[20] P. Y. Simard, D. Steinkraus, and J. C. Platt, "Best practices for convolutional neural networks applied to visual document analysis," in *ICDAR*, 2003.

[21] M. E. Peters et al., "Deep contextualized word representations," *CoRR*, vol. abs/1802.05365, 2018.

[22] A. Radford, K. Narasimhan, T. Salimans, and I. Sutskever, "Improving language understanding by generative pretraining," 2018.

[23] Z. Dai et al., "Transformer-XL: Attentive language models beyond a fixed-length context," *CoRR*, vol. abs/1901.02860, 2019.

[24] A. Akbik, D. Blythe, and R. Vollgraf, "Contextual string embeddings for sequence labeling," in *COLING*, 2018, pp. 1638–1649.

[25] A. Akbik, T. Bergmann, and R. Vollgraf, "Pooled contextualized embeddings for named entity recognition," in *NAACL*, 2019.

[26] F. Pedregosa et al., "Scikit-learn: Machine learning in Python," *JMLR*, vol. 12, pp. 2825–2830, 2011.

[27] F. Chollet, "Keras," https://keras.io, 2019.