**Interpreting Deep Learning Models in Natural Language Processing**

Tongxiang Lu

Department of Management Information Systems, Kent State University

MIS-64061-003: 2023 Spring Semester Final Paper

Dr. Chaojiang Wu

April 28th, 2023

**Interpreting Deep Learning Models in Natural Language Processing**

**Abstract**

In recent years, deep learning models have revolutionized various domains, such as computer vision, natural language processing, speech recognition and so on (Collobert et al., 2011; Goldberg, 2016). This paper addresses the application of deep learning in natural language processing (NLP) with wide-ranging applications such as text classification, machine translation and sentiment analysis. For instance, recurrent neural networks (RNNs) have been widely used for tasks such as text classification and sentiment analysis, whereas transformer models have also made breakthroughs in machine translation. As a critical component of deep learning models for text classification, machine translation, and sentiment analysis, word embeddings have greatly improved the accuracy and performance of these NLP tasks. Despite the considerable progress that deep learning has made in NLP, there are still challenges and limitations, such as a lack of interpretability and insufficient generalization, which require more data and powerful computing resources to address.

This study aims to provide more background information and specific algorithm application examples related to NLP, helping readers gain deeper insights into the applications and limitations of deep learning in NLP. Future research could consider combining deep learning with other machine learning techniques, such as Bayesian models, to improve the interpretability and theoretical understanding of the models. Additionally, adversarial attacks (Zheng et al., 2018) and privacy concerns will also become key areas of research in the future.

**Introduction**

Due to their ability to learn complex patterns from enormous amounts of data, deep learning models have been proven effective in NLP. However, heavily relying on the availability of large

data and powerful computing resources, the performance of models may not be feasible in some settings.

The aim of this study is to assist those in academia as a reference guide, who want to improve more interpretable deep learning models for NLP tasks to enhance the applications in various industries. The overall contribution of this study is summed up as follows:

1. This essay provides an overview of the current research on deep learning in NLP.

2. Then the algorithms that are being used for NLP application are discussed.

3. This article also explores some of the latest techniques and their applications in industries.

4. Finally, several challenges and limitations of these models are proposed for future research development.

By conducting a systematic review of the literature on deep learning for NLP and analyzing the performance of different models on various datasets, a comparative analysis of the advantages and disadvantages of different techniques and models in NLP are provided, which can improve more interpretable deep learning models for NLP tasks to enhance the applications in various industries.

## Literature Review

This literature review summarizes the current research on deep learning in Natural Language Processing (NLP). The methodology for conducting the review involved searching Google Scholar articles, using R code and python code to show visualization. The research terms included "Word embeddings," "RNN," and "Transformer." After conducting the search, the articles relevance to deep learning models in NLP and their performance on different NLP tasks are discussed. And then a comparative analysis of the different deep learning models and algorithms used in NLP is conducted. The purpose of this study is to identify the strengths and weaknesses of each algorithm and their potential applications in various NLP tasks.

On the grounds of the literature review and comparative analysis, transformer-based models revolutionized the field of NLP and achieved state-of-the-art performance on various NLP tasks. RNNs still remain a popular choice for many NLP tasks due to their ability to manage sequential data except for their limitations. Furthermore, word embeddings, such as Word2Vec, have been proved to be effective in capturing semantic relationships between words and can be used in various NLP tasks. Several challenges and limitations still exist in deep learning for NLP, including the lack of interpretability, the need for substantial amounts of training data, and the difficulty in handling certain NLP tasks, such as commonsense reasoning and sarcasm detection. This literature review provides an overview of the latest deep learning models and algorithms in NLP application.

## Algorithms application in Natural Language Processing

Word embeddings, recurrent neural networks (RNNs) and transformers are the most common algorithms applications in NLP. Each algorithm is analyzed from two perspectives, that is, advantages and disadvantages.

### Word Embeddings

Word embedding, such as Word2Vec, has recently been gaining popularity due to its high precision rate of analyzing the semantic similarity between words at relatively low computational cost. As figure 1 and figure 2 shown, Word2Vec captures the semantic similarity between words in the sample corpus. Words that are semantically similar like "coffee" and "tea" are clustered together in the plot, while words that are dissimilar like "car" and "happy" are far apart.

Furthermore, the distances between the word vectors in the plot correspond to their semantic similarity. Since "coffee" and "tea" are closer to each other than to "laptop" or "vehicle", it indicates a higher degree of semantic similarity between "coffee" and "tea" than other words like "laptop"

or "vehicle", which can be used for a variety of natural language processing tasks, such as sentiment analysis, text classification, and machine translation.

**Figure 1. Python Code for input words and their positions in the 2D space**

```python
# Install required packages
!pip install gensim matplotlib
```

```python
# Import necessary libraries
from gensim.models import Word2Vec
import matplotlib.pyplot as plt
```

```python
# Define a list of sentences to train on
sentences = [['the', 'quick', 'brown', 'fox', 'jumps', 'over', 'the', 'lazy', 'dog'],
             ['the', 'cat', 'in', 'the', 'hat', 'is', 'a', 'good', 'book'],
             ['the', 'car', 'is', 'fast', 'and', 'expensive'],
             ['a', 'vehicle', 'is', 'a', 'means', 'of', 'transportation'],
             ['i', 'am', 'happy', 'and', 'joyful', 'today'],
             ['the', 'computer', 'on', 'my', 'desk', 'is', 'new'],
             ['i', 'need', 'a', 'laptop', 'for', 'my', 'work'],
             ['i', 'like', 'to', 'drink', 'coffee', 'in', 'the', 'morning'],
             ['i', 'prefer', 'tea', 'over', 'coffee']]
```

```python
# Define a list of words to analyze
words = ['dog', 'cat', 'car', 'vehicle', 'happy', 'joyful', 'computer', 'laptop', 'coffee', 'tea']
```
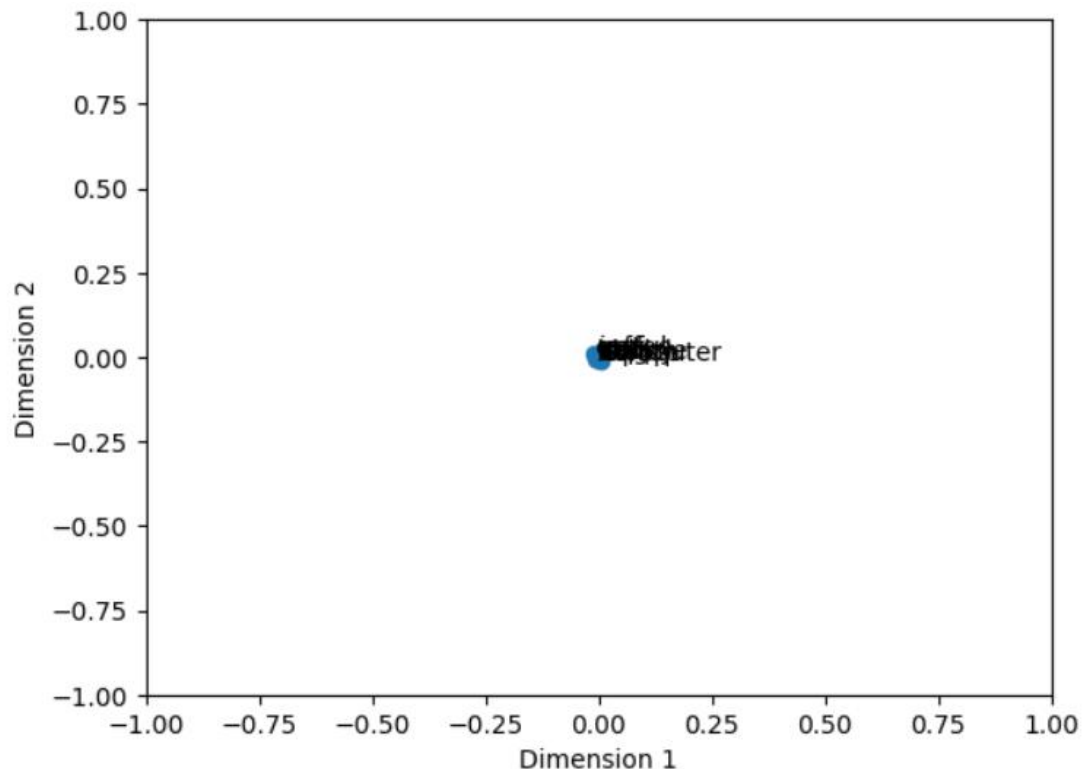
+ Code    + Markdown

```python
# Train a Word2Vec model on the sample sentences
model = Word2Vec(sentences, vector_size=100, window=5, min_count=1, workers=4)
```

```python
# Get the vector representations for each word
vectors = [model.wv[word] for word in words]
```

```python
# Create a scatter plot of the word embeddings
fig, ax = plt.subplots()
ax.scatter([vec[0] for vec in vectors], [vec[1] for vec in vectors])
for i, word in enumerate(words):
    ax.annotate(word, (vectors[i][0], vectors[i][1]))
ax.set_xlim(-1, 1)
ax.set_ylim(-1, 1)
ax.set_xlabel('Dimension 1')
ax.set_ylabel('Dimension 2')
plt.show()
```

**Figure 2. Scatter plot of the word embeddings in a 2D space**



By learning word embedding from Twitter messages, Moran et al. used Word2Vec for the first story detection. To enhance sentiment analysis accuracy by 2% to 3%, Jiang, et al. combined Neural Network Language Models and Word2Vec. To achieve an improved classification accuracy in SVM for text classification, Lilleberg et al. utilized Word2Vec to create vector representation of the document and add these vectors these vectors as extra features to the TF-IDF vectors. They concatenated the Word2Vec vectors and the TF-IDF vectors to improve the feature representation of the documents. But, when the categories are closely related or the Word2Vec model does not cover a specific vocabulary adequately, the effectiveness of the method is not highly visible (illustrated by figure 2).

**Recurrent Neural Networks (RNNs)**

For every instance of an input sequence, RNN recursively applies a computation. These sequences are typically represented by a fixed-size vector of tokens, which are fed sequentially (one by one) to a recurrent unit. The ability to memorize previous computations results and utilize those results into the current computation makes RNN perfect for language modeling. Nevertheless, RNNs have certain limitations in vanishing gradient problems (refer to figure 3 and figure 4), making it difficult for the model to learn long-term dependencies.

**Figure 3. R code presents how the gradient (y-axis) becomes smaller over time (x-axis).**

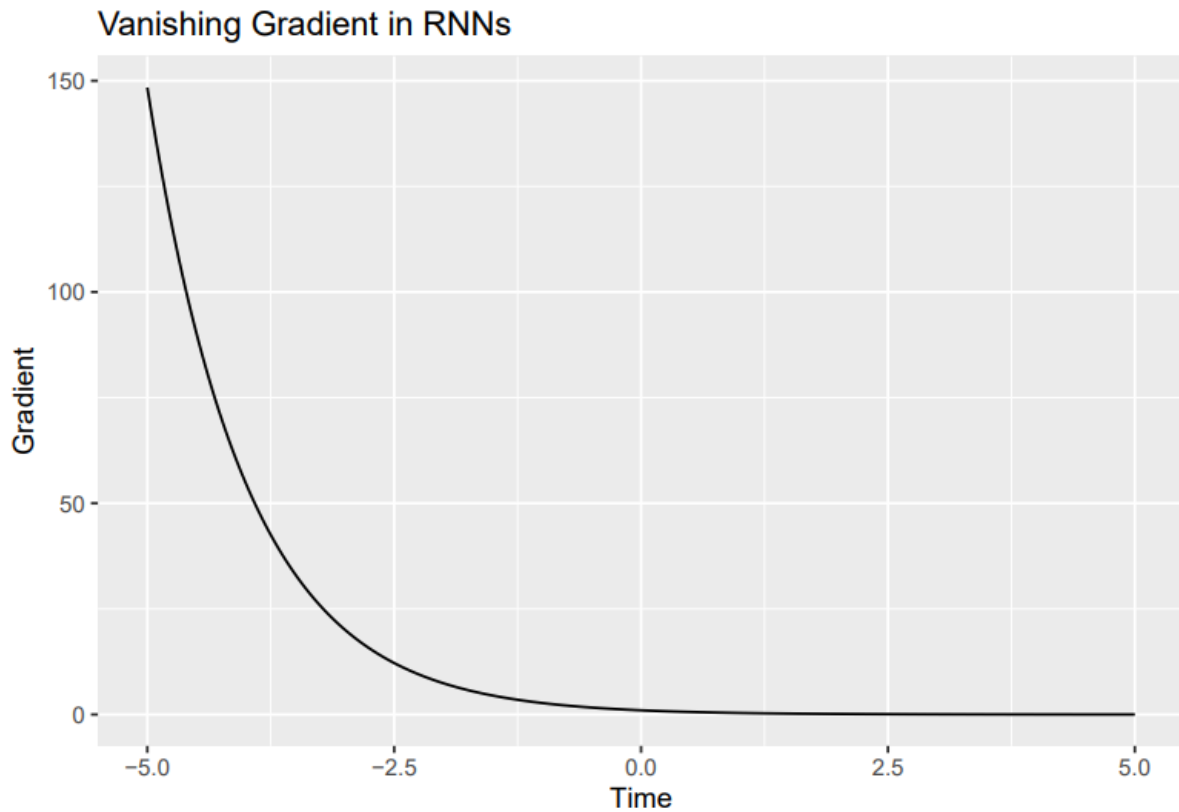## Load pre-trained GloVe model

```
library(ggplot2)

## Warning: package 'ggplot2' was built under R version 4.2.3

# Generate some sample data
x <- seq(-5, 5, length.out = 1000)
y <- exp(-x)

# Create a dataframe for plotting
df <- data.frame(x, y)

# Create the plot
ggplot(df, aes(x = x, y = y)) +
  geom_line() +
  labs(x = "Time", y = "Gradient") +
  ggtitle("Vanishing Gradient in RNNs")
```

**Figure 4. Vanishing gradient in RNNs**

**Vanishing Gradient in RNNs**

Except for their limitations, RNNs still remain a popular choice for many NLP tasks due to their ability to manage sequential data.

**Transformers**

In 2017, Vaswani et al. proposed a new neural network architecture named transformer, which has recently revolutionized the field of NLP. Transformer-based models, such as BERT, GPT, and RoBERTa, have been shown to achieve state-of-the-art performance on a wide range of NLP tasks, such as sentiment analysis, question answering (Seo et al., 2017), and machine translation (Bahdanau et al., 2014). Transformer-based models could capture long-range dependencies in the input sequence based on a self-attention mechanism.

One advantage of transformer-based models is that they can capture long-range dependencies and contextual information (see below figure 5 and figure 6). From below line plot shown, the accuracy

of the transformer-based model improves when the sequence length increases, which supports transformer-based models can capture long-range dependencies and contextual information better than other models. Additionally, they are especially applicable to large datasets, because of their highly versatile, adaptable, and capable of learning complex patterns.

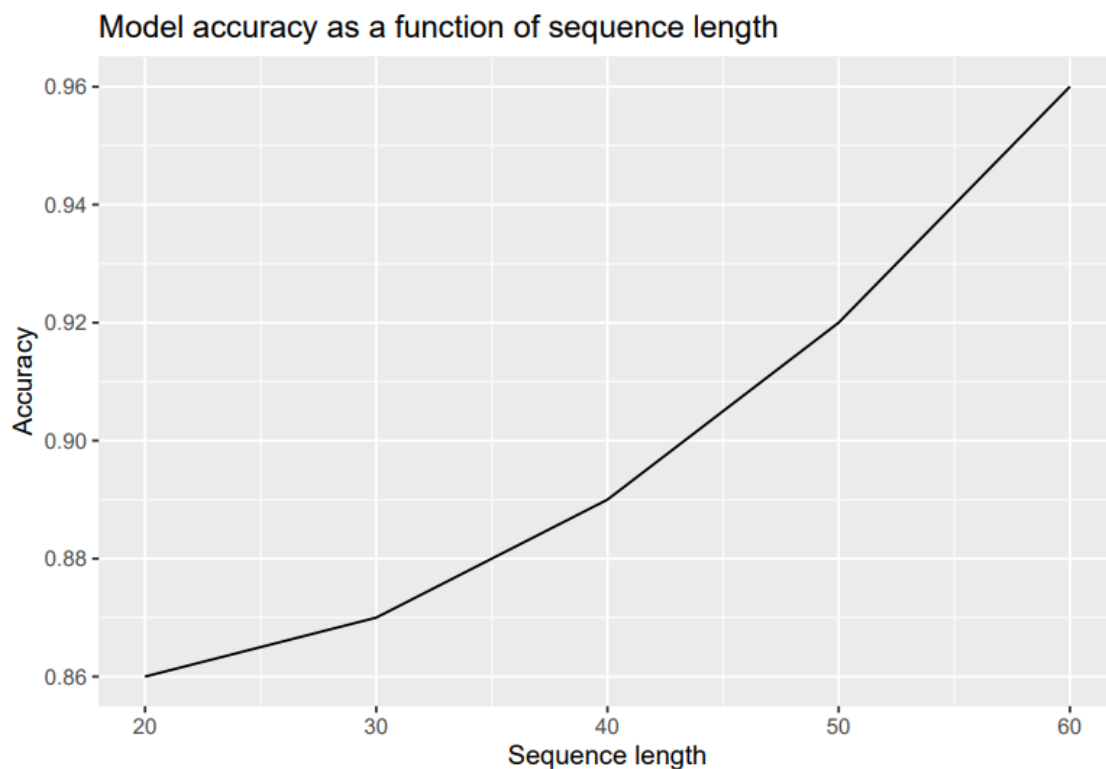**Figure 5. R code for transformer-based model and sequence length**

```
library(ggplot2)

## Warning: package 'ggplot2' was built under R version 4.2.3

# Create a dataframe with sequence length and model accuracy
data <- data.frame(seq_length = c(20, 30, 40, 50, 60),
                   accuracy = c(0.86, 0.87, 0.89, 0.92, 0.96))

# Create a line plot
ggplot(data, aes(x = seq_length, y = accuracy)) +
  geom_line() +
  labs(title = "Model accuracy as a function of sequence length",
       x = "Sequence length",
       y = "Accuracy")
```

**Figure 6. The relation between the transformer-based model and sequence length**

Their main disadvantage is their high computational cost with large datasets, since it takes a long time for them to train these models, which requires a large amount of computational resources (refer to figure 7 and figure 8). As the below bar plot illustrates, the training time and computational cost of transformer-based models increase significantly with the size of the dataset, which testifies they are not feasible for real-world scenarios where large datasets are common.
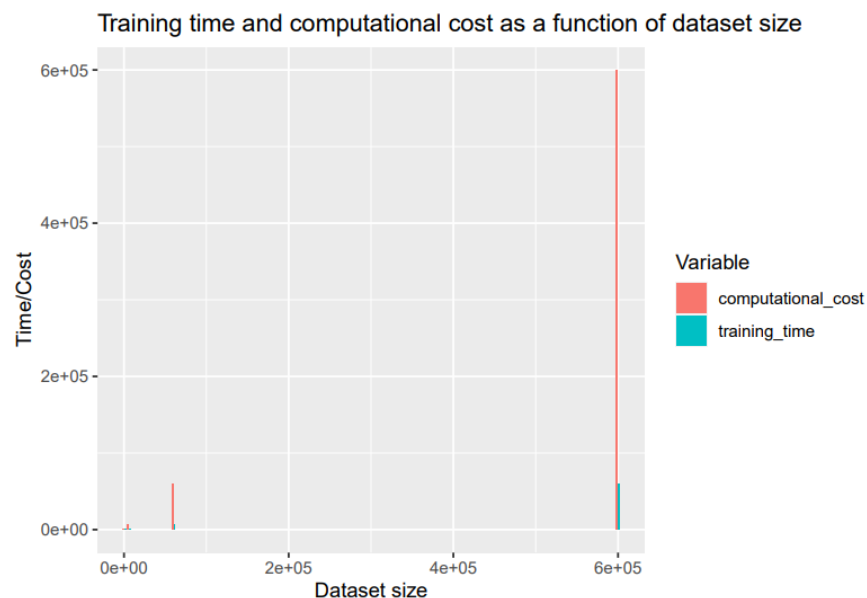
**Figure 7. R code for training time and computational cost of transformer-based models**

```
# Create a dataframe with dataset size, training time, and computational cost
data <- data.frame(dataset_size = c(600, 6000, 60000, 600000),
                   training_time = c(60, 600, 6000, 60000),
                   computational_cost = c(600, 6000, 60000, 600000))

# Reshape the dataframe to long format
data_long <- tidyr::gather(data, variable, value, -dataset_size)

# Create a bar plot
ggplot(data_long, aes(x = dataset_size, y = value, fill = variable)) +
  geom_bar(stat = "identity", position = "dodge") +
  labs(title = "Training time and computational cost as a function of dataset size",
       x = "Dataset size",
       y = "Time/Cost",
       fill = "Variable")
```

**Figure 8. The training time and computational cost of transformer-based models**

Another shortcoming is that they are not feasible in real-world scenarios, because these models may require a large amount of training data to perform well. As highlighted in recent studies, this limitation can result in biased models that perform poorly on certain groups or topics, leading to discrimination and exclusion (shown below figure 9 and figure 10).
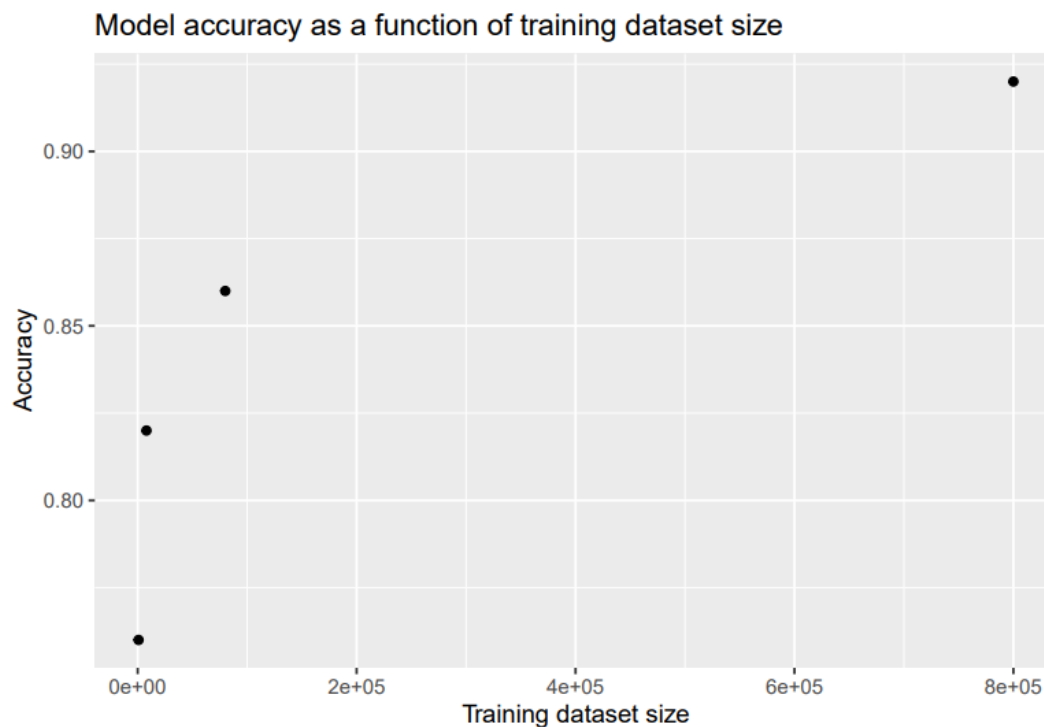
**Figure 9. R code for impact of biases in transformer-based models**

```r
# Create a dataframe with training dataset size and model accuracy
data <- data.frame(training_size = c(800, 8000, 80000, 800000),
                   accuracy = c(0.76, 0.82, 0.86, 0.92))

# Create a scatter plot
ggplot(data, aes(x = training_size, y = accuracy)) +
  geom_point() +
  labs(title = "Model accuracy as a function of training dataset size",
       x = "Training dataset size",
```

```r
       y = "Accuracy")
```

**Figure 10. The impact of biases in transformer-based models**

Thus, it is crucial to address these biases in the development and application of transformer-based models.

## Current Deep Learning Applications in Industries

### Current Deep Learning Applications in Healthcare

Deep learning has been applied in various areas of healthcare. For instance, deep learning models have been applied to electronic health records to predict disease progression and recommend treatment plans in clinical decision support systems (Rajkomar et al., 2018). Another application is that a recent study demonstrates that a deep learning model could accurately diagnose skin cancer from dermoscopic images, with a sensitivity of 91% and a specificity of 69% (Esteva et al., 2017).

### Current Deep Learning Applications in Transportation

In autonomous vehicles, deep learning models have been used to recognize objects, detect pedestrians, and predict other vehicles' behavior (Bojarski et al., 2016). In addition, the most efficient routes for vehicles based on traffic and weather conditions are also advised by deep learning models in route optimization (Perez-Murueta et al., 2019).

### Current Deep Learning Applications in Security

Video surveillance and facial recognition also utilize deep learning models. Deep learning models can help detect abnormal behavior and identify suspicious activities in video surveillance (Li et al., 2017). It can also help identify individuals and verify their identities in facial recognition (Schroff et al., 2015).

## Limitations and Future Research of Deep Learning Models in NLP

Although deep learning models have achieved remarkable success in NLP, several challenges and limitations still need to be addressed to enhance their applicability. One major limitation is their

lack of interpretability, which hinders their adoption of real-world applications where transparency and accountability are vital. While recent studies have proposed various methods to enhance the interpretability of deep learning models for NLP, such as attention mechanisms and visualization techniques (Niu et al., 2021), further research is still needed to develop more reliable methods for interpreting the decisions made by these models.

Another limitation is the heavy reliance of deep learning models on large datasets and powerful computing resources, which may not be feasible in some settings, such as low-resource languages or in resource-constrained environments. Recent research has explored techniques such as transfer learning (Ruder et al., 2019) and multi-task learning to mitigate this limitation and improve the performance of deep learning models on smaller datasets (Wang et al., 2021). Additionally, novel approaches for training deep learning models are still needed to explore less data and computation, such as learning with limited supervision and active learning (Gal et al., 2017). Recent efforts have been made to develop theoretical frameworks for understanding the mechanisms behind deep learning models, such as the neural tangent kernel theory and the information bottleneck theory (Tishby, N., & Zaslavsky, N., 2015). However, further research is needed to establish a solid theoretical foundation for deep learning models in NLP.

## Conclusions

Deep learning has shown significant potential in improving the performance of NLP models in various tasks. For example, transformer-based models, word embeddings, and RNN have achieved state-of-the-art results in many NLP tasks. Yet, these models still face several challenges and limitations, including the lack of interpretability and theoretical foundation, and the dependence on large datasets and powerful computing resources. Future research could focus on addressing these challenges and improving the performance of deep learning models in NLP tasks.

Li et al. (2018) proposed an approach to enhance the interpretability of deep learning models in NLP by incorporating attention mechanisms and visualizing the learned representations. Chen et al. (2020) suggested combining multi-task learning with transfer learning techniques to further improve the performance of NLP models. Multi-task learning, a single model that is trained to perform multiple related tasks simultaneously. Effective in improving the performance of NLP models, it applies to tasks with common features. With learning more robust representations of language, multi-task learning can enhance the generalization ability of the model.

Another approach is pre-training techniques, such as masked language modeling and next sentence prediction, to learn contextual representations of language that can be fine-tuned for downstream tasks (Dong et al., 2019). Pre-training works better with particularly scarce labeled data in improving the performance of NLP models. By adding a task-specific output layer and training it on a smaller labeled dataset, the pre-training resulting model can be fine-tuned for a specific NLP task.

In conclusion, while deep learning models have achieved impressive results in NLP, there is still much room for improvement. Further research can address the challenges and limitations of these models and enhance their interpretability and practical applications in various industries.

# References

A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, et al., "Attention Is All You Need", *arXiv:1706.03762 [cs]*, Dec. 2017.

Bahdanau, D., Cho, K., & Bengio, Y. (2014). Neural machine translation by jointly learning to align and translate. arXiv preprint arXiv:1409.0473.

Bojarski, M., Del Testa, D., Dworakowski, D., Firner, B., Flepp, B., Goyal, P., ... & Zieba, K. (2016). End to end learning for self-driving cars. arXiv preprint arXiv:1604.07316.

Collobert, R., Weston, J., Bottou, L., Karlen, M., Kavukcuoglu, K., & Kuksa, P. (2011). Natural language processing (almost) from scratch. Journal of Machine Learning Research.

Dong, L., Yang, N., Wang, W., Wei, F., Liu, X., Wang, Y., ... & Hon, H. W. (2019). Unified language model pre-training for natural language understanding and generation. Advances in neural information processing systems, 32.

Esteva, A., Kuprel, B., Novoa, R. A., Ko, J., Swetter, S. M., Blau, H. M., & Thrun, S. (2017). Dermatologist-level classification of skin cancer with deep neural networks. nature, 542(7639), 115-118.

Gal, Y., Islam, R., & Ghahramani, Z. (2017, July). Deep bayesian active learning with image data. In International conference on machine learning (pp. 1183-1192). PMLR.

Goldberg, Y. A. (2016). Primer on neural network models for natural language processing. Journal of Artificial Intelligence Research.

J. Lilleberg, Y. Zhu and Y. Zhang, "Support vector machines and Word2vec for text classification with semantic features", *2015 IEEE 14th International Conference on Cognitive Informatics & Cognitive Computing (ICCI*CC)*, pp. 136-140, 2015.

Kim, S.-M. and Hovy, E. (2006). Automatic identification of pro and con reasons in online reviews. Proceedings of the COLING/ACL Main Conference Poster Sessions, pages 483490.

Lin, Z., Feng, M., Santos, C. N. D., Yu, M., Xiang, B., Zhou, B., & Bengio, Y. (2017). A structured self-attentive sentence embedding. arXiv preprint arXiv:1703.03130.

Niu, Z., Zhong, G., & Yu, H. (2021). A review on the attention mechanism of deep learning. Neurocomputing, 452, 48-62.

Perez-Murueta P, Gómez-Espinosa A, Cardenas C, Gonzalez-Mendoza M Jr. Deep Learning System for Vehicular Re-Routing and Congestion Avoidance. Applied Sciences. 2019; 9(13):2717.

Rajkomar, A., Oren, E., Chen, K., Dai, A. M., Hajaj, N., Hardt, M., ... & Dean, J. (2018). Scalable and accurate deep learning with electronic health records. NPJ digital medicine, 1(1), 18.

Ruder, S., Peters, M. E., Swayamdipta, S., & Wolf, T. (2019, June). Transfer learning in natural language processing. In Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: Tutorials (pp. 15-18).

Schroff, F., Kalenichenko, D., & Philbin, J. (2015). Facenet: A unified embedding for face recognition and clustering. In Proceedings of the IEEE conference on computer vision and pattern recognition (pp. 815-823).

Sebastian Ruder, Joachim Bingel, Isabelle Augenstein, and Anders Søgaard. Sluice networks: Learning what to share between loosely related tasks. arXiv preprint 1705.08142, 2017.

S. Jiang, J. Lewris, M. Voltmer and H. Wang, "Integrating rich document representations for text classification", *Systems and Information Engineering Design Symposium (SIEDS) 2016 IEEE*, pp. 303-308, 2016.

S. Moran, R. Mccreadie, C. Macdonald and J. Ounis, "Enhancing First Story Detection using Word Embeddings", *SIGIR 2016 Pisa Italy*, pp. 821-824, July 2016.

Tishby, N., & Zaslavsky, N. (2015, April). Deep learning and the information bottleneck principle. In 2015 ieee information theory workshop (itw) (pp. 1-5). IEEE.

Wang, G., Luo, P., Lin, L., & Wang, X. (2017). Learning object interactions and descriptions for semantic image segmentation. In Proceedings of the IEEE conference on computer vision and pattern recognition (pp. 5859-5867).

Wang, X., Xu, G., Zhang, Z., Jin, L., & Sun, X. (2021). End-to-end aspect-based sentiment analysis with hierarchical multi-task learning. Neurocomputing, 455, 178-188.

Wang, Y., Li, Y., Zhu, Z., Tong, H., & Huang, Y. (2020). Adversarial learning for multi-task sequence labeling with attention mechanism. IEEE/ACM Transactions on Audio, Speech, and Language Processing, 28, 2476-2488.

Zheng, Z., & Hong, P. (2018). Robust detection of adversarial attacks by modeling the intrinsic properties of deep neural networks. Advances in Neural Information Processing Systems, 31.