

# Predicting the Category of a News Article from its Headline & Description

Abhinandan Padhi

apadhi@ucsd.edu

## 1 Introduction

In this study, we have worked on the task of classification of online news articles into various news categories using the headline and short description (or abstract) of the article. The dataset for this task was obtained from Kaggle (Misra, 2022), and contains nearly 210,000 news headlines as well as short descriptions (and some metadata) from 2012 to 2022, taken from news articles published by HuffPost. We were motivated to pursue this topic as news article classification/categorization could be quite useful and interesting, especially given that the dataset used has 41 different news categories, which adds a bit of challenge.

Besides, the task chosen appears to have some practical significance as it could be used by online news publishers to automatically assign a label or tag to articles just based on the headline and short description.

The following is a brief summary of the steps we took for this study:

- We obtained the dataset from Kaggle. The dataset was pre-processed and split into training, validation, and test sets.
- Before implementing our own models, we researched and reviewed the models used by prior works for the text classification.
- We implemented Recurrent Neural Network (RNN) and Long-Short Term Memory (LSTM) networks as baselines.
- We conducted experiments with BERT and DistilBERT, which are transformer-based models. These experiments accounted for the bulk of the time spent on this project, especially due to the immense compute requirements for these models.

## 2 Related work

### 2.1 Text Classification Algorithms: A Survey

The paper “Text Classification Algorithms: A Survey” (Kowsari, 2019) (Machine Learning on Scientific Data and Information, 2019) reviews a variety of text classification algorithms, categorizing them into traditional machine learning methods (e.g., Naïve Bayes, Support Vector Machines, and Logistic Regression), ensemble techniques (e.g., Boosting and Bagging), and deep learning models (e.g., RNNs, LSTMs, and Transformers). They observed that traditional ML algorithms rely on statistical approaches and are computationally efficient, but usually underperform on complex text relationships compared to deep learning models. Ensemble methods improve predictive performance by aggregating weaker models but are computationally intensive. Deep learning approaches leverage neural networks to model complex relationships in text data, with Transformers like BERT leading in state-of-the-art performance due to their contextual word representation capabilities.

Their results highlight that deep learning models, particularly LSTMs and Transformers, outperform traditional methods like RNNs by better capturing sequential and contextual information in text. BERT stands out with superior results due to its bidirectional architecture and contextualized embeddings. Traditional methods such as Naïve Bayes and SVMs still provide competitive baselines for simpler problems or when computational resources are limited. Ensemble methods, while powerful, require significant computational overhead and are less interpretable than standalone models.

Based on the survey of text classification algorithms by Kowsari et al., it was decided that this project would focus on: Recurrent Neu-

ral Networks (RNNs), Long Short-Term Memory (LSTMs), and Transformers like BERT and DistilBERT. These models are all reviewed below.

## 2.2 Recurrent Neural Networks (RNNs)

The Recurrent Neural Networks (RNN) is a classical neural network architecture for tasks like text classification. While research on RNN can be credited to many authors, one of the most influential papers on the RNN was written by Elman et al. (of UC San Diego) in 1990 (Elman, 1990). Essentially, the RNN is characterized by so-called “recurrent” connections that enable state propagation, where hidden states are recursively updated through a transformation function that processes each input in sequence. The main idea behind the RNN architecture is that makes it possible to “share weights across time steps”, allowing the network to handle variable-length sequences by maintaining a compressed representation of historical context.

However, the performance of classical RNNs is limited due to the problem of vanishing and exploding gradients during backpropagation, which severely limits their ability to capture long-term dependencies. Hence, in recent years, it has proved to be challenging for RNNs to achieve good performance on complex NLP tasks. Nevertheless, the relative ease of implementation and low computational costs make RNNs a suitable baseline for our task of news classification based on text features.

## 2.3 Long Short-Term Memory (LSTMs)

Long Short-Term Memory models or LSTMs (Sepp Hochreiter, 1997) have historically been immensely influential in the field of NLP. The LSTM network is effectively an evolution of Recurrent Neural Networks (RNNs). The issue with standard RNNs is that they struggle to preserve information over long sequences due to the problem of “exploding and vanishing gradients” during backpropagation. Hence, Hochreiter et al.’s LSTM architecture features a special type of RNN cell with a more sophisticated memory mechanism called a “gating mechanism” that can selectively remember or forget information.

An LSTM cell features three gates: the forget gate, input gate, and output gate. These gates provide fine-grained control over information flow. The input gate decides what new information to store, forget gate determines what infor-

mation to discard from the cell state, and output gate controls which part of the cell state is output. While the original paper was not specifically focused on text classification tasks, the architecture’s ability to capture long-range dependencies makes LSTMs particularly powerful for sequence-based tasks like text classification. LSTMs can learn meaningful representations of text that, at a certain time step, consider context from much further back in the sequence compared to traditional RNNs. Hence, LSTMs are well-suited for tasks like sentiment analysis, document categorization, and in our case, the task of news category classification.

## 2.4 Bidirectional Encoder Representations from Transformers (BERT)

Bidirectional Encoder Representations from Transformers, better known as BERT (Jacob Devlin, 2018), is a powerful language representation model designed to pre-train deep bi-directional representations. Unlike prior models, BERT employs a Masked Language Model (MLM) to learn context from both left and right directions in the text input, and incorporates a next sentence prediction task for text-pair understanding. Fine-tuning BERT for specific tasks involves minimal modifications, making it versatile for a broad range of applications in Natural Language Processing.

For general text classification tasks, BERT significantly outperformed previous methods. According to the original paper, it achieved a GLUE score of 82.1 with its larger configuration (BERT-Large), setting new benchmarks across tasks like natural language inference (MNLI), sentiment analysis (SST-2), and paraphrase detection (QQP). Such results indicate that BERT truly is powerful and versatile in handling classification tasks.

## 2.5 DistilBERT

DistilBERT (Victor Sanh, 2019) is a compact and efficient version of BERT, designed to retain most of the BERT’s performance while reducing size and computational requirements. It employs knowledge distillation during pre-training, using a so-called “triple loss approach”, which consists of combining language modeling, distillation, and cosine-distance losses to preserve the inductive biases of the larger teacher model. By reducing the number of layers by half, DistilBERT achieves a 40% reduction in size and a 60% increase in infer-

ence speed, with only a slight performance drop of 3%. This efficiency makes it suitable when the compute capacity and / or time available for training is limited.

In text classification tasks, DistilBERT performs surprisingly well, achieving 97% of BERT’s score on the GLUE benchmark (i.e. for many Natural Language Understanding tasks). Even on specific benchmarks like IMDB sentiment classification, the accuracy gap between BERT and DistilBERT appears to be negligible. Hence, DistilBERT has proven to be capable of handling text classification tasks effectively, while being highly efficient in terms of compute resources. Thus, we decided to fine-tune DistilBERT for this project, and compare it to BERT.

### 3 Dataset

The dataset (Misra, 2022) for this project was obtained from Kaggle, which is an online AI competitions and datasets platform. The dataset was compiled by by Risabh Misra in 2022, and is available for public use under the Attribution 4.0 International Creative Commons license. According to dataset creator, the data was based on news articles published by HuffPost between 2012 and 2022.

We were motivated to use this dataset as news article classification/categorization could be quite useful and interesting, especially given that the dataset contains 41 different news categories. While most articles are assigned labels from the top 15 categories (e.g. Politics, Wellness, Entertainment, etc.), this is still a large number of categories, which necessitates a powerful and complex model to attain high performance. This task and dataset combination is also interesting because there is an imbalance in the dataset in terms of class labels. For instance, the news articles from the 15 least-frequent categories are more likely to be misclassified. I discuss these issues in more detail in the sub-sections below.

#### 3.1 Data Statistics and Analysis

This dataset contains 209,527 rows and 6 columns: headline, short description, category, authors, date, and link. Of these features, we are considering the headline and short description as the main features; the category column is our target feature; and the authors, date, and link are metadata, which have not used for our predictive task for news classification. Figure 1 shows the main features for

five examples from the dataset. Figure 2 shows the metadata features for the same five examples.

	headline	short_description	category
0	over 4 million americans roll up sleeves for o...	over 4 million americans roll up sleeves for o...	U.S. NEWS
1	american airlines flyer charged, banned for li...	american airlines flyer charged, banned for li...	U.S. NEWS
2	23 of the funniest tweets about cats and dogs ...	23 of the funniest tweets about cats and dogs ...	COMEDY
3	the funniest tweets from parents this week (se...	the funniest tweets from parents this week (se...	PARENTING
4	woman who called cops on black bird-watcher lo...	woman who called cops on black bird-watcher lo...	U.S. NEWS

Figure 1: The main data features for five examples.

	authors	date	link
0	Carla K. Johnson, AP	2022-09-23	<a href="https://www.huffpost.com/entry/covid-boosters-...">https://www.huffpost.com/entry/covid-boosters-...</a>
1	Mary Papenfuss	2022-09-23	<a href="https://www.huffpost.com/entry/american-airlin...">https://www.huffpost.com/entry/american-airlin...</a>
2	Elyse Wanshel	2022-09-23	<a href="https://www.huffpost.com/entry/funniest-tweets...">https://www.huffpost.com/entry/funniest-tweets...</a>
3	Caroline Bologna	2022-09-23	<a href="https://www.huffpost.com/entry/funniest-parent...">https://www.huffpost.com/entry/funniest-parent...</a>
4	Nina Golgowski	2022-09-22	<a href="https://www.huffpost.com/entry/amy-cooper-lose...">https://www.huffpost.com/entry/amy-cooper-lose...</a>

Figure 2: The metadata features for the same five examples as Figure 1.

Since our models will only be taking text-based features (headline and short description) as input and the category as output, we shall explore these features in more detail below. But, first, here are some examples of the input/output pairs we shall be working with:

1. **Headline:** Biden Says U.S. Forces Would Defend Taiwan If China Invaded

**Short Description:** President issues vow as tensions with China rise.

**Category:** Politics

2. **Headline:** Over 4 Million Americans Roll Up Sleeves For Omicron-Targeted COVID Boosters

**Short Description:** Health experts said it is too early to predict whether demand would match up with the 171 million doses of the new boosters the U.S. ordered for the fall.

**Category:** U.S. News

#### 3.1.1 Distribution of Words over the Dataset

The dataset contains a total of 6,132,807 words across both text columns, where the headline column contains 2,011,615 words and the short description column contains a total of 4,121,192 words.

The histogram in Figure 3 visualizes the distribution of the number of words in the headline over the dataset’s rows. From this histogram, it is obvious that the word count in the headlines follows a normal distribution, although we did observe some

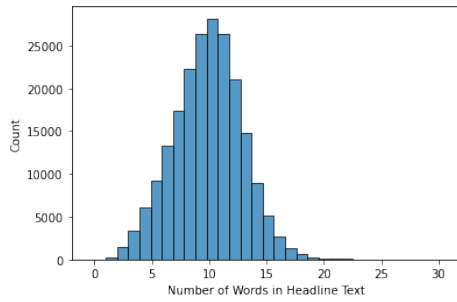


Figure 3: Histogram visualizing distribution of headline length over the dataset’s rows.

anomalies where the headline had more than 30 words, with a maximum of 44 words.

The histogram in Figure 4 visualizes the distribution of the number of words in the short description over the dataset’s rows. This distribution appears to be skewed towards the left, and similar to the headline, there were some anomalies observed where description length was over 200 words.

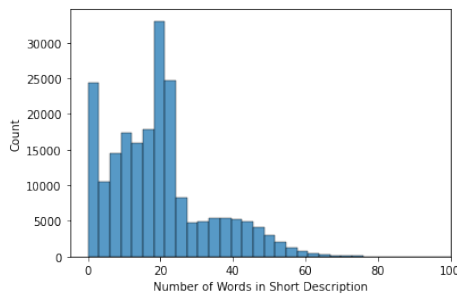


Figure 4: Histogram visualizing distribution of short description’s length over the dataset’s rows.

Figure 5 shows the count, mean, minimum, median (50%), and maximum of the number of words (length) of the headline and short description columns.

	Description Length	Headline Length
count	209527.00	209527.00
mean	19.67	9.60
std	14.15	3.07
min	0.00	0.00
50%	19.00	10.00
max	243.00	44.00

Figure 5: Histogram visualizing distribution of number of words per row (headline + description).

### 3.1.2 Categories

Figure 6 shows the number of rows (count) for the top 10 categories.

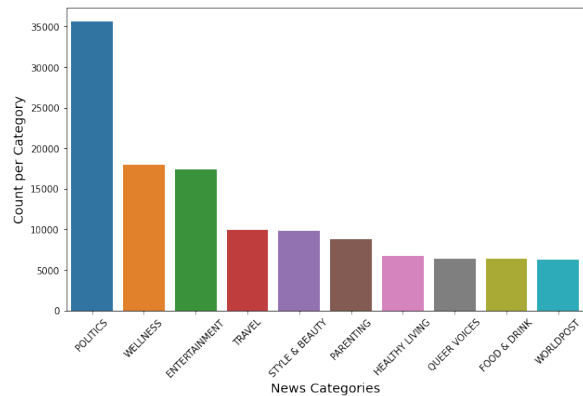


Figure 6: Number of rows for the Top 10 news categories.

The top 3 categories are: “Politics” (35,602 rows), “Wellness” (17,945 rows), and “Entertainment” (17,362 rows).

We also created a Word Cloud for further visualization/analysis of the text in our dataset, shown in Figure 7. It is interesting to see that the name “Donald Trump” appears so frequently in our dataset, and this is consistent with the fact that “Politics” is the most common category in our data.

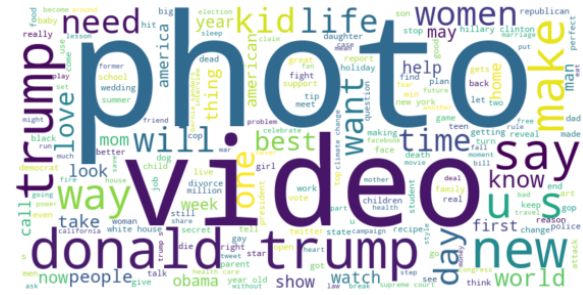


Figure 7: Word cloud showing the most frequent words in our dataset.

### 3.1.3 Challenges

From the dataset analysis above, it is clear that there is a class imbalance in the dataset. Since *politics* is the most common category, a naive model would inherently be biased towards *politics*, predicting this category more often. This would likely result in less common classes being predicted rarely (or not at all). This is especially problematic since there are 41 categories, which is

an immense number of classes for a classification task.

Moreover, I noticed it is fairly difficult to categorize some news articles that mention ideas related to multiple categories. For instance, an example article I listed earlier had the headline “Biden Says U.S. Forces Would Defend Taiwan If China Invaded”. - such an article could be categorized as both *politics* and *U.S. News*. In these situations, the model may perform poorly, with its bias towards *politics* further complicating the classification task.

### 3.2 Data Pre-Processing

The following are the steps we took to pre-process the dataset:

- There were two categories that appeared to be the same: *THEWORLDPOST* and *WORLDPOST*, so these two were merged, reducing the category count from 42 to 41.
- For some rows, the description was missing (usually for articles with short descriptions) - we handled such cases by adding the headline text to the empty description.
- Converted the two text columns to lower-case, and depending on the model being used, we did some text cleaning, such as removing non-alphanumeric symbols (including punctuation marks), stop-word removal, and lemmatization.
- Next, the text features were tokenized. The tokenizer is dependent on the model or baseline used. For example, for our BERT experiments, the tokenizer used was Hugging Face’s BertTokenizer. For some of our experiments, we also had to manually add padding at the end of token sequences to ensure uniform length for each input feature across the dataset.
- Lastly, the dataset was split into three subsets: train (80%), validation (10%), and test (10%) subsets.

## 4 Baselines

### 4.1 Recurrent Neural Network (RNN)

We briefly described the RNN architecture in the related works section. To implement and fine-tune hyperparameters efficiently for our RNN, we used

the TensorFlow framework. The following are the hyperparameters we experimented with:

1. **Embedding Dimensionality:** 16, 32, 64.
2. **Number of RNN Layers:** 1, 2, 3
3. **Neurons per RNN Layer:** 32, 64
4. **Dropout Regularization (Drop rate):** 0.1, 0.2, 0.4
5. **Gradient Descent Optimizer:** Adaptive Moment Estimation (Adam), RMSProp.
6. **Learning Rate:**  $10^{-3}$ ,  $10^{-2}$ .
7. **Number of Epochs:** 10, 20.
8. **Loss:** Categorical Cross-Entropy

### 4.2 Long Short-Term Memory (LSTM)

Based on our brief descriptions of LSTMs in our related works section, it is clear that the LSTM network is quite powerful; it also requires less compute resources than transformers. Hence, we chose LSTM to be a baseline. To implement and fine-tune hyperparameters efficiently for our LSTM, we used the TensorFlow framework. The following are the hyperparameters we experimented with:

1. **Embedding Dimensionality:** 16, 32, 64.
2. **Neurons in LSTM Layer:** 64, 128
3. **Dropout Regularization (Drop rate):** 0.1, 0.2, 0.4
4. **Gradient Descent Optimizer:** Adaptive Moment Estimation (Adam)
5. **Learning Rate:**  $10^{-3}$
6. **Number of Epochs:** 10, 20.
7. **Loss:** Categorical Cross-Entropy

Note: for both baselines, as mentioned in the dataset section, the train/val/test split was: 80/10/10. Also, as our dataset was quite large and due to the limited compute resources available, our hyperparameter tuning was not as thorough as we would have preferred. We tuned our hyperparameters based on a random search rather than a grid search.

## 5 Approach

### 5.1 Conceptual Approach

Based on our research and survey of past works in the area of text classification, we decided that the best models for this project would be BERT and DistilBERT, especially given our compute constraints, dataset size, and feasibility of implementation. Below are some specific notes for each model.

#### 5.1.1 BERT

As explained in the related works section, Bidirectional Encoder Representations from Transformers (BERT) is an extremely powerful (and popular) model for encoding-based NLP tasks (i.e. Natural Language Understanding or NLU tasks), like our task of text classification. Moreover, BERT is a pre-trained model, i.e. it has already been trained on large amounts of text data, making it ideal for fine-tuning for our specific downstream NLP task. Since our dataset involves text from news articles, which is unlikely to be overly domain-specific, we were confident that we would not need to fine-tune for as many epochs as would be necessary for more domain-specific data like for scientific, medical, or legal data.

We chose to fine-tune the “base” version (and not the “large” version) of BERT. Specifically, we used the “bert-base-uncased” variant of BERT, as it is more compact and efficient compared to the large variants, allowing us to lower the computational cost for fine-tuning.

#### 5.1.2 DistilBERT

DistilBERT is a so-called “distilled” version of the BERT base model that we have also experimented with (described above). Essentially, it is a smaller and faster version of BERT, and it was pre-trained on the same corpus as BERT base in a self-supervised fashion, using the BERT base model as a teacher.

We chose to use DistilBERT for this project as it allows us to evaluate just how well a distilled version of a contemporary model like BERT compares to its base model. Moreover, using DistilBERT allowed us to further reduce the load on our limited compute resources.

#### 5.1.3 How the Models Work

In our implementation, most of the internal working of BERT and DistilBERT is handled by the libraries used (Hugging Face transformers, Tensor-

Flow, etc.). However, since we’re discussing the conceptual approach in this subsection, we briefly describe the internal working of BERT (and DistilBERT, which is similar since it is based on BERT) below.

First, input text features are pre-processed and tokenized by BERT’s tokenizer class, which is based on the WordPiece tokenizer. It splits the words into sub-word tokens to handle out-of-vocabulary words, adds special “CLS” tokens to the beginning of each input sequence, and adds padding or truncates token sequences to a fixed length. Each token is converted into three embedded representations: token embedding, segment embedding, and position embedding, which are added up, resulting in a rich, contextual input representation. BERT uses multi-head self-attention mechanisms to capture complex contextual relationships across the entire input sequence, which is what makes it so effective.

For our specific task of news classification, we made use of transfer learning by adding a classification layer on top of the pre-trained BERT model. A simple feedforward neural network with a softmax activation was added as the classification head, which outputs the probabilities for each of the 41 news categories in our dataset.

The fine-tuning process primarily involves training the classification head; the pre-trained BERT parameters are not adjusted significantly. This allows the model to retain its existing language understanding ability while adapting to the specific nuances of the news classification task.

### 5.2 Implementation

The pre-trained models for BERT and DistilBERT were implemented mainly using the Hugging Face transformers and TensorFlow (with Keras) libraries. We also explored some variants of the model implementation in PyTorch. Regarding code submission, the code for dataset exploration, pre-processing, and all our models (including baselines) are provided in the IPython notebook called “CSE256.Project.Code.ipynb”. The IPython notebook is split into sections and subsections using Markdown and Python comments for each of the aforementioned code sections.

### 5.3 Compute

Our personal computer is an Apple MacBook Pro 16 with an M3 Max Processor with 36 gigabytes of system memory. However, our com-



Table 1: Best Results Achieved by Models

<i>Model</i>	<i>Train Accuracy</i>	<i>Test Accuracy</i>
RNN	0.5487	0.4920
LSTM	0.7089	0.6531
BERT	0.7837	0.6912
DistilBERT	0.7422	0.6753

puter lacks a compatible, dedicated accelerator like a CUDA-based Graphical Processing Unit (GPU) and has inadequate system memory for loading large datasets. Hence, we made use of the GPUs provided on online platforms like Kaggle and Google Colab.

We encountered several issues with both Kaggle and Colab. While Kaggle does provide GPUs such as Nvidia T4 and P100 GPUs to accelerate deep learning, there is a time limit on the GPUs of 30 hours per week. Besides, we can only use one GPU at a time across all our IPython notebooks. Kaggle also provides Tensor Processing Units (TPUs), but using these required waiting in long queues and we were unable to get the TPUs to work correctly. Often, the Kernel for our code notebooks would crash when using a TPU. Besides, it was tricky to determine how to adjust our regular code to actually take benefit of TPUs. Lastly, due to simultaneous work on other projects (academic and personal) that also required GPUs, the actual amount of allowed time with the GPUs was severely limited. To alleviate this issue, we used Google Colab in addition to Kaggle, which had similar compute resources available, but working on various notebooks simultaneously was tedious, and often Kaggle and Colab would terminate the runtime due to idle limits or GPU usage limits.

## 5.4 Results

To evaluate our models and baselines, we used test set accuracy. While it would be preferable to use more evaluation metrics like the F1 score, it is not particularly viable as the F1 score would have to be computed over 41 different categories, which is atypical for an evaluation metric. Nonetheless, the accuracy achieved by each of our baselines and models on both the training and test sets are shown in table 1.

A note on runtime for each model: it is difficult to say the combined time taken to tune, train, and evaluate the models due to the complicated way in

which the compute resources were being used for this project, and especially considering how Kaggle and Google Colab would automatically terminate runtime due to usage and/or idle time limits. But, the relative amount of time allocated to work with each model is clear. In ascending order: RNN, LSTM, DistilBERT, BERT.

## 6 Error analysis

While our BERT and DistilBERT models achieved reasonable accuracy on the dataset, there were several cases where they failed. Given below are three examples of erroneous test set predictions made by our models and baselines along with error analysis. Note that the categories for Examples 1 and 2 were predicted incorrectly by all 4 models (baselines and transformers). The category for Example 3 was predicted incorrectly only by the baselines, not the transformers.

1. **Headline:** Coming Out Of The Mental Health Closet

**Short Description:** This article first appeared on the blog of Intentional Insights, a nonprofit organization that empowers people.

*True Category:* Healthy Living

**Predicted Category:** Wellness

**Analysis:** the headline and description can be easily interpreted as belonging to either category of *wellness* or *healthy living*, even for a human. In such a case, it seems unfair to criticize the model for an incorrect prediction. Rather, the error occurred mainly because of how similar some of the labels in the dataset are to each other. Another example we pointed out earlier (in the dataset section) was the pair (*US News, Politics*) - since *politics* is the most common class in the dataset, and most political articles are also related to US news.

2. **Headline:** You Can Assess An Italian Restaurant Based On Its Wine

**Short Description:** Italian wine — I can assure you — is good. Very good. However, you need to know what to choose, especially at a restaurant. Outside Silvia Paoli, La Dolce Vita

**True Category:** Travel

**Predicted Category:** Food and Drink

**Analysis:** similar to the first example, predicting the news category correctly for this example is tricky even for humans, let alone an NLP model. Once again, the issue lies in the fact that the input can be classified as both *travel* and *food and drink*. This is again more of an issue with the dataset, than it is with our models. It is interesting to note, however, that all models predicted *food and drink* - this is likely because the terms “restaurant” and “wine” are present in the input, causing the models to favor the *food and drink* category.

### 3. **Headline:** Man Drove Drunk With 100 Chickens In Car, Many Dead

**Short Description:** [No description provided]

**True Category:** Weird News

**Predicted Category:** Crime

**Analysis:** this example did not include a description, only a headline. Thus, there is very little text for the models to analyze. Nonetheless, BERT / DistilBERT predicted correctly; only the baselines made an incorrect prediction, with both of them predicting *crime* as the predicted label. We believe the reason for the incorrect baseline prediction is that they were unable to detect the “weirdness” in the headline. Clearly, to us humans, it is unusual or “weird” for a man to drive with 100 chickens in their car. However, the RNN / LSTM could not account for this weirdness likely because these models are not complex enough. Instead, they may have noticed the word “dead”, and related this word to the *crime* category.

## 7 Conclusion

To conclude, this project was an excellent opportunity to explore various models for the NLP task of text classification. We were able to experiment with classical models like RNNs and LSTMs as well as state-of-the-art models like BERT and DistilBERT. We were especially impressed by how well LSTMs worked for our task and dataset, which we found surprising considering how ubiquitous transformer-based models like BERT are currently. Similarly, DistilBERT’s performance was impressive considering how it is supposed to

be a compact version of BERT - this goes to show how techniques like model distilling and pruning have immense potential for maximizing performance while minimizing compute cost and time.

Yet, despite the interesting results attained in this study, we were constantly frustrated by the compute constraints and accelerators like TPUs not working correctly. This difficulty was most apparent when fine-tuning BERT, which could easily take hours.

If we could continue working on the project in the future, we could explore how well some alternative transformer-based models might work on the same dataset and task. For instance, it would be interesting to work with RoBERTa, i.e. the model described in the paper: “A Robustly Optimized BERT Pretraining Approach”.

## 8 Acknowledgments

We would like to acknowledge and thank (Misra, 2022; Misra and Grover, 2021) for allowing us to use their dataset. We would also like to acknowledge that certain parts of the code for this project are based on the official documentation by HuggingFace, TensorFlow, Keras, and PyTorch.

## References

- Elman, J. L. (1990). Finding structure in time. *Cognitive Science* 14, 179-211.
- Jacob Devlin, Ming-Wei Chang, K. L. K. T. (2018). Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv.org*.
- Kowsari, K., J. M. K. H. M. M. S. B. L. . B. D. (2019). Text classification algorithms: A survey. *Machine Learning on Scientific Data and Information*.
- Misra, R. (2022). News category dataset. *arXiv preprint arXiv:2209.11429*.
- Misra, R. and Grover, J. (2021). *Sculpting Data for ML: The first act of Machine Learning*.
- Sepp Hochreiter, J. S. (1997). Long short-term memory. *Neural Computation* 9, 1735–1780.
- Victor Sanh, Lysandre Debut, J. C. T. W. (2019). Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. *EMC<sup>2</sup> : 5th Edition Co – located with NeurIPS’19*.