# Counterfeit: Sentiment Improved Text Summarization

**Piotr Parkitny**
University of California, Berkeley
pparkitny@berkeley.edu

## Abstract

I present Counterfeit, abstractive text summarization model that improves on existing approaches by adding an additional sentiment model used for text preprocessing. The sentiment model duplicates strong sentiment sentences to bias or deceive (Counterfeit) the second model into including those duplicated sentences in the summary. Counterfeit improves the performance of pretrained models ROUGE score for a new dataset. Counterfeit hyperparameters allow for precision-recall trade-off which results in the ability to create adjustable summaries that are more suitable in a practical setting.

## 1 Introduction

The Transformer(Vaswani et al., 2017) is the dominant architecture for natural language processing tasks. Its use in creating abstractive text summarizations has created new baseline benchmarks for state-of-the-art results.

In Counterfeit I use the Transformers library (Wolf et al., 2020) that is dedicated in supporting Transformer-based architectures and distribution of pretrained models to use two fully stacked pretrained models. The first model will be pretrained for sentence sentiment task and the second model on text summarization. Feeding results from the first models into the second to create the final text summarization. The problem that Counterfeit will solve is improving text summarization on datasets with no ground-truth available for training. In real-life application where ground-truth summaries are not available Counterfeit provides ROUGE score improvement.

The novelty of this approach is that it allows to transfer the knowledge from one model to another using the sentence duplication task. I have chosen to duplicate sentences with strong sentiment but it could be sentences with places, people, events, dates or any other sentence whose probability of making it into the summary needs to be increased.

## 2 Background

The problem of text summarization is well known with many existing solutions. I am going to use pretrained models that are available for download from the Transformers library. I will use 4 different text summarization models and one sentiment model.

I have decided to use two BART based models: facebook/bart-large-cnn and facebook/bart-large-xsum and two distillation models: sshleifer/distilbart-cnn-12-6 and sshleifer/distilbart-xsum-12-6. BART (Lewis et al., 2019) a denoising autoencoder for pretraining sequence-to-sequence models where the distillation (Shleifer and Rush, 2020) version is created by distilling BART and creating a smaller faster version while retaining strong performance which makes it better for practical use.

For sentiment analysis which is the first part of my stacked model I have chosen to also use a model from the Transformers library. The model is cardiffnlp/twitter-roberta-base-sentiment which is a roBERTa-base model trained on ~58M tweets and finetuned for sentiment analysis with the TweetEval (Barbieri et al., 2020) benchmark. The Roberta model is a text classification model that produces three scores for each sentence it is provided. The three scores are:

- LABEL0 - [NEG] negative sentiment score for the sentence

- LABEL1 - [NEU] neutral sentiment score for the sentence

- LABEL2 - [POS] positive sentiment score for the sentence

Having a scalar metric based on NEG,NEU and POS sentence sentiment is a key in identifying sentences that should be duplicated. In Counterfeit only the NEG and POS scores are used as they identify strong sentiment sentences. Provided below are three sentences along with a score that is produced by the sentiment model.

- A bomb blast on a bus kills 12 people: NEG=0.952 NEU=0.045 POS=0.003

- This is the best assignment ever: NEG=0.002 NEU=0.007 POS=0.991

- I am going home: NEG=0.079 NEU=0.729 POS=0.191

## 3 Objective

The purpose behind Counterfeit is to achieve better text summarization from pretrained models by stacking them together. The baseline is created using an existing pretrained summarization model and comparing that model against itself with a stacked sentiment model on top of it.

### 3.1 Evaluation Metric

Counterfeit evaluation metric is the ROUGE score (Lin, 2004). ROUGE stands for Recall Oriented Understudy for Gisting Evaluation. ROUGE is a set of metrics specifically designed for evaluation of automatic summarizations. ROUGE score is composed of three values

1. F1-Score - calculated as the harmonic mean of precision and recall.

2. Precision - fraction of relevant instances among the retrieved instances. It is the accuracy of positive predictions.

3. Recall - fraction of relevant instances that were retrieved. It is the ratio of positive instance that are correctly detected.

ROUGE score has different versions. I will be using ROUGE-1, ROUGE-2, and ROUGE-L scores as those are the most popular ROUGE versions used at this time. ROUGE-1 measures the overlap of unigram between the generated summary and ground-truth summary where ROUGE-2 refers to the overlap of bigrams. ROUGE-L is the longest common subsequence-based statistics. As it stands the ROUGE score is the most common evaluation used for text summarization. It is the right evaluation since at the time is the industry standard. Scoring text summarization using the ROUGE score is what other researchers are using currently. In practice ROUGE Recall metric is the most useful metric to use but it can be easily gamed which is the reason why the F1 metric is primarily used for comparing results.

The baseline for comparison will be the ROUGE score produced by the pretrained models compared to Counterfeit results that uses the same model. It should be noted that ROUGE score is not an ideal method for comparison based on reviewing ROUGE scored examples but is the best choice.

## 4 Methods

### 4.1 Architecture

Counterfeit architecture can be described as a stacked model pipeline where we feed information thought stages and obtain the summarized text. The overall model is composed of two tasks. First create the duplicate text using the sentiment model task and then create the summarization using the summarization model task.

### 4.2 Stage 1:Sentiment Model Task

The sentiment model will be used to score sentences so they can be duplicated if they have a strong sentiment. The cardiffnlp/twitter-roberta-base-sentiment model available from the Hugging Face website is used for this task. Sentence Boundary Detection which is the splitting of the article into individual sentences is done using the spaCy python package. spaCy brands itself as Industrial-Strength Natural Language Processing package that is free, open-source and written in Cython. It is a very well developed and documented package.

### 4.3 Stage 2: Summarization Models Task

I am using four different models for text summarization available from the Hugging Face website:

- facebook/bart-large-cnn

- facebook/bart-large-xsum

- sshleifer/distilbart-cnn-12-6

- sshleifer/distilbart-xsum-12-6

The four BART based models are pretrained for text summarization using CNN(Hermann et al., 2015) or XSUM(Narayan et al., 2018) datasets.

## 4.4 Hyperparameters

Counterfeit has four hyperparameters that are used to fine-tune the model. Based on the very large possible pool of hyperparameters I have chosen to use random selection for selecting the hyperparameter values along with some specific choices that are used to illustrate the impact of the hyperparameter on the ROUGE score. The four parameters are:

- TPOS - Top positive sentiment sentences to duplicate

- TNEG - Top negative sentiment sentences to duplicate

- DCNT - Duplication counter, the number of times to duplicate the sentence

- BEAM - The beam parameter used by the text summarization model

## 5 Datasets

Four datasets are used for scoring the model. Each dataset is composed of an article and a summary. A sample of 1000 observations from each dataset is used. Figure 1 demonstrates that the overall Counterfeit improvement stabilizes at around 1000 samples and there is no real need to score on the entire set as the results would not change. It provides evidence based on the near to zero slope of the lines that 1000 sample size is sufficient, and it will not bias the model ROUGE score. Using a sample of 1000 observations is imposed due to performance limitations as it takes about 1.5 hour to run the entire 1000 samples through the model using a 16vCPU with 15GB RAM and an NVIDIA Tesla V100 GPU. Increasing the sample size would be too costly. The following dataset are used for scoring Counterfeit:

**NEWSROOM**[NR] (Grusky et al., 2018) is a large dataset for training and evaluating summarization systems. It contains 1.3 million articles and summaries written by authors and editors in the newsrooms of 38 major publications.

**Multi-News**[MN] (Fabbri et al., 2019) consists of news articles and human-written summaries of these articles from the site newser.com. Each summary is professionally written by editors and includes links to the original articles cited.

**CNN-DailyMail**[CNN] (Hermann et al., 2015) is an English-language dataset containing just over 300k unique news articles as written by journalists

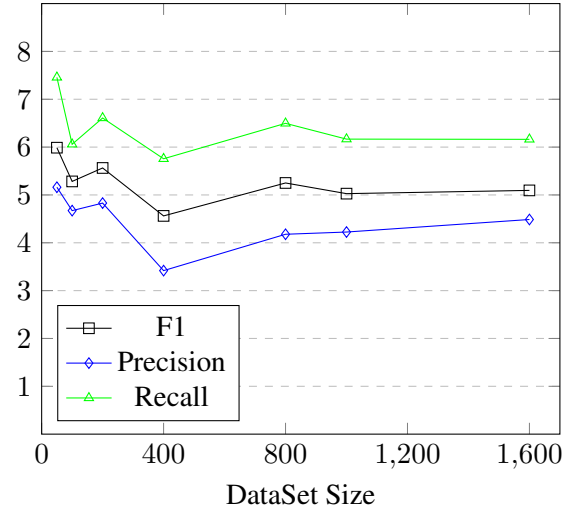Dataset Size VS ROUGE-1 Counterfeit Improvement



Figure 1: Counterfeit ROUGE-1 Improvement by Dataset size using NR with TPOS = TNEG = 10, DCNT = 5, BEAM = 4 and sshleifer-distilbart-xsum-12-6

at CNN and the Daily Mail. The current version supports both extractive and abstractive summarization, though the original version was created for machine reading and comprehension and abstractive question answering.

## 6 Hyperparameters Results

Changes to hyperparameters adjust the model ROUGE score. The values reported thought-out the report are the absolute difference between the two results such as F1 score improvement from 23% to 25% would be reported as improvement of 2.

### 6.1 TPOS and TNEG

A key hyperparameters in the model are the TPOS and TNEG which are always kept as the same value to minimize the total count of hyperparameters. From Figure 2 we can see that increasing the TPOS and TNEG value drastically increases Counterfeit improvements over the baseline model summary.

### 6.2 BEAM

Precision-recall trade-off in the ROUGE score can be achieved by adjusting the BEAM hyperparameter for the summarization model. Figure 3 demonstrates such relationship for BEAM hyperparameter value from 1 to 10. We can see that as we increase the BEAM hyperparameter recall increases but precision decreases. This relationship between

TPOS TNEG VS ROUGE-1 Counterfeit Improvement  BEAM VS ROUGE-1 Counterfeit Improvement
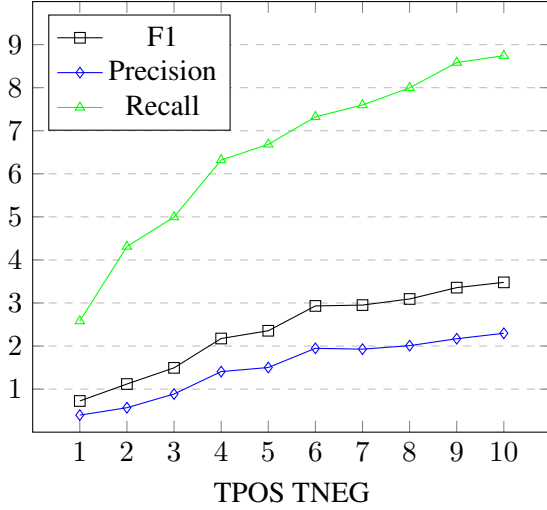


Figure 2: Coneterfeit ROUGE-1 Improvement by TPOS TNEG using NR and facebook-bart-large-cnn
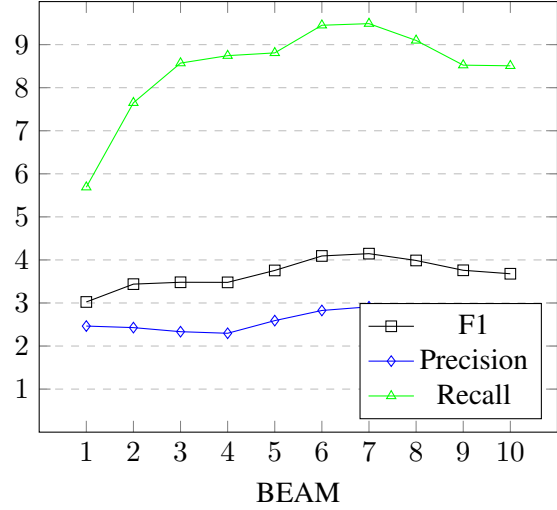


Figure 3: Counterfeit ROUGE-1 Improvement by BEAM using NR and facebook-bart-large

BEAM and recall can be adjusted based on what the type of summary is required. One major drawback of increasing the beam search value is that it increases processing time required to create the summary. When adjusting BEAM hyperparameter I have always kept the value the same in Counterfeit and baseline model and both models showed improvement however Counterfeit improved more as shown in Figure 3. Baseline model impact based on BEAM hyperparameter is shown in Figure 4. The baseline model Recall improves as we increase the BEAM hyperparameter however both F1 and precision decline.

### 6.2.1 DCNT

Increasing the DCNT hyperparameter increases Counterfeit ROUGE score. Figure 5 illustrates this behavior as all three of ROUGE-1 metric increase directly proportionally to DCNT increase. F1, Precision and Recall appear to stabilize at DCNT=10.

## 7 Results

Counterfeit ROUGE-1, ROUGE-2 and ROUGE-L metrics are higher when compared to the baseline model. Counterfeit summary does not end up with duplicated sentences in the summary which was a concern at the start of this project. Reviewing the results of the text summarization it becomes clear that generating a better summary does not mean scoring higher on the ROUGE score.
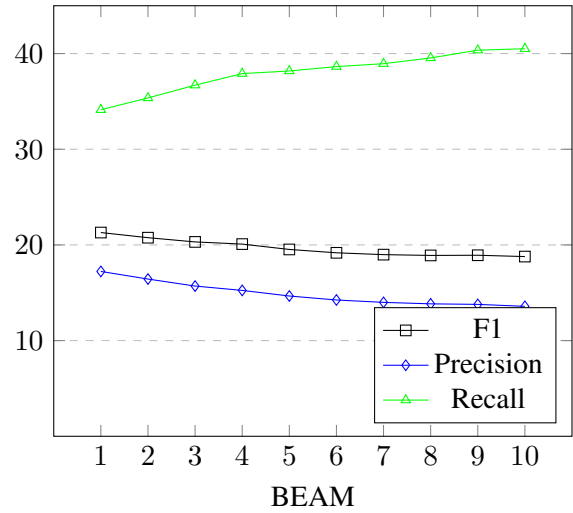
BEAM VS ROUGE-1 Baseline Improvement



Figure 4: Baseline ROUGE-1 Improvement by BEAM using NR and facebook-bart-large

### 7.1 Baseline Comparison

I have compared Counterfeit on each of the data sets and provided ROUGE score results. The complete results are provided in the appendix, in total 47 comparison are provided across the hyperparameters and the different pretrained models. The NR dataset results are provided in Table 1. Counterfeit outperforms the baseline on each of the ROUGE-1 metrics. The largest increase is produced on the Recall metric.

Counterfeit improvements for the CNN data set are in Table 2. The improvements are smaller than
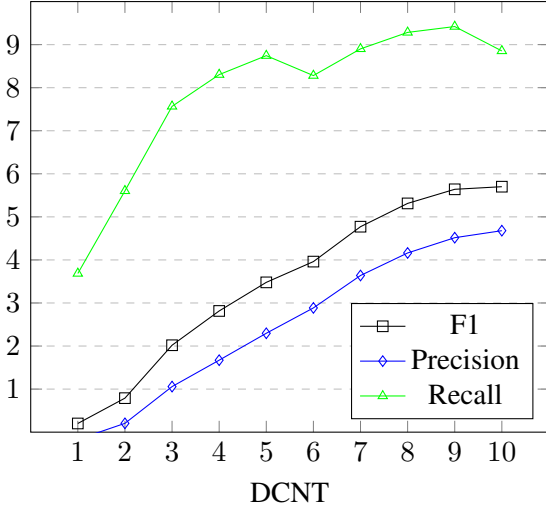
DCNT VS ROUGE-1 Counterfeit Improvement

Figure 5: Counterfeit ROUGE-1 Improvement by DCNT using NR and facebook-bart-large

|  | Score | Baseline | Counterfeit |
|---|---|---|---|
| **ROUGE-1** | **F1** | 18.98 | 23.13 |
|  | **Precision** | 13.99 | 16.91 |
|  | **Recall** | 38.94 | 48.43 |
| **ROUGE-2** | **F1** | 7.52 | 13.5 |
|  | **Precision** | 5.71 | 10.27 |
|  | **Recall** | 14.85 | 26.06 |
| **ROUGE-L** | **F1** | 18.59 | 23.52 |
|  | **Precision** | 13.93 | 17.56 |
|  | **Recall** | 34.85 | 44.39 |

Table 1: Counterfeit results for NR dataset on facebook-large-bart-cnn. Counterfeit hyperparameters TPOS=10, TNEG=10, DCNT=5, BEAM=7

the results obtained using the NR dataset but still impressive.

|  | Score | Baseline | Counterfeit |
|---|---|---|---|
| **ROUGE-1** | **F1** | 20.96 | 24.14 |
|  | **Precision** | 34.95 | 37.04 |
|  | **Recall** | 15.75 | 18.97 |
| **ROUGE-2** | **F1** | 5.9 | 8.33 |
|  | **Precision** | 10.02 | 12.9 |
|  | **Recall** | 4.4 | 6.53 |
| **ROUGE-L** | **F1** | 19.3 | 22.32 |
|  | **Precision** | 30.74 | 32.74 |
|  | **Recall** | 14.67 | 17.77 |

Table 2: Counterfeit results for CNN dataset on facebook-bart-large-xsum. Counterfeit hyperparameters TPOS=10, TNEG=10, DCNT=5, BEAM=4

Counterfeit produces lower improvement using the MN dataset. This can be seen in Table 3. The explanation appears to be related to the ground-truth summary and hyperparameter selection. MN dataset has the longest summary and a high value of the BEAM hyperparameter should have been used. The improvement here is small, in practical terms zero.

|  | Score | Baseline | Counterfeit |
|---|---|---|---|
| **ROUGE-1** | **F1** | 23.65 | 24.41 |
|  | **Precision** | 48.32 | 47.39 |
|  | **Recall** | 16.4 | 17.31 |
| **ROUGE-2** | **F1** | 6.73 | 6.8 |
|  | **Precision** | 13.81 | 13.32 |
|  | **Recall** | 4.64 | 4.78 |
| **ROUGE-L** | **F1** | 21.94 | 22.07 |
|  | **Precision** | 40.41 | 38.87 |
|  | **Recall** | 15.6 | 16.03 |

Table 3: Counterfeit results for MN dataset on facebook-bart-large-cnn. Counterfeit hyperparameters TPOS=10, TNEG=10, DCNT=1, BEAM=4

To be fully transparent Appendix-A contains the full list off all the tests that have been conducted. As I have presented some of the better results that Counterfeit has produced the results of all the scores are included in the appendix table. Overall, across all the tests Counterfeit produces improvement of about 2.9 for F1 ROUGE-1. The average improvement is included in Table 4. However, results for test summarization model that is matched to the dataset that it was trained on the same dataset are much lower or even negative.

|  | Score | Improvement |
|---|---|---|
| **ROUGE-1** | **F1** | 2.90 |
|  | **Precision** | 2.07 |
|  | **Recall** | 5.44 |
| **ROUGE-2** | **F1** | 3.74 |
|  | **Precision** | 3.27 |
|  | **Recall** | 5.93 |
| **ROUGE-L** | **F1** | 3.17 |
|  | **Precision** | 2.38 |
|  | **Recall** | 5.31 |

Table 4: Counterfeit Improvement across all the tests

# 8 Conclusion

I have introduced Counterfeit, a stacking pre-trained model approach to text summarization

which achieves higher ROUGE score when compared to baseline model on a dataset for which the model was not trained on. This approach is ideally suited to a production environment where you can't train the model due to lack of ground-truth summaries. The demonstrated model stacking approach can be expanded for text duplication to other stacked models depending on the priority of the text to be included in the summary. BEAM hyperparameter can be used for adjusting precision vs recall metric, it is extremely useful at making a longer summary as the trained model length hyperparameter is not as effective. Having a higher BEAM hyperparameter increases the length of the summary. Interestingly based on the results of the text duplication tasks the TPOS and TNEG hyperparameters can be used to adjust the summary into a more positive or negative version.

# References

Francesco Barbieri, Jose Camacho-Collados, Luis Espinosa Anke, and Leonardo Neves. 2020. TweetEval: Unified benchmark and comparative evaluation for tweet classification. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1644–1650, Online. Association for Computational Linguistics.

Alexander R. Fabbri, Irene Li, Tianwei She, Suyi Li, and Dragomir R. Radev. 2019. Multi-news: a large-scale multi-document summarization dataset and abstractive hierarchical model.

Max Grusky, Mor Naaman, and Yoav Artzi. 2018. Newsroom: A dataset of 1.3 million summaries with diverse extractive strategies. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 708–719, New Orleans, Louisiana. Association for Computational Linguistics.

Karl Moritz Hermann, Tomáš Kočiský, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, and Phil Blunsom. 2015. Teaching machines to read and comprehend.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. 2019. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension.

Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.

Shashi Narayan, Shay B. Cohen, and Mirella Lapata. 2018. Don't give me the details, just the summary! Topic-aware convolutional neural networks for extreme summarization. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, Brussels, Belgium.

Sam Shleifer and Alexander M. Rush. 2020. Pretrained summarization distillation.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.

# A Appendices

# B All Results

The table below is rotated 90° on the page. It reports summarization results across multiple models and datasets, with metric groups for **Baseline Model**, **Counterfeit**, and **Improvement**, each containing Rouge-1, Rouge-2, and Rouge-L (with F1, P, R sub-columns).

| # | DataSet | TPOS | TNEG | DCNT | SIZE | Summarization Model | BEAM |
|---|---------|------|------|------|------|---------------------|------|
| 1 | cnndailymail | 10 | 10 | 10 | 1000 | facebook-bart-large-cnn | 4 |
| 2 | newsroom | 10 | 10 | 4 | 1000 | facebook-bart-large-cnn | 4 |
| 3 | newsroom | 4 | 4 | 5 | 1000 | facebook-bart-large-cnn | 4 |
| 4 | newsroom | 9 | 9 | 5 | 1000 | facebook-bart-large-cnn | 4 |
| 5 | newsroom | 3 | 3 | 5 | 1000 | facebook-bart-large-cnn | 4 |
| 6 | multinews | 10 | 10 | 10 | 1000 | facebook-bart-large-cnn | 4 |
| 7 | cnndailymail | 10 | 10 | 10 | 1000 | sshleifer-distilbart-xsum-12-1 | 4 |
| 8 | newsroom | 1 | 1 | 5 | 1000 | facebook-bart-large-cnn | 4 |
| 9 | newsroom | 10 | 10 | 5 | 1000 | facebook-bart-large-cnn | 5 |
| 10 | newsroom | 10 | 10 | 5 | 1000 | sshleifer-distilbart-xsum-12-1 | 4 |
| 11 | multinews | 10 | 10 | 10 | 1000 | facebook-bart-large-cnn | 4 |
| 12 | newsroom | 7 | 7 | 5 | 1000 | facebook-bart-large-cnn | 4 |
| 13 | newsroom | 10 | 10 | 5 | 50 | sshleifer-distilbart-xsum-12-1 | 4 |
| 14 | newsroom | 10 | 10 | 5 | 1000 | facebook-bart-large-cnn | 3 |
| 15 | multinews | 10 | 10 | 10 | 1000 | sshleifer-distilbart-cnn-12-6 | 4 |
| 16 | newsroom | 5 | 5 | 5 | 1000 | facebook-bart-large-cnn | 4 |
| 17 | multinews | 10 | 10 | 10 | 1000 | sshleifer-distilbart-xsum-12-1 | 4 |
| 18 | newsroom | 10 | 10 | 5 | 1000 | facebook-bart-large-cnn | 10 |
| 19 | newsroom | 10 | 10 | 5 | 1000 | facebook-bart-large-cnn | 2 |
| 20 | multinews | 10 | 10 | 1 | 1000 | facebook-bart-large-xsum | 4 |
| 21 | newsroom | 10 | 10 | 3 | 1000 | facebook-bart-large-cnn | 4 |
| 22 | newsroom | 6 | 6 | 1 | 1000 | facebook-bart-large-cnn | 4 |
| 23 | newsroom | 10 | 10 | 1 | 1000 | facebook-bart-large-cnn | 7 |
| 24 | newsroom | 10 | 10 | 1 | 1000 | facebook-bart-large-cnn | 4 |
| 25 | newsroom | 10 | 10 | 5 | 1000 | facebook-bart-large-xsum | 4 |
| 26 | newsroom | 10 | 10 | 9 | 1000 | facebook-bart-large-cnn | 4 |
| 27 | newsroom | 10 | 10 | 6 | 1000 | facebook-bart-large-cnn | 4 |
| 28 | cnndailymail | 10 | 10 | 5 | 1000 | facebook-bart-large-xsum | 6 |
| 29 | newsroom | 10 | 10 | 5 | 200 | sshleifer-distilbart-xsum-12-1 | 4 |
| 30 | newsroom | 10 | 10 | 10 | 1000 | facebook-bart-large-cnn | 4 |
| 31 | newsroom | 10 | 10 | 5 | 1000 | facebook-bart-large-cnn | 4 |
| 32 | newsroom | 2 | 2 | 5 | 1000 | facebook-bart-large-cnn | 4 |
| 33 | cnndailymail | 10 | 10 | 10 | 1000 | sshleifer-distilbart-cnn-12-6 | 1 |
| 34 | newsroom | 10 | 10 | 5 | 1000 | facebook-bart-large-cnn | 4 |
| 35 | newsroom | 10 | 10 | 5 | 1000 | facebook-bart-large-cnn | 8 |
| 36 | cnndailymail | 10 | 10 | 5 | 1000 | facebook-bart-large-xsum | 4 |
| 37 | newsroom | 10 | 10 | 5 | 1000 | facebook-bart-large-cnn | 6 |
| 38 | newsroom | 10 | 10 | 8 | 800 | sshleifer-distilbart-xsum-12-1 | 4 |
| 39 | newsroom | 10 | 10 | 1 | 1000 | facebook-bart-large-cnn | 4 |
| 40 | multinews | 10 | 10 | 5 | 1000 | facebook-bart-large-cnn | 4 |
| 41 | newsroom | 10 | 10 | 5 | 100 | sshleifer-distilbart-xsum-12-1 | 4 |
| 42 | newsroom | 10 | 10 | 2 | 1000 | facebook-bart-large-cnn | 4 |
| 43 | newsroom | 10 | 10 | 5 | 1000 | facebook-bart-large-cnn | 9 |
| 44 | newsroom | 10 | 10 | 5 | 400 | sshleifer-distilbart-xsum-12-1 | 4 |
| 45 | newsroom | 10 | 10 | 7 | 1000 | facebook-bart-large-cnn | 4 |
| 46 | newsroom | 10 | 10 | 5 | 1600 | sshleifer-distilbart-xsum-12-1 | 4 |
| 47 | newsroom | 8 | 8 | 5 | 1000 | facebook-bart-large-cnn | 4 |
| | Total Mean | | | | | | |

## C  BEAM Hyperparameter

### C.1  NR, TPOS=10, TNEG=10, DCNT=5, facebook-bart-large-cnn

**BEAM=1** U.S. President Obama and U.N. Secretary-General Ban Ki-moon have put climate change in the spotlight. Australian Prime Minister Tony Abbott had said he wanted the G20 meeting to focus on economic and security issues. Obama ensured climate change was front and center before he even landed in Queensland.

**BEAM=2** Australian Prime Minister Tony Abbott had said he wanted the G20 meeting to focus on economic and security issues. But U.S. President Obama and U.N. Secretary-General Ban Ki-moon have put climate change squarely in the spotlight. On Wednesday, Obama announced landmark joint emissions commitments with Chinese President Xi Jinping. On Friday, in a speech at the University of Queensland in Brisbane, Obama said the United States would contribute 3 billion into the Green Climate Fund.

**BEAM=10** Australian Prime Minister Tony Abbott may have thought he left global climate change off the agenda for the G20 summit in Brisbane, Australia. But U.S. President Obama and U.N. Secretary-General Ban Ki-moon have put it squarely in the spotlight through a series of actions in the past few days. Climate change has been a traditional agenda item for G20, G8 and G7 meetings in recent years, making its absence from the Brisbane gathering of world leaders noteworthy. Obama ensured climate change was front and center before he even landed in Queensland's capital city on Friday. On Wednesday in Beijing, he announced landmark joint emissions commitments with Chinese President Xi Jinping. China agreed for the first time to peak its greenhouse gas emissions by 2030, and to dramatically scale up the use of renewable energy in its economy, to about 20 by 2020. Then on Friday, in a speech at the University of Queensland in Queensland, Obama announced that the U.s. would contribute 3 billion into the Green Climate Fund which aims to help developing countries improve their climate resilience, and transition away from fossil fuels.

## D  Summary Example

### D.1  NR, TPOS=10, TNEG=10, DCNT=1, BEAM=4, facebook-bart-large-xsum, ID-1

**Article:** Police at the scene of the Tunis blast. (AAP)bomb blast on a bus transporting Tunisiaś presidential guard in central Tunis has killed at least 12 people and wounded 16 more.security source at the site of the attack said "most of the agents who were on the bus are dead."presidency had said the blast on the bus killed at least 14 people, but this was revised down, with the interior ministry giving the latest toll as least 12 dead.explosion, described as an "attack" by presidential spokesman Moez Sinaoui, struck on the capitalś Mohamed V Avenue.AFP journalist reported seeing the partly burnt-out shell of the bus, with police, ambulances and fire trucks at the scene.people were in tears.bank employee working nearby reported hearing a large explosion and seeing the bus on fire.there was no immediate claim of responsibility, Tunisia has been plagued by Islamist violence since the 2011 overthrow of longtime dictator Zine El Abidine Ben Ali.attacks earlier this year claimed by ISIL terror group targeted foreigners - at the National Bardo Museum in March, killing 21 tourists and a policeman, and at a resort hotel in Sousse in June, killing 38 tourists.of members of the security forces have also been killed by extremists since 2011.

**Summary:** A bomb blast on a bus transporting Tunisia's presidential guard in central Tunis on Tuesday killed at least 12 people and wounded 16, the interior ministry said.

**Baseline Summary:** Tunisia has been plagued by Islamist violence since the 2011 overthrow of longtime dictator Zine El Abidine Ben Ali. Two attacks earlier this year claimed by ISIL terror group targeted foreigners - at the National Bardo Museum in March, killing 21 tourists and a policeman, and at a resort hotel in Sousse in June, killing 38 tourists. Dozens of members of the security forces have also been killed by extremists since 2011.'

**Counterfeit Summary:** A bomb blast on a bus transporting Tunisia's presidential guard in central Tunis has killed at least 12 people and wounded 16 more. A security source at the site of the attack said "most of the agents who were on the bus are dead" The explosion, described as an "attack" by presidential spokesman Moez Sinaoui, struck on the capital
ś Mohamed V Avenue. No immediate claim of responsibility, Tunisia has been plagued by Islamist violence since the 2011 overthrow of longtime dictator Zine El Abidine Ben Ali. Two attacks earlier this year claimed by ISIL terror group targeted foreigners - at the National Bardo Museum in March, killing 21 tourists and a policeman, and at a resort hotel in Sousse in June, killing 38 tourists.