

UCD-CS at TREC 2021 Incident Streams Track

Congcong Wang

School of Computer Science
University College Dublin
Dublin, Ireland
conggong.wang@ucdconnect.ie

David Lillis

School of Computer Science
University College Dublin
Dublin, Ireland
david.lillis@ucd.ie

Abstract

In recent years, the task of mining important information from social media posts during crises has become a focus of research for the purposes of assisting emergency response (ES). The TREC Incident Streams (IS) track is a research challenge organised for this purpose. The track asks participating systems to both classify a stream of crisis-related tweets into humanitarian aid related information types and estimate their importance regarding criticality. The former refers to a multi-label information type classification task and the latter refers to a priority estimation task. In this paper, we report on the participation of the University College Dublin School of Computer Science (UCD-CS) in TREC-IS 2021. We explored a variety of approaches, including simple machine learning algorithms, multi-task learning techniques, text augmentation, and ensemble approaches. The official evaluation results indicate that our runs achieve the highest scores in many metrics. To aid reproducibility, our code is publicly available¹.

1 Introduction

Unexpected Emergencies can cause substantial loss of both life and property if assistance is not available in a timely manner. Recent studies have sought solutions for more efficient emergency response (ES) using computational techniques (Caragea et al., 2011; Vieweg, 2012; Imran et al., 2015). Among these works, social media is acknowledged as a promising venue for mining important messages for ES given that some people do tend to seek help by posting messages on social media as a crisis situation unfolds, these messages may contain critical information of relevance to

¹<https://github.com/wangconggong123/crisis-mtl>

emergency responders (Imran et al., 2015; McCreadie et al., 2019, 2020).

This motivated the Incident streams (IS) track (McCreadie et al., 2019, 2020), which challenges the community to explore effective approaches for identifying important messages from user-posted streams on social media during crises. The IS track is a research challenge consisting of two main tasks. The first asks participating systems to classify a stream of crisis-related tweets into humanitarian aid related categories, known as the multi-label information types (ITs) classification task. IS comprises a total of 25 information types that are defined as the categories of possible aid needs in a crisis such as requesting donations, call for search and rescue, reporting weather, etc. The 25 ITs are further divided into two sub-categories; 6 are defined as “actionable” ITs (e.g., search and rescue) and the remaining 19 are “non-actionable” ones (e.g., reporting weather)². The second task is known as the priority estimation task. It requires participants to estimate the criticality of those tweets that have been classified into ITs. This criticality is represented by a numeric value from 0 to 1 indicating the least to the most importance.

Having participated in this track since 2019 (the second iteration of the IS track), our system has evolved based on the experience learnt from our prior participations in past TREC-IS editions³. Unlike previous IS editions (McCreadie et al., 2019, 2020), TREC-IS 2021 initiated an online leaderboard for participants⁴. It is noted that the leaderboard only reports the performance of par-

²For a full list of the ITs, see the official website at http://dcs.gla.ac.uk/~richardm/TREC_IS/

³The IS track normally runs two editions every year and a new test set is annotated and added to the training set after each edition.

⁴<https://trecis.github.io/>

ticipating runs in the 2021A Edition where the test set is partially annotated within events based on pooling by priority (the submitted test tweets are predicted by ITs and sorted by priority score within each event). In the 2021B Edition, the test set comprises the tweets of more annotated events and deeper pooling (new judgements). Hence, the 2021B Edition acted as an enhanced evaluation for the participating runs that had been submitted to the 2021A leaderboard. Given the timeliness of performance feedback from the leaderboard, we explored a wide range of approaches including a Naïve Bayes classifier using contextual sentence embeddings as the features, multi-task learning approaches with text augmentations, and an ensemble technique. We found our runs perform consistently well in both A and B editions and in particular our multi-task learning runs and ensemble runs perform the best in many metrics amongst all participating runs. However, the results did not show that text augmentations can bring overall improvements.

2 Related Work

Since the launch of TREC-IS, many works have been produced on the topic of crisis tweet classification and priority estimation (CTC-PE). Wang et al. (2019) applied Naïve Bayes, Support Vector Machine (SVM), Random Forest, and the ensemble of these models with hand-crafted features for CTC-PE. Choi et al. (2018) applied SVM and deep learning models which combine class activation mapping with one-shot learning in convolutional neural networks for CTC-PE. Miyazaki et al. (2019) applied a BiLSTM model for CTC-PE by incorporating the hierarchical structure of labels into the model. Wang and Lillis (2020) applied a BiLSTM model along with pre-trained ELMo embeddings and trainable embeddings as the input features for CTC-PE. Wang et al. (2021a) fine-tuned BERT (Devlin et al., 2019) in a multi-task learning manner for CTC-PE while Wang and Lillis (2021) extended the multi-task learning approach to a sequence-to-sequence transformer-based model T5 (Raffel et al., 2020). To alleviate the class imbalanced problem, SHARMA and BUNTAIN (2020) applied synonym replacements as well as crisis image labels to augment the original training data. Other techniques such as downsampling the training data or generating new examples via GPT-2 are also found in the lit-

erature (Wang and Lillis, 2020; Hepburn and McCreadie, 2020).

3 Methods and Experiments

Table 1 summarises the runs we submitted to TREC-IS 2021. The major techniques used in the runs are described as follows.

ML run: In this run, we convert each tweet to a representation via pre-trained sentence embeddings (SBERT) models (Reimers and Gurevych, 2019). Having tested multiple combinations of the available publicly-available pre-trained variants of SBERT⁵, we finally choose `all-mpnet-base-v2` and `paraphrase-xlm-r-multilingual-v1` to embed the tweets, where each tweet’s representation is the concatenation of outputs of the two models. Similarly, in choosing the downstream classifier, we exhaustively searched over a list of candidates including SVC, logistic regression, decision tree and random forest. We finally used GaussianNB as the downstream classifier as it brought the best result on the development set. Here the classifier is only for IT prediction whereas the priority is simply mapped from the predicted ITs (An IT’s mapped priority score is the average priority of all tweets belonging to this IT in the training set). In this approach, priority is assumed to be a function of the IT.

Multi-task and ensemble run. Similar to Wang et al. (2021a), we train a single model for both the downstream IT classification and the priority estimation tasks in a multi-task learning manner. In simple terms, we fine-tuned a pre-trained DeBERTa model (He et al., 2020) jointly on the two tasks through adding a multi-label classification head and a regression head on top of the model. The model is optimised on a linear combination of the cross entropy loss of classification and MSE regression loss. By doing so, the model is capable of making predictions on both tasks for the test tweets with only one input forward at inference time. Based on this idea, we train multiple individual models varying in model size and training data size. Ultimately the individual models consist of a fine-tuned `deberta-base`, `deberta-base` with Easy Data Augmentation (EDA) (Wei and Zou, 2019) and `deberta-large`. EDA is used in our

⁵<https://huggingface.co/sentence-transformers>

system to augment the training data in order to ensure that every IT has at least 500 examples. We apply this augmentation since the original training data is heavily class-imbalanced. Moreover, we adopt the ensemble approach from Wang et al. (2021a) to leverage the predictions of individual models for IT classification and priority estimation. The ensemble approach is simple, using the union of predicted ITs by individual models as the final IT prediction and the highest priority among individual priority predictions as the final priority score for test tweets.

Ensemble run with post-processing. Among the pre-defined 25 ITs⁶, there is an IT called “Irrelevant”. The multi-label ITs predicted by the above ensemble approach can contain this class along with other ITs. However, a tweet that is classified as “Irrelevant” cannot also be labelled with other ITs. We thus adopt a post-processing step to handle this issue. For any tweet with this type of prediction, we compare the prediction probability for “Irrelevant” with the probabilities of other ITs. The tweet is assigned “Irrelevant” if its probability for “Irrelevant” is greater than all the individual probabilities of the other ITs. Otherwise it is predicted to be one of the other ITs. As a result, the tweet’s priority score also becomes 0 if it is considered to be “Irrelevant”.

Direct-Generation Augmentation (DGA) and Noise Label Annealing (NLA). Aside from EDA augmentation, described above, we also explored other augmentation techniques. Inspired by Wang et al. (2021b) who applied large pre-trained language models to generate training data without any human annotation and model training but through carefully-crafted prompts, we utilise a similar approach using a small number of examples as the prompt. We choose the pre-trained checkpoint `gpt-neo-2.7B`⁷ as the generation model and the prompt template is formulated as follows:

Tweet for help in disaster

Title: {IT name}

Content: {Tweet text}

The template constructs a stream of natural language, starting with a task description⁸, followed

⁶http://dcs.gla.ac.uk/~richardm/TREC_IS/

⁷<https://huggingface.co/EleutherAI/gpt-neo-2.7B>

⁸The task description is carefully chosen based on our pre-

by the title and content fields, which are replaced by the IT name and the tweet text respectively. This is something we refer to Direct-Generation Augmentation (DGA). In our DGA-based runs, we sampled two examples of non-target ITs from the training data to construct the prompt. To generate a new example for a target IT, we omit the textual part of the content so that the model learns from the prompt (two sampled non-target examples) to complete the content part of the target IT. Finally, we used DGA to augment the training data, thus ensuring that every IT has at least 1000 examples. One challenge associated with this kind of augmentation is that the generated texts are likely not to be label-aligned with the label it should be and these generated texts are deemed to be noisy or label-incompatible data that is harmful to the downstream task performance. We adopt a strategy called Noisy Label Annealing (NLA) introduced in Wang et al. (2021b) to filter out noisy training signals as training progresses. The general idea is that we check the predictions of augmented training examples at the end of each epoch of downstream model training and remove an example if the model disagrees with its label with high confidence.

Regarding model training, we remove approximately 10% of the original training data to use as the development set. We fine-tune the multi-task learning model with 10 epochs and select the best checkpoint based on the IT macro-F1 score on the development set. The model’s parameters are tuned on batches (batch size = 16) of training data using Adam (Kingma and Ba, 2015) as the optimizer with a linear warm-up scheduler changing the learning rate from 0 to $5e - 5$ within the first 10% of total training steps and then linearly decays to 0. Apart from these, the rest of hyperparameters are set up the same as the default by the transformers library (Wolf et al., 2020).

4 Results and Discussions

In order to measure a system’s performance from different perspectives, the IS track defines multiple metrics. The metrics can be broadly divided into two categories: IT classification measurements and priority estimation measurements. They are described as follows:

- **IT classification measurements:** To measure the performance of ITs classification,

liminary experiments evaluated on the development set.

Run names	Description
ucdcs-strans.nb	ML run of using SBERT as the fixed features and GaussianNB as the downstream classifier
ucdcs-run1	Multi-task run using deberta-base
ucdcs-run2	Multi-task run using deberta-base with EDA augmentation
ucdcs-run3	Multi-task run using deberta-large
ucdcs-mtl.ens (run4)	Ensemble run of run 1, 2 and 3
ucdcs-mtl.ens.new	Ensemble run with post processing
ucdcs-mtl.fta	Multi-task run of deberta-base with direct-generation augmentation (DGA)
ucdcs-mtl.fta.nla	Multi-task run of deberta-base with DGA plus noise label annealing (NLA)
ucdcs-mtl.ens.fta	Ensemble run of run1, 3 and mtl.fta.nla

Table 1: Overview of UCD-CS runs at TREC-IS 2021. Details of the techniques in bold are elaborated in Section 3.

	nDCG	IT F1 [A]	IT F1 [All]	IT Acc.	Pri F1 [A]	Pri F1 [All]	Pri R [A]	Pri R [All]
ucdcs-strans.nb	0.4297	0.2423	0.2695	0.8294	0.1998	0.1988	0.147	0.1514
ucdcs-run1	0.6115	0.215	0.2951	0.8837	0.3032	0.3068	0.2592	0.297
ucdcs-run2	0.5848	0.2215	0.2984	0.8835	0.25	0.2781	0.2305	0.2748
ucdcs-run3	0.6051	0.2391	0.31	0.8852	0.272	0.3066	0.3112	0.3325
ucdcs-mtl.ens (run4)	0.5907	0.2579	0.3211	0.8646	0.3052	0.3125	0.325	0.3416
ucdcs-mtl.ens.new	0.5951	0.2627	0.3205	0.8686	0.305	0.3211	0.2892	0.3089
ucdcs-mtl.fta	0.589	0.1986	0.2793	0.8902	0.2769	0.2807	0.2471	0.3001
ucdcs-mtl.fta.nla	0.529	0.2007	0.2751	0.8815	0.262	0.281	0.1721	0.2193
ucdcs-mtl.ens.fta	0.5755	0.1592	0.2597	0.8034	0.306	0.3141	0.2786	0.2855
med	0.5695	0.206	0.2823	0.8827	0.2113	0.2175	0.1728	0.2099
max	0.6115	0.2815	0.3211	0.8902	0.306	0.3211	0.4349	0.3585

Table 2: The performance of UCD-CS runs at TREC-IS 2021 based on results using only the judgments in 2021A Edition. The figures in **bold** indicate the best scores across all participating runs. The med and max rows present the median and maximum scores of each metric respectively across all participating runs.

three metrics are defined. They are IT F1 [A], IT F1 [All] and IT Acc., referring to the F1 score of only actionable ITs classification, the F1 score and the accuracy of all 25 ITs classification respectively.

- **Priority estimation measurements:** There are five metrics related to the evaluation of priority estimation. Four of them are: Pri F1 [A], Pri F1 [All], Pri R [A] and Pri R [All], referring to the F1 scores and recall scores of only actionable and all ITs classification respectively. Besides these, nDCG is a ranking metric included in this category to measure a run’s average performance in ranking the top 100 test tweets per event by priority.

As TREC-IS 2021 has been run with two editions (A and B) that produce two sets of judgements, we report our runs’ performance separately in the judgements of each edition as well as in the combined judgements of both editions, as presented in Tables 2, 3 and 4.

First in overview, most of our runs perform well consistently in both editions across the participating runs. When compared to the median and maximum of each metric, we find that our multi-task and ensemble runs in particular achieve strong performance, hitting the best scores in many cases. To

examine the figures by task, we notice that some of our runs can perform well in one task while under-performing in the other. For example, the ML run achieves decent scores in IT classification but its scores for priority estimation are relatively poor. This is likely due to the simple mapping from IT predictions to priority estimation in that run. We expect the ML run to be improved upon by modelling not just the IT classification but also the priority estimation (a regression task).

In terms of our multi-task runs, we find that these runs tend to achieve strong scores in both tasks. This indicates that the joint learning on both tasks through fine-tuning pre-trained language models (DeBERTa in our case) can help achieve strong performance, which adds support to the results in (Wang et al., 2021a) in previous TREC-IS editions. Also unsurprisingly, the bigger fine-tuned model brings slightly-improved performance when comparing our run3 to our run1. Regarding our ensemble runs without augmentation, they outperform our other runs in almost every metric for both tasks as well as achieving highest scores in many metrics among the participating runs. It is noted the mtl.ens.new run with post-processing (to deal with the “Irrelevant” IT) further improves the performance in IT classification as compared to run4.

	nDCG	IT F1 [A]	IT F1 [All]	IT Acc.	Pri F1 [A]	Pri F1 [All]	Pri R [A]	Pri R [All]
ucdcs-strans.nb	0.338	0.1861	0.2395	0.8557	0.1733	0.1688	0.0425	0.1003
ucdcs-run1	0.4499	0.2177	0.247	0.8966	0.2376	0.2566	0.1547	0.2525
ucdcs-run2	0.4361	0.2087	0.2433	0.8947	0.242	0.2528	0.207	0.2622
ucdcs-run3	0.4583	0.2218	0.2539	0.8964	0.2543	0.2756	0.221	0.2571
ucdcs-mtl.ens (run4)	0.4521	0.2361	0.2591	0.8708	0.2753	0.2582	0.2302	0.2952
ucdcs-mtl.ens.new	0.4555	0.251	0.2623	0.8753	0.2783	0.2703	0.2116	0.2604
ucdcs-mtl.fta	0.4464	0.1958	0.2369	0.9067	0.2309	0.2333	0.2044	0.2454
ucdcs-mtl.fta.nla	0.4069	0.1687	0.2187	0.8945	0.2588	0.2635	0.1861	0.2122
ucdcs-mtl.ens.fta	0.4448	0.0946	0.1889	0.8113	0.2798	0.2724	0.2149	0.2629
med	0.4272	0.1842	0.233	0.8947	0.2107	0.2031	0.1495	0.1993
max	0.4791	0.251	0.2623	0.9067	0.2798	0.2756	0.2302	0.2952

Table 3: The performance of UCD-CS runs at TREC-IS 2021 based on results using only the judgments in 2021B Edition. The figures in **bold** indicate the best scores across all participating runs. The med and max rows present the median and maximum scores of each metric respectively across all participating runs.

	nDCG	IT F1 [A]	IT F1 [All]	IT Acc.	Pri F1 [A]	Pri F1 [All]	Pri R [A]	Pri R [All]
ucdcs-strans.nb	0.3368	0.2083	0.2575	0.8474	0.1959	0.1712	0.1096	0.1417
ucdcs-run1	0.4727	0.2433	0.2772	0.8926	0.2657	0.2632	0.259	0.2888
ucdcs-run2	0.4569	0.2326	0.2753	0.8911	0.2536	0.2524	0.1995	0.2686
ucdcs-run3	0.4707	0.2538	0.286	0.893	0.253	0.2694	0.2741	0.3053
ucdcs-mtl.ens (run4)	0.4617	0.267	0.2923	0.8685	0.2817	0.2623	0.2886	0.3182
ucdcs-mtl.ens.new	0.4643	0.2784	0.2946	0.8728	0.2864	0.2734	0.2748	0.2827
ucdcs-mtl.fta	0.463	0.2159	0.2647	0.9016	0.2497	0.2361	0.2779	0.2834
ucdcs-mtl.fta.nla	0.4193	0.1936	0.2488	0.8907	0.2727	0.2609	0.265	0.2339
ucdcs-mtl.ens.fta	0.4515	0.1131	0.217	0.8073	0.2852	0.2724	0.2479	0.2665
med	0.4381	0.2008	0.26	0.8911	0.2086	0.2044	0.2087	0.2431
max	0.4904	0.2784	0.2946	0.9016	0.2864	0.2734	0.3072	0.3182

Table 4: The performance of UCD-CS runs at TREC-IS 2021 based on results using the judgments in 2021A and 2021B Editions. The figures in **bold** indicate the best scores across all participating runs. The med and max rows present the median and maximum scores of each metric respectively across all participating runs.

To check the effect of EDA augmentation when comparing run2 to run1, we see marginal improvements in priority estimation in 2021B (Table 3) but not in the other two tables. Hence, it is difficult to conclude whether the EDA augmentation adds benefit in this scenario. We also see similar results for our DGA-based runs (see our `mtl.fta`, `mtl.fta.nla` and `mtl.ens.fta` runs). This seems somewhat inconsistent with the previous study (Wang and Lillis, 2020). We attribute these results to two possible reasons. First, the training data has grown from around 5,000 to 50,000 since then, so it is quite possible that the advantage of text augmentation in a low-data situation is more obvious. Second, since the downstream model used in our current runs is pre-trained on big general text data, the new examples generated by text augmentation may be noisy as well as be redundant (the model learns general language features at pre-training and is likely to augment similar examples itself implicitly during fine-tuning). Hence, better approaches for denoising and diversifying the augmented examples are avenues of research that we seek to explore in

the future.

5 Conclusion

In this paper, we report UCD-CS’s participation at the TREC 2021 Incident Streams track (TREC-IS). We submitted multiple runs and our approaches included machine learning algorithms, multi-task learning techniques and ensemble approaches. Among these runs, we find in particular that our multi-task and ensemble runs achieve strong performance in both the information type classification and priority estimation tasks through two rounds of evaluation: TREC-IS 2021A and B editions. Although we explored some text augmentation approaches with the intent of boosting the performance, the results did not indicate consistent performance improvements and thus we seek better augmentation techniques in the future.

References

Cornelia Caragea, Nathan J McNeese, Anuj R Jaiswal, Greg Traylor, Hyun-Woo Kim, Prasenjit Mitra, Dinghao Wu, Andrea H Tapia, C Lee Giles, Bernard J Jansen, et al. 2011. Classifying text

- messages for the haiti earthquake. In *Proceedings of the 8th International Conference on Information Systems for Crisis Response and Management (IS-CRAM 2011)*. Citeseer.
- Won-Gyu Choi, Seung-Hyeon Jo, and Kyung-Soon Lee. 2018. Cbnu at trec 2018 incident streams track.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. **BERT: Pre-training of deep bidirectional transformers for language understanding**. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2020. Deberta: Decoding-enhanced bert with disentangled attention. *arXiv preprint arXiv:2006.03654*.
- Alexander J Hepburn and Richard McCreadie. 2020. University of glasgow terrier team (uogtr) at the trec 2020 incident streams track. *Image*, 8:5.
- Muhammad Imran, Carlos Castillo, Fernando Diaz, and Sarah Vieweg. 2015. Processing social media messages in mass emergency: A survey. *ACM Computing Surveys (CSUR)*, 47(4):67.
- Diederik P. Kingma and Jimmy Ba. 2015. Adam: A method for stochastic optimization. In *Proceedings of the 3rd International Conference for Learning Representations*, San Diego.
- Richard McCreadie, Cody Buntain, and Ian Soboroff. 2019. TREC incident streams: Finding actionable information on social media. *Proceedings of the International ISCRAM Conference, 2019-May*(May):691–705.
- Richard McCreadie, Cody Buntain, and Ian Soboroff. 2020. Incident Streams 2019: Actionable Insights and How to Find Them. In *Proceedings of the 17th International Conference on Information Systems for Crisis Response and Management (IS-CRAM 2020)*.
- Taro Miyazaki, Kiminobu Makino, Yuka Takei, Hiroki Okamoto, and Jun Goto. 2019. Label embedding using hierarchical structure of labels for twitter classification. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 6318–6323.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67.
- Nils Reimers and Iryna Gurevych. 2019. **Sentencebert: Sentence embeddings using siamese bert-networks**. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics.
- SHIVAM SHARMA and CODY BUNTAIN. 2020. Improving classification of crisis-related social media content via text augmentation and image analysis.
- Sarah Elizabeth Vieweg. 2012. *Situational awareness in mass emergency: A behavioral and linguistic analysis of microblogged communications*. Ph.D. thesis, University of Colorado at Boulder.
- Congcong Wang and David Lillis. 2020. **Classification for Crisis-Related Tweets Leveraging Word Embeddings and Data Augmentation**. In *Proceedings of the Twenty-Eighth Text REtrieval Conference (TREC 2019)*, Gaithersburg, MD.
- Congcong Wang and David Lillis. 2021. **Multi-task transfer learning for finding actionable information from crisis-related messages on social media**. In *Proceedings of the Twenty-Ninth Text REtrieval Conference (TREC 2020)*, Gaithersburg, MD, USA.
- Congcong Wang, Paul Nulty, and David Lillis. 2021a. Transformer-based Multi-task Learning for Disaster Tweet Categorisation. *Proceedings of the International ISCRAM Conference, 2021-May*(May).
- Junpei Zhou Xinyu Wang, Po-yao Huang, and Alexander Hauptmann. 2019. CMU-Informedia at TREC 2019 Incident Streams Track. In *Proceedings of the Twenty-Eighth Text REtrieval Conference (TREC 2019)*, Gaithersburg, MD.
- Zirui Wang, Adams Wei Yu, Orhan Firat, and Yuan Cao. 2021b. Towards zero-label language learning. *arXiv preprint arXiv:2109.09193*.
- Jason Wei and Kai Zou. 2019. Eda: Easy data augmentation techniques for boosting performance on text classification tasks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 6382–6388.
- 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.