

# Adapters for Resource-Efficient Deployment of Intent Detection Models

Anonymous ACL submission

## Abstract

Modern Transformer-based language models such as BERT are huge and expensive to deploy in practical applications. In environments such as commercial chatbot-as-a-service platforms that deploy many NLP models in parallel, less powerful models with a smaller number of parameters are often used to keep deployment costs down. Also, in times of climate change and scarce resources, the deployment of many huge models is no longer adequate. This paper proposes the BERT+Adapter architecture for hosting many models and saving considerable amounts of (GPU) memory. We further demonstrate the effectiveness of this approach using the example of intent detection for dialog systems. Many task-specific adapters can share one large Transformer model with the adapter framework. To deploy 100 NLU models requires 1 GB of memory for the proposed BERT+Adapter architecture, compared to 41.78 GB for a BERT-only architecture. Also, we show that the training time of the BERT+Adapter architecture is on average 14.43 times shorter than that of vanilla BERT. Furthermore, we demonstrate that the accuracy of BERT+Adapter on intent detection tasks is on par with a vanilla BERT architecture.

## 1 Introduction

Natural Language Understanding (NLU) is an essential component of dialog systems (DS). The NLU converts an unstructured user utterance into structured information: It comprises a) Intent Detection (ID) where the dialog system classifies a user utterance into a predefined list of intents. The system can understand what the user wants to say based on this classification. The other part of NLU is b) slot filling / entity recognition (ER) in which the dialog system fills in specific slots that belong to an intent.

Transformer-based models (Vaswani et al.) currently show the best results in ID (Mehri and Eric,

2021). Transformers models are very large, a standard BERT model (Devlin et al., 2019) for ID on the HWU64 dataset (Liu et al., 2019) occupies 417.75 MB of data in the Hugging Face implementation (Wolf et al., 2020). Hosting these large models is expensive and uses large amounts of computational power which is no longer adequate in times of climate change and scarce resources (Strubell et al., 2019). Therefore, smaller architectures were proposed such as DIET (Bunk et al., 2020) or ConveRT (Henderson et al., 2020) which use substantially less memory (63 MB for DIET in our HWU64 example). While state-of-the-art Transformer-based NLU models reach performances up to 92% accuracy (Mehri and Eric, 2021), practical NLU models show a weaker performance of 86-89% (Liu et al., 2019).

In this paper, we apply adapters (Rebuffi et al., 2017; Houlsby et al., 2019) to ID. We use a general-purpose pre-trained BERT model and introduce a small number of additional parameters (5.87 MB for our HWU64 example in the AdapterHub implementation (Pfeiffer et al., 2020)). During training, we freeze the parameters of the original BERT model and train only the additional parameters. Using this approach, we propose a resource-efficient method to deploy multiple ID models: Instead of deploying multiple BERT models, we only need a single, shared BERT model and one adapter for each downstream NLU application. The deployment of, for example, 100 ID models using this new framework requires 1 GB (417.75 MB for the BERT model +  $100 \times 5.87$  MB for the adapters), instead of 41.78 GB ( $100 \times 417.75$  MB size of the BERT model).

Adapters show their strength in environments where many models are deployed in parallel. We demonstrate our approach in two areas of applications. First, we propose its use in chatbot-as-a-service platforms such as IBM Watson Assistant<sup>1</sup>,

<sup>1</sup><https://www.ibm.com/products/watson-assistant>

Google Dialogflow<sup>2</sup>, and SAP Conversational AI<sup>3</sup>. These systems offer chatbots as a service and therefore host many DS and a large number of NLU models. This industry can save significant amounts of costs, resources, and energy using our approach.

The second application area is modular dialog systems (MDS) (Nehring and Ahmed, 2021) in which a DS is composed of several sub-DS. Often, each sub-DS uses its own NLU model and can benefit from the adapter approach’s resource savings.

We perform experiments on three datasets to compare the performance of the BERT+Adapter model to a fully fine-tuned, vanilla BERT model. We determine the sizes of the models and show that for a large number of models, BERT+Adapter saves a significant amount of memory, compared to BERT and to the DIET (Bunk et al., 2020) model. The experiments show that the accuracy of ID of BERT+Adapter is comparable to BERT, despite the memory savings of BERT+Adapter. Also, we show that BERT+Adapter is trained considerably faster than BERT.

We publish the source code to reproduce our experiments on GitHub<sup>4</sup>.

## 2 Background

### 2.1 Natural Language Understanding

The recent increase in research on dialog systems was a catalyst for research in NLU. Systems like the Dual Intent and Entity Transformer (DIET) (Bunk et al., 2020) or the Dual Sentence Encoders (Henderson et al., 2020; Cer et al., 2018) focus on lightweight models. DIET model uses pre-trained word embeddings like BERT or ConVeRT as dense features. It combines these features with character n-grams and passes them through a Transformer with two instead of the usual 12 encoder layers. Smaller models are competitive with large-scale models in terms of performance and are much faster in training and inference. Full-size Transformer models are used for NLU also and achieve a stronger performance than the efficient architectures (Mehri et al., 2020; Mehri and Eric, 2021).

<sup>2</sup><https://cloud.google.com/dialogflow>

<sup>3</sup><https://cai.tools.sap>

<sup>4</sup>GitHub link available with camera-ready version. Until then, please download the source codes from <https://bit.ly/3ltjPd2>

### 2.2 Adapters

Instead of replacing the original model for a smaller one, adapters (Rebuffi et al., 2017; Houlsby et al., 2019) are a lightweight addition to Transformer models. An adapter is a small set of parameters inserted between the original model’s layers, usually a pre-trained transformer model. For a text classification task as in our paper, the adapter also adds a classification layer on top of the original model. During training, the parameters of the original Transformer model are frozen, and only the parameters of the adapter are modified. The performance measured in accuracy or F1-score is similar to full fine-tuning on most tasks (Peters et al., 2019). In dialog systems research, adapters were used for the tasks dialog state tracking, response retrieval (Hung et al., 2021) and neural end-to-end dialog (Madotto et al., 2020), but the resource efficiency was not investigated by these works.

### 2.3 Modular Dialog Systems

Dialog systems are often composed of multiple sub-dialog systems in practical applications. Nehring and Ahmed (2021) presented the modular dialog system (MDS) architecture to combine several dialog systems into one. Each dialog system is called a module and performs its NLU downstream task. The module selection (MS) component selects the appropriate module for each user utterance. MS is a text classification task very similar to ID. The deployment of large NLU models is also an issue in MDS. A MDS with 5 modules uses 6 models ( $5 \times \text{NLU}$  and  $1 \times \text{MS}$ ). Using our approach, one can reduce the entire MDS’s memory requirements considerably.

## 3 Experiments

### 3.1 Datasets

We use three datasets in our experiments. We compute the performance metric of the NLU models separately for each dataset.

The HWU64 dataset (Liu et al., 2019) is a dataset for NLU. It contains 25k user utterances from the domain of personal assistants, such as “set my alarm to 07:30” or “play next song”. HWU64 spans 64 intents and 21 domains. NLU models from Mehri and Eric (2021) were trained on the DialoGLUE dataset (Mehri et al., 2020) which comprises several datasets, including HWU64. However, the DialoGLUE version of HWU64 contains

only 11k user utterances<sup>5</sup>. For better comparability, we refer to this dataset as HWU64-DialoGLUE and compute the performance metrics of our models for HWU64 and HWU-DialoGLUE.

Another dataset for ID is CLINC150 (Larson et al., 2019). It contains 20k user utterances spanning 150 intents and ten domains. CLINC150 also includes out-of-scope utterances that do not belong to any intents, but are not used in our experiment. The original version and the DialoGLUE version of CLINC150 are identical.

In both datasets, each utterance has exactly one intent label. We split the datasets into train, valid, and test partitions. In CLINC150 and HWU64-DialoGLUE, we reuse the train / valid / test split given by DialoGLUE. In HWU64, we split the data into 80% train, 10% valid and 10% test.

### 3.2 Models

We conduct our experiments on two models. The BERT model is the standard BERT architecture (Devlin et al., 2019) with a sequence classification head. The output layer consists of one neuron for each target label. We use the implementation of the Hugging Face transformers library (Wolf et al., 2020).

For the BERT+Adapter model, we use the implementation from Pfeiffer et al. (2021).

Further, we compare our models to other models for ID. Mehri et al. (2020) used a BERT model similar to ours for NLU. Further, Mehri et al. (2020) proposed ConvBERT, a BERT model that is pre-trained on conversational text instead of on general-purpose text. The best performing model is, to the best of our knowledge, a ConvBERT model with additional *observers* and *examples* (Mehri and Eric, 2021). Observers are tokens that are not attended to and are an alternative to the [CLS] token as a semantic representation of utterances *examples* (Mehri and Eric, 2021). Example-driven training improves the prediction step by assigning inputs with similar semantic meanings a similar representation in the semantic space *examples* (Mehri and Eric, 2021).

### 3.3 Experimental setup

We measure the quality of ID using accuracy, in line with Mehri and Eric (2021); Mehri et al. (2020); Casanueva et al. (2020); Henderson et al. (2020);

<sup>5</sup>We contacted the author, and he confirmed that this is a bug.

Bunk et al. (2020), in both the modular and the non-modular scenario. The non-modular scenario is the standard ID setting in which a single DS processes the whole dataset. In the modular scenario, we split the dataset into multiple parts which we use to generate modules in a MDS. We split the dataset along the domains instead of the intents, i.e. all intents of one domain belong to a single module. We generated MDS with three and ten modules.

We use the same model architecture for ID and MS for the modular experiments. MS models are trained on the same training sets as the ID models, but they use the domain labels as target labels. Further, we perform a grid search to find optimal hyperparameters for learning rate and training batch size using the HWU64-DialoGLUE dataset. The grid search determined a learning rate of  $10^{-5}$  for BERT and  $10^{-3}$  for BERT+Adapter and a batch size of 256. We also apply the optimal hyperparameters obtained on this dataset to the experiments on the other datasets. To determine the number of training epochs, we measure the accuracy of ID on the validation set after each training epoch and stop training as soon as the validation accuracy increases by less than 0.2% compared to the last run. In order to measure the models' sizes, we serialized them to disk and measured their size in MB.

We train and evaluate the models and measure the execution time on the HWU64-DialoGLUE dataset. First, we measure the training and inference duration<sup>6</sup>. Then, we calculate how many samples are processed per second (SPS). SPS makes it easier to compare the training duration of differently sized datasets to each other. We repeat this process ten times and average the resulting durations. All experiments were run on a machine with an RTX6000 graphics card.

## 4 Results

### 4.1 Model size

While BERT requires 417.75 MB and DIET 63.00 MB for an ID model trained on the HWU64 dataset with 64 intent labels, an adapter only needs 5.87 MB of memory. Figure 1 shows how the frameworks scale in terms of memory when the number of models grows. DIET has the smallest memory requirements when deploying a single model, BERT is considerably higher, and BERT+Adapter

<sup>6</sup>We measured the time for processing samples only, excluding other components of training and inference, such as initially loading a model from disk to memory.

Model	Dataset	SPS train	SPS inference
BERT	CLINC150	8.67 (0.06)	792.17 (5.10)
BERT	HWU64	8.20 (0.06)	777.59 (20.45)
BERT	HWU64-DialoGLUE	16.44 (0.10)	791.26 (6.45)
<b>average BERT</b>		<b>11.10</b>	<b>787.01</b>
BERT+Adapter	CLINC150	132.35 (0.52)	752.39 (7.15)
BERT+Adapter	HWU64	138.96 (1.21)	750.19 (19.59)
BERT+Adapter	HWU64-DialoGLUE	209.28 (1.21)	762.72 (4.40)
<b>average BERT+Adapter</b>		<b>160.20</b>	<b>755.10</b>

Table 1: Processing speed of the different models and dataset in samples per second (SPS). The first number is the mean SPS over ten runs, the second number in brackets is the standard deviation over the ten runs. The table lists the average SPS for each model over all datasets.

is slightly higher than BERT. However, when the number of modules grows, the vanilla BERT setup does not scale well. While the slopes of both BERT+Adapter and Rasa are much smaller, BERT+Adapter outperforms Rasa after seven modules.

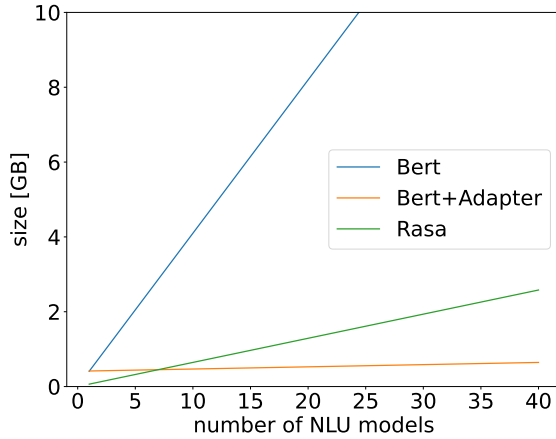


Figure 1: Comparison of the total amounts of memory needed to store the respective framework in GB.

## 4.2 Training and inference time

Table 1 shows SPS, the number of samples processed per second during training and inference. BERT+Adapter trains considerably faster. BERT+Adapter is slightly faster at inference.

## 4.3 Quality of intent detection

Table 2 shows the accuracy values of various models for the task of intent detection in both the non-modular scenario (number of modules 1) and the modular scenario (number of modules 3 and 10). Mehri and Eric (2021) achieve the best results

across both HWU64-DialoGLUE and CLINC150. We list two results of BERT models, one from our experiments and one from Mehri et al. (2020). We conclude that this difference is due to different hyperparameters used during training.

BERT+Adapter achieves higher scores than BERT in our implementation and similar results to BERT in the implementation of Mehri and Eric (2021). Moreover, in the modular setting, the BERT+Adapter model shows a better performance than the BERT model.

## 5 Discussion

The model sizes (section 4.1) show that the DIET architecture is most resource-efficient for the deployment of fewer than five models. On the other hand, when the number of models is above six, the BERT+Adapter architecture is more efficient. In addition, the memory saving increases when more models are deployed which is particularly important for chatbot-as-a-service platforms.

SPS improves significantly for training and slightly for inference. We hypothesize that the difference in training time results from the substantially smaller number of parameters that have to be trained in BERT+Adapter, compared to vanilla BERT. Shorter training times are generally favorable. In the chatbot-as-a-service scenario, a shorter training time is helpful for the designer of the DS, because he does not have to wait for training and has a better user experience.

The experiment shows that the accuracy of ID of BERT+Adapter is comparable to the performance of a BERT model. Still, both models do not meet the state of the art. We leave it for future research to find out if the ConvBERT+ model is compatible



Num. Modules	Model	HWU64	HWU64- DialoGLUE	CLINC150
1	BERT	87.65%	85.50%	91.00%
	BERT+Adapter	<b>89.60%</b>	89.03%	93.04%
	BERT (Mehri et al., 2020)	-	89.97%	95.93%
	ConvBERT (Mehri et al., 2020)	-	90.43%	97.07%
	ConvBERT+ (Mehri and Eric, 2021)	-	<b>93.03%</b>	<b>97.31%</b>
3	BERT	84.94%	83.92%	88.80%
	BERT+Adapter	<b>89.04%</b>	<b>86.25%</b>	<b>91.00%</b>
10	BERT	83.07%	80.67%	88.00%
	BERT+Adapter	<b>88.80%</b>	<b>83.55%</b>	<b>93.31%</b>

Table 2: Accuracy values as percentage of different models, datasets and number of modules. Number of modules=1 is the non-modular scenario. The best performing values are highlighted in bold.

with adapters while still performing just as well.

## 6 Practical considerations on the deployment of NLU models with adapters

Regarding the deployment of NLU models with adapters in a real-world settings, we assume a scenario in which a company offers many user-generated chatbots that are all different from each other in their training data and use case and potentially even come in many languages.

To host 10 languages with 100 chatbots each, we would count  $10 \times \text{Size BERT} + 10 \times 100 \times \text{Size of Adapter} \approx 11\text{GB}$ . 11 GB fit in the memory of a single commercial GPU. The server that hosts the 1000 NLU models would run at all times, keeping the models in memory. In the case of heavy usage, one can horizontally scale out this architecture across several GPUs or several servers.

One downside of our architecture is security and multi-tenant considerations. In multi-tenant environments, it is often a requirement to separate the data of different clients in the architecture, such that data of different clients never get mixed in the same pipeline. The sharing approach suggested in this paper would require a separate pre-trained BERT model for each tenant and language in the multi-tenant environment, decreasing resource savings.

## 7 Conclusion and Future Work

We examined if the BERT+Adapter architecture is favorable for deploying NLP models compared to BERT and DIET. BERT+Adapter can save a considerable amount of memory compared to BERT and

DIET. Also, BERT+Adapter is considerably faster than BERT. The accuracy of ID of BERT+Adapter is still comparable to vanilla BERT.

To the best of our knowledge, our work is the first to point out that adapters can host many NLU models more efficiently. We believe that this approach should be applicable in any environment that hosts multiple pre-trained models. It is especially useful in environments with a skewed usage scheme, meaning that only a few models are used often, and many models are used rarely, such as a chatbots-as-a-service environment. Further, we believe it is applicable in many settings where users can generate and fine-tune their models, e.g. named entity recognition as a service with user-generated content. Finally, adapters are not limited to models with the same task. They can also be used in environments where models for different tasks are deployed (Houlsby et al., 2019), e.g. one adapter for Named Entity Recognition and another adapter for Sentiment Analysis that share the same pre-trained Transformer model.

## References

- Tanja Bunk, Daksh Varshneya, Vladimir Vlasov, and Alan Nichol. 2020. *DIET: Lightweight language understanding for dialogue systems*. *arXiv*.
- Iñigo Casanueva, Tadas Temčinas, Daniela Gerz, Matthew Henderson, and Ivan Vulić. 2020. *Efficient intent detection with dual sentence encoders*. In *Proceedings of the 2nd Workshop on Natural Language Processing for Conversational AI*, pages 38–45, Online. Association for Computational Linguistics.
- Daniel Cer, Yinfei Yang, Sheng yi Kong, Nan Hua, Nicole Limtiaco, Rhomni St. John, Noah Constant,

386	Mario Guajardo-Cespedes, Steve Yuan, Chris Tar,	444
387	Yun-Hsuan Sung, Brian Strope, and Ray Kurzweil.	445
388	2018. <a href="#">Universal sentence encoder</a> . <i>arXiv</i> .	446
389	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and	447
390	Kristina Toutanova. 2019. <a href="#">BERT: Pre-training of</a>	448
391	<a href="#">deep bidirectional transformers for language under-</a>	449
392	<a href="#">standing</a> . In <i>Proceedings of the 2019 Conference of</i>	450
393	<i>the North American Chapter of the Association for</i>	451
394	<i>Computational Linguistics: Human Language Tech-</i>	452
395	<i>nologies, Volume 1 (Long and Short Papers)</i> , pages	
396	4171–4186, Minneapolis, Minnesota. Association for	
397	Computational Linguistics.	
398	Matthew Henderson, Iñigo Casanueva, Nikola Mrkšić,	
399	Pei-Hao Su, Tsung-Hsien Wen, and Ivan Vulić. 2020.	
400	<a href="#">ConveRT: Efficient and accurate conversational repre-</a>	
401	<a href="#">sentations from transformers</a> . In <i>Findings of the Asso-</i>	
402	<i>ciation for Computational Linguistics: EMNLP 2020</i> ,	
403	pages 2161–2174, Online. Association for Computa-	
404	tional Linguistics.	
405	Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski,	
406	Bruna Morroni, Quentin De Laroussilhe, Andrea	
407	Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019.	
408	<a href="#">Parameter-efficient transfer learning for NLP</a> . In	
409	<i>Proceedings of the 36th International Conference</i>	
410	<i>on Machine Learning</i> , volume 97 of <i>Proceedings</i>	
411	<i>of Machine Learning Research</i> , pages 2790–2799.	
412	PMLR.	
413	Chia-Chien Hung, Anne Lauscher, Simone Paolo	
414	Ponza, and Goran Glavaš. 2021. <a href="#">DS-TOD: Ef-</a>	
415	<a href="#">ficient domain specialization for task oriented dialog</a> .	
416	<i>arXiv</i> .	
417	Stefan Larson, Anish Mahendran, Joseph J. Peper,	
418	Christopher Clarke, Andrew Lee, Parker Hill,	
419	Jonathan K. Kummerfeld, Kevin Leach, Michael A.	
420	Laurenzano, Lingjia Tang, and Jason Mars. 2019. <a href="#">An</a>	
421	<a href="#">evaluation dataset for intent classification and out-of-</a>	
422	<a href="#">scope prediction</a> . In <i>Proceedings of the 2019 Confer-</i>	
423	<i>ence on Empirical Methods in Natural Language Pro-</i>	
424	<i>cessing and the 9th International Joint Conference</i>	
425	<i>on Natural Language Processing (EMNLP-IJCNLP)</i> ,	
426	pages 1311–1316, Hong Kong, China. Association	
427	for Computational Linguistics.	
428	Xingkun Liu, Arash Eshghi, Pawel Swietojanski, and	
429	Verena Rieser. 2019. <a href="#">Benchmarking Natural Lan-</a>	
430	<a href="#">guage Understanding Services for building Conversa-</a>	
431	<a href="#">tional Agents</a> . In <i>Proceedings of the Tenth Inter-</i>	
432	<i>national Workshop on Spoken Dialogue Systems</i>	
433	<i>Technology (IWSDS)</i> , Ortigia, Siracusa (SR), Italy.	
434	Springer.	
435	Andrea Madotto, Zhaojiang Lin, Yejin Bang, and Pas-	
436	cale Fung. 2020. <a href="#">The adapter-bot: All-in-one con-</a>	
437	<a href="#">trollable conversational model</a> . <i>arXiv</i> .	
438	Shikib Mehri and Mihail Eric. 2021. <a href="#">Example-driven</a>	
439	<a href="#">intent prediction with observers</a> . In <i>Proceedings of</i>	
440	<i>the 2021 Conference of the North American Chap-</i>	
441	<i>ter of the Association for Computational Linguistics:</i>	
442	<i>Human Language Technologies</i> , pages 2979–2992,	
443	Online. Association for Computational Linguistics.	
	Shikib Mehri, Mihail Eric, and Dilek Hakkani-Tur.	444
	2020. <a href="#">DialoGLUE: A natural language understand-</a>	445
	<a href="#">ing benchmark for task-oriented dialogue</a> . <i>ArXiv</i> .	446
	Jan Nehring and Akhyar Ahmed. 2021. <a href="#">Normal-</a>	447
	<a href="#">isierungsmethoden für Intent Erkennung Modu-</a>	448
	<a href="#">larer Dialogsysteme</a> . In <i>Tagungsband der 32.</i>	449
	<i>Konferenz. Elektronische Sprachsignalverarbeitung</i>	450
	<i>(ESSV-2021)</i> , March 3-5, Berlin, Germany. TUD-	451
	press.	452
	Matthew E. Peters, Sebastian Ruder, and Noah A. Smith.	453
	2019. <a href="#">To tune or not to tune? adapting pretrained</a>	454
	<a href="#">representations to diverse tasks</a> . In <i>Proceedings of</i>	455
	<i>the 4th Workshop on Representation Learning for</i>	456
	<i>NLP (RepL4NLP-2019)</i> , pages 7–14, Florence, Italy.	457
	Association for Computational Linguistics.	458
	Jonas Pfeiffer, Aishwarya Kamath, Andreas Rücklé,	459
	Kyunghyun Cho, and Iryna Gurevych. 2021.	460
	<a href="#">AdapterFusion: Non-destructive task composition</a>	461
	<a href="#">for transfer learning</a> . In <i>Proceedings of the 16th Con-</i>	462
	<i>ference of the European Chapter of the Association</i>	463
	<i>for Computational Linguistics: Main Volume</i> , pages	464
	487–503, Online. Association for Computational Lin-	465
	guistics.	466
	Jonas Pfeiffer, Andreas Rücklé, Clifton Poth, Aishwarya	467
	Kamath, Ivan Vulić, Sebastian Ruder, Kyunghyun	468
	Cho, and Iryna Gurevych. 2020. <a href="#">AdapterHub: A</a>	469
	<a href="#">framework for adapting transformers</a> . In <i>Proceedings</i>	470
	<i>of the 2020 Conference on Empirical Methods in Nat-</i>	471
	<i>ural Language Processing: System Demonstrations</i> ,	472
	pages 46–54, Online. Association for Computational	473
	Linguistics.	474
	Sylvestre-Alvise Rebuffi, Hakan Bilen, and Andrea	475
	Vedaldi. 2017. <a href="#">Learning multiple visual domains</a>	476
	<a href="#">with residual adapters</a> . In <i>Advances in Neural In-</i>	477
	<i>formation Processing Systems</i> , volume 30. Curran	478
	Associates, Inc.	479
	Emma Strubell, Ananya Ganesh, and Andrew McCal-	480
	lum. 2019. <a href="#">Energy and policy considerations for</a>	481
	<a href="#">deep learning in NLP</a> . In <i>Proceedings of the 57th</i>	482
	<i>Annual Meeting of the Association for Computational</i>	483
	<i>Linguistics</i> , pages 3645–3650, Florence, Italy. Asso-	484
	ciation for Computational Linguistics.	485
	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob	486
	Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz	487
	Kaiser, and Illia Polosukhin. <a href="#">Attention is all you</a>	488
	<a href="#">need</a> . In <i>Advances in Neural Information Processing</i>	489
	<i>Systems</i> , pages 5999–6009. Curran Associates, Inc.	490
	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien	491
	Chaumond, Clement Delangue, Anthony Moi, Pier-	492
	ric Cistac, Tim Rault, Remi Louf, Morgan Funtow-	493
	icz, Joe Davison, Sam Shleifer, Patrick von Platen,	494
	Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu,	495
	Teven Le Scao, Sylvain Gugger, Mariama Drame,	496
	Quentin Lhoest, and Alexander Rush. 2020. <a href="#">Trans-</a>	497
	<a href="#">formers: State-of-the-art natural language processing</a> .	498
	In <i>Proceedings of the 2020 Conference on Empirical</i>	499
	<i>Methods in Natural Language Processing: System</i>	500

501 *Demonstrations*, pages 38–45, Online. Association  
502 for Computational Linguistics.