# Effective Ensembling of Transformer based Language Models for Acronyms Identification

## Divesh Kubal, Apurva Nagvenkar

CRIMSON AI

divesh.kubal@crimsoni.ai, apurva.nagvenkar@crimsoni.ai

## Abstract

An acronym can be viewed as a word which is constructed by taking initial components from a phrase. This study deals with the problem of identification and extraction of an acronym's short and long-form. The proposed approach solves the Acronym Identification (AI) task mentioned in the scientific document understanding (SDU@AAAI-21) task. This paper model the Acronym identification task as the sentence level sequence-labeling problem. The proposed method is computed by an ensemble of various Language Models trained by hyper-parameter tuning. This ensembling technique is then coupled with post-processing steps to extract the best possible predictions. The trained model's performance is evaluated against standard evaluation metrics such as precision, recall, and F1-score. The final model achieves an F1 score of 95.60%, a precision of 93.97%, and a recall of 97.95% on the development dataset. On the test data, the proposed model achieves an F1 score of 92.08%, a precision of 89.70%, and the highest recall of 94.59%, compared to other participants results in the competition.

## Introduction

The acronyms in the technical/scientific domain are increasing at an exponential rate. This is due to a huge amount of research conducted in the technical and non-technical domains. The term "acronym" is defined as the name given to a particular word or a phrase by taking the first letters of each word of a phrase (Mack 2012). For example, 'ANN' is an acronym that stands for 'Artificial Neural Network'. The more common or general term used is "abbreviation". Abbreviations encompass acronyms and few abbreviations that use letters other than initial characters of phrases, such as 'Mr.' for 'Ministers'. Hence, there is a thin line that distinguishes acronyms from abbreviations. The acronyms serve a vital role in writing science or research-related technical documents, patents, etc., by preventing content repetition. This enables speeding the reading process and paving the way for an easier understanding of the content written inside a document.

Several techniques have been proposed to extract acronyms from a given input text corpus. These systems are

rule-based(Schwartz and Hearst 2002) or machine learning-based(Jacobs, Itai, and Wintner 2020; Kuo et al. 2009; Liu, Liu, and Huang 2017). Some techniques use word-embedding techniques (Kirchhoff and Turner 2016) to extract acronyms. There are several packages available in python (Cook 2019; Schwartz and Hearst 2002) which extract acronyms and their expansions. The understanding of acronyms and their expansions is an important task in the following use-cases:

- Text understanding: There can be multiple expansions of an acronym. Hence, identification of correct contextual meaning is important to understand the text in an unambiguous manner.

- Information retrieval: When a document is queried by inputting a query containing an acronym, the results should contain the relevant results.

- Machine translation: Acronyms posses a big challenge when translating a source language to its target language.

- Text Summarization: It is advisable to use an acronym counterpart of its expansion to summarize the text.



Figure 1: Examples of Acronym Identification

This paper presents an effective ensembling based approach to automatically extract acronyms along with their extended/long-forms. The proposed approach combines ensembling Language Models + hyper-parameter tuning + post-processing to get the best possible results. This paper is structured by giving a survey of related work where a quick brief about existing approaches used to solve the problem of Acronym Identification (Veyseh et al. 2020a) is given. After

the related work section, an in-depth explanation of the proposed system architecture is discussed. A thorough comparative analysis of results on different scenarios is explained in the section results and discussions. Finally, this paper concludes by giving a quick conclusion and future directions.

## Related Work

The existing Acronym Identification systems can be broadly classified into rule-based, features & machine learning-based and Deep Learning-based as depicted in figure 2. The rule-based systems mainly consist of techniques that use rules, patterns and regular expressions to extract acronyms and their expansions. The machine learning-based techniques first extract features required for Acronym Identification, and then a classifier (like Support Vector Machine, Conditional Random Fields) is trained over these features. The Deep Learning-based techniques use state-of-the-art algorithms like Recurrent Neural Networks (RNN), Long Short Term Memory (LSTM), transformer-based models, etc.
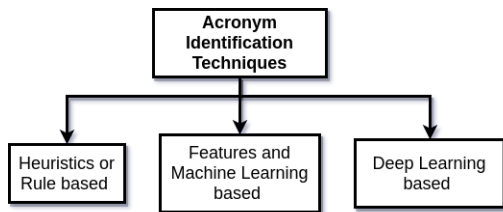


Figure 2: Acronym Identification Systems Broad Classification

The Acronym Identification techniques can be traced back to the year 1999 where (Taghva and Gilbreth 1999) proposed a system named as Acronym Finding Program (AFP). It is a simple regex-based system that identified candidate acronyms (upper-case words of three to ten grams). It attempts to find acronym expansion by scanning a 2n-window (n: number of letters in candidate acronyms). The final system was evaluated on 1328 files. Many regex-based Acronym Identification systems were proposed afterward. The issue with regex-based techniques was that it is impossible to incorporate all the possible rules, and it required more manual work to identify patterns and then to write rules.

A popular approach to detect acronyms was proposed by (Schwartz and Hearst 2002)[1]. This technique is capable of extracting complicated acronyms, and it's expansions. It works in two stages. The first stage identifies candidate acronyms by using predefined patterns ("acronym" and "expansion" and vice-versa). In the second stage, the overlapping characters in an acronym and expansion are counted, and finally, this count is compared to a specified threshold. As this technique cannot store contextual information, it failed to extract acronyms and their expansions with long-term dependencies.

---

[1]https://github.com/philgooch/abbreviation-extraction

One of the machine learning-based Acronym Identification systems proposed by (Kuo et al. 2009) extracts features and then uses various algorithms like Support Vector Machine, Naive Bayes, Logistic Regression, and Monte-Carlo Sampling Logistic Regression for training.

(Liu, Liu, and Huang 2017) proposed a Latent-state Neural Conditional Random Fields model (LNCRF) system which couples Conditional Random Fields with nonlinear hidden layers. This system models the task of Acronym Identification as a sequence labeling problem, and it surpasses many baseline models that were in existence at that time.

A machine learning-based approach proposed by (Jacobs, Itai, and Wintner 2020) aims to extract acronyms by automatically building an acronym-based dictionary from an unannotated dataset. One of the critical parts of this system is that it is capable of extracting non-local acronyms too. It means extracting the expanded form of an acronym even if the short form is not present in a given sentence.

Another approach that utilizes Long Short Term Memory - Conditional Random Field (LSTM-CRF) proposed by (Veyseh et al. 2020b) provides an in-depth explanation about overall Acronym identification and Disambiguation implementations.

The rule-based and machine learning-based models were not able to capture the contextual information. An acronym extraction system that uses machine learning combined with a neural network based-contextual model was proposed by (Kirchhoff and Turner 2016). This model can store contextual information and hence can also be used to solve the task of acronym disambiguation.

## Effective Ensembling of Language Models for Sequence Labeling Problem: EELM-SLP

The proposed system architecture is depicted in figure 3. The main components of Effective Ensembling of Language Models for Sequence Labeling Problem (EELM-SLP) System Architecture are as follows:

1. Data Acquisition/Collection.

2. Finetuning of transformer-based Language Models for sequence-labeling task.

3. Hyperparameter tuning and retraining of Language Models.

4. Ensembling and postprocessing to obtain final predictions.

5. Evaluation on development and test datasets.

### Data Acquisition/Collection

The entire data can be downloaded from https://github.com/amirveyseh/AAAI-21-SDU-shared-task-1-AI. A sample snapshot of the data is depicted in figure 1. The data is bifurcated into training, development, and test data. The training dataset consists of 14006 labeled sentences (Veyseh et al. 2020b). The development data consists of 1717 labeled sentences. Each datapoint in train and test data has 'id' representing train or development datapoint ID, 'tokens', a
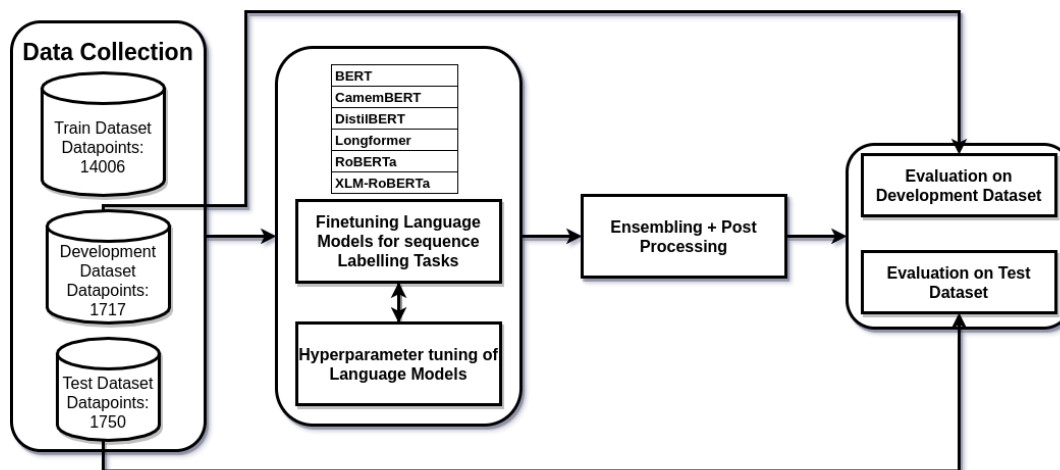
Figure 3: Effective Ensembling of Language Models for Sequence Labeling Problem (EELM-SLP) System Architecture

list containing the sentence split into individual tokens, and 'labels' are the token-wise annotations for the tokens associated with them. The test dataset consists of about 1750 datapoints where only 'id' and 'tokens' are provided. The 'labels' depict the short and long-form acronyms in BIO format (short for inside, outside, beginning). B-long, I-long, B-short, and I-short are the labels used.

## Finetuning of transformer-based Language Models for sequence-labeling task.

The proposed approach uses six transformer-based pre-trained Language Models, who are fined tuned on a downstream sequence labeling task. The following Language Models are used: BERT (Devlin et al. 2018), RoBERTa (Liu et al. 2019), XLM-RoBERTa (Conneau et al. 2019), CamemBERT (Martin et al. 2019), Longformer (Beltagy, Peters, and Cohan 2020) and DistilBERT (Sanh et al. 2019). BERT is trained with the objective of masked language modeling (MLM) and the next sentence prediction (NSP). For BERT, RoBERTa, XLM-RoBERTa, CamemBERT, and DistilBERT, the base-cased configurations are used. AllenAI-base-4096 configuration used for longformer.

The configuration used for finetuning Language Models are as follows:

- The BERT and CamemBERT model is 12-layered having 768-hidden layers, 12-heads, and 109M parameters.

- RoBERTa has the same configuration as that of BERT, except it has 125M parameters.

- XLM-RoBERTa has approximately 270M parameters with 12-layers, 768-hidden-state, 3072 feed-forward hidden-state, 8-heads, which is pretrained on Common-Crawl data in 100 languages.

- DistilBERT has less parameters of 65M, and it is 6-layered. DistilBERT model is distilled from the BERT-based configuration.

- Longformer has approximately 149M parameters. Here 4096 represents that the model is pretrained on documents of maximum length 4096.

All the above-mentioned Language Models are finetuned on the training dataset for sentence-level sequence labeling. Every model is trained separately, and the latest iteration is stored.

## Hyperparameter tuning of Language Models for Sequence Labeling Task

It was observed that during the inference phase, the output length was truncated to 128 tokens as it was the default 'max_seq_length' parameter. To preserve the entire token length, two parameters were used heavily. These two parameters are 'sliding_window' and 'max_seq_length'. The sliding_window prevents the truncation of sentences by splitting the input sequence into multiple windows if it exceeds the default maximum sequence value. The sliding window problem is that the contextual information is broken while predicting, and hence it was not used. Some experiments were carried on with max_seq_length parameter, and it was observed that the model was performing better when max_seq_length was kept to 350. Other max_seq_length values were 128, 256, 300, 400, 450, and 512. The language models were finetuned by using 'Simple Transformers' (Rajapakse 2020)[2] library and by using 'Hugging Face' (Wolf et al. 2019)[3] pretrained models. GEFORCE RTX 2080 Ti 11GB GPU was used to train all the models with a 32GB primary memory system.

Table 1 depicts the final list of hyperparameters used to finetune the Language Models on the sequence labeling task. The 'adam_epsilon' is the epsilon value to use for adam optimizer. 'best_model_dir' is the directory where the best model automatically gets stored after the completion of all epochs. 'cache_dir' is where the processed data is stored,

---

[2]https://simpletransformers.ai/
[3]https://huggingface.co/

| Hyperparameter | Value |
|---|---|
| adam_epsilon | 1e-08 |
| best_model_dir | outputs/best_model/ |
| cache_dir | cache/ |
| train_custom_parameters_only | False |
| dataloader_num_workers | 4 |
| do_lower_case | False |
| early_stopping_consider_epochs | False |
| early_stopping_delta | 0 |
| early_stopping_metric | eval_loss |
| early_stopping_metric_minimize | True |
| early_stopping_patience | 3 |
| encoding | null |
| eval_batch_size | 8 |
| evaluate_during_training | False |
| evaluate_during_training_silent | True |
| evaluate_during_training_steps | 2000 |
| evaluate_during_training_verbose | False |
| fp16 | True |
| gradient_accumulation_steps | 1 |
| learning_rate | 4e-05 |
| local_rank | -1 |
| logging_steps | 50 |
| max_grad_norm | 1.0 |
| max_seq_length | 350 |
| multiprocessing_chunksize | 500 |
| n_gpu | 1 |
| no_cache | False |
| no_save | False |
| num_train_epochs | 50 |
| output_dir | outputs/ |
| overwrite_output_dir | True |
| process_count | 4 |
| reprocess_input_data | True |
| save_best_model | True |
| save_eval_checkpoints | True |
| save_model_every_epoch | True |
| save_steps | 20000 |
| save_optimizer_and_scheduler | True |
| silent | False |
| train_batch_size | 16 |
| use_cached_eval_features | False |
| use_early_stopping | False |
| use_multiprocessing | True |
| warmup_ratio | 0.06 |
| warmup_steps | 2628 |
| weight_decay | 0 |
| classification_report | False |
| labels_list | B-long, I-long, B-short, I-short, O |
| lazy_loading | False |
| lazy_loading_start_line | 0 |

Table 1: Final hyperparameters list used to finetune Language models for Sequence Labeling Task)

which is consumable by PyTorch (Paszke et al. 2019). The 'train_custom_parameters_only' is kept False as we are utilizing all the hyperparameters available in Language Models. The 'dataloader_num_workers' specifies the number of CPUs which will be used for data processing. As this is a sequence labeling task, the 'cased' configurations of Language Models are used, and hence the 'do_lower_case' is kept to be False. The 'early_stopping_metric' is kept to be the evalua-

tion loss, and the 'patience' is kept to be 3. The training process avoids evaluation while training. Hence, the parameter 'evaluate_while_training' is kept to be False. 'fp16' corresponds to the 16-bits training and sometimes also referred to as mixed-precision training is kept to True. The 'gradient_accumulation_step' is kept to 1 and the 'learning_rate' to be 4e-05. The training was carried on a single Graphics Processing Unit (GPU) system, and so the 'n_gpu' was fixed to 1. As we don't want the model to save very frequently and save the secondary storage space, the 'save_steps' is kept to a higher value of 20,000. The training epochs were set to 50, but it was observed that the models were able to train around 20 epochs. Hence, the training process was stopped manually as soon as the loss was stagnant and at the lowest point. The 'labes_list' corresponds to the target labels and it was set to ["B-long", "I-long", "B-short", "I-short", "O"]. The 'max_seq_length' is one of the hyperparameters which was finetuned. These are some of the important hyperparameters known to have a huge impact on finetuning Language Models for sequence labeling task. After the final hyperparameters were computed, all the Language Models were again finetuned for Sentence Level Sequence Labeling task to identify acronyms from a given input text.

## Ensembling and postprocessing to obtain final predictions

After all the language models were finetuned on training data, the ensembling technique was applied. The advantage of performing ensembling is to construct a robust discriminator or model by using comparatively weaker models. In this paper, a novel-ensembling approach based on RoBERTa prioritizing was used. The table 2 shows that the F1-score of RoBERTa for the development dataset is higher than other algorithms. The steps are as follows:

1. Initially, all the predictions from six pre-trained language models were extracted.

2. For each sentence, token level predictions were extracted for all six models. Hence for each token for every sentence, six predictions were obtained. For simplicity, the predictions can be from any of the following labels - 'O', 'B-short', 'I-short', 'B-long' and 'I-long'.

3. If all the predictions are 'O' and any of the models gives any predictions from the bouquet of 'B-short', 'I-short', 'B-long' and 'I-long', then the prediction other than 'O' is considered.

4. If there is a conflict between 'B-short', 'I-short', 'B-long', and 'I-long', then the prediction suggested by RoBERTa is taken into consideration. But if, in this case, if RoBERTa's prediction is 'O', then the prediction with the highest frequency is considered.

The above ensembling technique gave rise to an increase in recall. The ensembling approach's limitation is that it predicted labels having 'I-short' and 'I-long' as the beginning tags, which lowered precision for some tokens. To overcome this problem, the following post-processing rules were applied:

| Algorithm (hyperparameters finetuned) | Precision | Recall | F-1 Score | LONG | | | SHORT | | |
|---|---|---|---|---|---|---|---|---|---|
| | | | | Precision | Recall | F1-score | Precision | Recall | F1-score |
| **Ensembling + Post Processing (EELM-SLP)** | **93.35** | **97.95** | **95.6** | **91.43** | **97.27** | **94.26** | **95.27** | **98.63** | **96.92** |
| Ensembling Technique | 92.36 | 97.2 | 94.72 | 90.9 | 96.65 | 93.69 | 93.81 | 97.74 | 95.74 |
| RoBERTa | 93.19 | 91.2 | 92.18 | 91.65 | 92.56 | 92.1 | 94.74 | 89.83 | 92.22 |
| XLM-RoBERTa | 93.05 | 90.89 | 91.95 | 91.49 | 92 | 91.75 | 94.61 | 89.77 | 92.12 |
| BERT | 92.34 | 90.88 | 91.6 | 90.43 | 90.27 | 90.35 | 94.24 | 91.48 | 92.84 |
| CamemBERT | 93.45 | 88.93 | 91.13 | 92.32 | 90.64 | 90.98 | 95.58 | 87.23 | 91.21 |
| Longformer | 93.14 | 88.59 | 90.81 | 90.49 | 90.2 | 90.34 | 95.8 | 86.97 | 91.17 |
| DistilBERT | 91.04 | 90.51 | 90.77 | 88.47 | 88.03 | 88.25 | 93.6 | 92.98 | 93.29 |

Table 2: Evaluation results for Development Dataset (all values are in percentage)

| Algorithm | Precision | Recall | F-1 Score |
|---|---|---|---|
| Ensembling + Post Processing (EELM-SLP) | 89.7 | 94.59 | **92.08** |
| Ensembling Technique | 89.93 | 93.87 | 91.86 |
| RoBERTa after hyperparameter tuning | 90.26 | 92.46 | 91.34 |
| RoBERTa before hyperparameter tuning | 90.85 | 91.73 | 91.29 |

Table 3: Evaluation results for Test Dataset (all values are in percentage)

- If at any point of time the prediction starts with 'I-short', then it is replaced as 'B-short'.
- The pattern 'O', 'I-long' was replaced with 'B-long', 'I-long' to maintain consistency.

By applying the above post-processing rules, the precision improved the proposed, and hence the proposed EELM-SLP system surpassed the other individual trained models on the development dataset.

All the models are trained on 11GB Nividia 2080 Ti GPU for 20-22 epochs, with each epoch taking around three to four minutes. Each model took around 90 minutes to train, and due to early stopping, the training was stopped as soon as the loss stopped reducing. The proposed architecture is efficient and can be integrated into an actual production environment. The prediction pipeline is designed in such a way that it takes advantage of CPU parallelism. So if the input contains a batch of sentences, the predictions will be computed in parallel independent batches and then later combined to form in order as they appear in the original paragraph or list of sentences.

## Results and Discussions

The table 2 depicts the evaluation metrics on the development dataset. It can be seen that the top-3 performers among individual language models based on F1-score are RoBERTa, followed by XLM-RoBERTa, which is followed by BERT. The RoBERTa achieves the highest F1-score of 92.18%, while XLM-RoBERTa achieves the F1-score of 91.95%. The third best performing algorithm BERT achieves the F1-score of 91.6%. The ensembling technique surpasses the hyperparameter finetuned individual Language Models and achieves a precision of 92.36%, recall of 97.2%, and F1-score of 95.6%. However, the proposed approach of Effective Ensembling of Language Models for Sequence Labeling Problem (EELM-SLP) achieves the highest F1-score

of 95.6% recorded on development data. It can be clearly seen that the proposed approach surpasses other individual language models by quite a considerable margin. For all the individual finetuned models and the proposed approach, the class scores are also computed, namely precision, recall, and F1-score for 'short' and 'long' labels.

In table 3, the evaluation on the test dataset is presented on standard evaluation metrics like precision, recall, and F1-score. Initially, the test dataset was evaluated on RoBERTa finetuned model before finetuning. This initial evaluation resulted in precision, recall, and F1-score to be 90.85%, 91.73%, and 91.29%, respectively. The RoBERTa was then finetuned after hyperparameter tuning, which further accelerated the metrics. The hyperparameter tuned RoBERTa achieved the precision of 90.26%, recall of 92.46%, and F1-score of 91.34%. One of the significant highlights is the results obtained after applying the ensembling technique, which improved the F1-score from 91.34% to 91.86%. Finally, the highest F1-score was achieved after the application of post-processing. The final proposed system achieved an F1-score of 92.08%, precision of 89.7%, and recall of 94.59%.

## Conclusion and Future Scope

In this paper, an Effective Ensembling of Language Models for Sequence Labeling (EELM-SLP) is proposed. The ensembling technique's importance can be clearly seen based on the development and test dataset results. The final model achieves a precision of 93.35%, recall of 97.95% and F1-score of 95.6% on the development dataset and precision of 89.7%, recall of 94.59% and F1-score of 92.08% on the test set. There is an improvement on all the three metrics of precision, recall, and F1-score for the proposed approach on the development dataset. Although the recall and F1-score are highest for the test dataset, the precision is lowered. This

might be due to the introduction of false-positives while ensembling and post-processing step. On the one hand, the ensembling and post-processing step increased the recall and F1-score, but on the other hand, it lowered the precision. Hence, the post-processing can be finetuned in-depth. Further, this paper uses all the language models in their 'base' configuration and not 'large' configuration. In the future, experiments can be carried on 'large' configuration based language models, which might improve the scores.

# References

Beltagy, I.; Peters, M. E.; and Cohan, A. 2020. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150* .

Conneau, A.; Khandelwal, K.; Goyal, N.; Chaudhary, V.; Wenzek, G.; Guzmán, F.; Grave, E.; Ott, M.; Zettlemoyer, L.; and Stoyanov, V. 2019. Unsupervised crosslingual representation learning at scale. *arXiv preprint arXiv:1911.02116* .

Cook, B. 2019. ACRONYM: Acronym CReatiON for You and Me. *arXiv preprint arXiv:1903.12180* .

Devlin, J.; Chang, M.-W.; Lee, K.; and Toutanova, K. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* .

Jacobs, K.; Itai, A.; and Wintner, S. 2020. Acronyms: identification, expansion and disambiguation. *Annals of Mathematics and Artificial Intelligence* 88(5): 517–532.

Kirchhoff, K.; and Turner, A. M. 2016. Unsupervised resolution of acronyms and abbreviations in nursing notes using document-level context models. In *Proceedings of the Seventh International Workshop on Health Text Mining and Information Analysis*, 52–60.

Kuo, C.-J.; Ling, M. H.; Lin, K.-T.; and Hsu, C.-N. 2009. BIOADI: a machine learning approach to identifying abbreviations and definitions in biological literature. In *BMC bioinformatics*, volume 10, S7. Springer.

Liu, J.; Liu, C.; and Huang, Y. 2017. Multi-granularity sequence labeling model for acronym expansion identification. *Information Sciences* 378: 462–474.

Liu, Y.; Ott, M.; Goyal, N.; Du, J.; Joshi, M.; Chen, D.; Levy, O.; Lewis, M.; Zettlemoyer, L.; and Stoyanov, V. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692* .

Mack, C. A. 2012. How to write a good scientific paper: acronyms. *Journal of micro/nanolithography, MEMS, and MOEMS* 11(4): 040102.

Martin, L.; Muller, B.; Suárez, P. J. O.; Dupont, Y.; Romary, L.; de la Clergerie, É. V.; Seddah, D.; and Sagot, B. 2019. Camembert: a tasty french language model. *arXiv preprint arXiv:1911.03894* .

Paszke, A.; Gross, S.; Massa, F.; Lerer, A.; Bradbury, J.; Chanan, G.; Killeen, T.; Lin, Z.; Gimelshein, N.; Antiga, L.; et al. 2019. Pytorch: An imperative style, high-performance deep learning library. In *Advances in neural information processing systems*, 8026–8037.

Rajapakse, T. 2020. Simple transformers.

Sanh, V.; Debut, L.; Chaumond, J.; and Wolf, T. 2019. DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108* .

Schwartz, A. S.; and Hearst, M. A. 2002. A simple algorithm for identifying abbreviation definitions in biomedical text. In *Biocomputing 2003*, 451–462. World Scientific.

Taghva, K.; and Gilbreth, J. 1999. Recognizing acronyms and their definitions. *International Journal on Document Analysis and Recognition* 1(4): 191–198.

Veyseh, A. P. B.; Dernoncourt, F.; Nguyen, T. H.; Chang, W.; and Celi, L. A. 2020a. Acronym Identification and Disambiguation shared tasksfor Scientific Document Understanding. *arXiv preprint arXiv:2012.11760* .

Veyseh, A. P. B.; Dernoncourt, F.; Tran, Q. H.; and Nguyen, T. H. 2020b. What Does This Acronym Mean? Introducing a New Dataset for Acronym Identification and Disambiguation. *arXiv preprint arXiv:2010.14678* .

Wolf, T.; Debut, L.; Sanh, V.; Chaumond, J.; Delangue, C.; Moi, A.; Cistac, P.; Rault, T.; Louf, R.; Funtowicz, M.; et al. 2019. HuggingFace's Transformers: State-of-the-art Natural Language Processing. *ArXiv* arXiv–1910.