# Tag and correct: high precision post-editing approach to correction of speech recognition errors

Tomasz Ziętkiewicz

Adam Mickiewicz University in Poznań
ul. Uniwersytetu Poznańskiego 4, 61-614 Poznań, Poland
Samsung R&D Institute Poland
Plac Europejski 1, 00-844 Warsaw, Poland
Email: t.zietkiewic@samsung.com, tomasz.zietkiewicz@amu.edu.pl

*Abstract*—**This paper presents a new approach to the problem of correcting speech recognition errors by means of post-editing. It consists of using a neural sequence tagger that learns how to correct an ASR (Automatic Speech Recognition) hypothesis word by word and a corrector module that applies corrections returned by the tagger. The proposed solution is applicable to any ASR system, regardless of its architecture, and provides high-precision control over errors being corrected. This is especially crucial in production environments, where avoiding the introduction of new mistakes by the error correction model may be more important than the net gain in overall results. The results show that the performance of the proposed error correction models is comparable with previous approaches while requiring much smaller resources to train, which makes it suitable for industrial applications, where both inference latency and training times are critical factors that limit the use of other techniques.**

**Index Terms**: speech recognition, error correction, post-processing, post-editing, natural language processing

## I. Introduction

Automatic Speech Recognition (ASR) models have been developed for over 70 years. During this time, they evolved from machines that could recognize single digits spoken by one person, created for demonstration purposes without production use back in the 1950s [1], to omnipresent voice assistants and speech transcription engines used every day by millions of people around the world. Although speech recognition technology has reached maturity and is production-ready, it is still far from perfect. Professional human transcribers do not reach 100% transcription accuracy, and although recent deep neural network-powered speech recognition systems are reported to slightly outperform humans [2], they still make mistakes. Furthermore, the near-perfect results of Word Error Rate (WER) below 2% are often reported on popular benchmark datasets such as the test subset of the Librispeech corpus [3]. When evaluated on corpora from different domains or recorded under different conditions, the results are far from perfect and require adaptation [4]. In real-world production settings, input to a speech recognition system may change frequently, caused by changes in topics that are interesting to users. Changes in the real world, such as the COVID-19 pandemic, influence the vocabulary used by speakers and may lead to out-of-vocabulary errors. The adaptation process of ASR models takes considerable time and resources. In some settings, direct adaptation of the model is impossible, for example, when using a cloud-based speech recognition service. One way to address the imperfections of speech recognition models mentioned above is to improve their results in a post-editor, a module operating on the textual output of a speech recognition model. Typically, in production environments, the post-editor provides a means of correcting ASR errors manually, for example, with hand-written regular expressions. The process of creating and maintaining such corrections is laborious and requires a lot of experience [5]. Therefore, the post-editor can include an error correction model which learns how to correct errors of a particular ASR system. This approach can be applied regardless of an ASR model architecture, also for systems where direct modification of the model is impossible. It requires only text data to train, and its training requires considerably less resources and time than performing the adaptation of a speech recognition model. This paper presents such an error correction model. In Section II we present an overview of previous work on the subject. Section III describes the data used for training and evaluation. Details of the proposed approach are presented in Section IV. Results and conclusions can be found in Sections V and VI.

## II. Related work

For a review of ASR error detection and correction systems together with a description of ASR evaluation metrics see [6]. Cucu et al. [7] propose error correction using SMT (Statistical Machine Translation) model trained on a relatively small parallel corpus of 2000 ASR transcripts and their manually corrected versions. At an evaluation time, the model is used to "translate" ASR hypothesis into corrected form. The system achieves 10.5% relative WER improvement by reducing WER of the baseline ASR system from 11.4 to 10.2. A similar approach, but using a neural LSTM sequence-to-sequence model and trained on a much larger dataset (40M utterances), is presented in [8]. To produce ASR hypotheses, the authors use a speech corpus generated from plain text data with a text-to-speech (TTS) system. In addition to the spelling correction model, authors experiment with improving the results of an end-to-end ASR system by incorporating an external language model and a combination of the two approaches. The proposed system achieves satisfactory results (19% relative

WER improvement and 29% relative WER improvement with additional LM re-scoring, with baseline ASR WER of 6.03) but requires a large speech corpus or high-quality TTS system to generate such corpus from a plain text. One of the recent works [9] presents an error correction model for Mandarin. The authors stress the importance of a low latency of ASR error correction model A realroduction environments and propose a non-autoregressive transformer model, faster than its autoregressive counterparts. The model is modeled on a large, artificially created parallel corpus of correct-incorrect sentence pairs, generated by randomly deleting, inserting, and replacing words in a text corpus. Real ASR corrections dataset is used to fine-tune the model to a specific ASR system. Relative WER reduction reported by authors on a publicly available testset is 13.87, which is slightly worse than an autoregressive model (15.53) while introducing over 6 times lower latency ($21ms$).

A similar approach to the one presented in this work is proposed in [10], but serves different tasks (grammatical error correction), operates on a poorer set of edit operations and uses different tagging models.

## III. DATA

We performed experiments for 3 European languages: Spanish, French, and German. The presented error correction models were trained and evaluated on pairs of ASR hypotheses and corresponding reference sentences. To create a corpus of such pairs, recordings from speech corpora for each language were processed using a corresponding model of an end-to-end speech recognition system. Reference transcriptions from speech corpora were paired with their corresponding hypotheses, creating parallel corpora of corrections for each language. To discard any differences caused by different normalization of transcriptions in speech corpora and in ASR output, we additionally performed automatic normalization of both reference and hypotheses sentences by lowercasing, removing punctuation characters and inverse-normalizing numbers. The data preparation pipeline is presented on Figure 1. For an example of a freely available corpus of ASR corrections for Polish prepared with the same pipeline, see [11].
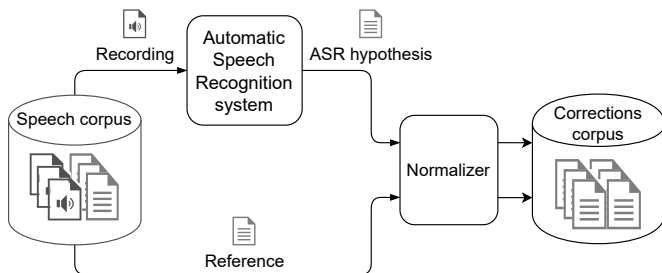


Fig. 1. Data preparation pipeline

Speech corpora, used to create corrections corpora, contain utterances used for virtual assistant development. They include commands and questions targeted at the assistant. Some of them are recorded in a studio for development purposes, and

|  | de-DE | es-ES | fr-FR |
|---|---|---|---|
| Sentences | 12242 | 16905 | 7180 |
| Tokens | 37955 | 55567 | 28004 |

some originate from usage data, after performing anonymization.

Data statistics are shown in Table I. Data sets were randomly divided into training, development, and test subsets in a proportion 8:1:1.

For a description of tagger training data preparation process, see section IV-A.

## IV. METHOD

The proposed error correction method is designed to be precise, easily controllable, and data-efficient. In contrast to methods inspired by machine translation, such as [8], our model does not have to model and reproduce all tokens of the output sequence. Instead, it learns only which tokens to modify to correct the sentence. To make it precisely controllable, the error correction mechanism works in two steps, depicted on figure 2. First, a sequence-tagging model assigns a tag to each token in the input sentence. The tag indicates whether the token is correctly recognized or requires correction. In the latter case, the tag specifies an edit operation that is needed to transform the incorrect sentence into a correct one. In the second step, the assigned tags are used to perform edit operations that correct ASR errors. This approach is especially suitable for production settings because it allows one to precisely control which edit operations to include in the model and which edit operations to perform on inference time. The control can be based on scores returned by the tagger (for example, by setting a global scorer threshold) or on rules excluding certain operations in certain contexts from being performed.

### A. Tagging data preparation

To train a sequence tagging model, a corpus of ASR hypotheses with assigned edit operation tags is needed. We create it from the parallel corrections corpus described in Section III. First, hypothesis and reference sentences are compared and aligned using Ratcliff-Obershelp algorithm [12] implemented in difflib library [13]. As a result, we get a list of operation codes describing how to turn corresponding parts of the hypothesis into reference. There are four operation codes: "replace", "delete", "insert" and "equal". For parts of a sentence which are not equal, we use a set of conditional rules involving recursively invoking the difflib alignment to tag each token with one of the edit operations. Examples of tags and corresponding edit operations are presented in table II. For an illustrated example of tags generation process, see Figure 3.

The set of available edit operation classes can be adjusted to the needs of a specific speech recognition system and a natural language. Ideally, the set should cover all errors and be as
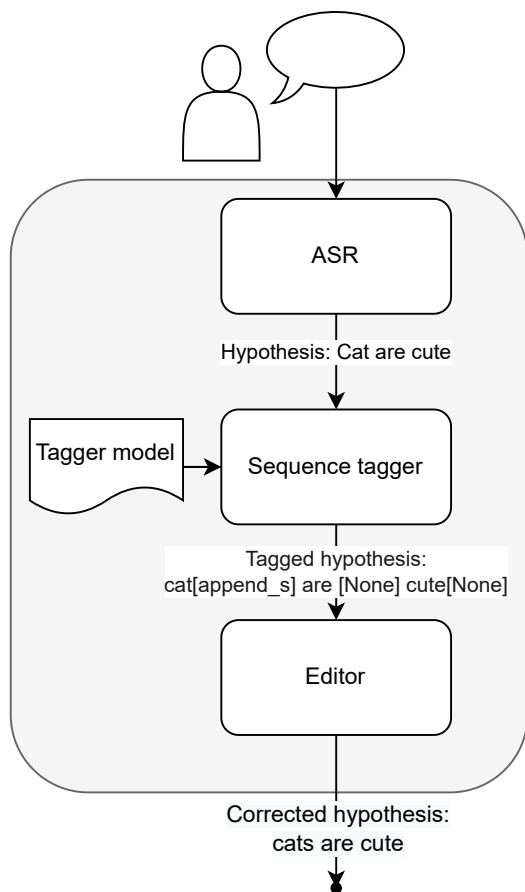
Fig. 2.   Error correction system architecture

the filtered edit operations with a special edit operation called "unsupported". This operation tag is present in the training set and the model learns to tag some tokens with it, but when correcting sentences on the inference time, there is no edit operation performed on them. A sequence of edit operations assigned to adjacent tokens can be interdependent - performing only some of them may deteriorate the results. Therefore, when one of the tokens is tagged with "unsupported", all surrounding, non-empty tags are replaced with the "unsupported" tag.
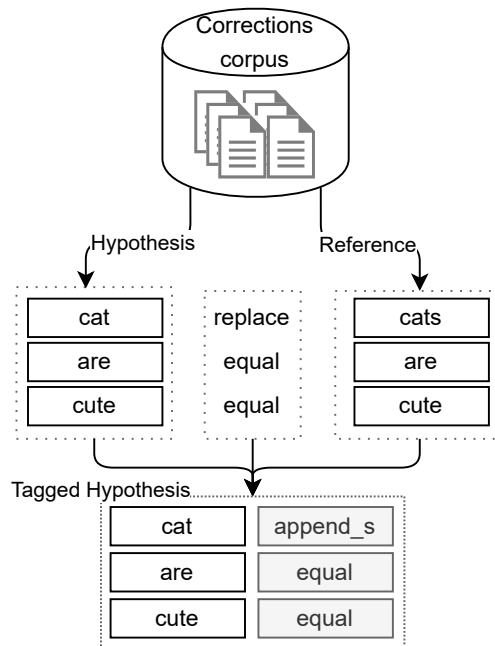


Fig. 3.   Tags generation example

small as possible. If the set is too big, edit operations become too sparse for the tagger to learn them effectively. Therefore, we chose to prioritize use o most expressive edit operations, like „append_s", which for example in English could cover most of the errors associated with a singular noun in place of a plural and missing "s" in third-person singular verbs. The same errors could be covered by more precise rules correcting only particular words, e.g. "replace_with_cats" (edit operation assigned to the word "cat"). By using more general operations, the model can generalize to unseen examples of errors. This is where our approach is different from [10] which uses only two main types of operations: KEEP ("None" in our approach), DELETE ("del" in our approach) and can insert any phrase or word from vocabulary $V$ before the current token. Such an approach cannot cover multiple errors with one tag.

Despite using edit operations that cover multiple errors, the edit operations dataset is still sparse. About $15\%$ of edit operations are supported only with one example. To make training more efficient and the model less overfitted to singular examples, we use a cut-off of 150 most frequent edit operations, also filtering-out all which are found only once in the dataset. To differentiate between tokens which are correct and those, whose errors are not frequent enough, we replace all

## B. Tagging models

We train and evaluate two tagging model types: BERT token classification model and contextual string embeddings model [14], referred later as "Flair". BERT tagger model uses locale variations of BERT transformer: Gbert for German [15], Camembert for French [16] and Beto for Spanish [17]. A single linear layer is added at the output, and the whole network is fine-tuned for the tagging task. Training was performed for 6 epochs, extending the training time did not improve results. The "Flair" model contains Bert models mentioned before, extended with contextual string embeddings [18], LSTM [19] and CRF [20] layers. Training was carried out for a maximum of 100 epochs.

## V. RESULTS

Both model types were evaluated using hold-out test sets by calculating WRR (Word Recognition Rate) of original ASR hypotheses and their corrected versions. Table III presents averaged results. As can be seen on an example of German dataset compared with other two - the worse the original ASR

TABLE II
EXAMPLES OF EDIT OPERATIONS

| name | description | example |
|------|-------------|---------|
| `del` | deletes a token | "a" → "" |
| `append_s` | appends given suffix to the token | "cat" → "cats" |
| `add_prefix_` | prepends given prefix to the token | "owl" → "howl" |
| `remove_suffix_1` | removes 1 character from the end of the token | "cats" → "cat" |
| `remove_prefix_1` | removes 1 character from the beginning of the token | "howl" → "owl" |
| `join` | joins token with previous one | "book store" → "bookstore" |
| `join_-` | joins token with previous one using given separator | "long term" → "long-term" |
| `replace_` | replaces token with given string | "cat" → "hat" |

TABLE III
RESULTS OF ERROR CORRECTION MODELS

|  |  | de-DE | es-ES | fr-FR |
|---|---|---|---|---|
| | base WRR | 78.07 | 90.70 | 93.51% |
| | corrected WRR | 83.48 | 92.65 | 94.97% |
| BERT | WRR gain | 5.41 | 1.95 | 1.46 |
| | Rel. WER reduction | 24.67% | 20.97% | 22.50% |
| | corrected WRR | 83.40 | 92.86 | 95.04 |
| Flair | WRR gain | 5.33 | 2.16 | 1.53 |
| | Rel. WER reduction | 24.30% | 23.23% | 23.57% |

TABLE IV
TRAINING (IN MINUTES) AND INFERENCE (IN MILLISECONDS) TIMES.

|  |  | de-DE | es-ES | fr-FR |
|---|---|---|---|---|
| | training time | 28 | 40 | 10 |
| BERT | inference latency | 14 | 14 | 13 |
| | training time | 180 | 154 | 67 |
| Flair | inference latency | 28 | 29 | 18 |

results are, the easier for the correction model to achieve higher absolute gains. Therefore, to compare correction models trained on datasets with different original results, we calculate a relative WER (Word Error Rate) reduction metric, which is calculated as

$$\frac{WER_{ASR} - WER_{Corrected}}{WER_{ASR}} \tag{1}$$

where $WER_{ASR}$ is WER of ASR system and $WER_{Corrected}$ is WER after applying the correction model. Relative improvements of our models vary between $21\%$ and $24.7\%$, making them comparable with state-of-the-art results reported in [8], while using much smaller training datasets, without the need of generating synthetic data with TTS engine. The Flair tagger offers slightly better results (except for German, where results are equal).

Table IV presents the time required for training the models and the average times needed to correct a single sentence from the test set. Both training and inference were performed using a machine with a single Tesla P40 GPU.

## VI. CONCLUSIONS

We presented a new approach to ASR errors correction problem. As demonstrated using three independent datasets, correction models trained using this approach are effective even for relatively small training datasets. The method allows to precisely control which errors should be included in the model and which of the included ones should be corrected at

the inference time. The evaluations performed on the models show that they can significantly improve the ASR results by reducing the WER by more than $20\%$. All of the models presented offer very good inference latency, making them suitable for use with streaming ASR systems.

The presented method is well suited for industrial applications where the ability to precisely control how the error correction model works, as well as small latency, are crucial.

## REFERENCES

[1] K. H. Davis, R. Biddulph, and S. Balashek, "Automatic recognition of spoken digits," *The Journal of the Acoustical Society of America*, vol. 24, no. 6, pp. 637–642, 1952. [Online]. Available: https://doi.org/10.1121/1.1906946

[2] W. Xiong, J. Droppo, X. Huang, F. Seide, M. L. Seltzer, A. Stolcke, D. Yu, and G. Zweig, "Achieving human parity in conversational speech recognition," *ArXiv*, vol. abs/1610.05256, 2016.

[3] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, "Librispeech: An asr corpus based on public domain audio books," *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 5206–5210, 2015.

[4] A. Narayanan, A. Misra, K. C. Sim, G. Pundak, A. Tripathi, M. Elfeky, P. Haghani, T. Strohman, and M. Bacchiani, "Toward domain-invariant speech recognition via large scale training," in *2018 IEEE Spoken Language Technology Workshop (SLT)*, 2018, pp. 441–447.

[5] A. Jeziorski, F. Sawicki, O. Solop, M. Junczyk, M. Sikora, and T. Zietkiewicz, "Industrial asr troubleshooting tool," in *Proceedings of the LREC2020 Industry Track*. Marseille, France: European Language Resources Association (ELRA), May 2020, pp. 10–14. [Online]. Available: https://www.aclweb.org/anthology/2020.lrec2020industrytrack-1.3

[6] R. Errattahi, A. El Hannani, and H. Ouahmane, "Automatic speech recognition errors detection and correction: A review," *Procedia Computer Science*, vol. 128, pp. 32–37, 01 2018.

[7] H. Cucu, A. Buzo, L. Besacier, and C. Burileanu, "Statistical Error Correction Methods for Domain-Specific ASR Systems," in *Statistical Language and Speech Processing*, A.-H. Dediu, C. Martín-Vide, R. Mitkov, and B. Truthe, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2013, pp. 83–92.

[8] J. Guo, T. N. Sainath, and R. J. Weiss, "A spelling correction model for end-to-end speech recognition," in *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2019, pp. 5651–5655. [Online]. Available: https://arxiv.org/pdf/1902.07178

[9] Y. Leng, X. Tan, L. Zhu, J. Xu, R. Luo, L. Liu, T. Qin, X.-Y. Li, E. Lin, and T.-Y. Liu, "Fastcorrect: Fast error correction with edit alignment for automatic speech recognition," 2021.

[10] E. Malmi, S. Krause, S. Rothe, D. Mirylenka, and A. Severyn, "Encode, tag, realize: High-precision text editing," 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)*. Hong Kong, China: Association for Computational Linguistics, Nov. 2019, pp. 5054–5065. [Online]. Available: https://aclanthology.org/D19-1510

[11] M. Kubis, Z. Vetulani, M. Wypych, and T. Ziętkiewicz, "Open challenge for correcting errors of speech recognition systems," in *Proceedings of the 9th Language and Technology Conference: Human Language Technologies as a Challenge for Computer Science and Linguistics,*

Z. Vetulani and P. Paroubek, Eds. Poznań, Poland: Wydawnictwo Nauka i Innowacje, 2019, pp. 219–223. [Online]. Available: https://arxiv.org/abs/2001.03041https://gonito.net/gitlist/asr-corrections.git/

[12] D. E. M. John W. Ratcliff, "Pattern matching: The gestalt approach," p. 46, 7 1988. [Online]. Available: https://www.drdobbs.com/database/pattern-matching-the-gestalt-approach/184407970

[13] "difflib — helpers for computing deltas," https://docs.python.org/3/library/difflib.html, accessed: 2022-03-20.

[14] A. Akbik, T. Bergmann, D. Blythe, K. Rasul, S. Schweter, and R. Vollgraf, "FLAIR: An easy-to-use framework for state-of-the-art NLP," in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations)*. Minneapolis, Minnesota: Association for Computational Linguistics, Jun. 2019, pp. 54–59. [Online]. Available: https://www.aclweb.org/anthology/N19-4010

[15] B. Chan, S. Schweter, and T. Möller, "German's next language model," in *Proceedings of the 28th International Conference on Computational Linguistics*. Barcelona, Spain (Online): International Committee on Computational Linguistics, Dec. 2020, pp. 6788–6796. [Online]. Available: https://aclanthology.org/2020.coling-main.598

[16] L. Martin, B. Muller, P. J. Ortiz Suárez, Y. Dupont, L. Romary, É. de la Clergerie, D. Seddah, and B. Sagot, "CamemBERT: a tasty French language model," in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Online: Association for Computational Linguistics, Jul. 2020, pp. 7203–7219. [Online]. Available: https://www.aclweb.org/anthology/2020.acl-main.645

[17] J. Cañete, G. Chaperon, R. Fuentes, J.-H. Ho, H. Kang, and J. Pérez, "Spanish pre-trained bert model and evaluation data," in *PML4DC at ICLR 2020*, 2020.

[18] A. Akbik, D. Blythe, and R. Vollgraf, "Contextual string embeddings for sequence labeling," in *Proceedings of the 27th International Conference on Computational Linguistics*. Santa Fe, New Mexico, USA: Association for Computational Linguistics, Aug. 2018, pp. 1638–1649. [Online]. Available: https://aclanthology.org/C18-1139

[19] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, pp. 1735–80, 12 1997.

[20] J. D. Lafferty, A. McCallum, and F. C. N. Pereira, "Conditional random fields: Probabilistic models for segmenting and labeling sequence data," in *Proceedings of the Eighteenth International Conference on Machine Learning*, ser. ICML '01. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc., 2001, pp. 282–289. [Online]. Available: http://dl.acm.org/citation.cfm?id=645530.655813