

# PATCorrect: Non-autoregressive Phoneme-augmented Transformer for ASR Error Correction

Ziji Zhang<sup>1</sup>, Zhehui Wang<sup>2</sup>, Rajesh Kamma<sup>2</sup>, Sharanya Eswaran<sup>2</sup>, Narayanan Sadagopan<sup>2</sup>

<sup>1</sup>Department of Applied Mathematics and Statistics, Stony Brook University, Stony Brook, NY, USA

<sup>2</sup>Amazon.com Inc, Seattle, WA, USA

ziji.zhang@stonybrook.edu, zhehwang@umich.edu, {srikamma, sharanye, sdgpn}@amazon.com

## Abstract

Speech-to-text errors made by automatic speech recognition (ASR) systems negatively impact downstream models. Error correction models as a post-processing text editing method have been recently developed for refining the ASR outputs. However, efficient models that meet the low latency requirements of industrial grade production systems have not been well studied. We propose PATCorrect—a novel non-autoregressive (NAR) approach based on multi-modal fusion leveraging representations from both text and phoneme modalities, to reduce word error rate (WER) and perform robustly with varying input transcription quality. We demonstrate that PATCorrect consistently outperforms state-of-the-art NAR method on English corpus across different upstream ASR systems, with an overall 11.62% WER reduction (WERR) compared to 9.46% WERR achieved by other methods using text only modality. Besides, its inference latency is at tens of milliseconds, making it ideal for systems with low latency requirements.

**Index Terms:** Non-autoregressive Transformer, phoneme augmented learning, ASR error correction

## 1. Introduction

Automatic speech recognition (ASR) model transcribes human speech into readable text, making it a critical component in large-scale natural language processing (NLP) systems like Amazon Alexa, Google Home and Apple Siri. Transcribed texts serve as input for downstream models such as intent detection in voice assistants and response generation in voice chatbots. Errors made in ASR transcriptions can severely impact the accuracy of downstream models and thus lower the performance of the entire NLP system. To improve the quality of ASR transcriptions, error correction models are applied to the outputs from ASR systems. ASR error correction can be formulated as a sequence-to-sequence (seq2seq) generation task, taking the ASR transcribed texts as input source sequences and the ground-truth transcriptions as target sequences. Previous studies [1, 2, 3] have proposed seq2seq models that decode the target sequence in an autoregressive (AR) manner. Although AR models have achieved great accuracy, high inference latency makes them infeasible for online production systems with low-latency constraints. For example, the end-to-end response time for voice digital assistants is of the order of tens of milliseconds for desired user experience. Under these latency constraints it is impractical to deploy AR models for error correction. Motivated by these considerations, we focus on non-autoregressive (NAR) models over AR models to meet the low latency requirement of industry production systems as NAR models are faster with parallel decoding during inference [4].

We propose PATCorrect (Phoneme Augmented Trans-

former for ASR error Correction) as shown in Figure 1, a novel NAR based upstream-model-agnostic ASR error correction model using both text and phoneme representations of the ASR transcribed sentences. PATCorrect applies multi-modal fusion to combine phoneme representation and text representation into joint feature embedding as input for the length tagging predictor. Both text and phoneme encoders interact with NAR decoder via encoder-decoder attention mechanism. PATCorrect improves the WER reduction (WERR) to 11.62% compared to FastCorrect which is state-of-the-art (SOTA) NAR method that solely uses text only representation of the input, with comparable inference latency at the scale of tens of milliseconds. Our main contributions are summarized as follows:

- We propose PATCorrect, a novel model based on the Transformer architecture for NAR ASR correction. This model uses a multi-modal fusion approach that augments the traditional input text encoding with an additional phoneme encoder to incorporate pronunciation information, which is one of the key characteristics of spoken utterances.
- Through extensive offline evaluations on English corpus, We demonstrate that PATCorrect outperforms the SOTA NAR ASR error correction model that uses text only modality. For example, PATCorrect improves WERR to 11.62% with an inference latency at the same tens of milliseconds scale, while still being about 4.2 - 6.7x times faster than AR models.
- To the best of our knowledge, we are the first to establish that multi-modal fusion is a promising direction for improving the accuracy of low latency NAR methods for ASR error correction, and comprehensively study the performance on English corpus across ASR systems with varying levels of quality.

## 2. Related Work

ASR error correction can be viewed as Neural Machine Translation (NMT) problem with erroneous sentences as source language, and corrected sentences as target language. [5, 6] applied the NMT approach to domain-specific ASR systems for error correction, and further utilized ontology learning to repair ASR outputs by a 4-step method. Recent NMT methods based on Transformers [7, 8] have become increasingly accurate and have inspired applications to ASR error correction [2, 3, 9]. Based on the intuition that phonetic information helps with understanding ASR errors [10, 11], [12] found that adding phoneme information for domain-agnostic ASR system could benefit entity retrieval task. Although achieving high accuracy, these encoder-decoder based autoregressive generation models suffer from slow inference speed (hundreds of milliseconds), error propagation and demand for a large amount of training data.

To address these issues in AR models, NAR sequence generation methods in NMT [13], which aim to speed up the in-



probability of  $\hat{\mathbf{w}}$  with a variational lower bound:

$$\mathcal{L}_{\text{PATCorrect}} = \log P_{\text{PATCorrect}}(\hat{\mathbf{w}}|\mathbf{w}, \mathbf{p}; \theta) \geq \mathbb{E}_{\mathbf{t} \sim q} \left( \sum_{i=1}^n \log P_{\text{TagP}}(t_i|\mathbf{w}, \mathbf{p}; \theta) + \sum_{i=1}^{\hat{n}} \log P(\hat{w}_i|\mathbf{w}'_i, \mathbf{p}; \theta) \right) \quad (3)$$

where  $\mathbf{w}'_i$  is the adjusted source inputs based on  $\mathbf{t}$ , the actual input to the NAR decoder. With the optimal edit alignment paths, we provide an approximate distribution  $q$  for the tag sequence with the most likely sample as the expectation. The first term within  $\mathbb{E}_{\mathbf{t} \sim q}(\cdot)$  trains the tag predictor *TagP*, and the second term trains the error correction model. This design enables the parallel decoding for every target token  $\hat{w}_i$  at the same time.

Meanwhile, we include the encoder-decoder attention mechanism for both encoders in the joint NAR decoder, by sequentially combining them together [12]. After getting the self-attention output from decoder  $H_{\hat{\mathbf{w}}} \in \mathbb{R}^{\hat{n} \times d_h}$  with causal mask removed to enable parallel calculation [13], we calculate the text and phoneme encoder-decoder attention by first using  $H_{\hat{\mathbf{w}}}$  as query, text encoder output  $H_w$  as key and value:

$$Z_w = \text{Softmax}\left(\frac{(H_{\hat{\mathbf{w}}}W_w^Q)(H_wW_w^K)^T}{\sqrt{d_h}}\right)(H_wW_w^V) \quad (4)$$

and then using the output  $Z_w$  as query, phoneme encoder output  $H_p$  as key and value:

$$Z_p = \text{Softmax}\left(\frac{(Z_wW_p^Q)(H_pW_p^K)^T}{\sqrt{d_h}}\right)(H_pW_p^V) \quad (5)$$

where  $\{W_w^Q, W_w^K, W_w^V\}, \{W_p^Q, W_p^K, W_p^V\} \in \mathbb{R}^{d_h \times d_h}$  are the parameter matrices in the text and phoneme encoder-decoder multi-head attention layers, respectively. The remaining of the decoder layers are the same as the vanilla Transformer model where  $Z_p \in \mathbb{R}^{\hat{n} \times d_h}$  will be input into fully-connected MLPs with residual connection and layer normalization.

## 4. Experiments

### 4.1. Datasets and ASR models

We create the ASR transcription dataset  $\mathbb{S}^{\text{trans}}$  specifically designed for correcting ASR errors. Multiple ASR systems are used to transcribe public English audio corpus, which has the corresponding ground truths from human transcriptions. The ASR transcriptions of the audio clips contain actual ASR errors, which are paired with golden texts. We combine two public English datasets, LibriSpeech [18] and Common Voice v9.0 [19], together to get a large and diverse training corpus. In order to test different ASR systems with different architecture and performances, we choose 3 pre-trained ASR systems implemented in NeMo [20]: The convolution-augmented Transformer **Conformer** [21] with top-tier performance; The convolution-based **Jasper** [22] with above average performance; A light-weight 5x5 **QuartzNet** [23] with subpar performance. All 3 ASR systems are trained with the Connectionist Temporal Classification (CTC) loss. We use the default splits in LibriSpeech and Common Voice with transcriptions from 3 ASR systems together to compose the dataset  $\mathbb{S}^{\text{trans}}$  as shown in Table 1, which in total has more than 3.5 million sentences pairs in training split. Note that LibriSpeech has been included in the training of all 3 ASR systems, so we use DEV and TEST splits from Common Voice as benchmarks in accuracy evaluations later.

To evaluate the performance of ASR error correction models, speed is measured by the latency of the whole inference process including encoding and decoding for different methods. For accuracy, similar to [4], we use Word Error Rate (WER), WER Reduction (WERR), *Precision* measures how many actual error tokens are edited among all of the edited tokens, *Recall* measures how many actual error tokens are edited among all of the error tokens. Since unnecessary edits may even lead to WER increase, we use  $F_{0.5}$  as an overall measurement to put more weight on *Precision* for error detection ability [2, 24, 25]. For error correction ability, *Correction*, is defined as the number of correctly edited error tokens divided by the number of edited error tokens, which measures the percentage of edited error tokens that match the ground truth.

Table 1: ASR transcription dataset and ASR original WER

Dataset $\mathbb{S}^{\text{trans}}$	TRAIN	DEV	TEST
# of sents	1,171,348	21,898	21,877
Avg. words/sent	15.8	12.1	11.8
Conformer WER	4.56	7.20	7.28
Jasper WER	7.80	14.12	15.29
QuartzNet WER	18.97	32.96	35.86

### 4.2. Model configurations

**PATCorrect** We use 6-layer text encoder, 6-layer phoneme encoder and 6-layer joint decoder with the hidden model dimension  $d_h = 512$ , and MLP dimension  $d_{\text{MLP}} = 2048$ . We use 8 attention heads for self-attention, encoder-decoder attention and cross attention respectively. In the cross attention setup for fusing two encoder outputs, we use 2 consecutive modules of cross attention with dropout, residual connection and layer normalization. For tag predictor *TagP* trained with MSE loss, we apply 5 layers of convolutional modules which consists of 1-D convolution layer with kernel size of 3, ReLU, layer normalization and dropout, followed by 2 layers of MLPs.

**Baseline models** We compare our model with both AR Transformer and popular NAR methods. For **AR Transformer**, we use the standard vanilla Transformer architecture with 6-layer encoder and 6-layer decoder and hidden size  $d_h = 512$ . For NAR methods Mask Predict (**CMLM**) [15] and Levenshtein Transformer (**LevT**) [14], we use the default training hyperparameters. For SOTA NAR ASR error correction method **FastCorrect** [4], we use the same 6-layer encoder-decoder architecture and the same architecture for length predictor.

**Training and inference details** All models are implemented using Fairseq [26], and trained using 4 NVIDIA Tesla V100 GPUs with maximum batch token size of 5000 and label smoothed cross entropy loss function. They are trained from scratch using the ASR transcription dataset  $\mathbb{S}^{\text{trans}}$  for 30 epochs. Source and target sentences are tokenized using sentencepiece [27], and the phoneme sequences are generated by English grapheme to phoneme conversion using the CMU pronouncing dictionary [28], an independent post-processing step of the 1-best ASR hypothesis as model inputs. We use Adam optimizer [29] and inverse square root for learning rate scheduling starting from  $5e^{-4}$ . During inference, we set the test batch size as 1 to simulate the online production environment with output from the upstream ASR system, and all NAR methods have the same

max decoding iteration of 1. Different hardware conditions are tested including single NVIDIA Tesla V100 GPU, 8 CPUs and 4 CPUs with Intel(R) Xeon(R) CPU E5-2686 v4 @ 2.30GHz.

## 5. Results

### 5.1. Accuracy

We compare the WER and WERR using different error correction models for all three ASR systems and their total combined transcriptions on Common Voice TEST split, as shown in Table 2. PATCorrect outperforms other NAR methods for all three ASR system transcriptions. Results show: 1) For averaged total results on three ASR systems, PATCorrect beats the SOTA FastCorrect method by improving the TEST set WERR from 9.46 to 11.62, which is more than 20% relative improvement; 2) Among all of the multi-modal fusion operations experimented, cross attention performs best across almost all datasets.

Table 2: WER (%) ↓ and WERR (%) ↑ using error correction models averaged over outputs from 3 upstream ASR systems

Models	WER ↓		WERR ↑	
	DEV	TEST	DEV	TEST
No error correction	25.55	27.96	-	-
AR Transformer	19.45	22.82	23.87	18.40
CMLM	23.92	26.49	6.36	5.26
LevT	23.19	25.70	9.24	8.11
FastCorrect	22.13	25.32	13.37	9.46
PATCorrect( <i>cat</i> )	22.09	25.24	13.53	9.75
PATCorrect( <i>add</i> )	22.14	25.28	13.33	9.61
PATCorrect( <i>max</i> )	21.93	25.07	14.16	10.35
PATCorrect( <i>cross_atten</i> )	<b>21.57</b>	<b>24.72</b>	<b>15.59</b>	<b>11.62</b>

### 5.2. Speed

We test inference speed on both GPU and CPUs as shown in Table 3 for Common Voice TEST split. Consistent with the observation from [13], AR model shows a linear latency trend with decoding lengths, while the latency of the NAR methods increases slightly. PATCorrect achieves an inference latency comparable with other NAR models, especially on GPU hardware, while still being about 4.2 - 6.7x times faster than AR models.

Table 3: Inference latency ↓ with unit ms/sentence

Models	1× GPU	8× CPUs	4× CPUs
AR Transformer	141.00	149.79	186.89
CMLM	35.54	43.78	47.90
LevT	47.86	45.68	55.33
FastCorrect	<b>20.81</b>	<b>23.77</b>	<b>29.14</b>
PATCorrect( <i>cross_atten</i> )	33.71	41.89	49.97

### 5.3. Sensitivity analysis

We conduct sensitivity analysis by comparing the error detection ability and error correction ability for different models with 3 ASR transcriptions on Common Voice TEST split. In Table 4,  $P$ ,  $R$ ,  $C$  denote *Precision*, *Recall*, *Correction* respectively. Results show that our PATCorrect not only has the highest  $F_{0.5}$

score with great *Precision* that is comparable to AR model, but also has better ability to edit error tokens to the correct targets indicated by higher *Correction*.

Table 4: Sensitivity metrics ↑ for error correction models

Models	$P$	$R$	$F_{0.5}$	$C$
AR Transformer	90.89	65.33	84.30	37.81
CMLM	84.87	61.45	78.86	27.37
LevT	86.49	<b>61.75</b>	80.07	28.42
FastCorrect	89.53	60.90	81.83	27.70
PATCorrect ( <i>cross_atten</i> )	<b>90.27</b>	59.64	<b>81.86</b>	<b>29.50</b>

### 5.4. Ablation study

We perform ablation study to understand the effectiveness of different components of PATCorrect by removing a component while retaining the others. No phoneme input for *TagP* means that we only use text information as input for predicting token tags while still using phoneme encoder-decoder attention in NAR decoder. No phoneme attention means that we remove phoneme encoder-decoder attention from the NAR decoder and only use text encoder-decoder attention, while still using cross-attention to fuse the text and phoneme encoder outputs for *TagP*. We also increase the amount of pre-training data by using a synthetic dataset  $\mathbb{S}^{\text{synth}}$  for data augmentation to pretrain the model for 20 epochs. We crawl 50M sentences from English Wiki pages [30], and add random editing operations like deletion, insertion, substitution with a homophone dictionary to produce erroneous sentences paired with the original correct texts, that mimics the ASR errors with a simulated error distribution. The results in Table 5 with WER and WERR on the Common Voice TEST split, equal-weight averaged from all 3 ASR systems, show that adding phoneme information in *TagP* and NAR decoder lead to better WER and WERR. In addition, using  $\mathbb{S}^{\text{synth}}$  for pretraining can also boost the model accuracy to further reduce WER.

Table 5: Ablation study for PATCorrect model

Models	WER ↓	WERR ↑
No error correction	27.96	-
FastCorrect	25.32	9.46
No phoneme input for <i>TagP</i>	25.13	10.14
No phoneme attention in decoder	25.06	10.39
PATCorrect( <i>cross_atten</i> )	24.72	11.62
PATCorrect( <i>cross_atten</i> ) + $\mathbb{S}^{\text{synth}}$	<b>24.33</b>	<b>13.00</b>

## 6. Conclusions

We propose PATCorrect, a novel NAR phoneme-augmented Transformer-based model with robust performance on different upstream ASR systems with varying speech-to-text transcription quality. Our model outperforms SOTA NAR ASR error correction models, and still 4.2 - 6.7x times faster than AR models, which makes it a great fit as industrial scale text editing method to refine ASR transcriptions. Our study establishes that multi-modal fusion is a promising direction for improving the accuracy of low latency NAR methods for ASR error correction.

## 7. References

- [1] L. F. D’Haro and R. E. Banchs, “Automatic correction of asr outputs by using machine translation,” in *Interspeech*, vol. 2016, 2016, pp. 3469–3473.
- [2] J. Liao, S. E. Eskimez, L. Lu, Y. Shi, M. Gong, L. Shou, H. Qu, and M. Zeng, “Improving readability for automatic speech recognition transcription,” *arXiv preprint arXiv:2004.04438*, 2020.
- [3] A. Mani, S. Palaskar, N. V. Meripo, S. Konam, and F. Metze, “Asr error correction and domain adaptation using machine translation,” in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 6344–6348.
- [4] Y. Leng, X. Tan, L. Zhu, J. Xu, R. Luo, L. Liu, T. Qin, X. Li, E. Lin, and T.-Y. Liu, “Fastcorrect: Fast error correction with edit alignment for automatic speech recognition,” *Advances in Neural Information Processing Systems*, vol. 34, pp. 21 708–21 719, 2021.
- [5] H. Cucu, A. Buzo, L. Besacier, and C. Burileanu, “Statistical error correction methods for domain-specific asr systems,” in *International Conference on Statistical Language and Speech Processing*. Springer, 2013, pp. 83–92.
- [6] C. Anantaram, A. Sangroya, M. Rawat, and A. Chhabra, “Repairing asr output by artificial development and ontology based learning,” in *IJCAI*, 2018, pp. 5799–5801.
- [7] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, “Attention is all you need,” *Advances in neural information processing systems*, vol. 30, 2017.
- [8] N. Ng, K. Yee, A. Baevski, M. Ott, M. Auli, and S. Edunov, “Facebook fair’s wmt19 news translation task submission,” *arXiv preprint arXiv:1907.06616*, 2019.
- [9] K. Hu, T. N. Sainath, R. Pang, and R. Prabhavalkar, “Deliberation model based two-pass end-to-end speech recognition,” in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 7799–7803.
- [10] A. Fang, S. Filice, N. Limsopatham, and O. Rokhlenko, “Using phoneme representations to build predictive models robust to asr errors,” in *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2020, pp. 699–708.
- [11] M. N. Sundararaman, A. Kumar, and J. Vepa, “Phoneme-bert: Joint language modelling of phoneme sequence and asr transcript,” *arXiv preprint arXiv:2102.00804*, 2021.
- [12] H. Wang, S. Dong, Y. Liu, J. Logan, A. K. Agrawal, and Y. Liu, “Asr error correction with augmented transformer for entity retrieval,” in *Interspeech*, 2020, pp. 1550–1554.
- [13] J. Gu, J. Bradbury, C. Xiong, V. O. Li, and R. Socher, “Non-autoregressive neural machine translation,” *arXiv preprint arXiv:1711.02281*, 2017.
- [14] J. Gu, C. Wang, and J. Zhao, “Levenshtein transformer,” *Advances in Neural Information Processing Systems*, vol. 32, 2019.
- [15] M. Ghazvininejad, O. Levy, Y. Liu, and L. Zettlemoyer, “Mask-predict: Parallel decoding of conditional masked language models,” *arXiv preprint arXiv:1904.09324*, 2019.
- [16] D. Kiela, E. Grave, A. Joulin, and T. Mikolov, “Efficient large-scale multi-modal classification,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32, 2018.
- [17] M. Ghazvininejad, C. Brockett, M.-W. Chang, B. Dolan, J. Gao, W.-t. Yih, and M. Galley, “A knowledge-grounded neural conversation model,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32, 2018.
- [18] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, “Librispeech: an asr corpus based on public domain audio books,” in *2015 IEEE international conference on acoustics, speech and signal processing (ICASSP)*. IEEE, 2015, pp. 5206–5210.
- [19] R. Ardila, M. Branson, K. Davis, M. Henretty, M. Kohler, J. Meyer, R. Morais, L. Saunders, F. M. Tyers, and G. Weber, “Common voice: A massively-multilingual speech corpus,” *arXiv preprint arXiv:1912.06670*, 2019.
- [20] O. Kuchaiev, J. Li, H. Nguyen, O. Hrinchuk, R. Leary, B. Ginsburg, S. Kriman, S. Beliaev, V. Lavrukhin, J. Cook *et al.*, “Nemo: a toolkit for building ai applications using neural modules,” *arXiv preprint arXiv:1909.09577*, 2019.
- [21] A. Gulati, J. Qin, C.-C. Chiu, N. Parmar, Y. Zhang, J. Yu, W. Han, S. Wang, Z. Zhang, Y. Wu *et al.*, “Conformer: Convolution-augmented transformer for speech recognition,” *arXiv preprint arXiv:2005.08100*, 2020.
- [22] J. Li, V. Lavrukhin, B. Ginsburg, R. Leary, O. Kuchaiev, J. M. Cohen, H. Nguyen, and R. T. Gadde, “Jasper: An end-to-end convolutional neural acoustic model,” *arXiv preprint arXiv:1904.03288*, 2019.
- [23] S. Kriman, S. Beliaev, B. Ginsburg, J. Huang, O. Kuchaiev, V. Lavrukhin, R. Leary, J. Li, and Y. Zhang, “Quartznet: Deep automatic speech recognition with 1d time-channel separable convolutions,” in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 6124–6128.
- [24] K. Omelianchuk, V. Atrasevych, A. Chernodub, and O. Skurzhan-skyi, “Gector-grammatical error correction: tag, not rewrite,” *arXiv preprint arXiv:2005.12592*, 2020.
- [25] S. Rothe, J. Mallinson, E. Malmi, S. Krause, and A. Severyn, “A simple recipe for multilingual grammatical error correction,” *arXiv preprint arXiv:2106.03830*, 2021.
- [26] M. Ott, S. Edunov, A. Baevski, A. Fan, S. Gross, N. Ng, D. Grangier, and M. Auli, “fairseq: A fast, extensible toolkit for sequence modeling,” in *Proceedings of NAACL-HLT 2019: Demonstrations*, 2019.
- [27] T. Kudo and J. Richardson, “Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing,” *arXiv preprint arXiv:1808.06226*, 2018.
- [28] K. Park and J. Kim, “g2pe,” <https://github.com/Kyubyong/g2p>, 2019.
- [29] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014.
- [30] G. Attardi, “Wikiextractor,” <https://github.com/attardi/wikiextractor>, 2015.