

Beyond Binary: Multiclass Paraphasia Detection with Generative Pretrained Transformers and End-to-End Models

Matthew Perez¹, Aneesha Sampath¹, Minxue Niu¹, Emily Mower Provost¹

¹University of Michigan, USA

mkperez@umich.edu, saneesha@umich.edu, sandymn@umich.edu, emilykmp@umich.edu

Abstract

Aphasia is a language disorder that can lead to speech errors known as paraphasias, which involve the misuse, substitution, or invention of words. Automatic paraphasia detection can help those with Aphasia by facilitating clinical assessment and treatment planning options. However, most automatic paraphasia detection works have focused solely on binary detection, which involves recognizing only the presence or absence of a paraphasia. Multiclass paraphasia detection represents an unexplored area of research that focuses on identifying multiple types of paraphasias and where they occur in a given speech segment. We present novel approaches that use a generative pretrained transformer (GPT) to identify paraphasias from transcripts as well as two end-to-end approaches that focus on modeling both automatic speech recognition (ASR) and paraphasia classification as multiple sequences vs. a single sequence. We demonstrate that a single sequence model outperforms GPT baselines for multiclass paraphasia detection.

Index Terms: paraphasia detection, disordered speech, aphasia speech analysis

1. Introduction

Aphasia is a common language disorder that occurs as a result of damage to the brain and can ultimately impair the communication abilities (both expressive and receptive) of an individual. Aphasia affects over two million people in the United States and nearly 225,000 acquire Aphasia each year following a medical event such as a traumatic brain injury or stroke [1]. Aphasia can manifest in a variety of ways, but generally, persons with Aphasia (PWAs) struggle with verbal communication and in some cases produce specific speech errors known as paraphasias.

There are several types of paraphasic errors. In this work, we focus specifically on phonemic, neologistic, and semantic paraphasias [2, 3].

- **phonemic** paraphasias involve substituting, omitting, or rearranging phonemes (i.e., ‘zut’ for ‘shut’)
- **neologistic** paraphasias involve substituting a non-sensical word (i.e., ‘flibber’ for ‘bottle’)
- **semantic** paraphasias involve substituting a semantically related word (i.e., ‘bed’ for ‘desk’)

Clinical research has highlighted the impact that accurate paraphasia detection plays in predicting recovery patterns and guiding treatment planning [4, 5]. In clinical settings, automated tools for detecting paraphasias in an individual’s speech can ultimately allow for more efficient and consistent assessment procedures. Additionally, for supplementary treatment

options such as remote, self-directed speech therapy (via smartphone), automatically identifying paraphasic errors is critical in providing constructive feedback to the user [6, 7].

Previous automatic paraphasia detection work has focused on identifying paraphasias from single-word elicitation tasks with manual transcriptions [4, 8, 9]. These works have limited applications, mainly in clinical settings. For applications with continuous or unsegmented speech, paraphasia detection includes identifying where in the given sequence a paraphasia occurs. Some previous works that have focused on continuous speech have treated paraphasia detection as a binary task [10, 11, 12]. However, these works are restricted to learning the presence or absence of paraphasic errors rather than learning to differentiate between paraphasia types. For remote speech therapy applications that process continuous speech, models that focus on multiclass classification are needed to characterize these different types of paraphasias and where they occur in a given utterance.

We present the first work into automatic multiclass paraphasia detection for continuous speech. We investigate several methods for automatic paraphasia classification, which include using a generative pretrained transformer to classify paraphasias from ASR transcripts (ASR+GPT), a single-sequence (single-seq) model where both ASR and paraphasia classification tasks are learned within the same sequence, and a multi-sequence (multi-seq) model where paraphasia classification and ASR are learned as separate sequences but jointly optimized with multi-task learning. The research goals of this work are:

- To investigate the utility of an off-the-shelf GPT model for paraphasia classification using imperfect (ASR) or perfect (oracle) transcripts.
- Explore single-seq and multi-seq models for word-level paraphasia classification.
- Analyze performance across different paraphasia classes.

Our findings demonstrate that GPT can be used to detect paraphasias with aphasic speech transcripts. However, we note that the single-seq model outperforms GPT for multiclass paraphasia detection, specifically for phonemic and neologistic paraphasias. Lastly, we discuss some limitations of the presented approaches for semantic paraphasia classification.

2. Related Work

One of the first works that explored statistical models for paraphasia detection was by Fergadotis et al., which used separate classification models to perform binary paraphasia detection in a one-vs-rest fashion [4]. This work was focused on the Moss Aphasia Psycholinguistics Project Database (MAPPD), which contains transcribed (text-input) single-word responses

Table 1: *Text Pre-Processing: CHAT transcriptions are processed to Oracle transcripts. Examples for each model output is also shown. Blue indicates paraphrastic words, red indicates paraphrastic labels*

CHAT Transcripts	aphasia fekts@u [: affects] [* p] my language not my ditikəlt@u [: intelligence] [* n]
Oracle	aphasia fekts [p] my language not my ditikalt [n]
Model Output	
ASR+GPT	aphasia [c] fekts [p] my [c] language [c] not [c] my [c] ditikalt [n]
single-seq	aphasia fekts [p] my language not my ditikalt [n]
multi-seq	ASR: aphasia fekts my language not my ditikalt Para: [c] [p] [c] [c] [c] [c] [n]

with paraphasia labels for each word [13]. The classifiers use linguistic features like word2vec or semantic similarity for paraphasia classification. One limitation of this work is that paraphasia classification is performed by separate classifiers, which lacks the specificity to differentiate between different paraphasias. More recent work [8, 9] has addressed these concerns by presenting a unified model for multiclass paraphasia detection on MAPPD where a multiclass decision tree based on the binary classifiers presented in [4] is used to perform multiclass paraphasia detection. A broader limitation of the works above is that MAPPD, is constrained to single-word responses that are transcribed (i.e., contain no audio data). This is useful in clinical applications where manual transcription for targeted tasks can be attained. However, for unconstrained speech applications such as remote speech therapy, automatic paraphasia detection models must be able to handle unsegmented or continuous speech data. In this work, we investigate multiclass paraphasia detection for continuous, read speech from the AphasiaBank corpus [14]. The paraphasia classification in continuous speech is complicated by temporal challenges, which involve not just identifying the different types of paraphasias but also the specific points in the utterance where they occur.

Some additional works have focused on paraphasia detection in continuous speech. Work by Le et al. showed that a fully automatic pipeline composed of ASR and a logistic regression model for binary paraphasia detection can detect phonemic and neologistic paraphasias in continuous speech [10]. Work by Pai et al. demonstrated that density-based clustering with manual segmentation can yield improved features for binary paraphasia detection [11]. Lastly, work by Perez et al. focused on automatic paraphasia detection using end-to-end models and showed that multitask learning improved performance over existing automated approaches [12]. These works are limited in that they focus on binary paraphasia detection in continuous speech and only detect phonemic and neologistic paraphasias. For remote applications, multiclass paraphasia detection models that can process continuous speech are needed to locate when paraphasias occur and accurately differentiate between paraphasia types. In this work, we present several methods for performing multiclass paraphasia detection in continuous speech that include phonemic, neologistic, and semantic paraphasias.

Outside of paraphasia detection, Omachi et al. has shown that a sequence-to-sequence model can learn to produce ASR and linguistic annotations in a single sequence [15]. In our work, we explore a similar framework that learns to produce ASR and paraphasia annotations from aphasic speech.

3. Data

We use two English datasets from the AphasiaBank corpus. The first is the Protocol dataset, which contains both aphasic

and healthy speakers performing a series of tasks following the Western Aphasia Battery [16]. The Protocol dataset is first used to train our models as this dataset is the larger of the two, containing roughly 100 hours of speech. The second dataset is the Fridriksson subset, which contains roughly three hours of speech, where PWAs were asked to read a series of scripts. The Fridriksson dataset contains a higher distribution of paraphasias compared to the Protocol dataset, with phonemic, neologistic, and semantic paraphasias representing roughly 13%, 7%, and 3% of the dataset, respectively. Similar to previous works, we use the Fridriksson dataset for finetuning and evaluation.

All utterances were manually transcribed in the CHAT format and included timestamps for both participant and interviewer speech segments [14]. We isolated participant speech and discarded utterances that were labeled unintelligible or containing overlapping speech between the participant and clinician. The CHAT transcription is processed by removing punctuation and reducing all characters to lowercase. For non-word phonological errors transcribed in the International Phonetic Alphabet (IPA) format, we convert these to pseudo-words following previous works [10] by mapping each IPA pronunciation to a sequence of phones and then converting phones to graphemes. Lastly, we label all phonemic, neologistic, and semantic paraphasias, indicated by [* p], [* n], and [* s] in the CHAT transcriptions with [p], [n], and [s] respectively. The final processed transcript can be seen in the Oracle section of Table 1.

4. Methods

4.1. Transcript+GPT

We explore a pipeline that consists of using speech transcripts and a generative pretrained transformer (GPT) model to classify paraphasias. We perform in-context learning on openAI’s GPT-3.5-turbo and GPT-4 by conditioning these models with task instructions and an example. In-context learning is an effective approach for tuning GPT output for a wide variety of NLP tasks without having to update model parameters [17, 18]. We experiment with two approaches for transcript generation, the first uses an in-domain ASR model and the second uses the manual (oracle) transcripts.

4.1.1. ASR+GPT

The ASR model shown in Figure 1a is inspired by previous work [12]. We use both connectionist temporal classification (CTC) loss [19] and cross-entropy (CE) loss. The model is optimized using a joint CTC-attention loss criterion [20]. We also use SpecAugment [21] to resample utterances at different time perturbations using rates of [0.8, 0.9, 0.95, 1.0, 1.05, 1.1, 1.2] which has been shown to be effective for disordered speech recognition [22, 12].

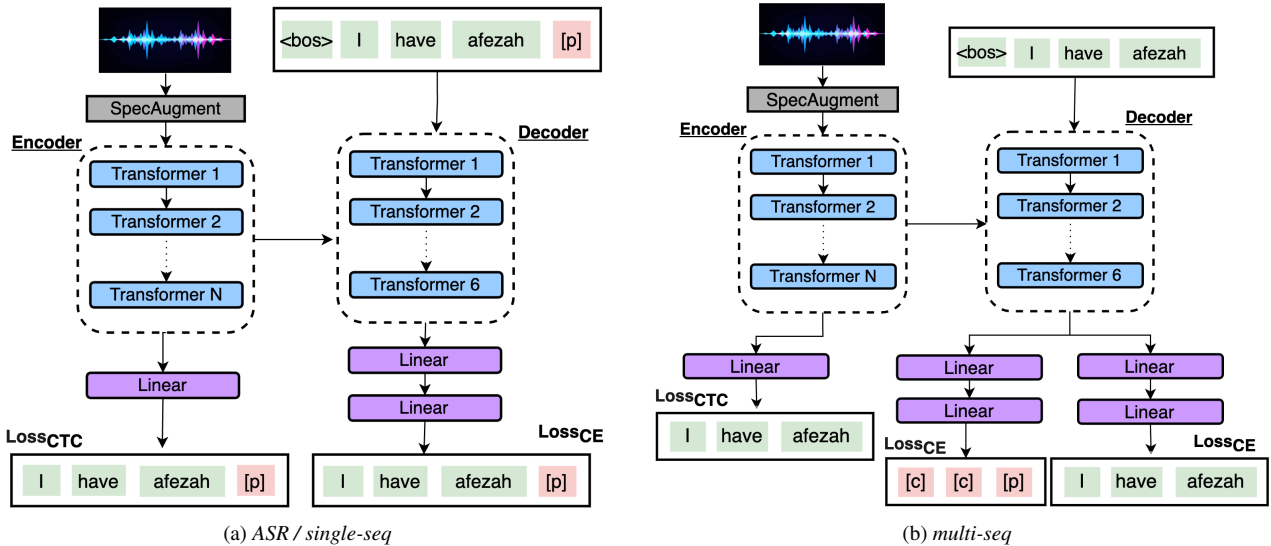


Figure 1: Paraphasia Classification Models.

4.1.2. Oracle+GPT

We explore the effect that ASR errors have on the performance of GPT’s paraphasia classification. We assume perfect ASR transcription (i.e., an oracle) by using the ground truth transcriptions as input to the GPT model. Comparing ASR+GPT and Oracle+GPT will highlight the performance gap in paraphasia classification introduced by inaccurate transcription.

4.2. Single-Seq

The single sequence model (Figure 1a) has the same architecture as the ASR model but the tokenizer has three added special tokens representing [p], [n], and [s] paraphasia labels. The seq2seq model learns to predict subword tokens (graphemes/morphemes) related to ASR and the special tokens related to paraphasia classification. The single-seq model learns to predict paraphasia labels after transcribing a paraphasic word. By treating paraphasias as a separate token, the model can learn contextual dependencies across subword and paraphasia tokens. The loss function for this model can be seen in equation 1, where z_t represents either a subword token or a paraphasia token at time t , M is the combined length of subword tokens and paraphasia labels in the utterance, and x is the audio representation from the encoder.

$$\mathcal{L} = - \sum_{t=0}^M \log P(z_t | z_0, \dots, z_{t-1}, x) \quad (1)$$

4.3. Multi-Seq

The multi-sequence model (shown in Figure 1b) has separate ASR and paraphasia classification heads to produce independent sequences. The paraphasia classification head consists of two fully connected layers with the final layer having an output size of 4 corresponding to [c] (non-paraphasia), [p], [n], and [s]. In this setup, paraphasia classification is performed at the subword level. We use a weighted paraphasia classification loss based on the inverse of the paraphasia class count. The total loss combines both the paraphasia classification loss and the ASR loss seen in equations 2 and 3, where y_t and p_t represent the ASR token and paraphasia class predictions at time t and T is the total length of the subword tokens in the utterance. We conduct a hyperparameter sweep for α ranging from 0.3 to 0.7 and select the optimal value based on development set performance.

$$\mathcal{L}_{asr} = - \sum_{t=0}^T \log P(y_t | y_0, \dots, y_{t-1}, x) \quad (2)$$

$$\mathcal{L}_{para} = - \sum_{t=0}^T \log P(p_t | y_0, \dots, y_{t-1}, x)$$

$$\mathcal{L} = \alpha \mathcal{L}_{asr} + (1 - \alpha) \mathcal{L}_{para} \quad (3)$$

The ASR and paraphasia classification heads use the same decoder representation, which ensures alignment at the subword level. For training, we map paraphasia labels to the subword level by assigning the paraphasia label of a given word to all its constituent subword tokens. At inference, we attain word-level paraphasia predictions by using the majority paraphasia class across all subwords in a given word. During decoding, beam pruning is based on ASR output scores only.

5. Experiment Setup

All our models were built using the SpeechBrain toolkit [23] and consist of an encoder and a decoder with 24- and 6-transformer layers respectively. We initialize the encoder parameters with a pretrained HuBERT model¹. We use a tokenizer with a vocabulary size of 500 that performs reductions using byte-pair encoding (BPE) and use learning rate annealing with a factor of 0.8 based on previous work [12]. Due to hardware constraints we use a batch size of 4 with a gradient accumulation of 4. Our code, model output samples, and additional supplementary information related to data splits, hyperparameters, and prompts for in-context learning can be found on github².

5.1. Standardizing Model Output

Prior to evaluation, we first standardize the output from each model to take the format of ASR+GPT in Table 1. Specifically, we pre-process all model outputs to the form $\hat{Y} = \hat{y}_0, \hat{p}_0, \dots, \hat{y}_n, \hat{p}_n$, where \hat{y}_n and \hat{p}_n represent the predicted word and paraphasia label at n -index.

¹<https://huggingface.co/facebook/hubert-large-ls960-ft>

²<https://github.com/chailab-umich/BeyondBinary-ParaphasiaDetection>

Table 2: Word-level results aggregated over all folds with best results in bold (oracle transcripts not included). † indicates statistical significance over the baseline GPT-4 approach. For all metrics, lower values indicate better performance.

Models	WER	AWER	TD-bin	TD-multiclass			
				[p]	[n]	[s]	all
<i>ASR Transcripts (Baseline)</i>							
GPT-3.5	38.6	32.9	1.17	1.02	0.67	0.30	1.99
GPT-4	38.6	32.7	0.68	0.86	0.73	0.29	1.88
<i>Proposed</i>							
Single-Seq	37.6	32.8	0.63 †	0.76 †	0.45 †	0.31	1.51 †
Multi-Seq	44.8	42.9	0.86	0.90	0.72	0.53	2.15
<i>Oracle Transcripts</i>							
GPT-3.5	–	10.6	0.82	0.97	0.54	0.30	1.81
GPT-4	–	7.9	0.25	0.67	0.44	0.27	1.38

5.2. WER Metrics

Word-error-rate (WER) measures ASR performance, where $\hat{Y} = \hat{y}_0, \dots, \hat{y}_n$ and $Y = y_0, \dots, y_m$ represent the predicted and ground truth sequences used in the evaluation. **Augmented-word-error-rate (AWER)** is used to measure both ASR and paraphasia classification performance, where WER is computed using $\hat{Y} = \hat{y}_0, \hat{p}_0, \dots, \hat{y}_n, \hat{p}_n$ and $Y = y_0, p_0, \dots, y_m, p_m$, which represent the predicted and ground truth sequences. This metric is based on previous work [10, 12].

We use the aggregated labels and predictions across all test folds and compute a single metric similar to [10, 11, 12]. We perform statistical significance testing for WER and AWER metrics using a bootstrap estimate, which is commonly used for determining statistical significance across ASR systems [24]. For each comparison, we perform 1000 iterations using a batch size of 100 and adopt a 95% confidence threshold for declaring statistical significance.

5.3. Distance Metrics

WER metrics have the limitation of measuring accuracy at a given index, which can be misleading for evaluating paraphasia classification if there is misalignment due to ASR. We use temporal distance (TD) to evaluate word-level paraphasia classification, which measures the proximity to the closest paraphasia class [25, 12]. We align the words of Y and \hat{Y} using the minimum edit operations, compute the TD metric using the paraphasia class labels, and normalize the resulting TD by the utterance length. We present several evaluations with TD. **TD-binary** measures the proximity of detecting any paraphasia instance. **TD-multiclass** measures the proximity between specific paraphasias, we compute individual metrics for each paraphasia class (**TD-[p,n,s]**) as well as an overall TD across all classes (**TD-all**). We present the average TD across all test utterances and perform statistical significance testing using a repeated measures ANOVA followed by a post-hoc tukey test, where statistical significance is determined by a p-value < 0.05 .

5.4. Utterance-level Metrics

We assess paraphasia classification at the utterance-level using a binary F1 score for each paraphasia type. Similar to previous work, we compute a single F1 score over all test utterances [10, 12]. For a given paraphasia class, an utterance is assigned a positive class label if there is any instance of that paraphasia present.

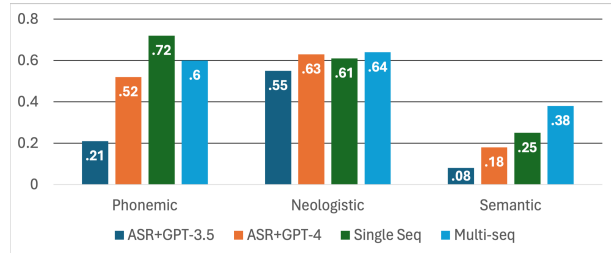


Figure 2: Utterance-level Binary F1-scores

6. Results

In Table 2 we see that single-seq has the lowest WER while GPT-4 achieves the lowest AWER. This indicates that although single-seq may produce more accurate transcriptions, GPT-4 is more accurate when factoring in paraphasia classification. We see that the single-seq model significantly outperforms the baseline GPT-4 for paraphasia detection (TD-binary) and multiclass classification (TD-multiclass) over most categories. Specifically, for phonemic paraphasias, we see that single-seq has the lowest TD-[p] and the highest F1 score (shown in Figure 2) indicating its effectiveness over other models for this paraphasia type. For neologistic paraphasias, we observe that single-seq has the lowest TD-[n], indicating that it is the most effective at locating where neologistic paraphasias occur within an utterance. However, the F1 score for single-seq is slightly below that of GPT-4 and multi-seq, suggesting that it may be slightly less effective (0.61 vs. 0.63 for single-sequence vs. GPT-4, respectively) at identifying the presence of a neologistic paraphasia in an utterance.

Across all approaches, classifying semantic paraphasias is a challenging task. Although TD-[s] is low compared to TD-[p] and TD-[n] this is due to the low representation of semantic paraphasias in the dataset, which ultimately results in low TD when averaged across all utterances. This is further evident when looking at the binary F1-scores for semantic paraphasias, where all approaches struggle at predicting semantic paraphasias at the utterance-level.

The use of oracle transcripts significantly improves GPT-4’s performance in detecting paraphasic instances, which is reflected by a 63% relative performance improvement for TD-binary. For classifying all paraphasias we note a 27% relative improvement for TD-multiclass, which highlights an existing gap for GPT-4 to differentiate between specific paraphasia classes despite having access to perfect transcription.

7. Conclusion

In this work, we present several methods for classifying paraphasias, including the use of transcripts and GPT, a seq2seq model with a single output (single-seq), and a seq2seq model with multi-sequence output (multi-seq). We show the efficacy of GPT with perfect and imperfect transcripts and its limitations for certain paraphasia types. However, within the context of differentiating between all paraphasia classes, we find that the single-seq approach provides statistically significant improvement over the baseline GPT approaches that use ASR for transcription. The performance of multi-seq warrants further investigation into methods where paraphasia predictions are made at the subword level. Future work for improving semantic paraphasia detection should explore including conditioning models with the target script in order to provide contextual information that is relevant for identifying associative word substitutions.

8. Acknowledgements

This paper was accepted to Interspeech 2024.

9. References

- [1] N. A. Association, <https://www.aphasia.org/>, [Online; accessed 01-March-2024].
- [2] N. Helm-Estabrooks and M. Albert, *Manual of Aphasia and Aphasia Therapy*. Pro-Ed, 2004. [Online]. Available: <https://books.google.com/books?id=adYLAQAAMAAJ>
- [3] M. M. Saling, "Chapter 3 - disorders of language," in *Neurology and Clinical Neuroscience*, A. H. Schapira, E. Byrne, S. DiMauro, R. S. Frackowiak, R. T. Johnson, Y. Mizuno, M. A. Samuels, S. D. Silberstein, and Z. K. Wszolek, Eds. Philadelphia: Mosby, 2007, pp. 31–42. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/B9780323033541500079>
- [4] G. Fergadiotis, K. Gorman, and S. Bedrick, "Algorithmic classification of five characteristic types of paraphasias," *American Journal of Speech-Language Pathology*, vol. 25, no. 4S, pp. S776–S787, 2016.
- [5] E. T. McKinnon, J. Fridriksson, A. Basilakos, G. Hickok, A. E. Hillis, M. V. Spampinato, E. Gleichgercht, C. Rorden, J. H. Jensen, J. A. Helpers *et al.*, "Types of naming errors in chronic post-stroke aphasia are dissociated by dual stream axonal loss," *Scientific reports*, vol. 8, no. 1, p. 14352, 2018.
- [6] K. J. Ballard, N. M. Etter, S. Shen, P. Monroe, and C. Tien Tan, "Feasibility of automatic speech recognition for providing feedback during tablet-based treatment for apraxia of speech plus aphasia," *American journal of speech-language pathology*, vol. 28, no. 2S, pp. 818–834, 2019.
- [7] J. Kurland, A. R. Wilkins, and P. Stokes, "ipractice: Piloting the effectiveness of a tablet-based home practice program in aphasia treatment," in *Seminars in speech and language*, vol. 35, no. 01. Thieme Medical Publishers, 2014, pp. 051–064.
- [8] M. Casilio, G. Fergadiotis, A. C. Salem, R. C. Gale, K. McKinney-Bock, and S. Bedrick, "Paralg: A paraphasia algorithm for multinomial classification of picture naming errors," *Journal of Speech, Language, and Hearing Research*, vol. 66, no. 3, pp. 966–986, 2023.
- [9] A. C. Salem, R. Gale, M. Casilio, M. Fleegle, G. Fergadiotis, and S. Bedrick, "Refining semantic similarity of paraphasias using a contextual language model," *Journal of Speech, Language, and Hearing Research*, vol. 66, no. 1, pp. 206–220, 2023.
- [10] D. Le, K. Licata, and E. M. Provost, "Automatic paraphasia detection from aphasic speech: A preliminary study," in *Interspeech*, 2017, pp. 294–298.
- [11] S. Pai, N. Sachdeva, P. Sachdeva, and R. Shah, "Unsupervised paraphasia classification in aphasic speech," in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: Student Research Workshop*, 2020, pp. 13–19.
- [12] M. Perez, D. Le, A. Romana, E. Jones, K. Licata, and E. M. Provost, "Seq2seq for automatic paraphasia detection in aphasic speech," *arXiv preprint arXiv:2312.10518*, 2023.
- [13] D. Mirman, T. J. Strauss, A. Brecher, G. M. Walker, P. Sobel, G. S. Dell, and M. F. Schwartz, "A large, searchable, web-based database of aphasic performance on picture naming and other tests of cognitive function," *Cognitive neuropsychology*, vol. 27, no. 6, pp. 495–504, 2010.
- [14] B. MacWhinney, *The CHILDES project: Tools for analyzing talk, Volume I: Transcription format and programs*. Psychology Press, 2014.
- [15] M. Omachi, Y. Fujita, S. Watanabe, and M. Wiesner, "End-to-end asr to jointly predict transcriptions and linguistic annotations," in *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2021, pp. 1861–1871.
- [16] A. Kertesz, "The western aphasia battery: A systematic review of research and clinical applications," *Aphasiology*, vol. 36, no. 1, pp. 21–50, 2022.
- [17] E. Perez, D. Kiela, and K. Cho, "True few-shot learning with language models," *Advances in neural information processing systems*, vol. 34, pp. 11 054–11 070, 2021.
- [18] Y. Lu, M. Bartolo, A. Moore, S. Riedel, and P. Stenetorp, "Fantastically ordered prompts and where to find them: Overcoming few-shot prompt order sensitivity," *arXiv preprint arXiv:2104.08786*, 2021.
- [19] A. Graves, S. Fernández, F. Gomez, and J. Schmidhuber, "Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks," in *Proceedings of the 23rd international conference on Machine learning*, 2006, pp. 369–376.
- [20] S. Kim, T. Hori, and S. Watanabe, "Joint ctc-attention based end-to-end speech recognition using multi-task learning," in *2017 IEEE international conference on acoustics, speech and signal processing (ICASSP)*. IEEE, 2017, pp. 4835–4839.
- [21] D. S. Park, W. Chan, Y. Zhang, C.-C. Chiu, B. Zoph, E. D. Cubuk, and Q. V. Le, "SpecAugment: A simple data augmentation method for automatic speech recognition," *Interspeech 2019*, 2019.
- [22] J. R. Green, R. L. MacDonald, P.-P. Jiang, J. Cattiau, R. Heywood, R. Cave, K. Seaver, M. A. Ladewig, J. Tobin, M. P. Brenner *et al.*, "Automatic speech recognition of disordered speech: Personalized models outperforming human listeners on short phrases," in *Interspeech*, 2021, pp. 4778–4782.
- [23] M. Ravanelli, T. Parcollet, P. Plantinga, A. Rouhe, S. Cornell, L. Lugosch, C. Subakan, N. Dawalatabad, A. Heba, J. Zhong *et al.*, "Speechbrain: A general-purpose speech toolkit," *arXiv preprint arXiv:2106.04624*, 2021.
- [24] M. Bisani and H. Ney, "Bootstrap estimates for confidence intervals in asr performance evaluation," in *2004 IEEE International Conference on Acoustics, Speech, and Signal Processing*, vol. 1. IEEE, 2004, pp. I–409.
- [25] G. Kovács, G. Sebestyen, and A. Hangan, "Evaluation metrics for anomaly detection algorithms in time-series," *Acta Universitatis Sapientiae, Informatica*, vol. 11, no. 2, pp. 113–130, 2019.