# Concurrent Validity of Automatic Speech and Pause Measures During Passage Reading in ALS

Saeid Alavi Naeini[1,2], Leif Simmatis[1], Yana Yunusova[1,3], and Babak Taati[1,2,4,5]

[1]Kite Research Institute, Toronto Rehabilitation Institute – University Health Network
[2]Institute of Biomedical Engineering, University of Toronto, Toronto, ON, Canada
[3]Department of Speech-Language Pathology, University of Toronto, Toronto, Canada
[4]Department of Computer Science, University of Toronto, Toronto, ON, Canada
[5]Vector Institute, Toronto, ON, Canada

*Abstract*—The analysis of speech measures in individuals with amyotrophic lateral sclerosis (ALS) can provide essential information for early diagnosis and tracking disease progression. However, current methods for extracting speech and pause features are manual or semi-automatic, which makes them time–consuming and labour–intensive. The advent of speech-text alignment algorithms provides an opportunity for inexpensive, automated, and accurate analysis of speech measures in individuals with ALS. There is a need to validate speech and pause features calculated by these algorithms against current gold standard methods. In this study, we extracted 8 speech/pause features from 646 audio files of individuals with ALS and healthy controls performing passage reading. Two pretrained forced alignment models – one using transformers and another using a Gaussian mixture / hidden Markov architecture – were used for automatic feature extraction. The results were then validated against semi-automatic speech/pause analysis software, with further subgroup analyses based on audio quality and disease severity. Features extracted using transformer-based forced alignment had the highest agreement with gold standards, including in terms of audio quality and disease severity. This study lays the groundwork for future intelligent diagnostic support systems for clinicians, and for novel methods of tracking disease progression remotely from home.

*Keywords*-Bulbar ALS; concurrent validity; speech; pause; forced alignment

## I. INTRODUCTION

Amyotrophic lateral sclerosis (ALS) is a rapidly progressing neurodegenerative disorder that destroys motor neurons in the cerebral cortex, brainstem, and spinal cord [1]. The emergence of bulbar signs is an essential milestone in the progression of ALS, as it causes the loss of speech intelligibility and, consequently, has a significant impact on the quality of life of people with ALS. The loss of communication, especially speech, has been reported as one of the worst aspects of the disease [2].

There are established clinical assessment techniques for diagnosing bulbar symptoms associated with ALS and tracking disease severity; however, these have important limitations. Clinical impressions of range of motion, speed, and symmetry of oral musculature are valuable, but these assessments are typically subjective [3]. The reliability of these clinician-based judgements have been questioned and lack of reliable assessment tools contributes to long diagnostic delay in ALS (up to 18 months) [4]. Earlier detection of bulbar symptoms in ALS could facilitate finding augmentative and alternative communication referrals, and assist early bulbar symptom management [5]. There is a substantial need to develop accessible, automated systems that can objectively and reliably assess bulbar ALS symptoms in nonclinical settings.

Speech timing can be broadly measured using the Speech Intelligibility Test (SIT) which provides a measure of speaking rate [6]. Speaking rate has been recommended as the preferred objective measure of bulbar decline by speech-language pathologists [7], [8]. Passage reading tasks are excellent alternatives to the SIT, since they can also provide measures of speech and pause duration [9]. Studies have shown that these features can distinguish bulbar disease at different stages, with or without respiratory deficits [10].

Current methods of obtaining speech and pause features are either manual or semi-automatic, which require time-intensive processing and training of professional staff [11]. Forced alignment techniques – which temporally align passage text and speech audio signals – could potentially overcome these limitations. Past attempts to automatically parse and align dysarthric speech have been error-prone [12], [13].

Transformers are a new class of deep learning models that have achieved state-of-the-art performance in a variety of tasks, e.g. in natural language processing or computer vision [14], [15]. The recent development of robust and accurate transformer-based forced alignment models [16] presents a new opportunity for fully automated speech/pause feature extraction in dysarthric speech. However, it is an open question as to whether novel methods could replace current gold standards for dysarthric speech alignment.

The aim of this study is to establish the concurrent validity of automatic speech-text alignment algorithms against

Table I: Spearman correlation ($\rho$) between features calculated from the Wav2Vec2 method and from the SPA software, divided based on audio quality. The number of audio files in each category is presented in parenthesis

| | Good (314) | | Fair (120) | | Poor (212) | |
|---|---|---|---|---|---|---|
| | $\rho$ | p | $\rho$ | p | $\rho$ | p |
| *Pause duration* | 0.82 | $< .001$ | 0.84 | $< .001$ | 0.62 | $< .001$ |
| *Total duration* | 0.98 | $< .001$ | 0.96 | $< .001$ | 0.93 | $< .001$ |
| *Speech duration* | 0.96 | $< .001$ | 0.91 | $< .001$ | 0.70 | $< .001$ |
| *Pause event* | 0.73 | $< .001$ | 0.80 | $< .001$ | 0.58 | $< .001$ |
| *% Pause* | 0.72 | $< .001$ | 0.75 | $< .001$ | 0.48 | $< .001$ |
| *Mean phrase* | 0.77 | $< .001$ | 0.70 | $< .001$ | 0.50 | $< .001$ |
| *CV phrase duration* | $-0.05$ | 0.39 | 0.12 | 0.18 | 0.01 | 0.86 |
| *CV pause duration* | 0.53 | $< .001$ | 0.38 | $< .001$ | 0.27 | $< .001$ |

Table II: Spearman correlation ($\rho$) between features calculated from the MFA model and from the SPA software, divided based on audio quality. The number of audio files in each category is presented in parenthesis

| | Good (314) | | Fair (120) | | Poor (212) | |
|---|---|---|---|---|---|---|
| | $\rho$ | p | $\rho$ | p | $\rho$ | p |
| *Pause duration* | 0.77 | $< .001$ | 0.74 | $< .001$ | 0.52 | $< .001$ |
| *Total duration* | 0.90 | $< .001$ | 0.86 | $< .001$ | 0.85 | $< .001$ |
| *Speech duration* | 0.87 | $< .001$ | 0.86 | $< .001$ | 0.76 | $< .001$ |
| *Pause event* | 0.74 | $< .001$ | 0.64 | $< .001$ | 0.56 | $< .001$ |
| *% Pause* | 0.65 | $< .001$ | 0.62 | $< .001$ | 0.40 | .04 |
| *Mean phrase* | 0.60 | $< .001$ | 0.50 | $< .001$ | 0.58 | $< .001$ |
| *CV phrase duration* | $-0.05$ | .51 | $-0.38$ | .05 | 0.02 | .86 |
| *CV pause duration* | 0.32 | $< .001$ | 0.31 | .11 | 0.27 | .11 |

the Speech and Pause Analysis (SPA) software, which is a validated semi-automated pause identification tool [11]. We investigate two open-source forced alignment models – including a new transformer-based model [16] – and report correlations between 8 speech and pause features calculated via SPA vs. via forced alignment. Our hypothesis is that features calculated via automatic text-speech alignment will have strong relationship with their SPA counterparts.

## II. METHODS

### A. Participants

We used a dataset of 646 speech recordings, obtained from 3 observational longitudinal studies of individuals with ALS and healthy controls (HC). The recordings were of the participants reading a short (60 word) passage – the bamboo passage [9], [10].

The total number of participants was 462 (ALS: 243, controls: 219). Additional demographic and clinical details are presented in [11]. The time interval between recordings varied from 3 to 6 months when more than one recording was available from a participant. Individuals with ALS were diagnosed with possible, probable, or definite ALS, as defined by the El Escorial Criteria [17].

Speaking rate was obtained for each speaker and session using the SIT and was calculated as the number of words per minute (WPM). Similar to [3], we classified the individuals with ALS into 3 groups of mild, moderate, and severe based on their speaking rate. The mild group comprised of individuals with speaking rates greater than 160 WPM; the moderate group included individuals with speaking rates between 120 and 160 WPM; and the severe group had a speaking rate less than 120 WPM.

We also categorized the quality of audio files into 3 groups of Good, Fair, and Poor based on signal-to-noise ratio (SNR) of each acoustic file. Following previous research [18], audio files with SNR below 15 dB were labeled as Poor, 15–20 dB were labeled as Fair, and above 20 dB were labeled as Good. The reason for having different quality groups were due to recordings being obtained in a clinical (in contrast to a laboratory) context, resulting in environmental noise and quiet/muffed recordings in some cases.

### B. Procedure

We applied 2 different forced alignment techniques for the automatic parsing and speech/text alignment of the Bamboo passage - a standard passage used for the analysis of bulbar disease severity in ALS [9]. All recordings were also

Table III: Spearman correlation ($\rho$) between features calculated from the Wav2Vec2 model on Good audio quality data and from the SPA software, divided based on disease severity. The number of audio files in each category is presented in parenthesis

| | HC (147) | | Mild (70) | | Moderate (59) | | Severe (38) | |
|---|---|---|---|---|---|---|---|---|
| | $\rho$ | p | $\rho$ | p | $\rho$ | p | $\rho$ | p |
| *Pause duration* | 0.82 | < .001 | 0.84 | < .001 | 0.66 | < .001 | 0.94 | < .001 |
| *Total duration* | 0.97 | < .001 | 0.95 | < .001 | 0.96 | < .001 | 0.99 | < .001 |
| *Speech duration* | 0.93 | < .001 | 0.93 | < .001 | 0.90 | < .001 | 0.98 | < .001 |
| *Pause event* | 0.53 | < .001 | 0.73 | < .001 | 0.60 | < .001 | 0.92 | < .001 |
| *% Pause* | 0.74 | < .001 | 0.82 | < .001 | 0.54 | < .001 | 0.83 | .001 |
| *Mean phrase* | 0.76 | < .001 | 0.80 | < .001 | 0.68 | < .001 | 0.91 | < .001 |
| *CV phrase duration* | 0.04 | .61 | −0.25 | .03 | 0.04 | .60 | −0.10 | .53 |
| *CV pause duration* | 0.32 | < .001 | 0.59 | < .001 | 0.67 | < .001 | 0.74 | < .001 |

Table IV: Spearman correlation ($\rho$) between features calculated from the MFA model on Good audio quality data and from the SPA software, divided based on disease severity. The number of audio files in each category is presented in parenthesis

| | HC (147) | | Mild (70) | | Moderate (59) | | Severe (38) | |
|---|---|---|---|---|---|---|---|---|
| | $\rho$ | p | $\rho$ | p | $\rho$ | p | $\rho$ | p |
| *Pause duration* | 0.64 | < .001 | 0.78 | < .001 | 0.51 | < .001 | 0.75 | < .001 |
| *Total duration* | 0.86 | < .001 | 0.92 | < .001 | 0.85 | < .001 | 0.44 | .07 |
| *Speech duration* | 0.92 | < .001 | 0.66 | .002 | 0.65 | .002 | 0.63 | .007 |
| *Pause event* | 0.52 | < .001 | 0.74 | < .001 | 0.77 | < .001 | 0.70 | .002 |
| *% Pause* | 0.68 | < .001 | 0.77 | < .001 | 0.41 | .07 | 0.41 | .09 |
| *Mean phrase* | 0.52 | < .001 | 0.74 | < .001 | 0.70 | < .001 | 0.39 | .12 |
| *CV phrase duration* | −0.02 | .83 | 0.13 | .60 | −0.08 | .73 | −0.10 | .71 |
| *CV pause duration* | 0.15 | .14 | 0.56 | .01 | 0.53 | .02 | 0.43 | .08 |

previously parsed using the semi-automated SPA software. The first forced alignment method was the Montreal Forced Aligner (MFA) [19]. The automatic speech recognition pipeline in MFA is based on a Gaussian mixture/hidden Markov (GMM/HMM) architecture from the open-source Kaldi speech recognition toolkit [19]. MFA uses Mel-frequency cepstral coefficients (MFCCs) as acoustic input features. For the second forced alignment method, we used a pretrained transformer-based Wav2Vec2 model [16] from the Charsiu library [20]. The model can perform both text-dependent and text-independent phone-to-audio alignment and maintains good performance in different settings. In all cases, we performed speech-text alignment in a supervised manner, i.e. by providing the text of bamboo passage as an input to the model.

*C. Feature Extraction*

The forced alignment output comprised the time stamps associated with each speech and pause event greater than 300 milliseconds. We then extracted the same 8 features used in [11] based on their ability to distinguish individuals with ALS from controls. They included:

1) *Pause duration* (sec): total duration of all pause events (excluding the start and end pauses).
2) *Total duration* (sec): recording time from the onset of the first sentence to the end of the last sentence.
3) *Speech duration* (sec): total duration of all speech events (pauses excluded).
4) *Pause events*: number of pauses while reading the passage.
5) *% Pause*: percentage of total reading time spent pausing (only pauses between speech events).
6) *Mean phrase* (sec): average duration of a phrase. Phrases were defined as sections of continuous speech without a pause.
7) *CV phrase duration*: Coefficient of variation of phrase durations (a normalized measure of variability).
8) *CV pause duration*: Coefficient of variation of pause durations.
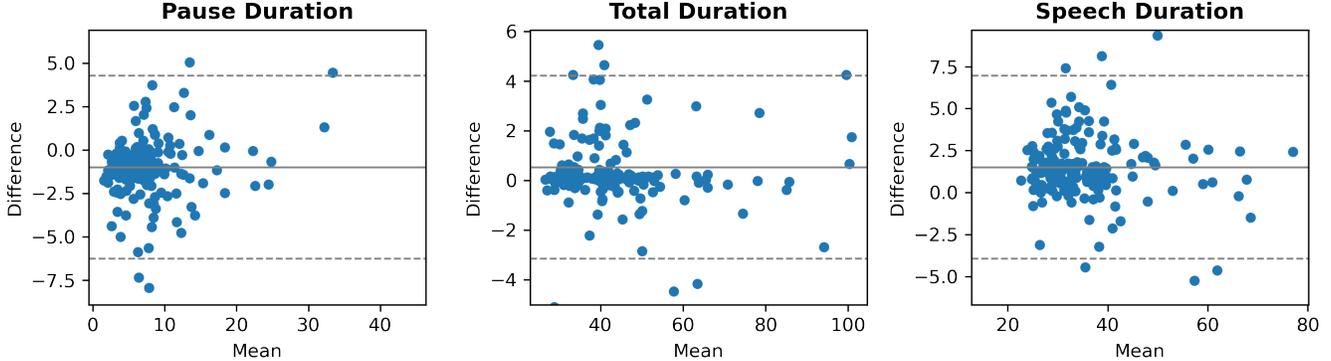
Figure 1: The Bland–Altman plots between Wav2Vec2 and SPA features of *pause duration*, *speech duration*, and *total duration* on Good audio quality data. The dashed lines are the lower and upper limits of agreement (1.96×standard deviation) and the solid line is zero.

## D. Evaluation

Spearman's correlation coefficients ($\rho$) were used to determine the agreement between the features extracted from forced alignment algorithms and from the SPA software. We compared these sets of 8 features across 3 different conditions: 1) two forced alignment algorithms (MFA, and Wav2Vec2), 2) four severity levels (HC, and ALS mild, moderate, or severe), and 3) three audio qualities (Good, Fair, Poor). The correlation coefficients ($\rho$) and p-values are reported (significance threshold was p<0.05). We also present the agreement between a subset of the Wav2Vec2 features and their SPA counterparts using Bland–Altman plots. The subset of 3 features was selected as those that were previously reported to have the highest area under the receiver operating characteristic curve (AUC-ROC) in detecting bulbar changes in otherwise bulbar presymptomatic individuals with ALS [11].

## III. RESULTS

The results of correlation analysis for the transformer-based (Wav2Vec2) and HMM-based (MFA) models are outlined in Tables I and II, respectively. Performance is separately reported for recordings based on their audio quality. Tables III and IV present the correlation results for Wav2Vec2 and MFA models based on disease severity. This analysis is performed on Good quality data to ensure disease severity is the only variable affecting the correlation.

The Bland–Altman plots for *pause duration*, *speech duration*, and *total duration* are shown in Figure 1. The plots illustrate the agreement between features extracted using the Wav2Vec2 model vs. using the SPA software. For all three features, there is similar agreement between Wav2Vec2 and SPA features as all 3 plots have relatively narrow limits of agreement band (±1.96×standard deviation) which contains 95% of the values.

## IV. DISCUSSION

The current gold standard for the parsing of clinical audio data and feature extraction is the manual or semi-automatic approach, which is time–consuming and labour–intensive. Here, we compared two novel, automated methods of audio-based parsing to SPA, which has been previously validated for parsing clinical audio data [9]. We demonstrated that, although both MFA and Wav2Vec2 performed well with respect to SPA, Wav2Vec2 generalized better across clinical severities. However, both methods were sensitive to audio data quality. The present work represents an important step toward developing automated tools for assessing neurodegenerative disorders, including ALS, in the home setting.

The general pattern than can be observed in the correlation results is that the Wav2Vec2 model performed better with most features being very strongly (or strongly) associated. We also observed that correlations between SPA and automated methods appeared to change as a function of audio sample quality (Tables I and II), with Good audio having the strongest correlations. This agrees with the results of previous studies; for example, the performance of forced alignment algorithms is reduced when sound to noise ratio is not optimized [19], [20]. Our results reinforce the need for high quality audio data in the assessment of neurological and neurodegenerative diseases.

The present results add to existing literature and support previous calls to incorporate more flexible methods into automated and semi-automated speech analysis. In [21], automated diadochokinetic analysis methods were explored, and it was found that automated methods' performance suffered as clinical severity increased. The results of the present study suggest that this could be alleviated by using flexible and robust methods, such as Wav2Vec2-based forced alignment in our current passage-reading paradigm. Importantly, we observed that correlations between SPA/Wave2Vec2 across clinical severities (Tables 3 and 4) were almost uniformly

better than they were in the SPA/MFA comparison. This suggests that transformer-based models are more robust to range of clinical severities and are viable replacements for current manual and semi-automatic procedures used for clinical audio data analysis.

In conclusion, our work has demonstrated that automated forced alignment methods perform well compared to a semi-automated ground truth, and that the Wav2Vec2-based approach in particular is robust across a range of clinical severity levels. Importantly, our results also highlight the need for good-quality audio data to be used for clinical speech analyses. These results have important implications for the development of new clinical tools for home-based speech assessment, and also for future work that will aim to prognose patients based on speech function. These are challenging problems [22], [23] that will require further future work. It will also be important in future work to extend our current findings to include other populations that experience speech impairments, such as stroke and Parkinson's disease.

## V. Acknowledgements

## References

[1] M. C. Kiernan, S. Vucic, B. C. Cheah, M. R. Turner, A. Eisen, O. Hardiman, J. R. Burrell, and M. C. Zoing, "Amyotrophic lateral sclerosis," *The lancet*, vol. 377, no. 9769, pp. 942–955, 2011.

[2] M. Hecht, T. Hillemacher, E. Gräsel, S. Tigges, M. Winterholler, D. Heuss, M.-J. Hilz, and B. Neundörfer, "Subjective experience and coping in als," *Amyotrophic Lateral Sclerosis and Other Motor Neuron Disorders*, vol. 3, no. 4, pp. 225–231, 2002.

[3] S. Shellikeri, J. R. Green, M. Kulkarni, P. Rong, R. Martino, L. Zinman, and Y. Yunusova, "Speech movement measures as markers of bulbar disease in amyotrophic lateral sclerosis," *Journal of Speech, Language, and Hearing Research*, vol. 59, no. 5, pp. 887–899, 2016.

[4] M. R. Turner, J. Scaber, J. A. Goodfellow, M. E. Lord, R. Marsden, and K. Talbot, "The diagnostic pathway and prognosis in bulbar-onset amyotrophic lateral sclerosis," *Journal of the neurological sciences*, vol. 294, no. 1-2, pp. 81–85, 2010.

[5] J. R. Green, Y. Yunusova, M. S. Kuruvilla, J. Wang, G. L. Pattee, L. Synhorst, L. Zinman, and J. D. Berry, "Bulbar and speech motor assessment in als: Challenges and future directions," *Amyotrophic Lateral Sclerosis and Frontotemporal Degeneration*, vol. 14, no. 7-8, pp. 494–500, 2013.

[6] M. Dorsey, K. Yorkston, D. Beukelman, and M. Hakel, "Speech intelligibility test for windows," *Institute for Rehabilitation Science and Engineering at Madonna*, 2007.

[7] K. M. Yorkston, "Speech deterioration in amyotrophic lateral sclerosis: Implications for the timing of intervention," *Jounal of Medical Speech-Language Pathology*, vol. 1, pp. 35–46, 1993.

[8] L. J. Ball, A. Willis, D. R. Beukelman, and G. L. Pattee, "A protocol for identification of early bulbar signs in amyotrophic lateral sclerosis," *Journal of the neurological sciences*, vol. 191, no. 1-2, pp. 43–53, 2001.

[9] J. R. Green, D. R. Beukelman, and L. J. Ball, "Algorithmic estimation of pauses in extended speech samples of dysarthric and typical speech," *Journal of medical speech-language pathology*, vol. 12, no. 4, p. 149, 2004.

[10] Y. Yunusova, N. L. Graham, S. Shellikeri, K. Phuong, M. Kulkarni, E. Rochon, D. F. Tang-Wai, T. W. Chow, S. E. Black, L. H. Zinman *et al.*, "Profiling speech and pausing in amyotrophic lateral sclerosis (als) and frontotemporal dementia (ftd)," *PloS one*, vol. 11, no. 1, p. e0147573, 2016.

[11] C. Barnett, J. R. Green, R. Marzouqah, K. L. Stipancic, J. D. Berry, L. Korngut, A. Genge, C. Shoesmith, H. Briemberg, A. Abrahao *et al.*, "Reliability and validity of speech & pause measures during passage reading in als," *Amyotrophic Lateral Sclerosis and Frontotemporal Degeneration*, vol. 21, no. 1-2, pp. 42–50, 2020.

[12] M. Pleva, J. Juhár, and A. S. Thiessen, "Automatic acoustic speech segmentation in praat using cloud based asr," in *2015 25th International Conference Radioelektronika (RADIOELEKTRONIKA)*. IEEE, 2015, pp. 172–175.

[13] Y. T. Yeung, K. H. Wong, and H. Meng, "Improving automatic forced alignment for dysarthric speech transcription," in *Sixteenth Annual Conference of the International Speech Communication Association*, 2015.

[14] 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.

[15] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," *arXiv preprint arXiv:2010.11929*, 2020.

[16] A. Baevski, Y. Zhou, A. Mohamed, and M. Auli, "wav2vec 2.0: A framework for self-supervised learning of speech representations," *Advances in Neural Information Processing Systems*, vol. 33, pp. 12 449–12 460, 2020.

[17] B. R. Brooks, R. G. Miller, M. Swash, and T. L. Munsat, "El escorial revisited: revised criteria for the diagnosis of amyotrophic lateral sclerosis," *Amyotrophic lateral sclerosis and other motor neuron disorders*, vol. 1, no. 5, pp. 293–299, 2000.

[18] H. Xu and H. Zheng, "The simple SNR estimation algorithms for mpsk signals," in *Proceedings 7th International Conference on Signal Processing, 2004. Proceedings. ICSP'04. 2004.*, vol. 2. IEEE, 2004, pp. 1781–1785.

[19] M. McAuliffe, M. Socolof, S. Mihuc, M. Wagner, and M. Sonderegger, "Montreal forced aligner: Trainable text-speech alignment using kaldi." in *Interspeech*, vol. 2017, 2017, pp. 498–502.

[20] J. Zhu, C. Zhang, and D. Jurgens, "Phone-to-audio alignment without text: A semi-supervised approach," *arXiv preprint arXiv:2110.03876*, 2021.

[21] C. Tanchip, D. L. Guarin, S. McKinlay, C. Barnett, S. Kalra, A. Genge, L. Korngut, J. R. Green, J. Berry, L. Zinman *et al.*, "Validating automatic diadochokinesis analysis methods across dysarthria severity and syllable task in amyotrophic lateral sclerosis," *Journal of Speech, Language, and Hearing Research*, vol. 65, no. 3, pp. 940–953, 2022.

[22] F. Kimura, C. Fujimura, S. Ishida, H. Nakajima, D. Furutama, H. Uehara, K. Shinoda, M. Sugino, and T. Hanafusa, "Progression rate of alsfrs-r at time of diagnosis predicts survival time in als," *Neurology*, vol. 66, no. 2, pp. 265–267, 2006.

[23] J. R.-J. Lee, J. F. Annegers, and S. H. Appel, "Prognosis of amyotrophic lateral sclerosis and the effect of referral selection," *Journal of the neurological sciences*, vol. 132, no. 2, pp. 207–215, 1995.