# LATENT REPRESENTATION ENCODING AND MULTI MODAL BIOMARKERS FOR POST-STROKE SPEECH AS SESSMENT

Anonymous authors

006

008 009 010

011

013

014

015

016

017

018

019

021

023

025

026

027

028

029

031

032

035

Paper under double-blind review

#### Abstract

Post-stroke language impairments arise from disruptions in neurophysiological pathways controlling speech production, affecting lexical, semantic, syntactic, and articulatory-prosodic functions. These deficits extend from impaired cognitivemotor planning to execution, manifesting as altered vocal fold dynamics that compromise speech fluency and intelligibility. The high-dimensional and multimodal nature of these impairments poses significant challenges to traditional assessment methods, necessitating automated solutions that can capture the heterogeneity of disfluencies. We present a multimodal framework that integrates foundation model embeddings with clinically relevant biomarkers for comprehensive speech assessment. Leveraging a purpose-built database of 600 poststroke patients, we fine-tune Whisper to extract encoder embeddings that capture pathological speech characteristics. These representations are integrated with linguistic complexity metrics, physiological glottal parameters, and acoustic features through neural networks, enabling biologically informed assessment. Our model achieves 92.4% classification accuracy in stroke detection, outperforming feature-based methods, with SHAP analysis validating modality-specific contributions. We further demonstrate clinical applicability through severity prediction on Comprehensive Aphasia Test (CAT) scores, achieving an N-RMSE of 0.1299. By combining speech-derived representations with domain-specific neurophysiological markers, this framework provides a scalable approach for automated diagnosis, severity tracking, and precision rehabilitation in post-stroke language disorders.

034 1 INTRODUCTION

Speech, as a biological signal, conveys information across multiple levels, from the physical vibrations of the vocal folds to the lexical and semantic structures of language. This multidimensional 037 nature makes speech a robust biomarker for the evaluation of physiological and cognitive health, particularly in neurological conditions (Ramanarayanan et al., 2022). However, speech analysis is inherently complex due to inter-individual variability in vocal tract anatomy, speaking patterns, and 040 disease-specific manifestations (Stefaniak et al., 2022; Olafson et al, 2024). Effective speech anal-041 ysis would require the integration of acoustic features (e.g., fundamental frequency, jitter, shimmer, 042 formant trajectories) and linguistic measures (e.g., lexical diversity, syntactic complexity, semantic 043 coherence). Traditional diagnostic tools, such as clinician-rated scales or standardised language tests 044 (e.g., Boston Diagnostic Aphasia Examination; Roth 2011), though valuable, are limited by their reliance on subjective interpretation, lack of granularity, and inability to capture subtle, real-time changes in speech. Additionally, these methods are often costly and resource-intensive, requiring 046 specialised training and significant time investment, which restricts their scalability and utility in 047 large-scale or remote clinical settings (Mahmoud et al., 2023; Le et al., 2018; Palmer & Enderby, 048 2012). These limitations underscore the need for objective, automated, and data-driven speech anal-049 ysis tools that can provide quantifiable, reproducible, and clinically actionable insights. 050

Recently, large-scale foundation models have emerged as a promising approach for medical analysis,
 particularly following the introduction of BioBERT (Lee et al., 2020) and ClinicalBERT (Alsentzer
 et al., 2019) in early 2019, both of which built upon BERT (Devlin, 2018). The advancements of these Large Language Models (LLMs) have been instrumental in enhancing foundation models,

particularly in their capacity to process and generate domain-specific, context-rich language-an es-055 sential factor for accurate language analysis in healthcare settings. More recently, multimodal LLMs 056 such as GPT-4V (Achiam et al., 2023) have demonstrated potential in supporting clinical decision-057 making and management, though challenges in regulation and validation remain (Qiu et al., 2024). Similarly, Me-LLaMA (Xie et al., 2024), a newly developed medical LLM family, has exhibited 058 superior performance across general and domain-specific medical tasks compared to existing opensource medical LLMs. Building on these advancements, recent research has explored the application 060 of Whisper (Radford et al., 2022) specifically on speech disorders, where its pre-trained representa-061 tions have demonstrated strong generalisability across diverse clinical contexts (Leung et al., 2024; 062 Jiang et al., 2024; Li & Zhang, 2024; Sanguedolce et al., 2023; 2024; Lee et al., 2024; Rathod et al., 063 2023; Best et al., 2024). Unlike traditional machine learning methods that rely on task-specific 064 training and extensive labeled datasets, foundation models leverage pre-training on diverse speech 065 corpora to learn robust representations that generalize well in low-resource scenarios. This is partic-066 ularly valuable for medical applications where labeled data is scarce and speech variability is high 067 due to disease severity and individual differences. However, while these models provide a strong 068 foundation, they still require domain-specific fine-tuning to capture the nuanced characteristics of 069 post-stroke speech disorders. The primary challenge remains the scarcity of large, clinically annotated datasets suitable for fine-tuning, limiting generalization across diverse patient populations and 070 clinical settings. 071

072 To address this, we introduce a robust in-house post-stroke speech dataset of  $\approx 1000$  patients, devel-073 oped over the years through a collaboration between different hospitals and our research facilities. 074 This corpus, specifically designed to capture the inherent variability and complexity of post-stroke 075 speech patterns, represents a significant step toward bridging the gap between general-purpose foundation models and the demands of clinical applications in speech pathology. Our approach lever-076 ages OpenAI Whisper (Radford et al., 2022), a state-of-the-art speech recognition model, to analyse 077 pathological speech. Whisper's architecture is particularly well-suited for this task because of its demonstrated near-human-level accuracy in low-resource settings and robust performance across di-079 verse speaking conditions (Radford et al., 2022). The model's effectiveness stems from extensive 080 pre-training on 680 000 hours of speech data encompassing varied acoustic conditions, speakers, 081 and languages. For clinical deployment, we implement a secure fine-tuning pipeline that addresses the unique challenges of medical data handling. All model adaptations are performed within an 083 isolated, on-site computing environment, ensuring compliance with privacy regulations while min-084 imising external network dependencies. Such approach ensures sensitive patient data throughout the 085 fine-tuning process. After adaptation, we extract embeddings from the model's weighted encoder, 086 capturing high-level representations of pathological speech patterns. These learned representations are then integrated into our multimodal classification framework alongside clinically relevant fea-087 tures - including glottal, linguistic, and acoustic parameters - to enhance diagnostic precision while 880 maintaining interpretability. This combination leverages both the strengths of foundation models 089 and the domain expertise encoded in more established clinical metrics. 090

091 To validate our framework's potential as an automated clinical assessment tool, we evaluate its 092 ability to replicate expert clinical judgment through regression analysis. Specifically, we predict the scores from Comprehensive Aphasia Test (CAT; Swinburn et al. 2004), a standardised metric 093 traditionally assigned by speech therapists through time-intensive manual assessment. This regres-094 sion task represents a step towards automated clinical decision support, as accurately predicting 095 CAT scores would enable rapid, objective, and consistent evaluation of speech impairment severity. 096 This capability, combined with our multimodal classification approach, provides a foundation for healthcare AI systems that could support the assessment and management of diverse speech and 098 communication disorders in clinical settings.

100 101

102

103

# 2 Methods

#### 2.1 POST-STROKE SPEECH DATABASE

The database used in this study is an in-house <sup>1</sup> comprehensive corpus of post-stroke speech, developed for clinical and scientific research aimed at improving Automatic Speech Recognition (ASR)

<sup>&</sup>lt;sup>1</sup>To ensure author anonymity, the name, link, ethics and database references will be added after the review process.

108 systems for disordered speech. It includes speech recordings from approximately 6000 healthy con-109 trols and 1000 individuals with a history of stroke, collected as part of two longitudinal studies 110 conducted in collaboration with a leading research institution and an associated national healthcare 111 network. All participants provided informed consent in accordance with ethical approval prior to 112 data collection. The speech recordings include picture description tasks based on standard clinical assessments of the Comprehensive Aphasia Test (CAT) and a beach scene picture stimulus we 113 designed. Additionally, the dataset features detailed orthographic English transcriptions performed 114 by trained speech pathologists, as well as phonetic transcriptions using the International Phonetic 115 Alphabet. The labeled dataset used for this study consists of 794 audio recordings from 578 unique 116 individuals, some of whom participated in multiple recording sessions to capture longitudinal recov-117 ery and individual speech pattern variations. In total, the dataset comprises approximately 15 hours 118 of speech data. The dataset is predominantly male (70.79%), with an average age of 61.96 years at 119 the time of testing, while female speakers had an average age of 58.52 years. This gender imbal-120 ance aligns with global trends indicating a higher incidence of stroke among males (Appelros et al., 121 2009). Furthermore, as expected in stroke populations, the age distributions exhibit a left-skewed 122 pattern.

123

# 124 2.2 DATA PRE-PROCESSING

126 Before fine-tuning Whisper, data preparation was conducted to support the extraction of embeddings 127 and linguistic features. To ensure consistency, trained speech therapists handled the transcription 128 process, resulting in an inter-rater reliability of 73%. The transcriptions followed the standard-129 ised Codes for the Human Analysis of Transcripts (CHAT; MacWhinney 2014) and were further processed using the Computerised Language ANalysis software (CLAN;Conti-Ramsden 1996). A 130 pre-processing step was applied to remove special symbols used for annotating linguistic errors, in-131 cluding those indicating semantic, phonological, and dysfluency-related issues. Additionally, when 132 neologisms or vocalisations appeared, transcribers provided phonetic representations, which were 133 later converted into Latin-alphabet phonemes while preserving their original sequence (Perez et al., 134 2020). 135

136 The transcriptions also included annotations for false starts, filler words (e.g., er, erm), and other interjections commonly found in dysfluent speech. Since filler words exhibit spelling variations 137 between American and British English, which could affect Word Error Rate (WER) in automatic 138 processing, these were normalised to match the conventions used in the Whisper training dataset. To 139 align transcriptions with corresponding audio files, utterance segmentation was manually marked 140 by expert transcribers following our established previous methodologies. This ensured that each 141 segment corresponded to a complete sentence rather than arbitrary fixed-length chunks, reducing the 142 risk of overfitting during model training. Instances of assessor speech present in the recordings were 143 also transcribed but subsequently excluded from the training dataset to ensure the model learned 144 only patient speech. 145

Since Whisper by default does not process audio recordings exceeding 30 seconds, longer files were 146 segmented into smaller units while maintaining alignment with their transcriptions. Additionally, 147 files shorter than 3 seconds posed computational challenges in Fourier transform calculations for 148 spectrogram generation (Torre & Romero, 2021). Instead of discarding these short recordings, they 149 were merged with adjacent utterances from the same speaker to preserve valuable data and optimise 150 training efficiency. All recordings were converted to WAV format, resampled to 16 kHz, encoded 151 with 16-bit resolution, and downmixed to mono. After pre-processing, the final curated dataset 152 contained approximately 13 hours of speech data. The data processing pipeline utilised several specialised tools, including SpeechBrain (Ravanelli et al., 2024), Pydub (Robert et al., 2018), FFmpeg 153 (Tomar, 2006), and SoX (Barras, 2012). 154

- 155
- 156 2.3 FINE-TUNING

Our fine-tuning approach builds on Whisper's encoder-decoder transformer architecture. The Ope nAI model processes acoustic input through an encoder that converts 80-channel log-Mel spectro grams into rich representations via two convolutional layers and sinusoidal positional encoding.
 These initial features are then refined through transformer blocks that capture long-range dependencies, while the decoder architecture mirrors this structure using learned positional embeddings

162 (Radford et al., 2022). To adapt this architecture for pathological speech recognition, we imple-163 mented a systematic fine-tuning protocol. The dataset was divided 70% (551 minutes) for training, 164 18% (141 minutes) for validation, and 12% (94 minutes) for testing. To prevent data leakage and 165 ensure unbiased evaluation, the test set comprised recordings exclusively from speakers not repre-166 sented in the training or validation sets. Whisper's medium-sized has been selected and all trainable parameters were tuned, allowing both acoustic and linguistic layers to adapt uniformly to patholog-167 ical speech patterns. This full model adaptation approach, with no frozen layers, maximised the 168 model's capacity to learn disorder-specific features. Training was conducted with a batch size of 16 per device, employing gradient accumulation for memory efficiency. The optimisation process 170 utilised the AdamW algorithm with cross-entropy loss minimisation, following a cosine learning 171 rate schedule initialised at  $1 \times 10^{-5}$  and incorporating a 1000-step warm-up phase. The model 172 was assessed at intervals of 1000 steps, with the checkpoint corresponding to the lowest WER on 173 the validation set retained as the best-performing version. Training continued up to a maximum of 174 6000 steps. To improve computational efficiency, mixed-precision arithmetic (fp16) and gradient 175 checkpointing were implemented, reducing memory overhead. The fine-tuning process required 176 approximately 8 hours and was executed using PyTorch (Paszke & Gross, 2019) alongside the Hug-177 ging Face Transformers library (Wolf & Debut, 2019), utilizing a single NVIDIA RTX 6000 GPU. The final evaluation was conducted using the WER, quantifying the discrepancy between Whisper's 178 transcriptions and human-annotated ground truth. To establish a performance benchmark, we first 179 measured WER using the pre-trained Whisper model before comparing it against the fine-tuned 180 variant<sup>2</sup>. Fine-tuning on the dataset significantly improved performance, reducing WER from the 181 baseline 39.60% to 21.51% on the validation set and from 43.62% to 21.93% on the test set. This 182 represents a substantial relative reduction, demonstrating the effectiveness of domain-specific adap-183 tation. Such fine-tuned model is then used to extract embeddings and linguistic features.

184 185

# 2.4 FEATURE EXTRACTION

187 To ensure unbiased evaluation, Whisper-derived embeddings and linguistic metrics were extracted 188 only from the unseen test set used during fine-tuning (94 minutes). This prevents as well data leak-189 age, as using the same speech samples from training could artificially inflate performance. Since em-190 beddings encode speech characteristics learned during fine-tuning, reusing them from the training set 191 would exploit prior exposure. Likewise, linguistic features from Whisper's ASR output might reflect learned transcription patterns rather than actual speech impairments. To maintain consistency, we 192 applied the same test-set-only extraction to glottal and acoustic features across all modalities. The 193 test set was complemented with age-matched healthy speakers ( $\mu = 61.51$  years,  $\sigma = 10.55$  years) 194 with no history of neurological impairments, which counts 90 minutes of recordings. 195

195 196

**Embeddings** To extract embeddings, we leverage the fine-tuned Whisper encoder, utilizing the fi-197 nal hidden state of its last transformer layer as a compact and informative representation of each 198 input utterance. The extracted embeddings consist of a 1024-dimensional latent feature vector, 199 encapsulating both acoustic and linguistic attributes of speech. In the earlier layers, the encoder 200 captures low-level acoustic features (e.g., pitch, formants, energy), while deeper layers encode in-201 creasingly abstract linguistic structures, including phonetic, lexical, and semantic information. The 202 HuggingFace implementation of Whisper was employed for this process (Wolf & Debut, 2019), ensuring reproducibility and consistency across feature extraction. These embeddings serve as a 203 high-dimensional latent representation, providing a data-driven alternative to already established 204 features. By preserving key speech characteristics in a structured feature space integrating acous-205 tic and linguistic cues within a single representation, these embeddings constitute a core modality 206 within our multimodal classification framework. 207

208

Linguistic Features Linguistic features were extracted from Whisper's ASR output and validated against manual transcriptions to ensure accuracy and reliability, as established in a previous work. Using the NLTK library (Bird et al., 2009) for tokenization, part-of-speech tagging, and stopword identification, the extracted 10 features capture lexical diversity, grammatical complexity, disfluencies, filler word usage, part-of-speech patterns and content richness. Lexical diversity was assessed using metrics such as total word count, unique words, type-token ratio (TTR), lexical density, and noun-to-verb ratio, capturing vocabulary richness and accessibility. Grammatical complexity was

<sup>215</sup> 

<sup>&</sup>lt;sup>2</sup>Papers with extensive comparison with benchmarks will later be referenced after acceptance.

216 measured through the ratio of function words to total words, indicating syntactic sophistication and 217 cohesion, while part-of-speech (POS) transitions helped identify syntactic patterns and irregularities. 218 Disfluencies were quantified by analysing filler words and their frequency, reflecting speech fluency. 219 Content richness, evaluated through the diversity of content words (e.g., nouns, verbs, adjectives), 220 provided insights into meaningful communication.

221 222 223

**Glottal Features** Glottal feature extraction provides a time-domain signal independent of vocal tract resonances, allowing precise assessment of laryngeal function. The process involved two steps: (1) isolating voiced segments using the PEFAC algorithm in MATLAB (Gonzalez & Brookes, 2014), 224 ensuring that only phonation-related portions were analysed, and (2) detecting glottal closure (GCI) 225 and opening instants (GOI) using the YAGA algorithm (Thomas et al., 2011). Once GCI and GOI 226 were identified, a set of 80 well established (Kadiri & Alku, 2019; Narendra & Alku, 2018; 2019; 227 2020; Corcoran et al., 2019) glottal parameters were extracted, categorised into time-domain and 228 frequency-domain features, through their summary statistics (mean, standard deviation, minimum, 229 maximum, kurtosis, median, skewness, and range). 230

The Time-Domain features quantify vocal fold motion over time. The Opening Quotient (OQ) and 231 Closing Quotient (CQ) measure the proportion of the glottal cycle in which the glottis remains open 232 or closed, respectively. Speed Quotient (SQ) captures the asymmetry of glottal pulses with devi-233 ations indicating irregular vocal fold vibrations, while Amplitude Quotient (AQ) and Normalised 234 Amplitude Quotient (NAQ; normalised by the fundamental frequency) characterise glottal closure 235 intensity, aiding in the detection of breathy or pressed phonation. Frequency-Domain features as-236 sess spectral properties of glottal vibrations. Harmonic Richness Factor (HRF) evaluates harmonic 237 energy relative to the fundamental frequency, reflecting the richness of voiced sounds. H1H2 and 238 H2H4 quantify harmonic amplitude differences, offering insight into vocal fold tension and phona-239 tion type. Parabolic Spectrum Parameter (PSP) measures how closely the power spectrum around each GCI resembles a parabolic shape, providing insights into spectral energy distribution. Peak 240 Slope (PS) assesses the sharpness of glottal closure, with steeper slopes indicating more abrupt 241 vocal fold contact, often linked to voice disorders. 242

- 243 Acoustic Features Acoustic parameters were extracted using the openSMILE Python library 244 (v. 3.0.1), employing the extended Geneva Minimalistic Acoustic Parameter Set (eGeMAPSv02; 245 Eyben et al. 2015). This standardised feature set comprises 88 parameters, selected for their rele-246 vance in clinical and paralinguistic speech analysis already widely clinically proven for assessing 247 pathological speech (Shahin et al., 2019; Barche et al., 2020; Liu et al., 2022; Mawalim et al., 2023; 248 Kumar et al., 2024). The acoustic feature set encompasses prosodic measures (e.g.,  $F_0$  stability, jit-249 ter, shimmer) to assess pitch and phonatory control, spectral features (e.g., MFCCs, flux) for vocal 250 resonance and articulatory precision, and energy-related metrics (e.g., loudness, harmonics-to-noise 251 ratio, pauses) to evaluate fluency and rhythmic disturbances.
- 252 253

254

# 2.5 MULTIMODAL ASSESSMENT AND FEATURE ATTRIBUTION

To systematically classify speech patterns associated with neurological impairments, we implement 255 a multimodal classification neural network integrating acoustic, linguistic, and glottal biomarkers 256 alongside learned fine-tuned Whisper embeddings. The proposed neural network is a feedforward 257 model comprising three layers: (1) a fully connected layer with 64 units, Layer Normalization, 258 LeakyReLU activation, and Dropout (0.4); (2) a second fully connected layer with 32 units, Batch 259 Normalization, LeakyReLU activation, and Dropout (0.4); and (3) a final output layer with a single 260 neuron and Sigmoid activation for binary classification between patient and healthy speech samples. 261 For the loss, we incorporated a weighted variant of the Binary Cross-Entropy (BCE) loss, integrating 262 Focal Loss and L1 regularization as follows:

263 264 265

$$\mathcal{L} = \alpha_t (1 - p_t)^{\gamma} \mathcal{L}_{\text{BCE}} + \lambda \|\theta\|_1 \tag{1}$$

266 where  $p_t$  represents the predicted probability for the true class,  $\alpha_t$  is a dynamically computed 267 weighting factor that compensates for class imbalance, and  $\gamma$  controls the focusing mechanism, reducing the loss contribution of easy-to-classify samples and emphasizing harder ones. The term 268  $(1-p_t)^{\gamma}$  adjusts the impact of each sample based on prediction confidence, effectively prioritiz-269 ing difficult cases. To further regularize the model, we introduce an L1 norm penalty  $\|\theta\|_1$ , which 270 promotes sparsity in the learned parameters, preventing overfitting and improving generalization. 271 Data is trained using 5-fold stratified group cross-validation on the training data before defined, 272 ensuring that speaker-level dependencies are maintained, followed by feature standardisation via 273 StandardScaler. Training employs also here the AdamW optimiser ( $lr = 1 \times 10^{-5}$ , weight de-274 cay = 0.01), gradient clipping ( $\|\nabla \theta\| \le 1.0$ ), and a ReduceLROnPlateau scheduler that adjusts the learning rate dynamically based on validation loss. Early stopping (patience = 12) is used to 275 prevent overfitting. Performance is measured via accuracy, precision, recall and F1-score. Next, 276 SHapley Additive exPlanations (SHAP) has been employed to assess the contribution of different 277 speech-derived features to patient classification. KernelSHAP was used to estimate feature attribu-278 tions by measuring the impact of perturbing input variables on the model's predictions. SHAP values 279 were computed across four feature modalities-acoustic, linguistic, glottal, and embeddings-to de-280 termine their relative importance on the best fold of the full multimodal model. We then calculated 281 the mean absolute SHAP values for each modality, providing a quantitative measure of their influ-282 ence on classification outcomes.

283 To assess the relationship between speech-derived biomarkers and the severity of impairment, we 284 performed a regression analysis on the patient cohort, using the Comprehensive Aphasia Test (CAT) 285 score as a measure of speech dysfunction tested by clinicians in hospitals. We trained a severity 286 regression model, namely a feedforward neural network consisting of a fully connected layer with 287 32 units and LeakyReLU activation, followed by an output layer for continuous score prediction. 288 The model was optimised using Huber loss with  $\delta = 2.0$ , which provides robustness to outliers 289 while preserving sensitivity to small deviations. Training employed the AdamW optimiser with a 290 learning rate of 0.02 and weight decay of 0.01, along with a ReduceLROnPlateau scheduler that adjusted the learning rate dynamically based on validation performance, with a patience of 11 291 epochs and a minimum learning rate of  $5 \times 10^{-5}$ . Performance was assessed using range normalised 292 root mean squared error (N-RMSE) to account for variations in score distribution, as well as absolute 293 error distribution. The model was trained via 5-fold cross-validation and the average over folds is 294 reported. By creating a severity estimation through the different feature modalities, this analysis 295 complements the classification framework in order to provide a more fine-grained assessment of 296 speech deficits. 297

298 299

300

# 3 Results

301 Tab. 1 show the results of the classification of individual and combined modalities. The multimodal 302 model that combined all features (including embeddings) achieved the highest performance, with 303 an accuracy of 92.4% and an F1-score of 0.924. Removing embeddings led to a drop in accuracy 304 (88.6%) and F1-score (0.885), highlighting the importance of learned representations in improving 305 classification accuracy. Among single modalities, acoustic features proved the most discriminative, achieving an accuracy of 90.4% and an F1-score of 0.903, followed by glottal modality with 84.6% 306 of accuracy and a F1 score of 0.846. The embeddings-based classifier showed comparable results 307 to glottal features (accuracy: 80.5%, F1: 0.803), indicating that even alone learned speech repre-308 sentations are able to decently capture key impairments. On the other hand, the linguistic modality 309 exhibited the most limited performance (accuracy: 66.3%, F1: 0.666), with precision (0.662) and 310 recall (0.657) suggesting a low power in distinguishing between post-stroke and healthy speech. 311 Post-classification, the bar chart (Fig. 1) illustrates the relative importance of each feature modality 312 in the Multimodal with embeddings model, as measured by the mean absolute SHAP values. The 313 plot reveals that embeddings have the highest mean absolute SHAP values, indicating they play a 314 crucial role in model predictions. This aligns with their contribution to the performance boost ob-315 served in the classification setup, highlighting their ability to encode complex speech patterns and pathologies effectively. Following embeddings, acoustic features exhibit the second-highest SHAP 316 values, underscoring their critical role in capturing temporal and spectral characteristics that are 317 highly indicative of neurological impairments. 318

Regression analysis of CAT scores demonstrated predictive errors across modalities ranging from
 12.99% to 20.23% relative to the CAT points range, showing high overall performance in severity
 prediction (Tab. 1). The box plot in Fig. 1 illustrates the distribution of absolute prediction er rors across modalities, with individual models showing more variability in errors. The multimodal
 approach without embeddings achieved the best performance (Tab. 1) with the lowest N-RMSE (0.1299) and a compact error distribution, suggesting that established clinical features remain cru-

326 Comprehensive Aphasia Test (CAT) total scores (Regression). 327 328 Classification Regression Modality Precision N-RMSE Accuracy F1-score Recall 330 Embeddings 0.805 0.803 0.804 0.802 0.2023 331 0.846 0.846 0.844 0.1883 Glottal 0.846 332 0.662 Linguistic 0.663 0.666 0.657 0.1608 333 Acoustic 0 904 0.903 0.902 0.902 0.1540 334 0.886 0.885 0.882 0.883 0.1299 Multimodal w/o Emb 335 0.925 0.923 Multimodal w/ Emb. 0.925 0.923 0.1801 336 337 338 Acoustic Error 339 Linguistic Prediction 340 341 Glottal 342 Absolute Embeddings 343

Table 1: Results for different feature modalities and multimodal combinations with (w/ Emb.) and 325 without (w/o Emb.) embeddings, for classifying controls and patients (Classification) and predicting



Embeddings

Glotta

Linguistic

Acoustic

Multimodal Multimodal w/ Emb. w/o Emb.

0.0010

cial for precise severity assessment. Among individual modalities, acoustic features exhibited the narrowest error range (N-RMSE: 0.1540), reflecting their reliability, followed by linguistic features (N-RMSE: 0.1608). Interestingly, while embeddings significantly enhanced classification performance, their integration in the regression task led to slightly reduced accuracy (N-RMSE: 0.1801 vs 0.1299 without embeddings - t = 4.76, p = 0.008). This underscores the importance of tailoring feature selection to the specific clinical objective, whether optimizing for diagnostic accuracy or enhancing severity estimation.

- 4
- 361 362

344

345

352 353

354

355

356

357

358

359 360

324

DISCUSSION

0.00000

0.00002

0.00004

0.00006

Mean Absolute SHAP Value

0.00008

Our analysis demonstrates the potential of foundation model-derived representations in conjunc-364 tion with clinical metrics in post-stroke speech assessment. By combining embeddings that capture 365 high-level representations, acoustic features that represent fine-grained spectral details, and glot-366 tal features that contribute unique motor-related insights, we achieved high classification accuracy. 367 The importance of embeddings in classification are in keeping with recent work (Syed et al., 2020; 368 Venugopalan et al., 2021; Bartelds et al., 2022; Neumann et al., 2024), who reported that learned 369 representations can capture subtle speech biomarkers overlooked by traditional features. Interestingly, while embeddings enhanced classification performance, their reduced accuracy in severity 370 prediction indicates that foundation models might excel at detecting pathological patterns but re-371 quire complementary clinical features for precise severity assessment. This suggests that, though 372 effective for capturing discriminative patterns in pathology detection, embeddings may struggle to 373 preserve the fine-grained relationships needed for severity scoring. 374

375 Compared to prior advancements in automated speech assessment, our work demonstrates notable improvements in performance and methodology in classification. While Venugopalan et al. (2021) 376 achieved 82% accuracy using ASR embeddings for disordered speech classification, we advance 377 this approach by integrating clinical features for better pathology detection. In the context of stroke

378 detection, Ou et al. (2025) validated the superiority of multimodal approaches, reaching 82.6% 379 accuracy and providing early evidence that integrated features outperform single modalities. Sim-380 ilarly, Soltau et al. (2023) achieved 83% accuracy using a Perceiver-based sequence classifier for 381 neurological speech abnormalities. Particularly relevant to our approach, Zusag et al. (2023) lever-382 aged Whisper with the AphasiaBank database to differentiate various types of aphasia from healthy speech, reaching an F1 score of 90.6%. Our framework's superior performance (92.4%) builds on 383 their findings while demonstrating the additional value of integrating clinical features with founda-384 tion model representations. 385

386 To the best of our knowledge, the prediction of CAT scores from speech, as well as the use of 387 our metrics, has not been widely explored yet, precluding direct performance comparisons with 388 existing approaches. Nevertheless, such a multimodal analysis with clinically established features shows high predictive ability, achieving a 12% error rate across the full CAT score range (0-216 389 points). This result is particularly noteworthy given that the CAT assessment encompasses multiple 390 linguistic domains including comprehension, expression, reading, naming, repetition and writing. 391 Indeed, by predicting such a score we were able to derive predictions for this comprehensive score 392 using solely audio recordings from the picture description task, suggesting that this focused speech 393 sample contains rich information about overall language abilities. These results suggest potential 394 clinical utility in supporting patient monitoring decisions, though further validation studies would 395 be valuable. 396

390 397 398

399

# 5 FUTURE WORK AND LIMITATIONS

Since our dataset includes additional neurobiological data such as blood samples and MRI scans, fu-400 ture work will focus on expanding the multimodal framework to incorporate these modalities. This 401 integration will enable a deeper investigation into how speech-derived features correlate with neu-402 roinflammatory markers, white matter integrity, and structural brain atrophy, offering a more com-403 prehensive understanding of neurological impairment. Additionally, we plan to assess the model's 404 predictive capabilities also for cognitive decline using the Montreal Cognitive Assessment (MoCA; 405 Nasreddine et al. 2005), a metric collected alongside CAT scores, to evaluate its ability to capture 406 broader neurological deficits beyond speech impairment. Given that stroke-induced speech impair-407 ments often overlap with symptoms of other neurological conditions (e.g., Parkinson's disease, mul-408 tiple sclerosis), our framework could potentially be adapted to additional disorder profiles. Indeed, 409 a key limitation of our study is that the model was trained and tested on a single dataset, making its generalisability to different clinical environments uncertain. Expanding validation to publicly avail-410 able datasets such as AphasiaBank (MacWhinney et al., 2011) or DementiaBank (Lanzi et al., 2023) 411 would enable cross-pathology evaluation and strengthen the model's robustness. Apart from data, 412 technical improvements could also enhance model performance; for example, speech enhancement 413 techniques would have improved a possible resilience to background noise and recording distortions. 414 Hyperparameter optimisation, alongside alternative architectures and dimensionality reduction ap-415 proaches could have improved predictive accuracy and efficiency. Long-term clinical validation 416 remains essential to ensure model outputs provide meaningful insights for assessment and monitor-417 ing. Our ongoing collaboration with an interdisciplinary team-including speech-language therapist 418 and neurologists-remains central to this process, ensuring technical advancements maintain clini-419 cal relevance.

420 421

422

# 6 CONCLUSION

423 Our work demonstrates that accurately assessing post-stroke speech impairments requires represen-424 tations that go beyond surface-level features, capturing the neurophysiological mechanisms underly-425 ing disordered speech while bridging computational modeling with clinically interpretable insights. 426 The strong performance of this tailored approach underscores its potential to transform clinical as-427 sessment, shifting from traditional, labor-intensive methods toward an intelligent fusion of mul-428 timodal data streams and foundation models. Just as BERT revolutionized clinical text analysis, 429 Whisper-based models may mark a significant advancement in speech-based diagnostics, particularly for early screening and severity assessment. By integrating computational models with clinical 430 expertise in a targeted manner, this research advances the adoption of automated speech analysis as 431 a viable tool for diagnosing and managing language disorders post-stroke

# 432 MEANINGFULNESS STATEMENT

434 This work explores the meaningful representation of life through modeling fundamental biological signals that reflect human cognition, communication, and motor control. Speech serves as a prime 435 example of such signals, emerging from the interplay of neurological, respiratory, and articulatory 436 systems. The heterogeneity of post-stroke lesions leads to diverse speech impairments, affecting lan-437 guage processing, articulation, and prosody in unique ways. Our multimodal framework captures 438 this complexity by combining Whisper's high-level speech representations with clinically grounded 439 features, enabling the detection of subtle biomarkers that could transform early diagnosis of neuro-440 logical conditions. 441

442 443

444

445

446

447

453

473

474

475

480

481

482

### References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Emily Alsentzer, John R Murphy, Willie Boag, Wei-Hung Weng, Di Jin, Tristan Naumann, and Matthew McDermott. Publicly available clinical bert embeddings. *arXiv preprint arXiv:1904.03323*, 2019.
- Peter Appelros, Birgitta Stegmayr, and Andreas Terént. Sex differences in stroke epidemiology: a systematic review. *Stroke*, 40(4):1082–1090, 2009.
- Purva Barche, Krishna Gurugubelli, and Anil Kumar Vuppala. Towards automatic assessment of voice disorders: A clinical approach. In *INTERSPEECH*, pp. 2537–2541, 2020.
- 456 Benjamin Barras. Sox : Sound exchange. 01 2012.
- Martijn Bartelds, Wietse de Vries, Faraz Sanal, Caitlin Richter, Mark Liberman, and Martijn Wiel ing. Neural representations for modeling variation in speech. *Journal of Phonetics*, 92:101137, 2022.
- Paul Best, Santiago Cuervo, and Ricard Marxer. Transfer learning from Whisper for microscopic intelligibility prediction. In *Interspeech 2024*, pp. 3839–3843, 2024. doi: 10.21437/Interspeech. 2024-2258.
- 465 Steven Bird, Ewan Klein, and Edward Loper. *Natural language processing with Python: analyzing* 466 *text with the natural language toolkit.* "O'Reilly Media, Inc.", 2009.
- Gina Conti-Ramsden. CLAN (Computerized Language Analysis). *Child Language Teaching and Therapy*, 12(3):345–349, 1996.
- Patrick Corcoran, Arnold Hensman, and Barry Kirkpatrick. Glottal flow analysis in parkinsonian
  speech. In *Proc. of the 12th Int. Joint Conf. on Biomedical Eng. Syst. and Technol. (BIOSTEC)*,
  pp. 116–123, 2019.
  - Jacob Devlin. Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805, 2018.
- Florian Eyben, Klaus R Scherer, Björn W Schuller, Johan Sundberg, Elisabeth André, Carlos Busso,
  Laurence Y Devillers, Julien Epps, Petri Laukka, Shrikanth S Narayanan, et al. The geneva
  minimalistic acoustic parameter set (gemaps) for voice research and affective computing. *IEEE transactions on affective computing*, 7(2):190–202, 2015.
  - Sira Gonzalez and Mike Brookes. Pefac-a pitch estimation algorithm robust to high levels of noise. *IEEE/ACM Trans. on Audio, Speech, and Lang. Process.*, 22(2):518–530, 2014.
- Yicong Jiang, Tianzi Wang, Xurong Xie, Juan Liu, Wei Sun, Nan Yan, Hui Chen, Lan Wang, Xunying Liu, and Feng Tian. Perceiver-prompt: Flexible speaker adaptation in Whisper for chinese disordered speech recognition. In *Interspeech 2024*, pp. 2025–2029, 2024. doi: 10.21437/Interspeech.2024-852.

486 Sudarsana Reddy Kadiri and Paavo Alku. Analysis and detection of pathological voice using glottal 487 source features. IEEE J. of Selected Topics in Signal Process., 14(2):367–379, 2019. 488 Deepak Kumar, Udit Satija, and Preetam Kumar. Pathological speech and electroglottography sig-489 nals analysis using invariance scattering network. *Circuits, Systems, and Signal Processing*, pp. 490 1-18, 2024.491 492 Alyssa M Lanzi, Anna K Saylor, Davida Fromm, Houjun Liu, Brian MacWhinney, and Matthew L 493 Cohen. Dementiabank: Theoretical rationale, protocol, and illustrative analyses. American Journal of Speech-Language Pathology, 32(2):426–438, 2023. 494 495 Duc Le, Keli Licata, and Emily Mower Provost. Automatic quantitative analysis of spontaneous 496 aphasic speech. Speech Communication, 100:1-12, 2018. 497 Jeehyun Lee, Yerin Choi, Tae-Jin Song, and Myoung-Wan Koo. Inappropriate pause detection 498 in dysarthric speech using large-scale speech recognition. In ICASSP 2024-2024 IEEE Inter-499 national Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 12486–12490. 500 IEEE, 2024. 501 502 Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. Biobert: a pre-trained biomedical language representation model for biomedical text 504 mining. Bioinformatics, 36(4):1234-1240, 2020. 505 Wing-Zin Leung, Mattias Cross, Anton Ragni, and Stefan Goetze. Training data augmentation for 506 dysarthric automatic speech recognition by text-to-dysarthric-speech synthesis. arXiv preprint 507 arXiv:2406.08568, 2024. 508 Jinpeng Li and Wei-Qiang Zhang. Whisper-based transfer learning for alzheimer disease classifica-509 tion: Leveraging speech segments with full transcripts as prompts. In ICASSP 2024-2024 IEEE 510 International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 11211– 511 11215. IEEE, 2024. 512 513 Yuanyuan Liu, Mittapalle Kiran Reddy, Nelly Penttilä, Tiina Ihalainen, Paavo Alku, and Okko 514 Räsänen. Automatic assessment of parkinson's disease using speech representations of phona-515 tion and articulation. IEEE/ACM Transactions on Audio, Speech, and Language Processing, 31: 242-255, 2022. 516 517 Brian MacWhinney. The CHILDES project: Tools for analyzing talk, Volume II: The database. 518 Psychology Press, 2014. 519 Brian MacWhinney, Davida Fromm, Margaret Forbes, and Audrey Holland. Aphasiabank: Methods 520 for studying discourse. Aphasiology, 25(11):1286–1307, 2011. 521 522 Seedahmed S Mahmoud, Raphael F Pallaud, Akshay Kumar, Serri Faisal, Yin Wang, and Qiang 523 Fang. A comparative investigation of automatic speech recognition platforms for aphasia assessment batteries. Sensors, 23(2):857, 2023. 524 525 Candy Olivia Mawalim, Benita Angela Titalim, Shogo Okada, and Masashi Unoki. Non-intrusive 526 speech intelligibility prediction using an auditory periphery model with hearing loss. Applied 527 Acoustics, 214:109663, 2023. 528 NP Narendra and Paavo Alku. Dysarthric speech classification using glottal features computed from 529 non-words, words and sentences. In Interspeech, pp. 3403–3407. Int. Speech Commun. Assoc. 530 (ISCA), 2018. 531 532 NP Narendra and Paavo Alku. Dysarthric speech classification from coded telephone speech using 533 glottal features. Speech Commun., 110:47-55, 2019. 534 NP Narendra and Paavo Alku. Glottal source information for pathological voice detection. IEEE 535 Access, 8:67745-67755, 2020. 536 Ziad S Nasreddine, Natalie A Phillips, Valérie Bédirian, Simon Charbonneau, Victor Whitehead, Isabelle Collin, Jeffrey L Cummings, and Howard Chertkow. The montreal cognitive assessment, 538 moca: a brief screening tool for mild cognitive impairment. Journal of the American Geriatrics Society, 53(4):695-699, 2005.

540	Michael Neumann, Hardik Kothare, and Vikram Ram	anarayanan. Multimodal speech biomarkers
541	for remote monitoring of als disease progression.	Computers in Biology and Medicine, 180:
542	108949, 2024.	
543		

- 544 Emily R Olafson et al. Data-driven biomarkers better associate with stroke motor outcomes than 545 theory-based biomarkers. *Brain Commun.*, 2024.
- Zijun Ou, Haitao Wang, Bin Zhang, Haobang Liang, Bei Hu, Longlong Ren, Yanjuan Liu, Yuhu Zhang, Chengbo Dai, Hejun Wu, et al. Early identification of stroke through deep learning with multi-modal human speech and movement data. *Neural Regeneration Research*, 20(1):234–241, 2025.
- Rebecca Palmer and Pam et al. Enderby. Computer therapy compared with usual care for people with long-standing aphasia poststroke: a pilot randomized controlled trial. *Stroke*, 43(7):1904–1911, 2012.
- Adam Paszke and Sam et al. Gross. Pytorch: An imperative style, high-performance deep learn ing library. In *Advances in Neural Information Processing Systems 32*, pp. 8024–8035. Curran Associates, Inc., 2019.
- Matthew Perez, Zakaria Aldeneh, and Emily Mower Provost. Aphasic speech recognition using a mixture of speech intelligibility experts. *arXiv preprint arXiv:2008.10788*, 2020.
- Jianing Qiu, Wu Yuan, and Kyle Lam. The application of multimodal large language models in
   medicine. *The Lancet Regional Health–Western Pacific*, 45, 2024.
- Alec Radford, Jong Wook Kim, and Tao et al. Xu. Robust speech recognition via large-scale weak supervision. *arXiv preprint arXiv:2212.04356*, 2022.
- Vikram Ramanarayanan, Adam C Lammert, Hannah P Rowe, Thomas F Quatieri, and Jordan R
   Green. Speech as a biomarker: opportunities, interpretability, and challenges. *Perspectives of the ASHA Special Interest Groups*, 7(1):276–283, 2022.
- Siddharth Rathod, Monil Charola, Akshat Vora, Yash Jogi, and Hemant A. Patil. Whisper features for dysarthric severity-level classification. In *Interspeech 2023*, pp. 1523–1527, 2023. doi: 10. 21437/Interspeech.2023-1891.
- Mirco Ravanelli, Titouan Parcollet, Adel Moumen, Sylvain de Langen, Cem Subakan, Peter 572 Plantinga, Yingzhi Wang, Pooneh Mousavi, Luca Della Libera, Artem Ploujnikov, Francesco 573 Paissan, Davide Borra, Salah Zaiem, Zeyu Zhao, Shucong Zhang, Georgios Karakasidis, Sung-574 Lin Yeh, Pierre Champion, Aku Rouhe, Rudolf Braun, Florian Mai, Juan Zuluaga-Gomez, 575 Seyed Mahed Mousavi, Andreas Nautsch, Xuechen Liu, Sangeet Sagar, Jarod Duret, Salima 576 Mdhaffar, Gaelle Laperriere, Mickael Rouvier, Renato De Mori, and Yannick Esteve. Open-577 source conversational ai with SpeechBrain 1.0, 2024. URL https://arxiv.org/abs/ 578 2407.00463. 579
- James Robert, Marc Webbie, et al. Pydub, 2018. URL http://pydub.com/.

581

582

583

584

588

589

590

- Carole Roth. *Boston Diagnostic Aphasia Examination*, pp. 428–430. Springer New York, New York, NY, 2011. ISBN 978-0-387-79948-3. doi: 10.1007/978-0-387-79948-3\_868. URL https://doi.org/10.1007/978-0-387-79948-3\_868.
- Giulia Sanguedolce, Patrick A Naylor, and Fatemeh Geranmayeh. Uncovering the potential for a
  weakly supervised end-to-end model in recognising speech from patient with post-stroke aphasia.
  In *Proceedings of the 5th Clinical Natural Language Processing Workshop*, pp. 182–190, 2023.
  - Giulia Sanguedolce, Sophie Brook, Dragos C Gruia, Patrick A Naylor, and Fatemeh Geranmayeh. When whisper listens to aphasia: Advancing robust post-stroke speech recognition. In *Interspeech*, 2024.
- Mostafa Shahin, Beena Ahmed, Daniel V Smith, Andreas Duenser, and Julien Epps. Automatic
   screening of children with speech sound disorders using paralinguistic features. In 2019 ieee 29th
   *international workshop on machine learning for signal processing (mlsp)*, pp. 1–5. IEEE, 2019.

- Hagen Soltau, Izhak Shafran, Alex Ottenwess, R Joseph Jr, Rene L Utianski, Leland R Barnard, John L Stricker, Daniela Wiepert, David T Jones, and Hugo Botha. Detecting speech abnormalities with a perceiver-based sequence classifier that leverages a universal speech model. In 2023 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU), pp. 1–7. IEEE, 2023.
- James D Stefaniak, Fatemeh Geranmayeh, and Matthew A Lambon Ralph. The multidimensional nature of aphasia recovery post-stroke. *Brain*, 145(4):1354–1367, 2022.
- Kate Swinburn, Gillian Porter, and David Howard. Comprehensive aphasia test. APA PsycTests, 2004.
  - Zafi Sherhan Syed, Sajjad Ali Memon, and Abdul Latif Memon. Deep acoustic embeddings for identifying parkinsonian speech. *International Journal of Advanced Computer Science and Applications*, 11(10):726–734, 2020.
- Mark RP Thomas, Jon Gudnason, and Patrick A Naylor. Estimation of glottal closing and opening
   instants in voiced speech using the yaga algorithm. *IEEE Trans. on Audio, Speech, and Lang. Process.*, 20(1):82–91, 2011.
- Suramya Tomar. Converting video formats with ffmpeg. *Linux Journal*, 2006(146):10, 2006.
- Iván G Torre and Mónica et al. Romero. Improving aphasic speech recognition by using novel semi-supervised learning methods on Aphasiabank for English and Spanish. *Applied Sciences*, 11(19): 8872, 2021.
- Subhashini Venugopalan, Joel Shor, Manoj Plakal, Jimmy Tobin, Katrin Tomanek, Jordan R
  Green, and Michael P Brenner. Comparing supervised models and learned speech representations for classifying intelligibility of disordered speech on selected phrases. *arXiv preprint arXiv:2107.03985*, 2021.
- Thomas Wolf and Lysandre et al. Debut. Huggingface's transformers: State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771*, 2019.
- Qianqian Xie, Qingyu Chen, Aokun Chen, Cheng Peng, Yan Hu, Fongci Lin, Xueqing Peng, Jimin Huang, Jeffrey Zhang, Vipina Keloth, et al. Me Ilama: Foundation large language models for medical applications. *arXiv preprint arXiv:2402.12749*, 2024.
  - Mario Zusag, Laurin Wagner, and Theresa Bloder. Careful Whisper leveraging advances in automatic speech recognition for robust and interpretable aphasia subtype classification. In *Interspeech 2023*, pp. 3013–3017, 2023. doi: 10.21437/Interspeech.2023-1653.