

Fantastic Features and Where to Find Them: Detecting Cognitive Impairment with a Subsequence Classification Guided Approach

Benjamin Eyre

Winterlight Labs

Toronto, Canada

{benjamin, aparna, jekaterina}@winterlightlabs.com

Aparna Balagopalan

Winterlight Labs

Toronto, Canada

Jekaterina Novikova

Winterlight Labs

Toronto, Canada

Abstract

Despite the widely reported success of embedding-based machine learning methods on natural language processing tasks, the use of more easily interpreted engineered features remains common in fields such as cognitive impairment (CI) detection. Manually engineering features from noisy text is time and resource consuming, and can potentially result in features that do not enhance model performance. To combat this, we describe a new approach to feature engineering that leverages sequential machine learning models and domain knowledge to predict which features help enhance performance. We provide a concrete example of this method on a standard data set of CI speech and demonstrate that CI classification accuracy improves by 2.3% over a strong baseline when using features produced by this method. This demonstration provides an example of how this method can be used to assist classification in fields where interpretability is important, such as health care.

1 Introduction

In recent years, word and sentence embedding-based methods have had a significant impact on the field of NLP (Devlin et al., 2019; Mikolov et al., 2013; Pennington et al., 2014; Di Palo and Parde, 2019). These approaches stand as an alternative to classical feature engineering approaches, where carefully crafted features, such as word length or part of speech tag, are extracted from text and used as input. Despite the promise of embedding-based methods, there are still several advantages to feature engineering. Most notably, using embeddings as input can lead to issues with interpretability (Heimerl and Gleicher, 2018; Hooker et al., 2019; Kindermans et al., 2017), which is especially important in a healthcare domain (Balagopalan et al., 2020). Meanwhile, feature engineering approaches directly lend themselves to easily inter-

pretable models (Ribeiro et al., 2016). As such, feature engineering remains an important practice for fields such as health care, where interpretability is imperative. An extensive body of work has been produced where ML methods and engineered features have been applied to cognitive impairment (CI) detection (Balagopalan et al., 2018; Karlekar et al., 2018; Zhu et al., 2019).

In this work, we present a new feature engineering method that is guided by classifying subsets of a pause-centred speech sequence (subsequences), and inspired by literature suggesting that CI could be indicated by the words that subjects pause before (Calley et al., 2010; Mack et al., 2013; Seifart et al., 2018). This approach aims to extract pause-related information while minimizing the noise added from unrelated factors. This method generates interpretable and effective features, potentially saving time and resources spent on excess feature engineering. We validate this method by presenting a 2.3% accuracy increase over a strong baseline on CI vs healthy (HC) classification, matching the state of the art (Hernández-Domínguez et al., 2018).

In summary, our major contributions are:

- A method of classifying speech using only a token of interest and a small context around it, i.e. *subsequence classification* (Sec. 3.3).
- A novel *feature engineering approach* guided by subsequence classification (Sec. 4).
- Validating this approach by showing that it aids in achieving classification results comparable to the state of the art (Sec. 5.2).

2 Related work

Several authors report increases in CI detection performance by extracting acoustic features such as filled and unfilled pause counts, as well as average pause duration (Tóth et al., 2015, 2018; Pistono

Data Subset	HC	CI	Total
DB (transcripts)	229 (42%)	321 (58%)	550
DB-C1	317 (33%)	645 (67%)	962
DB-C2	511 (35%)	963 (65%)	1,474
DB-C3	529 (35%)	980 (65%)	1,509
DB-Utt	755 (42%)	1,059 (58%)	1,814

Table 1: Overview of the number of samples (subsequences or transcripts) in different subsets of DB.

et al., 2016). However, we believe further performance increases can be achieved if we focus not only on the pauses themselves, but also the linguistic context in which the pauses occur in speech.

Recently, Hernández-Domínguez et al. (2018) achieved an accuracy of 78% when classifying Dementiabank (Sec. 3.1) transcripts as CI or HC using an extended set of lexical features, which is to the best of our knowledge the state of the art (SOTA). We use their recorded performance as a benchmark when validating our approach.

Several authors have reported performance gains by using subsequences to aid with classification. These authors use subsequences only as a means to process full sequences (Phan et al., 2017), or they use the presence of common subsequences as a feature for longer text sequences (Iglesias et al., 2007; Kumar et al., 2005). To the best of our knowledge, no prior work describes using subsequence classification to guide feature engineering.

3 Experimental Method

In this section, we define the data sets and methodology used in our experimental framework.

3.1 Data Sets

Dementiabank (DB): Dementiabank¹ is a large public data set of pathological speech (Becker et al., 1994), containing audio files and transcripts of participants describing the ‘Cookie Theft’ image. Transcripts are created manually by trained transcriptionists following the CHAT protocol (MacWhinney, 2014). Out of the 286 participants, 193 are diagnosed with some form of CI (CI; N = 321 transcripts) and 93 are healthy controls (HC; N = 229 transcripts). Transcripts receive a CI or HC label corresponding to whether the participant who produced the transcribed speech was cognitively impaired or not.

Subsequence-based Data Subsets: To conduct subsequence classification, we extract subsequences of varying length from each transcript.

¹<https://dementia.talkbank.org>

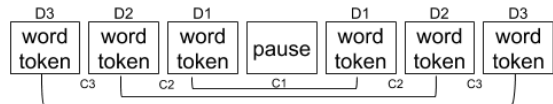


Figure 1: Visualization of the difference between contexts and distances in a pause-focused subsequence.

For each transcript in DB, each utterance containing a pause was extracted and labelled as positive if the sample contains CI speech, or negative otherwise. Subsequences were extracted from these utterances by taking the first one, two, or three speech tokens before and after each pause.² We created three data subsets by including subsequences of at most one, two, or three tokens around the pause: Context 1 (DB-C1), Context 2 (DB-C2), and Context 3 (DB-C3), respectively. We also included one data subset including full utterances that include pauses, DB-Utt (Tab.1). Identical subsequences found in both classes were removed. Furthermore, subsequences extracted from HC transcripts are labeled as HC, and subsequences extracted from CI transcripts are labeled as CI. We refer to the tokens that are next to the pause as Distance 1 (D1), the tokens that are one token away from the pause as Distance 2 (D2), and the tokens that are two tokens away from the pause as Distance 3 (D3). The differences between *context* and *distance* are shown in Fig.1. For example, for the pause sequence “The boy is *uh* stealing a cookie”, only the tokens “boy” and “a” would be considered the Distance 2 tokens for this sequence, while the tokens “boy”, “is”, “*uh*”, “stealing”, and “a” would be considered the Context 2 tokens for this sequence.

3.2 Feature Extraction

In this section, we describe how features are extracted on the transcript-level (for transcript classification) and on the token-level (for subsequence classification).

Transcript-Level Features: We extract over 500 linguistic and acoustic features from each transcript, such as part of speech counts and average word length (App.A). These features, referred to as the *Original* feature set, are used to provide a baseline to benchmark transcript-level classification performance. We also use the *Original* feature set as a base that we extend with newly engineered transcript-level features (Sec. 4). To produce an ad-

²If there were less than two or three tokens before or after a pause, the largest possible sequence of tokens was extracted.

ditional baseline, we perform feature selection on the *Original* feature set, and found $k = 85$ features led to the greatest performance.

Token-Level Features: In order to conduct subsequence classification, we extract features on the token-level for each of the subsequence-based data subsets. Of the *Original* feature set, we select a subset of features that have a clear token-level analogue (App.A). For instance, the transcript-level feature of average word length has the token-level analogue of individual word length. After feature extraction, each token is represented by a 23-dimensional input vector. Consequently, each subsequence in the DB-C1, DB-C2, DB-C3, and DB-Utt data subsets is represented as a T by 23 matrix, where T is the length of the sequence in tokens.

3.3 Classification

In this section, we describe the methodology used for subsequence and transcript classification. Subsequence classification is used to guide the engineering of new features, while transcript classification validates the new features' effectiveness.

Subsequence Classification: Our subsequence classification experiment involves performing 5-fold cross validation with each of the DB-C1, DB-C2, DB-C3, and DB-Utt data subsets. Subsequences are classified as either HC or CI. We conduct classification using GRU-based (Cho et al., 2014) models with an attention mechanism designed for document classification (Yang et al., 2016), with model parameters tuned for each of the data subsets. We report accuracy for M-C1, the model that achieved the highest accuracy on DB-C1, M-C2, the model that achieved the highest accuracy on DB-C2, M-C3, the model that achieved the highest accuracy on DB-C3, and M-Utt, the model that achieved the highest accuracy on DB-Utt (additional training details provided in App.B).

Transcript Classification: We evaluate the efficacy of our feature engineering approach (Sec.4) by performing transcript-level 10-fold cross validation with a variety of feature sets. Transcripts are classified as either HC or CI. We use the *Original* feature set, as well as the top 85 of the *Original* features, based on their ANOVA F-values, as baselines. Additionally, we extend the *Original* feature set with the k best features, based on ANOVA F-values, from each of the feature sets generated using our novel feature engineering approach (Sec. 4), separately. k is optimized for accuracy for each

extending feature set separately. To classify DB transcripts, we use 5 ML models: an SVM, a gradient boosting ensemble, a 2-layer neural network (NN), a random forest, and an ensemble of the previous four models (Ens). We report the accuracy (Acc), precision (Prec), sensitivity (Sens), and specificity (Spec) for the model that achieved the greatest cross validated accuracy for each feature set, separately (training details provided in App.B).

4 Proposed Feature Engineering Approach

Our approach to engineering new transcript-level features involves three major steps:

1) Subsequences of varying length centred around a token of interest, in our case a pause, must be extracted from each of the input transcripts and grouped into subsets based on maximum length. Each of the tokens in these subsequences must have token-level features extracted. The token-level features, as well as the central token, should be chosen based on in-domain knowledge.

2) A sequential ML model must be cross validated on each of the subsequence data subsets from the previous step in a subsequence classification experiment. Here, we are specifically attempting to exploit the ability for sequential machine learning models to uncover patterns in sequential data. The mean cross validated accuracy on each of these length-based data subsets should be used as an indicator of how much distinguishing information can be extracted from tokens within the specified range of the pause.

3) Based on the recorded cross validated accuracies from the previous step, transcript-level aggregations of the token-level features must be created at various distances from the pause. We propose two methods of aggregating token-level features (DB-specific details provided in App.A):

- **Continuous features** can be aggregated simply by taking the average of a feature across each of the tokens. An example of this would be calculating the average word length for each of the tokens found at a specified distance from a pause.

- **Categorical features** can be aggregated using counts or ratios, such as the number of nouns occurring at a specified distance from a pause.

These transcript-level aggregates should only be extracted for the distances that produced the greatest cross validated accuracy during subsequence classification, as the subsequence classification per-

Feature Set	Model	Acc	Prec	Sens	Spec
Original	Ens	74.77±0.6*	82.08±0.4*	73.1±0.5*	79.74±0.4*
Original w/ feat.sel.	Ens	75.18±1.4	83.37±1.2	72.21±1.3*	81.67±1.6
Original + F-D1	NN	74.41±1.9*	78.59±1.1*	77.15±3.3	72.63±1.9*
Original + F-D2	Ens	77.09±1.0	84.40±0.8	75.21±1.0	82.32±0.9
Original + F-D3	Ens	76.05±0.8	84.02±0.9	71.92±0.7*	83.33±1.3
Original + F-C2	NN	75.14±1.4	79.85±1.5*	76.93±0.9*	74.02±2.5*
Original + F-C3	Ens	74.82±1.2*	83.68±1.2	70.62±1.9*	82.63±1.4

Table 2: Transcript classification performance for each feature set’s best performing classification model, averaged across four random seeds. Bold indicates best performance, and * indicates significance ($p < 0.05$) when compared to the model using F-D2 features.

Model	M-C1	M-C2	M-C3	M-Utt
Accuracy	59.6±2.6	60.7±2.5	59.8±0.9	60.3±1.0

Table 3: Subsequence classification performance. Accuracy is averaged across four random seeds.

formance indicates that the features found in that range are the most distinguishing. For instance, if subsequences of up to two tokens around a pause produced the most accurate subsequence classifier, transcript-level aggregates should only be extracted for tokens at the D1 and D2 positions in reference to the pause, and not the D3 position.

To validate this method, we create five transcript-level feature sets: features aggregated from tokens at the D1 position in reference to a pause (*F-D1*), features aggregated from the D2 position (*F-D2*), features aggregated from the D3 position (*F-D3*), the combination of F-D1, F-D2, and F-D3 (*F-C3*), and the combination of F-D1 and F-D2 (*F-C2*).

5 Results

In this section, we report the results for the subsequence and transcript classification experiments.

5.1 Subsequence Classification

After averaging across four random seeds, M-C2 was able to achieve an accuracy of 60.7%, higher than M-C1, M-C3, or M-Utt (Tab.3). This leads us to the conclusion that using features from the two tokens preceding and succeeding a pause could enhance transcript-level classification performance.

5.2 Transcript Classification

We create the F-D1, F-D2, and F-C2 aggregate feature sets, as the highest subsequence classification accuracy was achieved by a model trained on DB-C2. Additionally, in order to validate our claims, we create F-D3 and F-C3. The highest accuracy of 77.09% on transcript classification is achieved by

an ensemble model that used the *Original + F-D2* feature set (Tab. 2).

Using one of the four random seeds used to produce the average performance metrics presented in Tab. 2, the model using F-D2 features was able to achieve an accuracy of 78.36%, the same as the single-seed SOTA accuracy of 78% (Hernández-Domínguez et al., 2018).

6 Discussion

As shown in Tab. 3 and Tab. 2, features from tokens within 2 tokens of a pause were the most effective in enhancing both subsequence and transcript classification. To determine how these two tasks are connected, we conduct a statistical analysis on the token-level and transcript-level features. Two sided t-tests between features extracted from tokens found at D1, D2, and D3 from different classes show similar patterns for features that are significantly different between classes for both the token and transcript-level. Larger concentrations of distinguishing features are found at D1 and D2 than at D3. This could explain the effectiveness of features from the D2 position in both tasks (Tab. 4).

However, this pattern congruity does not explain why F-D3 features on their own are more effective than F-D1 features on their own. The trend that both the F-D3 and F-C3 feature sets produced greater transcript-level accuracy than the F-D1 feature set, and lower transcript-level accuracy than the F-D2 and F-C2 feature sets, is the same as the trend for the subsequence classification results reported in Tab. 3. This indicates that subsequence classification may be able to provide better insight into potential transcript classification performance than traditional statistical testing.

It is important to consider the implications of producing a model with the F-D2 feature set that achieved significantly higher accuracy than the most accurate model produced with the F-C3 fea-

Distance	Token-Level	Transcript-Level
D1	18	12
D2	21	12
D3	7	6

Table 4: Number of features that are significantly different between classes according to two sided t-tests for each distance.

ture set. As described in Section 3.3, we perform feature selection using ANOVA F-values for each of the aggregate feature sets. Since F-D2 is a subset of F-C3, this implies that this more traditional feature selection method did not select a group of features from F-C3 that was more effective than the features from F-D2, even though it was able to select any and all of the features in F-D2. This serves as a testament to our feature engineering method, as it demonstrates that even popular feature selection methods are not able to completely remove the negative effects of engineering an excessive amount of ineffective features.

Following several other works that used the DB data set (Hernández-Domínguez et al., 2018; Pou-Prom and Rudzicz, 2018; Sarawgi et al., 2020), all of our experiments are conducted with K-fold cross validation. While the small size of the DB data set helps to justify this as a validation procedure, optimizing a cross validated performance metric (accuracy, F1, etc.) may lead results using K-fold cross validation to be an overestimate of generalization performance.

DB-C2 produced a more accurate subsequence classifier than any other data subset of DB. This suggests that the class distinguishing signal from the pause is strongest within a two token radius around the pause. Beyond that radius, the signal may be obstructed by noise from other patterns in speech. However, in different data sets, a different subsequence length may present the strongest, least noisy signal. New aggregate features should be created for tokens within whichever range produces the best subsequence classification performance.

However, our results do indicate that there is a strong link between how well features from certain token positions contribute to both subsequence and transcript classification. This may relate to the effect of noise on those token positions, which we use subsequence classification to identify.

7 Conclusion and Future Work

In this work, we present two principle contributions. First, we describe a novel method for speech classification - subsequence classification - in which speech is modelled as a token of interest, such as a pause, along with surrounding tokens of context. Secondly, we demonstrate how subsequence classification can be used to engineer features that extract distinguishing information while minimizing added noise, and consequently match SOTA performance on a standard data set of CI speech.

Future work should be done to understand why certain context lengths are more conducive for subsequence classification than others, and when that performance can transfer to effective transcript-level classification. Finally, additional work should be done to develop techniques for finding tokens of interest, such as pauses, that can be exploited using our feature engineering technique.

References

- Aparna Balagopalan, Benjamin Eyre, and Jekaterina Novikova. 2020. To BERT or Not To BERT: Comparing Speech and Language-based Approaches for Alzheimer's Disease Detection. In *Proceedings of INTERSPEECH*.
- Aparna Balagopalan, Jekaterina Novikova, Frank Rudzicz, and Marzyeh Ghassemi. 2018. [The Effect of Heterogeneous Data for Alzheimer's Disease Detection from Speech](#). In *Proceedings of the Machine Learning for Health (MLAH) Workshop at NeurIPS 2018*.
- James T Becker, François Boiler, Oscar L Lopez, Judith Saxton, and Karen L McGonigle. 1994. The natural history of Alzheimer's disease: description of study cohort and accuracy of diagnosis. *Archives of Neurology*, 51(6):585–594.
- Marc Brysbaert, Amy Beth Warriner, and Victor Kuperman. 2014. Concreteness ratings for 40 thousand generally known english word lemmas. *Behavior Research Methods*, 46(3):904–911.
- Clifford S Calley, Gail D Tillman, Kyle Womack, Patricia Moore, John Hart Jr, and Michael A Kraut. 2010. Subjective report of word-finding and memory deficits in normal aging and dementia. *Cognitive and behavioral neurology: official journal of the Society for Behavioral and Cognitive Neurology*, 23(3):185.
- Nitesh V Chawla, Kevin W Bowyer, Lawrence O Hall, and W Philip Kegelmeyer. 2002. Smote: synthetic minority over-sampling technique. *Journal of artificial intelligence research*, 16:321–357.
- Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*.
- Mark Davies. 2009. The 385+ million word corpus of contemporary american english (1990–2008+): Design, architecture, and linguistic insights. *International journal of corpus linguistics*, 14(2):159–190.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186.
- Flavio Di Palo and Natalie Parde. 2019. Enriching neural models with targeted features for dementia detection. *arXiv preprint arXiv:1906.05483*.
- Kathleen C Fraser, Jed A Meltzer, and Frank Rudzicz. 2016. Linguistic features identify alzheimer's disease in narrative speech. *Journal of Alzheimer's Disease*, 49(2):407–422.
- Florian Heimerl and Michael Gleicher. 2018. Interactive analysis of word vector embeddings. In *Computer Graphics Forum*, volume 37, pages 253–265. Wiley Online Library.
- Laura Hernández-Domínguez, Sylvie Ratté, Gerardo Sierra-Martínez, and Andrés Roche-Bergua. 2018. Computer-based evaluation of alzheimer's disease and mild cognitive impairment patients during a picture description task. *Alzheimer's & Dementia: Diagnosis, Assessment & Disease Monitoring*, 10:260–268.
- Matthew Honnibal and Ines Montani. 2017. spaCy 2: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing. To appear.
- Sara Hooker, Dumitru Erhan, Pieter-Jan Kindermans, and Been Kim. 2019. A benchmark for interpretability methods in deep neural networks. In *Advances in Neural Information Processing Systems*, pages 9737–9748.
- José Antonio Iglesias, Agapito Ledezma, and Araceli Sanchis. 2007. Sequence classification using statistical pattern recognition. In *International Symposium on Intelligent Data Analysis*, pages 207–218. Springer.
- Sweta Karlekar, Tong Niu, and Mohit Bansal. 2018. Detecting linguistic characteristics of alzheimer's dementia by interpreting neural models. In *Proceedings of NAACL-HLT*, pages 701–707.
- Pieter-Jan Kindermans, Sara Hooker, Julius Adebayo, Maximilian Alber, Kristof T Schütt, Sven Dähne, Dumitru Erhan, and Been Kim. 2017. The (un) reliability of saliency methods. *arXiv preprint arXiv:1711.00867*.
- N Pradeep Kumar, M Venkateswara Rao, P Radha Krishna, and Raju S Bapi. 2005. Using sub-sequence information with knn for classification of sequential data. In *International Conference on Distributed Computing and Internet Technology*, pages 536–546. Springer.
- Victor Kuperman, Hans Stadthagen-Gonzalez, and Marc Brysbaert. 2012. Age-of-acquisition ratings for 30,000 english words. *Behavior research methods*, 44(4):978–990.
- Jennifer E Mack, Aya Meltzer-Asscher, Sarah D. Chandler, Sandra Weintraub, Marek Marsel Mesulam, and Cynthia K Thompson. 2013. [Word-finding pauses in primary progressive aphasia \(ppa\): Effects of lexical category](#). *Procedia - Social and Behavioral Sciences*, 94:129–130.

- Brian MacWhinney. 2014. *The CHILDES project: Tools for analyzing talk, Volume I: Transcription format and programs*. Psychology Press.
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*.
- Natalia B Mota, Nivaldo AP Vasconcelos, Nathalia Lemos, Ana C Pieretti, Osame Kinouchi, Guillermo A Cecchi, Mauro Copelli, and Sidarta Ribeiro. 2012. Speech graphs provide a quantitative measure of thought disorder in psychosis. *PLoS one*, 7(4):e34928.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. *PyTorch: An imperative style, high-performance deep learning library*. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc.
- Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. 2011. Scikit-learn: Machine learning in python. *the Journal of machine Learning research*, 12:2825–2830.
- Jeffrey Pennington, Richard Socher, and Christopher Manning. 2014. *GloVe: Global Vectors for Word Representation*. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, Doha, Qatar. Association for Computational Linguistics.
- Huy Phan, Philipp Koch, Fabrice Katzberg, Marco Maass, Radoslaw Mazur, and Alfred Mertins. 2017. Audio scene classification with deep recurrent neural networks. *arXiv preprint arXiv:1703.04770*.
- Aurélien Pistono, Mélanie Jucla, Emmanuel J Barbeau, Laure Saint-Aubert, Béatrice Lemesle, Benjamin Calvet, Barbara Köpke, Michèle Puel, and Jérémie Pariente. 2016. Pauses during autobiographical discourse reflect episodic memory processes in early alzheimer’s disease. *Journal of Alzheimer’s Disease*, 50(3):687–698.
- Chloé Pou-Prom and Frank Rudzicz. 2018. Learning multiview embeddings for assessing dementia. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2812–2817.
- Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. ” Why should I trust you?” Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1135–1144.
- Utkarsh Sarawgi, Wazeer Zulfikar, Nouran Soliman, and Pattie Maes. 2020. Multimodal inductive transfer learning for detection of alzheimer’s dementia and its severity. *arXiv preprint arXiv:2009.00700*.
- Frank Seifart, Jan Strunk, Swintha Danielsen, Iren Hartmann, Brigitte Pakendorf, Søren Wichmann, Alena Witzlack-Makarevich, Nivja H de Jong, and Balthasar Bickel. 2018. Nouns slow down speech across structurally and culturally diverse languages. *Proceedings of the National Academy of Sciences*, 115(22):5720–5725.
- Hans Stadthagen-Gonzalez and Colin J. Davis. 2006. The bristol norms for age of acquisition, imageability, and familiarity. *Behavior Research Methods*, 38(4):598–605.
- Laszlo Tóth, Gábor Gosztolya, Veronika Vincze, Ildikó Hoffmann, Gréta Szatlóczi, Edit Biró, Fruzsina Zsura, Magdolna Pákáski, and János Kálmán. 2015. Automatic detection of mild cognitive impairment from spontaneous speech using asr. In *Sixteenth Annual Conference of the International Speech Communication Association*.
- László Tóth, Ildikó Hoffmann, Gábor Gosztolya, Veronika Vincze, Gréta Szatlóczi, Zoltán Bánréti, Magdolna Pákáski, and János Kálmán. 2018. A speech recognition-based solution for the automatic detection of mild cognitive impairment from spontaneous speech. *Current Alzheimer Research*, 15(2):130–138.
- Amy Beth Warriner, Victor Kuperman, and Marc Brysbaert. 2013. Norms of valence, arousal, and dominance for 13,915 english lemmas. *Behavior Research Methods*, 45(4):1191–1207.
- Zichao Yang, Diyi Yang, Chris Dyer, Xiodong He, Alex Smola, and Hovy. 2016. Hierarchical attention networks for document classification. In *Proceedings of NAACL-HLT*, pages 1480–1489.
- Zining Zhu, Jekaterina Novikova, and Frank Rudzicz. 2019. Detecting cognitive impairments by agreeing on interpretations of linguistic features. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 1431–1441.

A Feature list

A.1 Token-Level features

For each token, we extract each of the features listed in Tab. 5. To ensure that each pause has at least one token preceding and succeeding it, *start* and *end* tokens are added to each utterance in DB (Sec. 3.1). For tokens that did not have these features, such as pauses and start/end tokens, all of the features were given a value of zero with the exception of the part of speech. For each feature, with the exception of word length and part of speech, we also extracted the same feature value from the lemmatized token. These features, along with a 5-dimensional, randomly initialized embedding for the part-of-speech, make up the 23-dimensions of each token input vector. Part-of-speech tagging is performed using Spacy (Honnibal and Montani, 2017). All word tokens then have missing values imputed with feature means, and are then normalized with respect to the feature means and standard deviations of the word tokens (i.e. excluding pauses, start/end tokens).

Identical subsequences found in both classes were removed from each data subset. Additionally, if multiple identical subsequences were found in only one of the classes, all but one of them are removed. Furthermore, we do not include utterances that include only a single pause as the final token, as we assume that this pause is occurring between consecutive sentences, rather than within a single sentence.

A.2 Transcript-Level Features

We classify transcripts using 500+ extracted features based on previous literature (Fraser et al., 2016; Tóth et al., 2018), which we refer to as the *Original* feature set. Each of these features come from one of 8 categories:

- **Information Units:** Semantic measures, pertaining to the ability to describe concepts and objects in the picture.
- **Discourse Mapping:** Features that help identify cohesion in speech using a visual representation of message organization in speech. We represent each word as a node to build a ‘speech graph’ (Mota et al., 2012), for the whole transcript. Examples of features extracted include number of edges in this graph, number of self-loops etc.

- **Coherence:** Semantic continuity that listeners perceive between utterances (locally or globally).
- **Lexical Complexity and Richness :** Different measures of lexical qualities and variation. Examples of features include average age of acquisition, number of occurrence of various POS tags etc.
- **Sentiment:** Sentiment lexical norms from Warriner et al. 2013. Examples include average sentiment valence over verbs, average sentiment dominance over nouns etc.
- **Syntactic Complexity:** Different measures to analyze the syntactic complexity of speech including features such as number of occurrence of various production rules, mean length of clause (in words) etc.
- **Word finding Difficulty:** Features quantifying difficulty in finding the right words. These include various pause features such as number of filled pauses, pause word ratio etc.
- **Acoustic:** Voice markers such as MFCC coefficients and Zero Crossing Rate (ZCR) related features.

All transcript-level features, including the aggregates described in Sec. A.3, have missing values imputed with feature medians. Features are then standardized by removing the feature median, and scaled according to the range between the first and third quartile (Pedregosa et al., 2011).

A.3 Transcript-Level Aggregates

In Sec. 4, we extend the *Original* set of features with transcript-level aggregates of token-level features. These include averages for each of the features described in Sec. A.1, with the exception of part of speech. Pauses are not considered when calculating the mean feature values. Additionally, for each part of speech (POS), we include the total amount of times that POS occurs at a certain distance from the pause, divided by the total number of pauses in the transcript multiplied by the percent of words in the transcript that are that POS.

B Classification

B.1 Subsequence Classification

In this work, we perform five-fold cross validation with each of the subsequence data subsets on sev-

Feature	# of features	Description
Word length	2	Length of the word, both in syllables and letters.
Sentiment	6	Three measures of the type and intensity of reaction a word produces.(Warriner et al., 2013)
Concreteness	2	Measure of the degree to which a word refers to a perceptible entity. (Brysaert et al., 2014)
Imageability	2	How easy it is for a word to elicit a mental image. (Stadthagen-Gonzalez and Davis, 2006)
Age of acquisition	2	Average age that the word is learned. (Kuperman et al., 2012)
Frequency	2	Word counts in a corpus of over 385 million words. (Davies, 2009)
Familiarity	2	The perceived popularity of a word.(Stadthagen-Gonzalez and Davis, 2006)
Part of Speech	5	Grammatical category of the word.

Table 5: Token-level linguistic features used as input to the models during subsequence classification and their description.

Model	Bidirectional	# of Layers	Dropout	Epochs	Learning Rate	Momentum	λ	Batch Size	Approx. Time
M-C1	False(12)	2(10, 5)	False	600	0.01	0.9	0.0001	20	6 Min.
M-C2	True(12)	2(10, 5)	True(p=0.5)	600	0.01	0.9	0.0001	20	17 Min.
M-C3	False(50)	1(40)	True(p=0.5)	600	0.01	0.9	0.0001	20	27 Min.
M-Utt	False(50)	2(40, 20)	True(p=0.5)	600	0.01	0.9	0.0001	20	90 Min.

Table 6: Hyperparameters used by the best performing models for each data subset in the subsequence classification task, as well as the approximate time required to complete cross validation. The number of hidden units in the GRU is indicated in the ‘‘Bidirectional’’ column, and the number of hidden units in each layer of the predicting network is indicated in the ‘‘# of Layers’’ column.

Feature Set	Model	# of Features Selected	SMOTE
Original	Ens	-	False
Original w/ feat.sel	Ens	85	False
Original + F-D1	NN	11	False
Original + F-D2	Ens	15	False
Original + F-D3	Ens	5	True
Original + F-C2	NN	20	False
Original + F-C3	Ens	5	True

Table 7: Parameters used by the best performing models for each feature set in the transcript classification task.

eral different GRU based models (Cho et al., 2014) with attention (Yang et al., 2016). Each model consists of a GRU that takes the subsequences as input, and outputs to a feed forward neural network which then makes predictions. The attention mechanism uses a linear layer that has as many hidden units as the GRU, as well as a context vector that has as many dimensions as the GRU has hidden units. An extensive search was conducted in terms of finding the most effective model parameters for each data subset. Each model was tested with variations on the number of intermediate layers in the predicting feed-forward network (1, 2 or 3), the addition of dropout, whether the GRU was bidirectional, and the number of hidden units in each layer (*large* or *small*, where *large* has approximately 4 times as many hidden units in each layer as *small*). This creates a total of 24 trials per data-subset. All the models were created with Pytorch (Paszke et al., 2019), and each model was trained for 600 epochs using SGD as an optimizer, learning rate = 0.01, momentum = 0.9, L2 regularization with $\lambda = 0.0001$, batch size of 20,

a Cosine Annealing learning rate scheduler, and cross entropy loss. Each layer with the exception of the final layer uses the ReLU activation function. Additionally, training scripts were run using the CPU of a p2.xlarge Elastic Compute Cloud instance provided by Amazon Web Services³.

A summary of the hyperparameters used for each of the best performing models reported in Sec. 5.1 are provided in Tab. 6. The number of hidden units in the GRU is indicated in the column indicating whether the GRU was bidirectional, and the number of hidden units in each layer of the predicting network is indicated next to the column indicating the number of layers.

While selecting the best performing model for DB-C1, DB-C2, DB-C3, and DB-Utt, a model was only considered if it was able to meet or exceed the specificity achieved by a model that was once SOTA, 28.8% (Di Palo and Parde, 2019).

B.2 Transcript Classification

We use a Random Forest (100 trees), Gradient Boosting Estimator (with 150 estimators), SVM (with RBF kernel), a 2-layer neural network (NN, 10 units, Adam optimizer, 200 epochs with learning rate initialized to 0.01), and an ensemble of all 4 aforementioned classifiers (Ens) (Pedregosa et al., 2011). For each extending feature set described in Sec. 3.2, we jointly optimize the number of features selected from the extending feature set (feature selection with k=3, 5, 7, 9, 11, 13, 15, 20,

³<https://aws.amazon.com/ec2/>

25, 30, and all features), whether or not we use the oversampling method SMOTE([Chawla et al., 2002](#)), and the model type used. For feature selection on the *Original* feature set, we attempted feature selection with $k=20, 25, 30, 35, 40, 45, 50, 55, 60, 65, 70, 75, 80, 85, 90, 95, 100, 150, 200, 250, 300,$ and 350. The best performing configuration for each extending feature set is recorded in [Tab. 7](#).