

# Too Big to Fail: Larger Language Models are Disproportionately Resilient to Induction of Dementia-Related Linguistic Anomalies

Changye Li<sup>1</sup>, Zhecheng Sheng<sup>1</sup>, Trevor Cohen<sup>2</sup>, and Serguei Pakhomov<sup>1</sup>

<sup>1</sup>University of Minnesota

<sup>2</sup>University of Washington

<sup>1</sup>{lix3013, sheng136, pakh0002}@umn.edu

<sup>2</sup>{cohenta}@uw.edu

## Abstract

As artificial neural networks grow in complexity, understanding their inner workings becomes increasingly challenging, which is particularly important in healthcare applications. The intrinsic evaluation metrics of autoregressive neural language models (NLMs), perplexity (PPL), can reflect how “surprised” an NLM model is at novel input. PPL has been widely used to understand the behavior of NLMs. Previous findings show that changes in PPL when masking attention layers in pre-trained transformer-based NLMs reflect linguistic anomalies associated with Alzheimer’s disease dementia. Building upon this, we explore a novel bidirectional attention head ablation method that exhibits properties attributed to the concepts of cognitive and brain reserve in human brain studies, which postulate that people with more neurons in the brain and more efficient processing are more resilient to neurodegeneration. Our results show that larger GPT-2 models require a disproportionately larger share of attention heads to be masked/ablated to display degradation of similar magnitude to masking in smaller models. These results suggest that the attention mechanism in transformer models may present an analogue to the notions of cognitive and brain reserve and could potentially be used to model certain aspects of the progression of neurodegenerative disorders and aging.

## 1 Introduction

Alzheimer’s disease (AD) dementia is a currently incurable neurodegenerative condition that leads to a progressive and irreversible decline in cognitive function. Due to the challenging nature of early diagnosis of this condition, there is a pressing need for efficient and cost-effective screening tools (Bradford et al., 2009) to mitigate the negative consequences of delayed or absent diagnosis (Stokes et al., 2015). Previous studies have demonstrated that changes in cognitive status can be reflected in

spoken language and spontaneous speech (Giles et al., 1996; Almor et al., 1999; Hier et al., 1985). Automated analysis of such speech, employing supervised machine learning models, shows its potential as an early screening tool. These models can be trained to identify subtle linguistic anomalies associated with dementia from transcripts of both healthy individuals and those with dementia. Recent advances in machine learning, such as deep learning models and the transformer with attention architecture (Vaswani et al., 2017), have mediated remarkable performance on this downstream task (for a review, see Shi et al. (2023)). Deep learning models, inspired by the human brain, are artificial neural networks (ANNs) that process vast amounts of data and learn complicated patterns, making them well-suited for analyzing subtle linguistic patterns. The transformer architecture, in particular, has advanced performance on natural language processing (NLP) tasks by enabling models to capture long-range dependencies more effectively via the attention mechanism (Vaswani et al., 2017).

As ANNs get larger and more complicated, it becomes even harder to interpret their inner workings. The performance of autoregressive neural language models (NLMs) (e.g., predicting the next word given the context) is frequently estimated with a single somewhat interpretable feature, perplexity (PPL), which has shown to be a suitable measurement for evaluating cognitive impairment from spontaneous speech (Fritsch et al., 2019; Cohen and Pakhomov, 2020). As the name “perplexity” suggests, it can be considered as an indicator of how “surprised” a model is by novel (i.e., not used in model’s training) input. The more different the input is from a particular model’s training data, the “harder” it is for the model to predict, resulting in higher PPL. Therefore, it is reasonable to hypothesize that PPL may have some degree of diagnostic utility, as an indicator of patterns of language use that fall outside the scope of the typical

language used to train a model. In the context of AD, changes in language and cognitive function often manifest as differences in language complexity, with individuals experiencing difficulty in forming coherent sentences and selecting appropriate words. As AD progresses, the language used by patients with dementia becomes more unpredictable and less coherent, leading to higher PPL with models trained on language from individuals presumed to be cognitively healthy.

While training data from cognitively healthy individuals is plentiful, language data produced by patients with dementia is much more impractical to obtain in sufficient quantity to train a large NLM. In hyperdimensional computing (Kanerva, 2009), high-dimensional vector representations are manipulated using operators that alter their distance from other learned representations. A prior work inspired by this concept (Li et al., 2022) demonstrates that masking the attention sub-modules of pre-trained transformer-based NLMs and thereby artificially increasing PPL on text from cognitively healthy individuals, can provide an effective solution to the challenge of limited data availability. By strategically altering these sub-modules and introducing controlled perturbations in the NLMs' attention layers, the degraded NLMs induce the linguistic anomalies and unpredictability associated with dementia.

Recent work in neuroscience using functional magnetic resonance imaging (fMRI) and electrocorticography (ECoG) has demonstrated that NLM's PPL is associated with predicting neural activation patterns during language comprehension tasks in the human brain (Schrimpf et al., 2021; Hosseini et al., 2024). This suggests a potential connection between the predictive capabilities of these models and understanding human information processing. In particular, one of the less well-understood phenomena in how neurodegeneration affects the human brain is the notion of cognitive and brain reserve. This notion is hypothesized to be responsible for findings that indicate individuals with higher innate abilities and/or aspects of life experience, such as educational and professional attainment, are able to mask the effects of dementia longer than those without these characteristics (Stern, 2002, 2009, 2012; Scarmeas and Stern, 2004, 2003; Snowdon et al., 1996). In some cases, the notion of cognitive and brain reserve may even allow individuals to revert from initial signs of cognitive impairment to normal function (Iraniparast et al., 2022).

Building upon these findings, our study seeks to further explore the potential of probing pre-trained GPT-2 family models (Radford et al., 2019) to simulate cognitive impairment observed in patients with dementia with a specific focus on the cognitive reserve hypothesis. Using a set of transcripts from a widely-used "Cookie Theft" picture description cognitive task, we propose that the impaired information processing as the disease progresses can be simulated by masking a certain share of attention heads in a pre-trained GPT-2 model. Specifically, we follow the previously established paired-perplexity paradigm (Li et al., 2022) using a pair of unmasked ("control") and masked ("dementia") NLMs. In this approach, the difference between PPLs produced by these two NLMs is used to discriminate between picture descriptions by patients with dementia and healthy controls. We hypothesize that larger GPT-2 models with more attention heads will exhibit greater resilience to masking (i.e., a proxy for neural degeneration), necessitating a larger share of attention heads to be masked to achieve comparable classification performance to smaller models. We evaluate this hypothesis by targeting two subsets of attention heads that are a) *most* important, and b) *least* important to representation of the content of the "Cookie Theft" task, in which the degree of importance is ranked by the gradient changes in each attention head during fine-tuning of a pre-trained GPT-2 model to the content of the "Cookie Theft" transcripts.

The contributions of this work can be summarized as follows: a) we provide preliminary evidence suggesting that the concept of cognitive reserve observed in human cognition appears to have an analog in ANNs; and b) our attention masking approach achieves comparable classification performance to another approach developed in prior work that directly artificially degrades NLM parameters (Li et al., 2022), and the state-of-the-art (SOTA) model trained from scratch (TaghiBeyglou and Rudzicz, 2024) with *significantly fewer* trainable parameter masking/fitting.<sup>1</sup>

## 2 Background

### 2.1 Cognitive Reserve

The notions of brain plasticity in the human brain and "graceful degradation" in ANNs have been

<sup>1</sup>The code to reproduce the results presented in this paper is available at [GitHub](#). The data are also publicly available but cannot be redistributed and must be obtained directly from Dementia Bank.

extensively investigated in the neuroscientific literature demonstrating, for example, that a large proportion (over 80%) of the connections in an ANN trained to simulate the motor cortex to generate signals directing body movement have to be ablated before the model's performance begins to collapse (Lukashin et al., 1994; Lukashin and Georgopoulos, 1994). The concepts of cognitive and brain reserves are closely related to brain plasticity applied to observations in neurodegenerative diseases as illustrated in Figure 1. One of the earlier observations of this phenomenon comes from the Nun Study which found that low linguistic ability early in life (possibly due to innate abilities or educational attainment) is predictive of poor cognitive function and AD later in life (Snowdon et al., 1996). The concept of cognitive reserve was further developed based on observations of the individual differences in effects of brain damage or pathology on clinical manifestations of cognitive function (Stern, 2002, 2009, 2012). A multi-site study (Esiri et al., 2001) reported that up to 25% older adults without signs of cognitive impairment during neuropsychological testing meet all the histopathological criteria for AD (amyloid plaques and tau protein tangles) prior to their death. While this study did not assess brain volume, another similar study did find that a subgroup of 10 study participants who had both AD pathology and preserved mental status had greater brain weights and number of neurons in their brains (Katzman et al., 1988).

A distinction can be made between the closely related notions of cognitive reserve and brain reserve. Cognitive reserve refers to the efficiency of brain networks, which manifests as greater educational and professional attainment. Brain reserve, on the other hand, refers to the physical properties of the brain, such as a larger number of neurons in biological neural network(s). This can manifest, for example, as a higher intelligence quotient. These two notions are difficult to disentangle due to their significant interdependence (Steffener and Stern, 2012). The properties of these notions have also been described using passive or active models that correspond to the notions of brain and cognitive reserves, respectively. Passive models (Katzman, 1993; Satz, 1993) measure the cognitive reserve by the size of the brain or the count of neurons in the brain. Passive models hypothesize that there is a threshold for brain reserve capacity - once an individual passes the "point of no return", the manifestation of neurodegenerative disease, such

as AD, will occur regardless. Contrary to passive models, active models (Stern, 2002) hypothesize that there is a neural compensatory effect for brain damage. This effect consists of the brain compensating for the damage by activating other biological neural network(s) to perform cognitive task-related activities. In this case, patients of similar brain impairment but with more cognitive reserve may be more resilient to the disease's progression before the clinical manifestations of neurodegeneration become apparent. Quantitatively, there is no clear difference between the passive and active models of cognitive reserve, as both of them rely on the physiologic basis of biological neural networks in the brain. This provides an opportunity to evaluate the underlying mechanisms that contribute to cognitive reserve across various neurological conditions computationally.

To avoid any potential confusion between these terms referring to different types of resilience and to avoid any inadvertent conflation between artificial and human brain networks, in the remainder of this paper we will refer to the phenomenon of resilience to damage that we observe in ANNs specifically as "artificial neural reserve" and use the terms "cognitive/brain reserve" to refer exclusively to human brain networks.

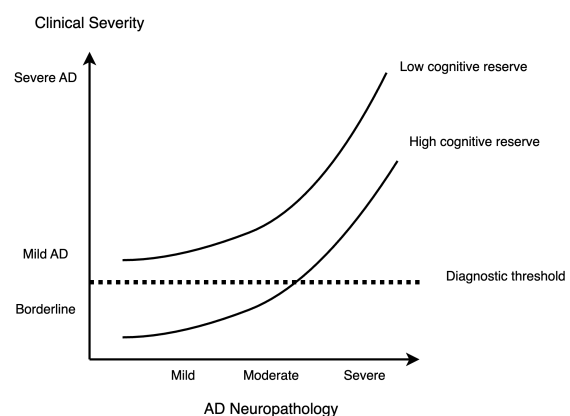


Figure 1: A theoretical illustration of cognitive reserve and its mediation effect between AD neuropathology (x-axis) and clinical outcome (y-axis). Illustration derived from Stern (2002, 2009). As the disease progresses (i.e., with more impairment), individuals with higher cognitive/brain reserve would be more resilient to the effects, resulting in a lower level of clinical severity.

## 2.2 Probing the Neural Network

The ablation of connections in ANNs is also referred to as probing in NLMs. This is a growing field aimed at understanding the inner workings

of large-scale transformer-based NLMs by probing the mechanism (i.e., attention weights, hidden states) to better understand the linguistic structure and representations encoded by such models. Similarly to the early findings of Lukashin et al. (1994) and Lukashin and Georgopoulos (1994), more recent work on transformers (i.e., Michel et al. (2019); Prasanna et al. (2020); Sanh et al. (2020)) demonstrates that a large percentage of attention heads or sub-modules can be removed at inference time without significantly impacting performance.

### 2.3 Linguistic Anomalies in AD

AD is a neurodegenerative disease, and progressively worsening linguistic deficits often accompany its progression (Kempler and Goral, 2008; Altmann and McClung, 2008). A widely-used diagnostic task to capture such linguistic anomalies is the “Cookie Theft” picture description task from the Boston Diagnostic Aphasia Examination (Goodglass and Kaplan, 1983). In this task, participants are asked to describe everything they see going on in Figure 2. Previous studies have demonstrated that dementia patients tend to overuse pronouns (Almor et al., 1999) and tend to perseverate (Hier et al., 1985) when describing the “Cookie Theft” picture.

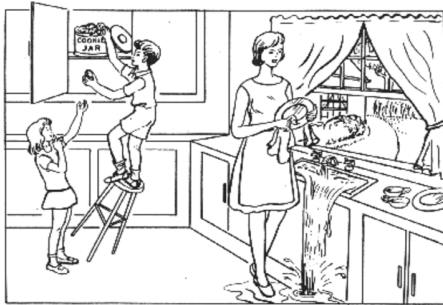


Figure 2: The “Cookie Theft” picture description stimuli.

There is a rich body of evidence that supervised machine learning and deep learning methods can learn to distinguish the subtle linguistic characteristics between healthy individuals and people with dementia. However, such models present a danger of overfitting, and hinder interpretability of model predictions, which are both critical concerns for clinical artificial intelligence (AI) applications (Graham et al., 2020). Alternatively, PPL is an easily interpretable measure used to evaluate model performance. With dementia, the difference of the *paired-perplexity* paradigm from a “healthy

control” NLM and a “dementia” NLM provides a diagnostically useful summary value that distinguishes language samples produced by dementia patients (Fritsch et al., 2019; Cohen and Pakhomov, 2020). Prior work (Li et al., 2022) has shown that the difference of PPLs from a pre-trained GPT-2 paired with an artificially degraded version of itself approximates SOTA classification performance *without* requiring a data set from dementia patients of comparable size to its comprehensive training data. However, this approach requires evaluating thousands of masking patterns in order to investigate the effects of masking various combinations of attention heads exhaustively. In the current work we obviate this requirement for extensive experimentation by using targeted masking (guided by the changes in gradients during training) of two subsets of attention heads that are a) *most* “important”, and *least* “important” with respect to the content of the “Cookie Theft” picture. We show that the resulting masked models can effectively identify transcripts from dementia patients with *significantly fewer* trainable parameters while exhibiting comparable classification performance to previous studies (Li et al., 2022; TaghiBeyglou and Rudzicz, 2024).

## 3 Methods

### 3.1 Data

We use two publicly available datasets that contain responses to the “Cookie Theft” picture description task: a) AD Recognition through Spontaneous Speech (ADReSS) Challenge<sup>2</sup> (Luz et al., 2020), and b) the Wisconsin Longitudinal Study (WLS)<sup>3</sup> (Herd et al., 2014). Table 1 shows basic characteristics of datasets used in this study. ADReSS is a subset of the Pitt corpus (Becker et al., 1994) designed to address the absence of a standardized train/test split in prior work. It is specifically matched on age and gender to reduce potential confounding effects. The WLS is a longitudinal study of 694 men and 675 women who graduated from Wisconsin high schools in 1957. The participants were interviewed up to 6 times between 1957 and 2011. The “Cookie Theft” picture description task was administered in the later round of interviews. In particular, we restricted the original WLS dataset to a total of 102 participants who a) agreed to partici-

<sup>2</sup><https://dementia.talkbank.org/ADReSS-2020/>

<sup>3</sup><https://dementia.talkbank.org/access/English/WLS.html>

pate in the “Cookie Theft” picture description task, and b) had either a clinical diagnosis of dementia or were deemed healthy in follow-up interviews conducted in 2020. This information was obtained through phone interviews and assessments by advanced practice providers. Subsequently, the collected data was presented to a panel of clinicians to obtain the diagnosis.

Dataset	# of participants ( $n$ )		
	Dementia	Healthy Controls	
ADReSS	Train	54	54
	Test	24	24
	Total	78	78
WLS		29	73

Table 1: The characteristics of ADReSS and WLS.

We perform verbatim transcripts pre-processing using TRESTLE (Toolkit for Reproducible Execution of Speech Text and Language Experiments) (Li et al., 2023) by removing utterances that do not belong to the participants, unintelligible words, and speech and non-speech artifacts event descriptions (i.e., “laughs”, “clear throat”).

### 3.2 Modeling and Evaluation

We follow a similar masking strategy to that proposed by Michel et al. (2019) to mask attention heads of the GPT-2 small, medium, large, and XL models via the rank of their importance to the task. We focus on the GPT-2 family models to minimize the variability that would result from multiple modeling architectures. The task-importance of attention heads in each model is determined by the gradient changes during the fine-tuning for subsequent word prediction task using transcripts of “Cookie Theft” picture descriptions in the training portion of the ADReSS dataset. Intuitively, if the gradient change of an attention head is large, this attention head is likely important with respect to predicting the language to which the model is being fine-tuned, and vice versa.

In contrast to the approach by Michel et al. (2019), which prunes the *least* important attention heads during testing, we anticipate that the *most* important attention heads are those relevant for predicting the text of the “Cookie Theft” task. This idea is supported by Yorkston and Beukel-

man (1980) and Berube et al. (2019), who found that the number of content units represented – a measure of how much relevant information is conveyed in the description – is sensitive to linguistic deficits often observed in individuals with neurodegenerative disease. However, we also reason that the *least* important attention heads may represent subtle differences in linguistic structure and representations that may distinguish between dementia patients and healthy controls. We also test the possibility that the semantic impairment observed in AD (Huff et al., 1986; Giffard et al., 2001; Hodges and Patterson, 1995) could be potentially simulated by masking a certain share of the columns in the pre-trained NLMs’ token embedding matrix, where each column contributes to the representation of the meaning of each token in the model’s vocabulary. Thus, masking columns in the embedding matrix leads to degrading the representation of *all* vocabulary items vs. degrading or deleting specific tokens from the otherwise intact vocabulary by operating on the rows of the embedding matrix.

Following these considerations, we design the masking strategies as follows: a) we fine-tune each of the GPT-2 models with a language model head layer as the top layer on the ADReSS training set to get the corresponding ranking of importance for each of the attention heads; b) we iteratively mask a small share ( $n\%$ ) of ranked attention heads *bidirectionally*, which consists of the  $\frac{n}{2}\%$  *most* important attention heads and the  $\frac{n}{2}\%$  *least* important attention heads, then gradually increase the percentage of attention heads for masking, and c) we iteratively mask columns of the word embedding matrix in reverse order, moving from right to left, and gradually increase the percentage of word embedding columns for masking<sup>4</sup>.

We examine the artificial neural reserve hypothesis using two evaluation approaches. The first approach consists of simply estimating the PPL of the progressively degraded NLMs based on healthy individuals’ transcripts from an independent dataset containing the same type of picture descriptions as the dataset that was used to rank attention heads by their importance to the dementia classification task. We use the WLS dataset and select only those WLS participants that remained cognitively healthy over the entire study period as the independently collected dataset for log PPL estimation. Using this

<sup>4</sup>All experiments in this study are done with HuggingFace’s transformers package (Wolf et al., 2020) on one A100 GPU.

approach, in addition to masking attention heads, we also experiment with masking model weights in the token embedding matrix to see if any observed effects are specific to the attention mechanism.

The second approach consists of evaluating the classification performance of ablated/degraded models paired with the original versions of the same GPT-2 model using the paired-perplexity paradigm (Fritsch et al., 2019; Cohen and Pakhomov, 2020; Li et al., 2022). These evaluations are conducted on the testing portion of the ADReSS dataset, with accuracy (ACC) and area under the receiver-operator characteristic (ROC) curve (AUC) as the evaluation metrics. Specifically, for the paired-perplexity paradigm, we estimate the ratio of PPLs  $\frac{PPL_{control}}{PPL_{dementia}}$  of each transcript from the test set. The ACC measure is calculated as accuracy at the equal error rate (EER), where the false acceptance rate is equal to false rejection rate on the ROC curve. The intuition behind this approach is based on the expectation that successful masking of a portion of attention heads in a pre-trained NLM will result in the NLM exhibiting dementia-like behavior, which would in turn result in high AUC and ACC values of the paired-perplexity classification.

## 4 Results

### 4.1 Effects of Masking on Perplexity

As illustrated in Figure 3a, the predictive ability of smaller GPT-2 models degrades linearly with the degree of damage inflicted on the attention mechanism by masking progressively larger proportion of attention heads. The predictive ability of the larger GPT-2 models, on the other hand, degraded in a non-linear fashion where increases in log PPL were relatively flat up to 40-50% of the attention heads being masked and then began to increase exponentially. Fitting the GPT-2 small, medium, large and XL model log PPL to a linear regression line resulted in  $r^2$  goodness-of-fit values of 0.99, 0.89, 0.91 and 0.83, respectively, whereas fitting to an exponential regression line failed to converge for the small and medium models and yielded  $r^2$  values of 0.97 and 0.99 for the large and XL models, respectively. The results of Dunn’s test further confirmed our observations, showing that the differences between log PPLs estimated by GPT-2 small and GPT-2 XL (adjusted p-value < 0.01), and GPT-2 medium and GPT-2 XL (adjusted p-value < 0.05) when masking attention heads are statistically significant. In contrast, *all* combina-

tions of log PPLs were not significantly different from each other for all GPT-2 models when masking the word embedding matrix (adjusted p-value > 0.05).

Compared to masking attention heads, with GPT-2 small, medium, large and XL model we needed to mask 93% (714 out of 768), 66% (675 out of 1024), 87% (1113 out of 1280), and 66% (1050 out of 1600) columns in the word embedding matrix to achieve ACCs of 0.75, 0.85, 0.79, and 0.81 respectively, on the ADReSS test set. Figure 3b can further support this claim, as estimated log PPLs of masking word embedding matrix show no significant statistical differences across various GPT-2 models.

### 4.2 Effects of Masking on Dementia Classification

As shown in Table 2, impairing 9% of attention heads ( $n=12$ ) of the GPT-2 small model (the “dementia” model) achieved an ACC of 0.83 and AUC of 0.86 when paired with the original unmasked version of itself (the “control” model) on the ADReSS test set. This is comparable to the prior work (Li et al., 2022) (ACC = 0.85, AUC = 0.89) but the masking approach uses *significantly fewer* masked parameters. Our results also show that a larger share of attention heads in the larger models must be masked to approximate a “dementia” model with the same level of performance in the paired-perplexity classification than with smaller models. Notably, masking of the word embedding matrix did not result in comparable observations. As anticipated by the results shown in Figure 3b, with GPT-2 models we needed to mask a majority portion (e.g., > 50%) of the word embedding matrix to obtain similar level of classification performance regardless of model size.<sup>5</sup>

As illustrated in Figure 4, we observed that once the best-performing masking pattern, marked by the highest ACC, was reached, the classification performance of all GPT-2 models started to fluctuate. However, this observation did not occur with the word embedding matrix masking. As illustrated in Figure 5 in Appendix A, the classification performance exhibited fluctuations prior to the emergence of the best-performing masking pattern, indicating that masking the columns of the word embedding matrix has less impact on identifying the signs of

<sup>5</sup>The importance of attention heads for each model can be found in Table 3, Table 4, Table 5, and Table 6 in the Appendix A.

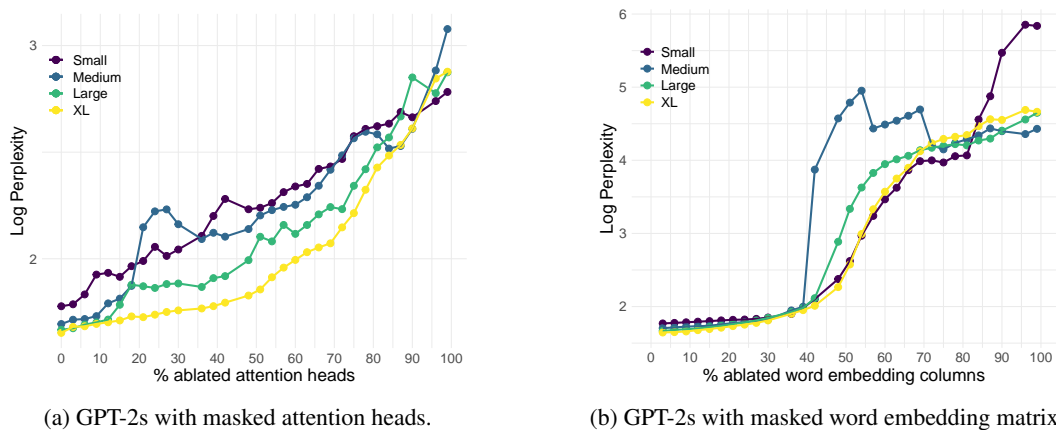


Figure 3: Changes in model log PPL as a function of the proportion of masked attention heads across GPT-2 models of various sizes. Note: the curves in panel (a) show that GPT-2 XL model has the most non-linear/concave shape indicating that the model starts to degrade rapidly only after masking of about 50% of its attention heads, followed by the curve for the GPT-2 large model. The smaller GPT-2 models begin to degrade with proportionally less masking, and exhibit a monotonic relationship between the magnitude of attention heads masking and model performance. The curves in panel (b) show almost completely preserved model performance without differences between models up to the point at which 40% - 50% of the columns in their embedding matrices have been masked. After that point, the performance of all models collapses “catastrophically”

Model	GPT-2 small	GPT-2 medium	GPT-2 large	GPT-2 XL
# of parameters	124M	355M	774M	1.5B
# of masked attention heads	12	92	388	1080
% of masked attention heads	9	24	54	90
ACC	0.83	0.83	0.81	0.81
AUC	0.86	0.85	0.80	0.82

Table 2: Classification performance of the paired-perplexity approach based on pre-trained and masked GPT-2 models on the ADRess test set.

cognitive impairment from text as it probably does not result in a good dementia-like model for the paired-perplexity classification task.

## 5 Discussion

The results of experiments presented in this paper suggest that the notion of cognitive reserve in the brain may have an analogue in transformer-based ANNs that is localized to the attention mechanism. Recent neuroscientific evidence shows that NLMs’ PPL is predictive of human behavioral responses and neural responses in functional MRI studies (Schrimpf et al., 2021; Hosseini et al., 2024). Based on this evidence, we interpret our findings of the differences in log PPL changes as a result of masking attention heads in NLMs of variable size

as at least suggestive that the resilience to damage is non-linear to the number of attention heads in NLMs. In other words, it takes disproportionately more masking to damage larger NLMs to elicit the same level of degradation in performance, as compared with smaller NLMs. Furthermore, the dissociation in performance as a result of damaging attention heads vs. the token embedding weights suggests that the NLM’s artificial neural reserve effects are localized to the attention mechanism.

Our results also suggest that masking attention heads within the paired-perplexity paradigm using the ratio of unmasked (“control”) and masked (“dementia”) pre-trained GPT-2 models results in good classification performance *without* requiring a corresponding large dataset produced by demen-

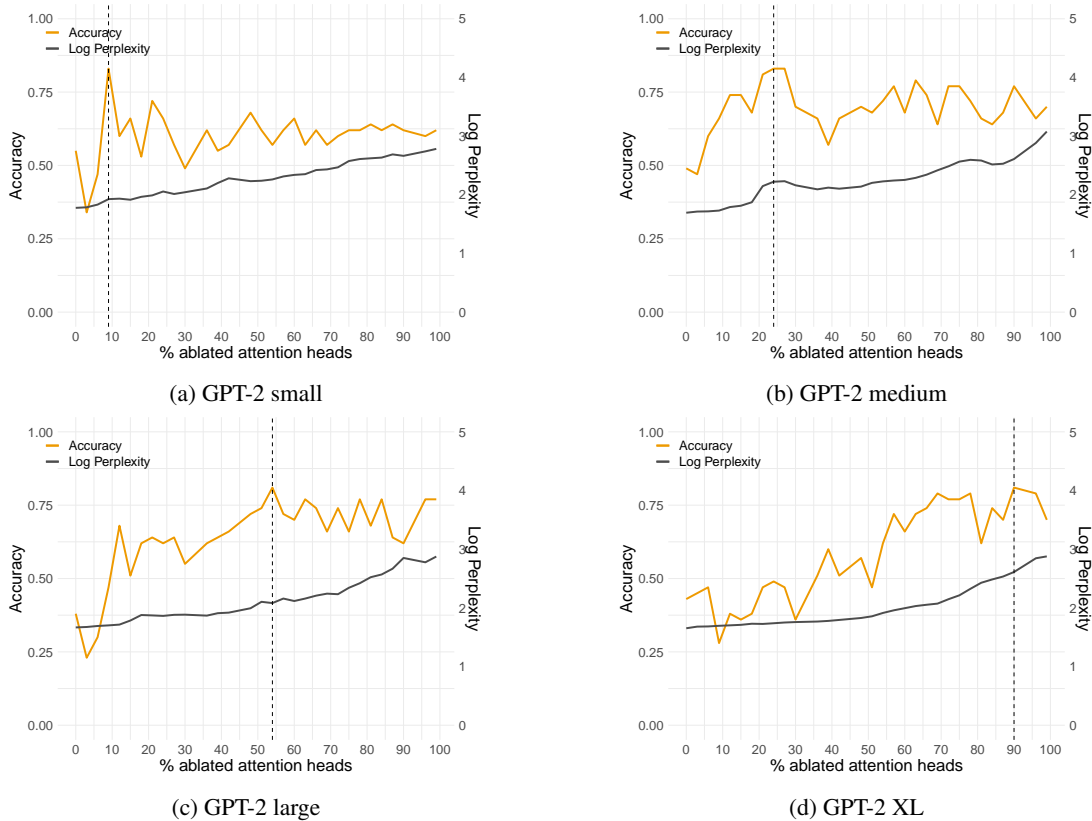


Figure 4: Comparison of GPT-2 models with masked attention heads on paired-perplexity classification performance. The left y-axis denotes classification performance using both masked and unmasked GPT-2 models on the ADReSS test set. The right y-axis indicates log PPL estimated from transcripts of WLS healthy individuals. The x-axis represents the percentage of attention heads getting masked. The vertical dashed line indicates the best-performing masking pattern, achieving the highest ACC.

tia patients and extensive parameter tuning. This can be achieved with as little as masking only 9% of attention heads of a pre-trained GPT-2 small model.

In contrast to previous studies, which typically involved purging attention heads determined to be the *least* important, our bidirectional masking method adds supporting evidence of content units (Yorkston and Beukelman, 1980; Berube et al., 2019), suggesting the importance of these contextual features in addition to the predominant emphasis on linguistic structure and representation modeling in previous research (e.g., Orimaye et al. (2014), Fraser et al. (2016)). The results of bidirectional masking also offers an interpretable explanation for transfer learning’s remarkable performance using pre-trained NLMs. It suggests that during fine-tuning, pre-trained NLMs use a combination of both task-specific (the *most* important) and task-agnostic (the *least* important) heads to achieve remarkable performance on various downstream tasks. Those task-agnostic attention heads

may play an important role in transfer learning. This also may explain why distilled NLMs in which the “nonvolitional cues” that fall outside of common NLP benchmarks are purged during the distillation, generalize less-than-ideally to other types of data produced by individuals with dementia (Li et al., 2022). With larger models, there are considerably more attention heads that can serve as those “nonvolitional cues,” helping a larger NLM perform better (Agbavor and Liang, 2022).

As the columns of the token embedding matrix in a pre-trained NLMs represent the global semantics of tokens in the vocabulary, the observations that the best-performing masking pattern appears in the later stage of the token embedding matrix masking are consistent with previously published findings that semantic impairment often occurs in the later stage (i.e., moderate) of the disease (Huff et al., 1986; Giffard et al., 2001; Hodges and Patterson, 1995). As illustrated in Figure 5, when masking the later 66% columns (675 out of 1024) of the word embedding matrix, the paired unmasked and



masked GPT-2 medium achieves an ACC of 0.85 on the ADReSS test set. This finding is consistent with a previous work (Hewitt and Manning, 2019), suggesting that some syntactic information is embedded implicitly in the word embedding matrix. This also provides an empirical support of our findings that masking word embedding matrix of a pre-trained NLM can provide some degree of discriminating effect on this downstream task. However, masking the word embedding matrix is far less effective than masking attention heads to simulate dementia-related cognitive impairment.

Our results suggest that similar mechanisms of resilience may exist in both human cognition and computational models. This could lead to more nuanced strategies in response to develop early screening tools for the delayed onset of cognitive impairment in the population with high risk. Our study holds promise for deploying such early-screening method in resource-constrained clinical settings to improve early intervention and patient management for AD.

## 6 Conclusion

We presented experimental findings suggesting the presence of artificial neural reserve in transformer-based NLMs, analogous to the concepts of brain/cognitive reserve in studies of human cognition. In addition, we introduced a novel bidirectional attention head ablation method that enables using unmasked and masked GPT-2 models in the paired-perplexity paradigm for detecting linguistic anomalies with significantly less parameter masking or fitting.

## Acknowledgement

Research reported in this publication was supported by the National Library of Medicine of the National Institutes of Health under award number R01 LM014056-02S1, and the National Institute on Aging of the National Institutes of Health under award number AG069792.

## Limitations

The work presented here has several limitations. First, the size of datasets used in this study is relatively small compared to datasets typically analyzed in the open-domain NLP tasks, therefore the results may not be readily generalizable. Second, all datasets used in our study are in American English, and many participants of these two studies

are representative of White, non-Hispanic American men and women located at the north part of the United States, which certainly limit their applicability to other languages. Third, while we propose that the findings presented in this paper may be interpreted as an analogue of the notions of cognitive or brain reserve, we do not suggest that GPT-2 models are accurate models of the human brain. Rather, our interpretation of these findings is that experimenting with masking of attention heads in models of various sizes and architectures may be useful in helping us understand cognitive processes that take place in the human brain. The observed effects of attention masking on the model's performance and behavior, while suggestive of an analog to cognitive reserve in the human brain, should not establish a direct causal link to human cognitive processes. Additionally, the ranking of attention heads by their relative importance is specific to the ADReSS dataset as it was derived in the training portion of the dataset and may not readily generalize to other datasets and types of data. Lastly, in this paper we did not address the distinction between the notions of cognitive and brain reserves. It would be important to investigate in future work if NLMs of the same size and architecture but different quantities and quality of the training data (i.e., as a simulation of educational attainment) exhibit differential resilience to damage independently of the effects observed in models of variable size.

## References

- Felix Agbavor and Hualou Liang. 2022. [Predicting dementia from spontaneous speech using large language models](#). *PLOS Digital Health*, 1(12):e0000168.
- Amit Almor, Daniel Kempler, Maryellen C MacDonald, Elaine S Andersen, and Lorraine K Tyler. 1999. [Why do alzheimer patients have difficulty with pronouns? working memory, semantics, and reference in comprehension and production in alzheimer's disease](#). *Brain and Language*, 67(3):202–227.
- Lori JP Altmann and Jill S McClung. 2008. [Effects of semantic impairment on language use in alzheimer's disease](#). In *Seminars in speech and language*, volume 29, pages 018–031. © Thieme Medical Publishers.
- 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.

- Shauna Berube, Jodi Nonnemacher, Cornelia Demsky, Shenly Glenn, Sadhvi Saxena, Amy Wright, Donna C Tippett, and Argye E Hillis. 2019. [Stealing cookies in the twenty-first century: Measures of spoken narrative in healthy versus speakers with aphasia](#). *American journal of speech-language pathology*, 28(1S):321–329.
- Andrea Bradford, Mark E Kunik, Paul Schulz, Susan P Williams, and Hardeep Singh. 2009. [Missed and delayed diagnosis of dementia in primary care: prevalence and contributing factors](#). *Alzheimer disease and associated disorders*, 23(4):306.
- Trevor Cohen and Serguei Pakhomov. 2020. [A tale of two perplexities: Sensitivity of neural language models to lexical retrieval deficits in dementia of the Alzheimer’s type](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1946–1957, Online. Association for Computational Linguistics.
- MM Esiri, F Matthews, C Brayne, PG Ince, FE Matthews, JH Xuereb, JC Broome, J McKenzie, M Rossi, IG McKeith, et al. 2001. [Pathological correlates of late-onset dementia in a multicentre, community-based population in england and wales](#). *Lancet*.
- 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.
- Julian Fritsch, Sebastian Wankerl, and Elmar Nöth. 2019. [Automatic diagnosis of alzheimer’s disease using neural network language models](#). In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 5841–5845. IEEE.
- Bénédicte Giffard, Béatrice Desgranges, Florence Nore-Mary, Catherine Lalevée, Vincent de la Sayette, Florence Pasquier, and Francis Eustache. 2001. [The nature of semantic memory deficits in Alzheimer’s disease: New insights from hyperpriming effects](#). *Brain*, 124(8):1522–1532.
- Elaine Giles, Karalyn Patterson, and John R Hodges. 1996. [Performance on the boston cookie theft picture description task in patients with early dementia of the alzheimer’s type: Missing information](#). *Aphasiology*, 10(4):395–408.
- Harold Goodglass and Edith Kaplan. 1983. *Boston diagnostic aphasia examination booklet*. Lea & Febiger.
- Sarah A Graham, Ellen E Lee, Dilip V Jeste, Ryan Van Patten, Elizabeth W Twamley, Camille Nebeker, Yasunori Yamada, Ho-Cheol Kim, and Colin A Depp. 2020. [Artificial intelligence approaches to predicting and detecting cognitive decline in older adults: A conceptual review](#). *Psychiatry research*, 284:112732.
- Pamela Herd, Deborah Carr, and Carol Roan. 2014. [Cohort profile: Wisconsin longitudinal study \(wls\)](#). *International journal of epidemiology*, 43(1):34–41.
- John Hewitt and Christopher D. Manning. 2019. [A structural probe for finding syntax in word representations](#). 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 4129–4138, Minneapolis, Minnesota. Association for Computational Linguistics.
- Daniel B. Hier, Karen Hagenlocker, and Andrea Gellin Shindler. 1985. [Language disintegration in dementia: Effects of etiology and severity](#). *Brain and Language*, 25(1):117–133.
- John R. Hodges and Karalyn Patterson. 1995. [Is semantic memory consistently impaired early in the course of alzheimer’s disease? neuroanatomical and diagnostic implications](#). *Neuropsychologia*, 33(4):441–459.
- Eghbal A Hosseini, Martin Schrimpf, Yian Zhang, Samuel R Bowman, Noga Zaslavsky, and Evelina Fedorenko. 2024. [Artificial neural network language models predict human brain responses to language even after a developmentally realistic amount of training](#). *Neurobiology of Language*, pages 1–50.
- F.Jacob Huff, Suzanne Corkin, and John H. Growdon. 1986. [Semantic impairment and anomia in alzheimer’s disease](#). *Brain and Language*, 28(2):235–249.
- Maryam Iraniparast, Yidan Shi, Ying Wu, Leilei Zeng, Colleen J Maxwell, Richard J Kryscio, Philip D St John, Karen S SantaCruz, and Suzanne L Tyas. 2022. [Cognitive reserve and mild cognitive impairment: predictors and rates of reversion to intact cognition vs progression to dementia](#). *Neurology*, 98(11):e1114–e1123.
- Pentti Kanerva. 2009. [Hyperdimensional computing: An introduction to computing in distributed representation with high-dimensional random vectors](#). *Cognitive computation*, 1:139–159.
- Robert Katzman. 1993. [Education and the prevalence of dementia and alzheimer’s disease](#). *Neurology*, 43(1\_part\_1):13–13.
- Robert Katzman, Robert Terry, Richard DeTeresa, Theodore Brown, Peter Davies, Paula Fuld, Xiong Renbing, and Arthur Peck. 1988. [Clinical, pathological, and neurochemical changes in dementia: a subgroup with preserved mental status and numerous neocortical plaques](#). *Annals of Neurology: Official Journal of the American Neurological Association and the Child Neurology Society*, 23(2):138–144.
- Daniel Kempler and Mira Goral. 2008. [Language and dementia: Neuropsychological aspects](#). *Annual review of applied linguistics*, 28:73–90.
- Changye Li, David Knopman, Weizhe Xu, Trevor Cohen, and Serguei Pakhomov. 2022. [GPT-D: Inducing dementia-related linguistic anomalies by deliberate degradation of artificial neural language models](#). In

- Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1866–1877, Dublin, Ireland. Association for Computational Linguistics.
- Changye Li, Weizhe Xu, Trevor Cohen, Martin Michalowski, and Serguei Pakhomov. 2023. Trestle: Toolkit for reproducible execution of speech, text and language experiments. *AMIA Summits on Translational Science Proceedings*, 2023:360.
- Alexander V Lukashin and Apostolos P Georgopoulos. 1994. A neural network for coding of trajectories by time series of neuronal population vectors. *Neural Computation*, 6(1):19–28.
- Alexander V Lukashin, George L Wilcox, and Apostolos P Georgopoulos. 1994. Overlapping neural networks for multiple motor engrams. *Proceedings of the National Academy of Sciences*, 91(18):8651–8654.
- Saturnino Luz, Fasih Haider, Sofia de la Fuente, Davida Fromm, and Brian MacWhinney. 2020. Alzheimer’s Dementia Recognition Through Spontaneous Speech: The ADReSS Challenge. In *Proc. Interspeech 2020*, pages 2172–2176.
- Paul Michel, Omer Levy, and Graham Neubig. 2019. Are sixteen heads really better than one? In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.
- Sylvester Olubolu Orimaye, Jojo Sze-Meng Wong, and Karen Jennifer Golden. 2014. Learning predictive linguistic features for Alzheimer’s disease and related dementias using verbal utterances. In *Proceedings of the Workshop on Computational Linguistics and Clinical Psychology: From Linguistic Signal to Clinical Reality*, pages 78–87, Baltimore, Maryland, USA. Association for Computational Linguistics.
- Sai Prasanna, Anna Rogers, and Anna Rumshisky. 2020. When BERT Plays the Lottery, All Tickets Are Winning. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 3208–3229, Online. Association for Computational Linguistics.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- Victor Sanh, Thomas Wolf, and Alexander Rush. 2020. Movement pruning: Adaptive sparsity by fine-tuning. In *Advances in Neural Information Processing Systems*, volume 33, pages 20378–20389. Curran Associates, Inc.
- Paul Satz. 1993. Brain reserve capacity on symptom onset after brain injury: a formulation and review of evidence for threshold theory. *Neuropsychology*, 7(3):273.
- Nikolaos Scarmeas and Yaakov Stern. 2003. Cognitive reserve and lifestyle. *Journal of clinical and experimental neuropsychology*, 25(5):625–633.
- Nikolaos Scarmeas and Yaakov Stern. 2004. Cognitive reserve: implications for diagnosis and prevention of alzheimer’s disease. *Current neurology and neuroscience reports*, 4:374–380.
- Martin Schrimpf, Idan Asher Blank, Greta Tuckute, Carina Kauf, Eghbal A. Hosseini, Nancy Kanwisher, Joshua B. Tenenbaum, and Evelina Fedorenko. 2021. The neural architecture of language: Integrative modeling converges on predictive processing. *Proceedings of the National Academy of Sciences*, 118(45):e2105646118.
- Mengke Shi, Gary Cheung, and Seyed Reza Shahamiri. 2023. Speech and language processing with deep learning for dementia diagnosis: A systematic review. *Psychiatry Research*, 329:115538.
- David A Snowdon, Susan J Kemper, James A Mortimer, Lydia H Greiner, David R Wekstein, and William R Markesbery. 1996. Linguistic ability in early life and cognitive function and alzheimer’s disease in late life: Findings from the nun study. *Jama*, 275(7):528–532.
- Jason Steffener and Yaakov Stern. 2012. Exploring the neural basis of cognitive reserve in aging. *Biochimica et Biophysica Acta (BBA)-Molecular Basis of Disease*, 1822(3):467–473.
- Yaakov Stern. 2002. What is cognitive reserve? theory and research application of the reserve concept. *Journal of the international neuropsychological society*, 8(3):448–460.
- Yaakov Stern. 2009. Cognitive reserve. *Neuropsychologia*, 47(10):2015–2028.
- Yaakov Stern. 2012. Cognitive reserve in ageing and alzheimer’s disease. *The Lancet Neurology*, 11(11):1006–1012.
- Laura Stokes, Helen Combes, and Graham Stokes. 2015. The dementia diagnosis: a literature review of information, understanding, and attributions. *Psychogeriatrics*, 15(3):218–225.
- Behrad TaghiBeyglou and Frank Rudzicz. 2024. Context is not key: Detecting alzheimer’s disease with both classical and transformer-based neural language models. *Natural Language Processing Journal*, 6:100046.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen,

Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. [Transformers: State-of-the-art natural language processing](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.

Kathryn M. Yorkston and David R. Beukelman. 1980. [An analysis of connected speech samples of aphasic and normal speakers](#). *Journal of Speech and Hearing Disorders*, 45(1):27–36.

## **A Appendix**

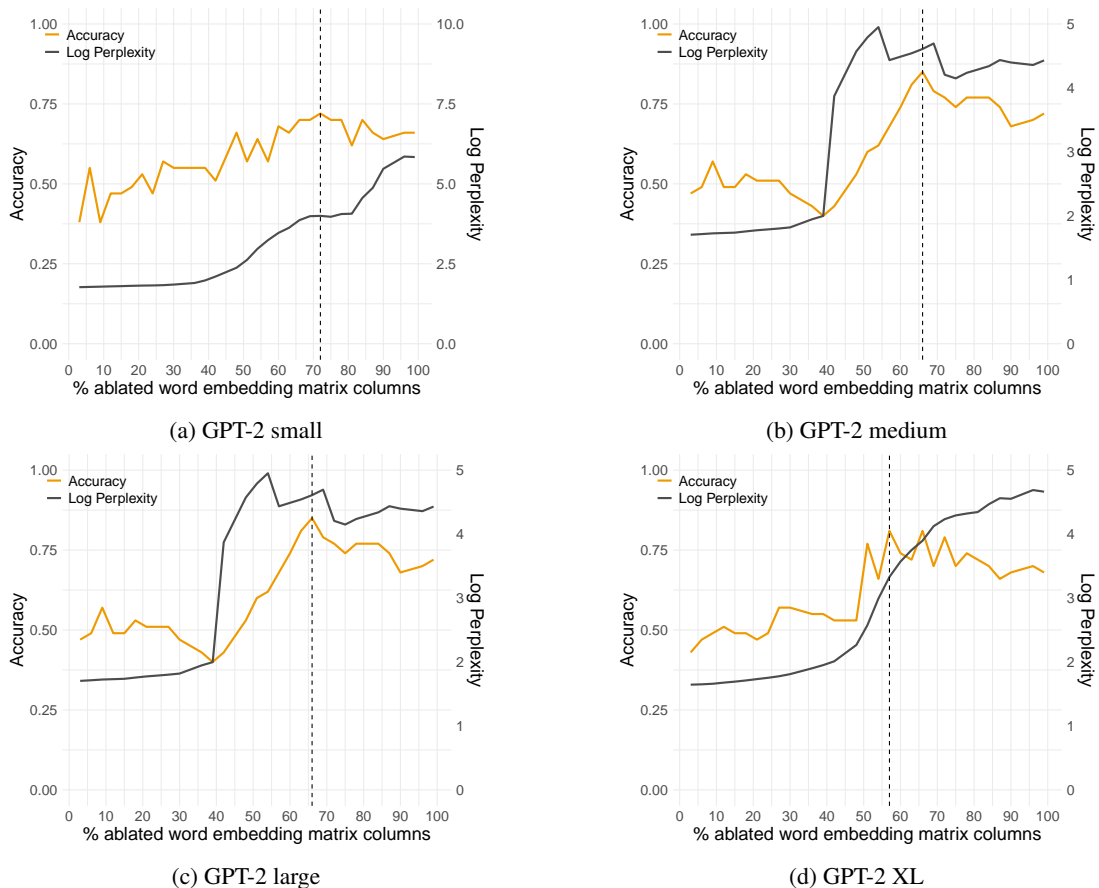


Figure 5: Comparison of GPT-2 models with masked columns of word embedding matrix on classification performance and cognitive reserve manifestation. The left y-axis denotes classification performance using both masked and unmasked GPT-2 models on the ADReSS test set. The right y-axis indicates log PPL estimated from transcripts of WLS healthy individuals. The x-axis represents the percentage of attention heads getting masked. The vertical dashed line indicates the best-performing masking pattern, achieving the highest ACC.

0	1	2	3	4	5	6	7	8	9	10	11
121	5	7	65	19	46	58	56	9	1	0	60
16	34	94	71	13	90	42	4	113	86	47	14
74	106	107	32	89	18	17	87	75	11	8	128
48	15	82	78	93	68	129	79	77	72	54	23
40	67	22	115	36	131	100	83	140	61	141	135
41	29	134	137	119	12	92	21	31	3	69	28
76	95	101	125	91	130	37	99	80	98	10	122
117	45	124	81	116	49	25	26	62	97	143	136
64	123	30	43	88	38	27	55	73	114	118	142
111	53	102	70	50	57	105	84	120	138	139	20
132	110	66	103	44	52	126	108	109	59	96	85
6	24	127	39	33	133	104	51	2	112	63	35

Table 3: The rank of importance for each attention head in the GPT-2 small model. The rows represent the layer of attention blocks in the model whereas the columns represent attention heads per layer.



