

# 脳内状態推定のための汎用言語モデル構築への取り組み

羅 桜 小林一郎  
お茶の水女子大学  
{luo.ying, koba}@is.ocha.ac.jp

## 概要

先行研究 [1] において、言語刺激下の脳内状態推定精度の向上を目的として、BERT で表現された言語特徴量と fMRI で観測された脳活動データとの対応関係をマルチモーダル深層学習の枠組みにより捉えた BrainBERT が提案されている。BrainBERT は、言語刺激と脳活動データの対応関係を捉えた被験者に依存しない汎用言語モデルであることが確認されている。しかし、BrainBERT を用いて脳活動を推定することにおいて、推定精度向上の課題がいくつか残っている。本研究では、それらの課題を克服するためモデルの改良を行い、得られた結果を客観的に検証するための実験を行った。実験結果から、新たに構築したモデルである BrainBERT 2.0 の推定精度が従来のモデルより向上したことを確認した。

## 1 はじめに

近年、深層学習を作業モデルとして用いることでヒト脳内の情報処理機構の解明を目指す研究が盛んになってきており、脳神経科学の分野において多くの新しい知見が得られている。また、言語モデルを用いた研究では、脳内の意味情報表現と実際の意味知覚の間に有意な相関があること [2] や Transformer [3] ベースの汎用言語モデルを用いて言語の特徴量を表現することで推定精度の向上が見られることが確認されている [4]。このように、言語刺激に対する脳内状態を予測することにより、様々な脳内情報処理機構が解明されてきている。そのような背景を受けて、言語刺激と脳活動状態との対応関係を捉える汎用言語モデル BrainBERT [1] が提案されている。

本研究では、言語刺激下のヒト脳内状態を予測する際に、より精度良く予測が可能になるように BrainBERT を改良し、その性能を検証することを目的とする。

## 2 BrainBERT

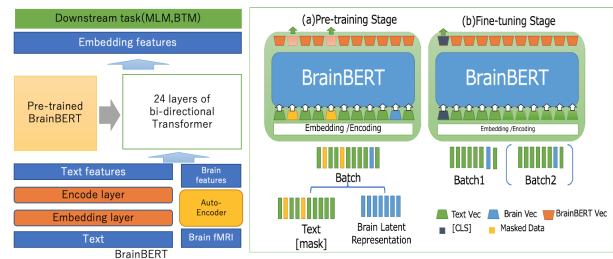


図 1 BrainBERT：脳活動データとテキストの対応関係を捉えた汎用言語モデル

BrainBERT [1] は、マルチモーダル深層学習モデルである UNITER [5] や SpeechBERT [6] を参考に、脳活動データとテキストの対応関係を捉えた汎用言語モデルとして構築された。図 1 に示すように、BrainBERT のアーキテクチャは、24 層の双方向 Transformer [3] を持つエンコーダブロックから構成される。脳活動データとして Alice Dataset [7] を用い、またプレーンテキストとして、Book コーパスと Wiki コーパスを用いて事前学習を行った。事前学習課題は、モデル学習のためにマスク言語タスク（以下、MLM タスクと略す）と脳活動データに対する対応文章の一致を捉える 2 値のマッチング自然言語処理タスク（以下、BTM タスクと略す）を用いた。

また、モデルの訓練と有効性の検証のために、[1] では 4 つの実験を段階的に実施している。

1. L1 損失による脳特徴量抽出効果検証実験：Autoencoder を使い、脳の特徴量を抽出。
2. BrainBERT 訓練での NLP タスク実験：BrainBERT の事前学習として fine-tuning の実施。
3. MLM タスク下でのアテンションの変化観察：アテンションを可視化し、言語モデルにおける脳特徴量の役割調査。
4. リッジ回帰を用いた脳活動信号予測における汎用言語モデルの比較実験：テキストによる脳内状態推定実験を行った。

実験結果から、事前学習済み BrainBERT モデルの精度は大規模汎用言語モデルには及ばずとも、脳活動状態の予測には他の言語モデルよりも優位性があることがわかった。とくに、BrainBERT のベースとなった bert-large-uncased における予測精度よりも対象としたすべての関心領域 (ROI) に対して精度が向上されていることから、脳活動データとテキストデータの対応関係が取れた汎用言語モデルが構築されたこと (以下、BrainBERT1.0 と略する) を確認した。

### 3 BrainBERT 2.0 へ向けて

マルチモーダル深層学習モデルの学習手続きを詳細に分析し、BrainBERT の 3 つの改善策を提案する。

#### 3.1 改善策 1：潜在表現の意味性

BrainBERT を学習する際、入力には ‘two people sitting on a bench watching a sailboat’ などの文と、それに対応する磁気共鳴機能画像法 (fMRI) による脳活動データを使用した。モデルの学習条件の整合性を考慮し、ボクセルの次元がおよそ 50,000~75,000 くらいの大脳皮質データを直接入力の次元とはせず、脳活動データを 1024 または 768 次元の特徴量に AutoEncoder を用いて変更した。AutoEncoder モデルの損失関数は SpeechBERT [6] を参考にし、式 1 を採用した。

$$L_{recon+MAE} = \frac{1}{N} \sum_{n=1}^N \|x_n - x'_n\|^2 + \frac{1}{N} \sum_{n=1}^N \|z_n - y_{tn}\| \quad (1)$$

式 1 中、L2 損失は入力と出力間の再構築損失、L1 損失は潜在表現ベクトルとテキストベクトル間の類似度損失である。テキストのベクトル化には、huggingface で提供されている事前学習済み bert-large-uncased モデル<sup>1)</sup>を用いている。この学習において、L1+L2 から得られる損失の合計が小さいほど良いということであるが、学習済みの bert-large-uncased モデルから得られるテキストベクトルをファインチューニングなど無しにそのまま該当テキストのベクトルであると仮定している点に問題がある。さらには、bert-large-uncased モデルから得られる CLS トークンは、エンコーダを通していないため実際にはランダムな埋め込みベクトルになり、文脈の関係性など、文の特徴をうまく表現することができないと考えられる。また、SimCSE モデル [8] は、単純な対比学習によって文の埋め込みを行

い、より質の高い文ベクトルを生成することができる。また、SentenceBERT[9] は、転移学習を用いて BERT モデルが CLS トークンを文の埋め込み表現として使うことの問題を改善している。このことから、SentenceBERT と SimCSE の拡張モデルから 6 つの派生バージョン (モデル名は表 1 の行に記載) を選択し、新たな比較実験を行った。単一変数規則を用いて、BrainBERT モデルを AutoEncoder で、上記実験の 1 段階目から再訓練した。学習結果を比較し、6 つのモデルの中からテキストのベクトル化において脳活動データと対応させるための最適モデルを見出した。

#### 3.2 改善策 2：単変数実験の制御

BrainBERT 構築の実験 4 段階目において、リッジ回帰モデルを用いて BrainBERT によるテキストの特徴量から脳活動状態を予測した際の精度を、他の 20 の汎用言語モデルを用いた際の予測精度と比較し、単一変数規則における BrainBERT の性能の検証を行った。実験の評価データには、訓練で用いていない全く新しい脳データセットである BOLD5000 [10] を使用した。このデータセットの収集実験では、3 つの大規模画像データセットの画像を用いた画像刺激を用い、脳活動の酸素飽和濃度 (BOLD) 信号を収集しており、視覚に関連する 5 つの脳部位のデータが公開されている。大規模画像データセットの 1 つである MSCoCo2014 には、画像のキャプションが付属しており、実験では脳活動データとのペアデータとして使用した。しかし、これまで BrainBERT への入力は他のモデルと異なり、脳特徴量も入力されてしまっていた。脳のデータをモデル計算の入力として誤って取り込んでいたため、客観性に欠ける結果となった。これは、単一変数の原則に沿わないものになっており、改善対象となった。このことから、4 段階目の実験を再度実施した。予測モデルとなるリッジ回帰モデル構築の際には、4-folds クロスバリデーションを使用し、20 のモデルが同じ入力変数に基づき、同じ正則化項を使用して比較された。

#### 3.3 改善策 3：公知に反する結論の妥当性検証

また、先行研究 [1] における実験 4 段階目においても、改善策 2 の 1 点を除き、もう一つの問題として、脳状態予測の結果から、roberta-large モデルの結果が roberta-base モデルの結果よりも推定精度が悪く、多くの先行研究における殆どの自然言語処理タスク

1) <https://huggingface.co/bert-base-uncased>

において、通常は roberta-large モデルが roberta-base モデルを上回った<sup>2)</sup>という実験結果と矛盾していることが挙げられる。この点において、結果の客観性・独立性を証明するために再度実験結果を精査をする必要があり、順列検定実験を行う。

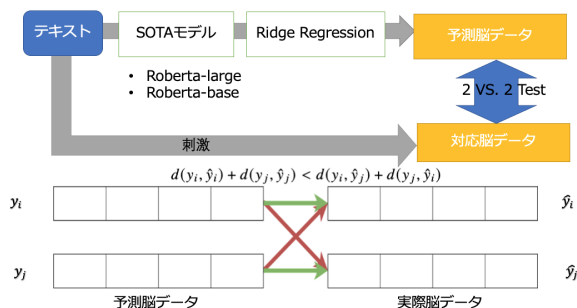


図2 2 VS. 2 検定手法流れ

この順列検定実験では、Foster ら [11] の論文で紹介されている 2 VS. 2 検定手法を用いた。

図2のように、2 VS. 2 検定は、ベクトルの比較を2値分類タスクに還元する相関テストである。このテストは、2つのデータペアが相関していることを証明するツールとして使用することができる。2 VS. 2 検定では、正しいペアの距離の和（図2の青い矢印）と、正しくないペアの距離の和（図2の赤い矢印）の比較を、以下の式で評価することにより行う。

$$d(y_i, \hat{y}_i) + d(y_j, \hat{y}_j) < d(y_i, \hat{y}_j) + d(y_j, \hat{y}_i) \quad (2)$$

式(2)が成り立つ場合、比較は成功したとみなされる。2 VS. 2 検定の精度は、2 VS. 2 の比較成功の割合となる。roberta-large モデルと roberta-base モデルについて別々に計算を行い、二つのモデルのうちどちらが正しいかを検証した。

## 4 実験

前章では、[1]における3つの問題点を説明し、それぞれの改善策を提案した。この章では、改善策に取り組む実験と、実験の具体的なパラメータについて説明する。

### 4.1 潜在表現の意味性検証実験

テキストの意味と脳活動データとのより良い対応関係を捉えるために、テキスト(文)の意味を sentenceBERT モデルおよび SimCSE モデルなどの合計6つのモデルを用い、実験を通じてそれらのモ

2) 注：本研究で使用「state-of-the-art」状態の言語モデルは、すべて hugging face(<https://huggingface.co/>) からダウンロードしたモデルを使用している。

デルを比較し、良い対応関係がとれるモデルを検証した。検証のための学習過程は以下の通りである。脳特徴量の抽出に AutoEncoder モデルを採用し、Pereira データセット [12] を学習データセットとして利用する。AutoEncoder の構造は、エンコーダ部が2つの Linear 層と2つの Batch Normalization 層からなり、エンコーダとデコーダは対称構造となる。また、ドロップアウト率は0.2とし、モデルの学習率は  $1e-5$ 、勾配法として AdamW を用い、バッチサイズは32、エポック数は75とした。

モデルの学習を容易にするため、異なる被験者に対して共通したボクセルサイズの46,840次元として大脳皮質のデータを抽出し、さらに中間変数のサイズを調整した。例えば、bert-large-uncased モデルの出力ベクトルの次元が1024の場合、中間変数の次元も1024に設定する。表1の1行目と2行目から見ると、今回の新たに条件を変更して再学習された BrainBERT の事前学習結果では、先行研究 [1] から訓練した精度15.7%の BrainBERT1.0 から258%向上し、最大55.56%の精度を達成した。その結果から、潜在表現の意味が最も優れているのは、bert-large-uncased モデルであることが確認された。

表1 改善策1に対する訓練結果

model type	Latent representation size	Autoencoder		BrainBERT		
		Loss	2VS2 Test(%)	MLM text Acc(%)	MLM Mix Acc(%)	
BERT	bert-large-uncased(前回)	1024	0.089	52.32	35.4	15.79
	bert-large-uncased(今回)	1024	<b>0.061</b>	<b>53.81</b>	10.08	<b>55.56</b>
SentenceBERT	sentence-camembert-large	1024	0.1074	52.42	10.08	20.14
SimCSE	unsup-simcse-roberta-large	1024	0.072	51.53	10.08	20.14
	sup-simcse-roberta-large	1024	0.078	51.35	10.08	20.35
	sup-simcse-bert-base-uncased	768	0.076	52.59	10.08	20.54
	unsup-simcse-bert-large-uncased	1024	0.1	50.73	-	-
	sup-simcse-bert-large-uncased	1024	0.09	49.83	-	-

### 4.2 単制御変量実験

24の汎用言語モデルをテキストの特徴量表現に用い、リッジ回帰を使って脳活動状態を予測した。予測精度には Pearson 相関係数を使用した。本実験で使用する BrainBERT は、先行研究 [1] での実験設定を一部踏襲しつつ新しく改善された設定の元、再学習されたモデルであり、以降、それを BrainBERT 2.0 と呼ぶ。比較対象となる他の単語埋め込みでは GloVe [13]、word2vec [14] の2つのモデルを採用し、汎用言語モデルでは、BERT [15]、RoBERTa [16]、ALBERT [17]、GPT [18] などを採用し、合計24の事前学習済み汎用言語モデルを使用した。符号化モデルを構築するためのデータセットとして、BOLD5000の一部とそれに対応する MSCoCo2014 データセットを使用した。脳活動データは4人の被験者のも



表 2 各相関脳領域の Pearson 相関係数の計算結果 (略表)

Models/ROIs	脳関心領域 (ROIs)					Average	検証結果 2 VS. 2 Test PC ACC
	PPA	OPA	EARLYVIS	RSC	LOC		
roberta-base	3.89	17.71	27	15.43	26.43	18.09	32.09
roberta-large	4.16	20.99	30.78	14.88	31.36	20.43	34.25
albert-xlarge-v1	50.03	<b>51.14</b>	42.55	50.17	46.19	48.02	48.81
albert-xlarge-v2	42.54	47.04	<b>53.33</b>	45.7	51.54	48.03	47.2
sentence-camembert-large	3.19	16.71	28.51	17.03	25.48	18.18	32.43
sup-simcse-bert-base-uncased	2.18	13.92	26.9	12.96	24.78	16.15	31.56
unsup-simcse-roberta-large	2.59	14.32	27.47	13.57	28.42	17.27	31.77
sup-simcse-roberta-large	2.64	13.41	26.6	13.46	25	16.22	30.95
BrainBERT1.0	49.44	44.86	51.29	44.51	48.9	47.80	50.29
BrainBERT2.0	<b>50.58</b>	47.12	50.8	<b>51.54</b>	<b>51.71</b>	<b>50.35</b>	<b>51.28</b>

のが使用され、合計 7,677 の画像刺激が与えられた際の脳活動データ (train, test, validation の 3 グループに分け、8:1:1 の割合) が用いられた。提供されている脳活動データは、PPA (Parahippocampal Place Area) など視覚に関連する 5 つの脳領域が対象となり、合計 3,566 のボリュームが処理された。比較のための予測精度として、多重補正検定に基づく Pearson 相関係数 ( $p < 0.05$ ) を計算した。表 2<sup>3)</sup> より、BrainBERT 1.0, BrainBERT 2.0 とともにテキストから脳活動を予測するタスクにおいて、他の汎用言語モデルと比べて良好な結果を示し、バージョン 2.0 は 1.0 より 5.58% 向上し、全モデルの中で最も良い精度を示した。

### 4.3 RoBERTa サイズ妥当性検証実験

RoBERTa のサイズにおける予測精度結果の客観性については、直接的に結論を導き出す術がないため、検定による実験を実施した。実験の詳細を以下に示す。MSCoCo2014 データセットを用いて RoBERTa への入力テキストとした。“leave 2 out”手法に基づき、770 文章のデータからランダムに 2 ペアを抽出し、合計 296,065 訓練データを構成した。予測脳データと実際脳データの Pearson 相関係数を計算し精度を求めた。RoBERTa-base の精度は 32.09%、RoBERTa-large の精度は 34.24% となった。この結果は、表中の各脳機能領域下でも一貫性を示し、先行

3) 注：表 2 は結果の一部を掲載した表で、主要なモデルの結果のみを掲載している。詳細については付録 (Appendix) ?? を参照。

研究 [1] の結果を是正する結果となった。これにより、RoBERTa-large が RoBERTa-base モデルよりも、脳内状態を予測する上で優れていることが確認された。

## 5 考察

実験で取り上げた汎用言語モデルの内、albert-xlarge 群は、albert モデルだけでなく、全体のモデルの中でも優れていた。本研究において採用した BERT を albert-xlarge に変更し、新たに学習しなおした脳特徴量を反映する汎用言語モデルによって、BrainBERT 2.0 より良い結果を得られるかどうか、さらに検証する価値があると考えられる。

## 6 おわりに

本研究は、先行研究 [1] での取り組みについて 3 つの大きな改善点を示し、実験を通じて改善点の効果を検証した。2 VS. 2 検定手法を導入し、テキストの意味と脳活動データの対応関係を捉える上で bert-large-uncased が有利であることを確認した。また、先行研究の結果において未解明であった RoBERTa のサイズによる精度の問題も検証実験を通じて疑問を解消することができた。さらに、改良されたモデルである BrainBERT 2.0 は、事前学習結果の精度を向上させ、脳内状態予測において他のモデルを上回る優位性を実現したことを確認した。

今後の課題として、albert を含む様々な汎用言語モデルを用いて再度実験を行うつもりである。

## 謝辞

本研究は、科研費（18H05521）の支援を受けた。  
ここに謝意を表す。

## 参考文献

- [1] 羅桜, 小林一郎. Brainbert: 脳活動とテキストの対応関係を捉えた言語モデル構築への取り組み. 日本知能情報ファジィ学会 ファジィ システム シンポジウム 講演論文集, Vol. 37, pp. 376–381, 2021.
- [2] Alexander G. Huth, Wendy A. de Heer, Thomas L. Griffiths, Frédéric E. Theunissen, and Jack L. Gallant. Natural speech reveals the semantic maps that tile human cerebral cortex. **Nat.**, Vol. 532, No. 7600, pp. 453–458, 2016.
- [3] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In **Proceedings of the 31st International Conference on Neural Information Processing Systems, NIPS'17**, p. 6000–6010, Red Hook, NY, USA, 2017. Curran Associates Inc.
- [4] Satoshi Nishida, Antoine Blanc, Naoya Maeda, Masataka Kado, and Shinji Nishimoto. Behavioral correlates of cortical semantic representations modeled by word vectors. **PLoS Computational Biology**, Vol. 17, No. 6, p. e1009138, 2021.
- [5] Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. Uniter: Universal image-text representation learning. In **Computer Vision – ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXX**, p. 104–120, Berlin, Heidelberg, 2020. Springer-Verlag.
- [6] Yung-Sung Chuang, Chi-Liang Liu, Hung-yi Lee, and Lin-Shan Lee. Speechbert: An audio-and-text jointly learned language model for end-to-end spoken question answering. In Helen Meng, Bo Xu, and Thomas Fang Zheng, editors, **Interspeech 2020, 21st Annual Conference of the International Speech Communication Association, Virtual Event, Shanghai, China, 25-29 October 2020**, pp. 4168–4172. ISCA, 2020.
- [7] Shohini Bhattasali, Jonathan Brennan, Wen-Ming Luh, Berta Franzluebbers, and John Hale. The alice datasets: fMRI & EEG observations of natural language comprehension. In **Proceedings of the 12th Language Resources and Evaluation Conference**, pp. 120–125, Marseille, France, May 2020. European Language Resources Association.
- [8] Tianyu Gao, Xingcheng Yao, and Danqi Chen. SimCSE: Simple contrastive learning of sentence embeddings. In **Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing**, pp. 6894–6910, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics.
- [9] Iryna Gurevych Nils Reimers. Sentence-bert: Sentence embeddings using siamese bert-networks. <https://arxiv.org/abs/1908.10084>, 2019.
- [10] Nadine Chang, John A. Pyles, Austin Marcus, Abhinav Gupta, Michael J. Tarr, and Elissa M. Aminoff. Bold5000, a public fmri dataset while viewing 5000 visual images. **Scientific Data**, Vol. 6, No. 1, May 2019.
- [11] Chris Foster, Dhanush Dharmaretnam, Haoyan Xu, Alona Fyshe, and George Tzanetakis. Decoding music in the human brain using eeg data. In **2018 IEEE 20th International Workshop on Multimedia Signal Processing (MMSP)**, pp. 1–6, 2018.
- [12] Francisco Pereira, Bin Lou, Brianna Pritchett, Samuel Ritter, Samuel Gershman, Nancy Kanwisher, Matthew Botvinick, and Evelina Fedorenko. Toward a universal decoder of linguistic meaning from brain activation. **Nature Communications**, Vol. 9, , 03 2018.
- [13] Jeffrey Pennington, Richard Socher, and Christopher Manning. GloVe: Global vectors for word representation. In **Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)**, pp. 1532–1543, Doha, Qatar, October 2014. Association for Computational Linguistics.
- [14] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space, 2013.
- [15] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding, 2019.
- [16] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach, 2019.
- [17] Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. Albert: A lite bert for self-supervised learning of language representations, 2020.
- [18] Xianrui Zheng, Chao Zhang, and Philip C. Woodland. Adapting gpt, gpt-2 and bert language models for speech recognition, 2021.

## 付録 (Appendix)

表 3 各相関脳領域の Pearson 相関係数の計算結果 (全表)

Models/ROIs	脳関心領域 (ROIs)					検証結果	
	PPA	OPA	EARLYVIS	RSC	LOC	Average	2 VS. 2 Test PC ACC
GloVe	2.92%	18.31%	27.80%	15.11%	28.10%	18.45%	32.99%
word2vec	12.90%	22.29%	37.55%	36.60%	37.10%	29.29%	40.70%
bert-base-uncased	3.61%	22.29%	31.42%	20.06%	31.34%	21.74%	34.83%
bert-large-uncased	22.22%	39.30%	47.55%	47.17%	46.32%	40.51%	42.48%
bert-base-multilingual-cased	4.70%	28.62%	35.60%	19.82%	30.23%	23.79%	34.67%
bert-large-uncased-whole-word-masking	4.40%	22.75%	31.75%	20.84%	31.08%	22.16%	34.38%
roberta-base	3.89%	17.71%	27.00%	15.43%	26.43%	18.09%	32.09%
roberta-large	4.16%	20.99%	30.78%	14.88%	31.36%	20.43%	34.25%
albert-base-v1	6.36%	29.29%	35.31%	24.27%	34.36%	25.92%	35.33%
albert-large-v1	12.60%	35.75%	39.33%	31.90%	40.37%	31.99%	37.99%
albert-xlarge-v1	50.03%	<b>51.14%</b>	42.55%	50.17%	46.19%	48.02%	48.81%
albert-xxlarge-v1	26.55%	44.55%	46.16%	43.90%	48.37%	41.91%	42.74%
albert-base-v2	9.14%	36.18%	38.85%	27.57%	37.35%	29.82%	36.82%
albert-large-v2	17.60%	38.84%	41.30%	28.14%	39.72%	33.12%	39.79%
albert-xlarge-v2	42.54%	47.04%	<b>53.33%</b>	45.70%	51.54%	48.03%	47.20%
albert-xxlarge-v2	36.39%	49.66%	45.95%	45.61%	45.17%	44.55%	43.86%
gpt2	44.45%	46.21%	52.69%	45.52%	49.23%	47.62%	59.49%
gpt2-medium	45.19%	47.84%	51.05%	43.53%	49.84%	47.49%	58.69%
gpt2-large	42.34%	46.79%	52.71%	44.38%	49.17%	47.08%	58.67%
gpt2-xl	43.25%	46.63%	52.65%	44.45%	49.11%	47.22%	58.28%
sentence-camembert-large	3.19%	16.71%	28.51%	17.03%	25.48%	18.18%	32.43%
sup-simcse-bert-base-uncased	2.18%	13.92%	26.90%	12.96%	24.78%	16.15%	31.56%
unsup-simcse-roberta-large	2.59%	14.32%	27.47%	13.57%	28.42%	17.27%	31.77%
sup-simcse-roberta-large	2.64%	13.41%	26.60%	13.46%	25.00%	16.22%	30.95%
BrainBERT1.0	49.44%	44.86%	51.29%	44.51%	48.90%	47.80%	50.29%
BrainBERT2.0	<b>50.58%</b>	47.12%	50.80%	<b>51.54%</b>	<b>51.71%</b>	<b>50.35%</b>	<b>51.28%</b>