

# Foundational GPT Model for MEG

**Richard Csaky<sup>\*,1,2,5</sup>, Mats W.J. van Es<sup>1,2</sup>, Oiwi Parker Jones<sup>2,3,4</sup>, and Mark Woolrich<sup>1,2</sup>**

<sup>1</sup>Oxford Centre for Human Brain Activity, Department of Psychiatry, University of Oxford, OX3 7JX, Oxford, UK

<sup>2</sup>Wellcome Centre for Integrative Neuroimaging, OX3 9DU, Oxford, UK

<sup>3</sup>Department of Engineering Science, University of Oxford, OX1 3PJ, Oxford, UK

<sup>4</sup>Jesus College, OX1 3DW, Oxford, UK

<sup>5</sup>Christ Church, OX1 1DP, Oxford, UK

{richard.csaky@psych, mats.vanes@psych,  
oiwi.parkerjones@eng, mark.woolrich@ohba}.ox.ac.uk

## Abstract

Deep learning techniques can be used to first training unsupervised models on large amounts of unlabelled data, before fine-tuning the models on specific tasks. This approach has seen massive success for various kinds of data, e.g. images, language, audio, and holds the promise of improving performance in various downstream tasks (e.g. encoding or decoding brain data). However, there has been limited progress taking this approach for modelling brain signals, such as Magneto-/electroencephalography (M/EEG). Here we propose two classes of deep learning foundational models that can be trained using forecasting of unlabelled MEG. First, we consider a modified Wavenet; and second, we consider a modified Transformer-based (GPT2) model. The modified GPT2 includes a novel application of tokenisation and embedding methods, allowing a model developed initially for the discrete domain of language to be applied to continuous multichannel time series data. We also extend the forecasting framework to include condition labels as inputs, enabling better modelling (encoding) of task data. We compare the performance of these deep learning models with standard linear autoregressive (AR) modelling on MEG data. This shows that GPT2-based models provide better modelling capabilities than Wavenet and linear AR models, by better reproducing the temporal, spatial and spectral characteristics of real data and evoked activity in task data. We show how the GPT2 model scales well to multiple subjects, while adapting its model to each subject through subject embedding. Finally, we show how such a model can be useful in downstream decoding tasks through data simulation. All code is available on GitHub<sup>1</sup>.

## 1 Introduction

Unsupervised learning provides a dataset-agnostic method for learning shared representations. Unsupervised learning techniques can be further differentiated between those that aim to learn interpretable representations and those that are purely data-driven. In functional neuroimaging, interpretable models can provide neuroscientific insights that are especially useful for rest or spontaneous data, where there is no known external stimuli or behaviour

---

<sup>\*</sup>Corresponding author.

<sup>1</sup><https://github.com/ricsinaruto/MEG-transfer-decoding>

associated with the brain activity (Baker et al., 2014). Models designed without focusing on interpretability can be used to generalise over multiple heterogeneous datasets and provide a pretrained foundation model (Yuan et al., 2024). Such foundation models can then be fine-tuned for downstream tasks on smaller amounts of data for which relevant labels are available (e.g. to do encoding or decoding). By leveraging large amounts of data, the hope is that the foundation model will be capable of generalising to new data types and provide improvement over just training the model on the labelled dataset. This is especially useful for brain-computer interface (BCI) applications.

The concept of using vast amounts of data to boost performance in downstream tasks originates from deep learning. Perhaps the most successful recent example is that of large language models, trained on diverse data sources and demonstrating enhanced capabilities over task-specific models in a multitude of language-related tasks (e.g., translation, summarisation) (Brown et al., 2020). This can also be viewed as a form of transfer learning. Zero-shot performance is obtained when no fine-tuning is done for the downstream task. Several factors enabled the success of large language models, including data scale, model size, fast GPUs, and effective neural network architectures (Kaplan et al., 2020; Fedus et al., 2022; Sutton, 2019). To adopt this paradigm for electrophysiology data, the primary obstacles are the model architecture and data size. In this paper, we focus on the former.

We aim to design general models well-suited to multichannel timeseries that can scale effectively. We also focus on foundational models that can be trained using forecasting (as opposed to masked prediction), which are causal and can generate data recursively. This allows for the interrogation of learned spatio-temporal dynamics over long temporal horizons. Self-supervised learning (SSL) has emerged as a promising approach for learning useful representations from unlabelled electrophysiological data. SSL reformulates an unsupervised learning problem as a supervised one by exploiting inherent structure in the data to generate "pseudo-labels". Forecasting is one such SLL task, as is reconstruction of masked sequence segments (Wang et al., 2023). In the context of electrophysiology, recent works have proposed SSL tasks tailored to the temporal and multivariate nature of neural time series data (Gohil et al., 2022; Banville et al., 2021; Kostas et al., 2021; Wang et al., 2023; Cui et al., 2023).

In our quest for designing expressive foundational models of MEG data, we can look to artificial intelligence domains with similar characteristics, such as audio or natural language processing. These domains share some similarities with MEG data, like the sequential nature of the modality. However, while audio data is also a continuous timeseries, it only contains a single channel and comes at a much higher sampling rate compared to M/EEG data. Language data is perhaps even more different as its timeseries are comprised of distinct units (words) from a finite vocabulary set. As such, starting with models developed for these domains and adapting them to handle the nuances of M/EEG data is a promising approach. In this work, we adapt two such approaches. First, we adapt Wavenet, originally developed for forecasting audio data (van den Oord et al., 2016); and second, we adapt a Transformer architecture (Vaswani et al., 2017) in the form of GPT2, originally developed for forecasting language (Radford et al., 2019).

Wavenet has been shown to be an effective model for forecasting time series, through its dilated convolutional architecture, which is fast and parameter-efficient (van den Oord et al., 2016). Here, we modify Wavenet to handle multichannel MEG data by first doing a channel-dependent quantisation (WavenetFullChannel) and then introducing a mixing layer

across the channel dimension (`WavenetFullChannelMix`).

In recent years, the Transformer architecture has driven a second deep learning revolution (Vaswani et al., 2017). Transformers use attention to model complex dependencies in long sequences, providing a more flexible inductive bias well-suited to language modelling and other tasks involving highly structured sequential data (Devlin et al., 2019; Brown et al., 2020). For example, self-attention provides direct connectivity between any two time steps, capturing long-range dependencies. The parallelisable architecture allows for more efficient computation compared to recurrent models. Transformers have also been applied to time series data with promising results (Wen et al., 2022). Zhou et al. (2021) adapted the self-attention mechanism for long-range forecasting and demonstrated state-of-the-art performance on multiple public datasets.

Here, we modify the Transformer architecture used in GPT2 (Radford et al., 2019) to handle multichannel MEG data. Specifically, we do a channel-independent quantisation, but augment the input with channel-embeddings to signify to the model which timeseries corresponds to which channel. We refer to this model as `ChannelGPT2`. While previous work has employed various modifications to the Transformer architecture we posit that keeping the categorical nature of the sequences and doing next time-step forecasting are essential for effective modelling. Our methods are further detailed in Section 3.

Both the modified Wavenet and `ChannelGPT2` models also include subject embeddings (Csaky et al., 2023b) and task condition embeddings. The latter provides the model with information about external stimuli synchronised to the MEG time-course. This means that the model can be treated as both a forecasting and an encoding model.

We evaluate how well the modified Wavenet and `ChannelGPT2` models can perform as foundational models. First, we assess how well the fitted models can recursively generate data with the same spatial, temporal and spectral characteristics as real MEG data. We evaluate this both at the subject- and group-level. Second, we see how well the foundation models can be used on downstream tasks, through data simulation. Finally, we show through several ablation studies how the best performing model, `ChannelGPT2`, leverages condition and channel embeddings.

## 2 Results

### 2.1 Next time-step prediction does not capture performance

We wanted to evaluate how well the modified Wavenet and `ChannelGPT2` models can perform as foundational models. One way to assess this is through the forecasting performance, i.e. the prediction accuracy of the label at the next time point. Here, we evaluated the performance of our two modified versions of Wavenet (`WavenetFullChannel` and `WavenetFullChannelMix`) alongside `ChannelGPT2`, as summarised in Table 1. For comparison, we also evaluated the performance of a linear autoregressive (AR) model of order 255. All deep learning models were trained on tokenised and embedded inputs, and cross-entropy loss was used to predict categorical probability distributions over tokens. AR models were trained on the continuous data using the mean-squared error loss, and we simply binned the predicted continuous output to compute accuracy and compare with other models.

The next-timestep forecasting accuracy for different models on a sample subject is shown in

Model	Univariate	Tokenised	Linear
AR(255)	yes	no	yes
WavenetFullChannel	yes	yes	no
WavenetFullChannelMix	no	yes	no
ChannelGPT2	yes	yes	no

Table 1: The 4 main models presented in the results section. Univariate means that the channel dimension is treated as a batch dimension during training. Non-tokenised models predict continuous values directly instead of categorical distributions.

Model	MSE	Top-1 Accuracy	Top-5 Accuracy
Repeat baseline	0.024	1.5	7.6
AR(255)	0.016	1.5	7.5
WavenetFullChannel	0.026	2.0	9.8
WavenetFullChannelMix	0.022	2.2	10.8
ChannelGPT2	0.023	2.2	10.9

Table 2: Test data next-timestep prediction performance across various models. Accuracy values are given in percentages. Chance-level is 1/256, however predicting the majority class (quantised value) is somewhat higher, since the distribution over bins is not uniform.

Tabel 2. Beyond standard accuracy (the number of true positives divided by the number of all examples), we also evaluated top-5 accuracy, counting a prediction as correct if the true bin was within the 5 most probable bins. Surprisingly, all models performed only slightly better than a naive baseline of repeating the previous timestep’s value.

As expected, the linear AR model had lower MSE but worse accuracy than the nonlinear models. This can be because MSE measures the distance of the prediction to the target, while accuracy is only 1 if the prediction is in the target bin. Thus, it may be that the AR model always predicts values that are slightly closer to the target, but never quite falling in the target bin. While WavenetFullChannel appears to be worse, WavenetFullChannelMix and ChannelGPT2 have nearly identical performance.

The choice of sampling rate can affect forecasting performance. A higher sampling rate makes the task easier as consecutive timesteps are more correlated, however this might make the model focus on very short-range temporal dependencies and overfit to noise. We analysed forecasting performance in relation to sampling rate in Supplementary Section A.2.

## 2.2 PSD of generated data matches real data

We have seen how forecasting accuracy does not differentiate well between the candidate foundation models. Despite this, the models may perform differently when evaluated using other criteria. For example, a good foundational model should be expected to be able to recursively generate data that looks like the real data. Here, we first assess the models’ ability to do this using the power spectra.

For deep learning models we used top-p sampling with  $p = 80\%$  (unless otherwise noted in the figure caption) to recursively generate data. We generated 3600 seconds with all models. For models that have task-conditioning (all except AR(255)) we use the task label timeseries

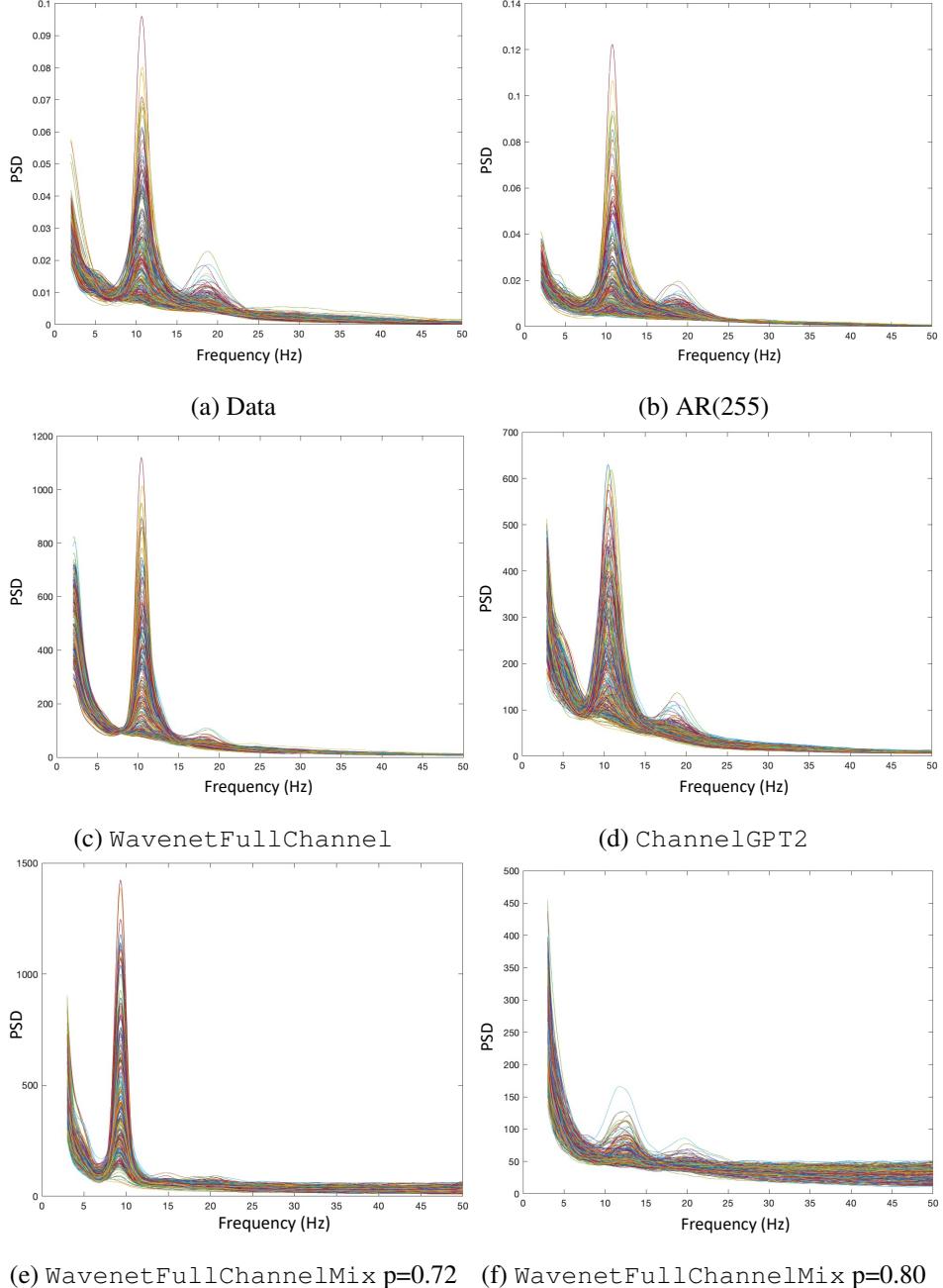


Figure 1: Comparison of generated data power spectral density (PSD) across models. Each line represents a different MEG channel. Note that in e) and f) the value of  $p$  corresponds to the probability mass used for top- $p$  sampling.

from the training set. The models in this section were trained on a single sample subject.

Generated token sequences are first de-tokenised and then the power spectral density (PSD) is computed on the continuous data. Figure 1 compares the PSD of the generated data across our models. Qualitatively, it is clear that AR(255) reproduces PSDs that match best with those computed directly on the MEG data, while WavenetFullChannel and ChannelGPT2 are not far behind. All models capture the characteristic  $1/f$  shape, and peaks at 10 and 19

Hz, likely related to alpha and beta band activity. Notably, `WavenetFullChannel` has reduced power at the 19 Hz peak, which could indicate issues in capturing higher frequency dynamics.

The performance of `WavenetFullChannelMix` was found to depend on the choice of the probability mass  $p$ , used in top- $p$  sampling from the model to generate the data.  $p$  prohibits the generation of any tokens that are not within the top- $p$  cumulative probability mass. Even slight modifications (e.g.,  $p = 0.72$  vs.  $p = 0.8$  for `WavenetFullChannelMix`) result in large differences in the frequency of the two peaks in the PSDs computed on the generated data, and also the width of the peaks. Ultimately both top- $p$  values provide subpar PSDs compared to channel-independent models. This is likely to be due to overfitting, as channel-mixing models lack the implicit regularisation of modelling each channel separately. The implicit regularisation is a consequence of having much more examples and much fewer input features when treating channels as a batch dimension.

### 2.3 HMM statistics of generated data match real data

Next, we assessed how well the fitted models can recursively generate data with the same spatial, temporal and spectral multi-channel characteristics as real MEG data. HMMs are an established way for doing unsupervised discovery of multi-channel dynamics in real neuroimaging data, and have been used to characterise the spatial, temporal and spectral characteristics of brain networks in MEG data (Rabiner, 1989; Vidaurre et al., 2018a).

Separately, we trained each model on a single sample subject (with condition embeddings), generated multi-channel data, and then inferred a 12-state HMM, with each state modelling the multi-channel data as a multi-variate Normal (MVN) distribution. The number of states was chosen based on previous work (Vidaurre et al., 2018a). Example state time-courses generated from all models are plotted in Figure 2 to qualitatively illustrate the differences in the generated dynamics. Note that since a separate HMM is trained each time, the states are not matched between models or with the real data. States could be matched post-HMM training by looking for similarity in the observation models, but here we opted to plot unmatched state time-courses. Since each generated timeseries is different, comparison across state timecourses would not be useful anyway.

We extracted four summary statistics on the inferred state timecourses and compared their distribution over states. These statistics are shown across models in Figure 3, alongside those for an HMM trained on the real multi-channel MEG data. Across the four summary statistics we can see that the real data has high variance in the distribution over states. AR(255) and `WavenetFullChannelMix` fail to produce data with variable state statistics, and even the mean over states is not captured well. `WavenetFullChannel` does a great job at capturing the mean of the state distributions, but still produces data with relatively invariant states. `ChannelGPT2` seems to best capture the distributions across all four statistics, especially for the mean interval and switching rate. This shows that Transformer-based models can generate data that better matches the HMM-inferred dynamics of real MEG data.

In addition to state statistics, we also computed the power spectra of each state across the timeseries. In MEG data different states might capture oscillatory activity with specific frequencies. The extracted power spectra from the inferred state time courses is shown in Figure 4. We can see that the HMM trained on the MEG data contains many states that capture the 10 Hz peak, with fewer states having a 20 Hz peak. It is also clear that the

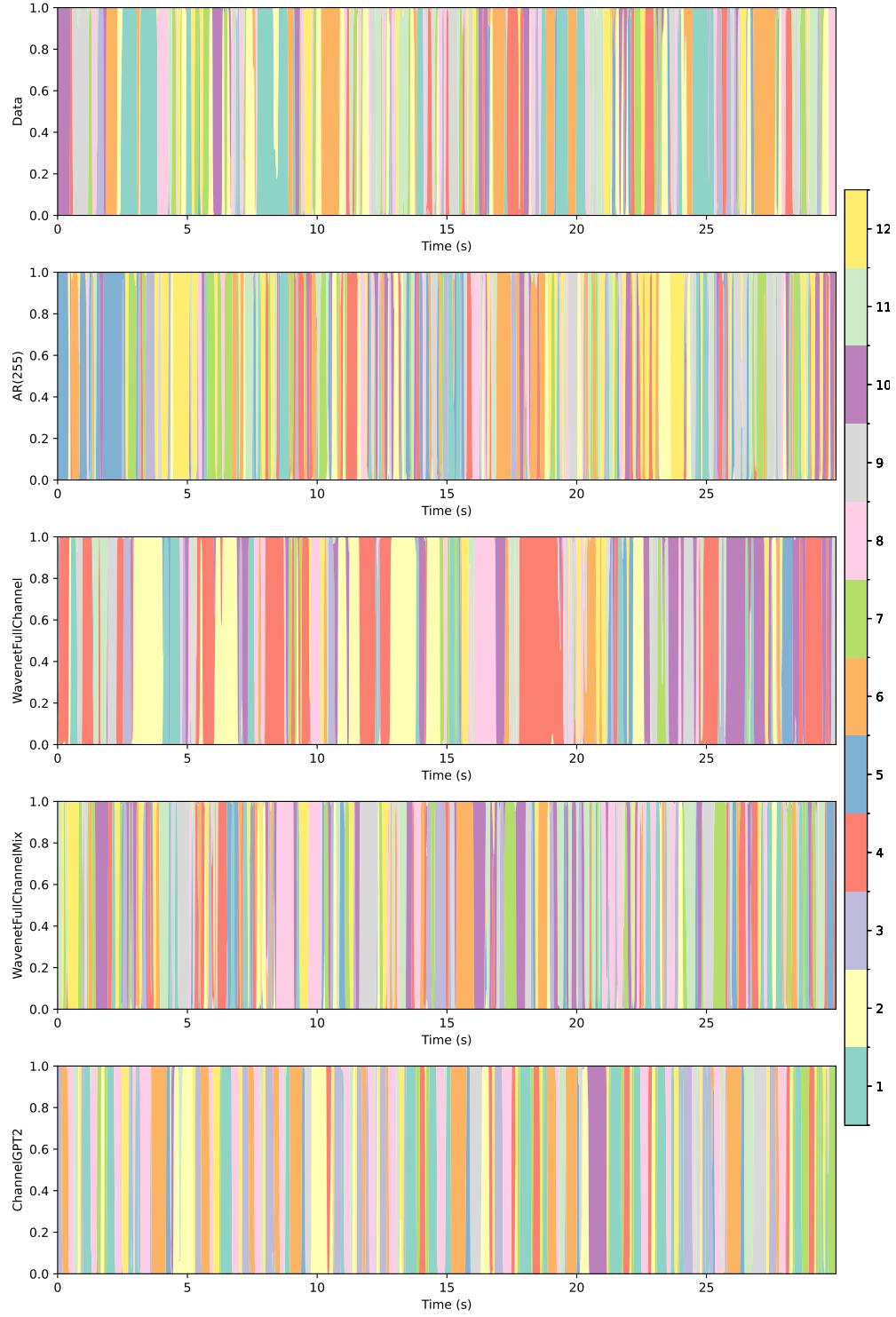


Figure 2: Example state timecourses from the HMMs trained on each model’s generated data (rows). Each state is represented by a different colour. Note that state indices and timecourses are not matched across models.

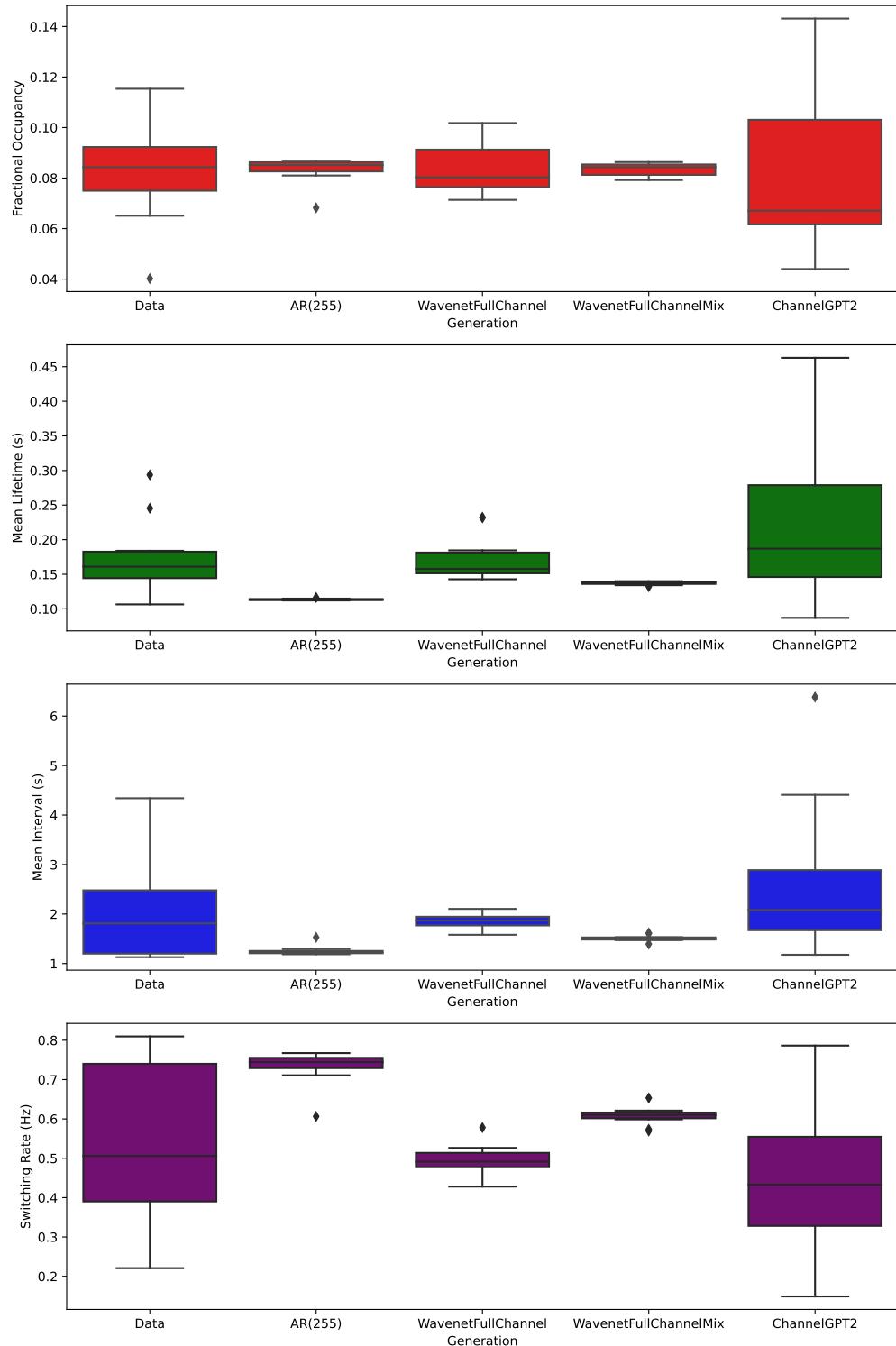


Figure 3: Distributions across the 12 states inferred by an HMM using multi-channel data generated from the different models, or using the real sensor space MEG data (columns). Four different summary statistics are shown describing the state dynamics (rows).

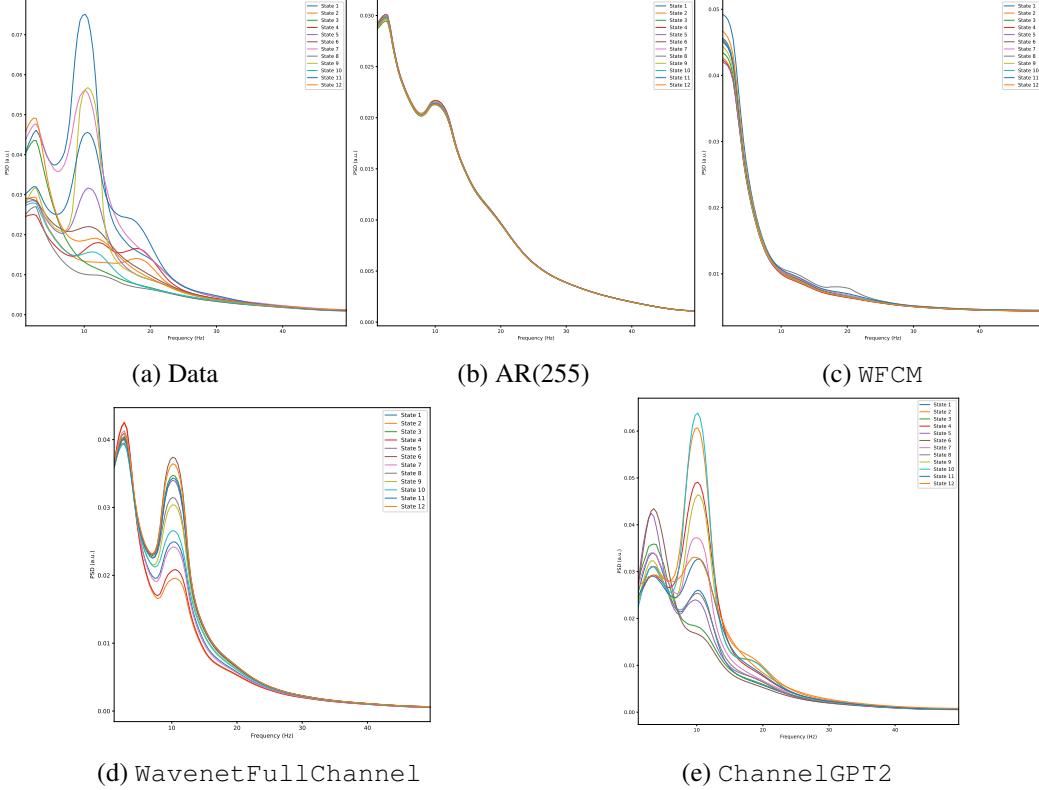


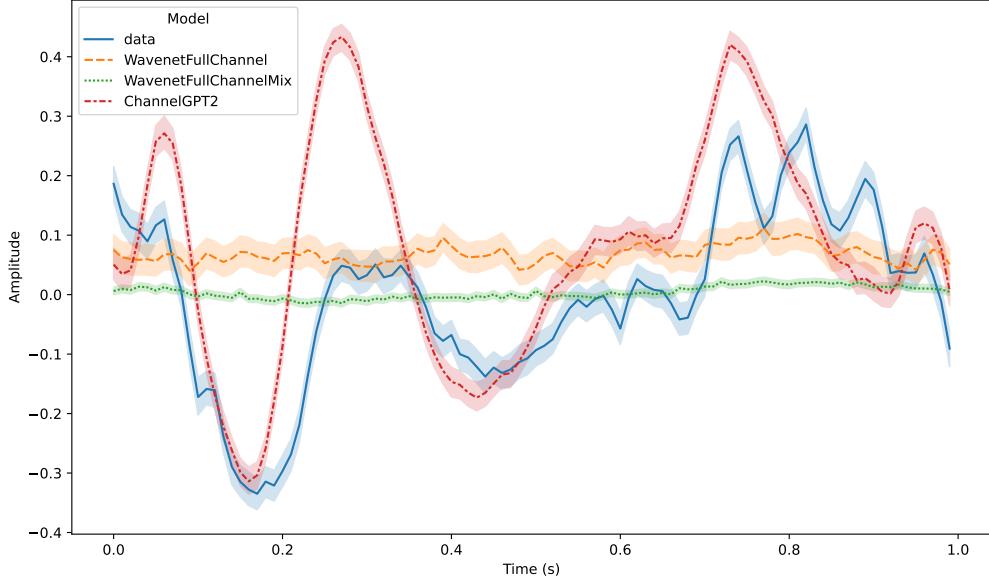
Figure 4: Power spectral density of HMM states inferred on the generated data of each model and on real MEG data. WFCM refers to WavenetFullChannelMix. Each line is the PSD of a different state. Note that states are not matched across models. Horizontal axis represents frequency in Hz.

states of the HMM fitted to the WavenetFullChannelMix generated timeseries do not contain these spectral peaks. While the AR(255) does contain states with a 10 Hz peak, the shape does not match the data well, and also states do not show the same variability as in real data.

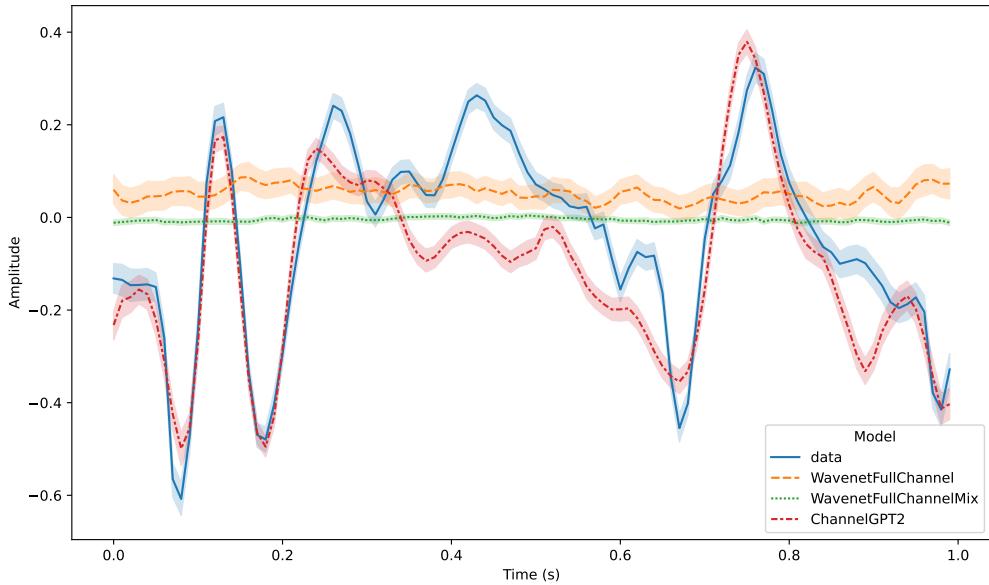
In contrast, ChannelGPT2, matches the state PSDs of the real data very well, further demonstrating the superiority of Transformer models in capturing complex neural dynamics. While WavenetFullChannel also improves substantially over the AR(255) power spectra, it falls short in capturing the 20 Hz peak and the heterogeneity between states observed in the real data and the generated data of ChannelGPT2. This and previous analyses show that the combination of channel-independence and a Transformer-based architecture are critical for matching the dynamics of real data.

## 2.4 Evoked activity in generated data matches real data

The analyses in the previous section considered metrics for assessing the ability of the candidate foundation models to generate timeseries without requiring any a priori knowledge of the timing of brain activity. Here, we use knowledge of the experimental task timings in the Cichy et al. (2016) data, to provide insight into the ability of the foundational models to generate realistic task data.



(a) Frontal channel (MEG0111)



(b) Visual channel (MEG2332)

Figure 5: Comparison of evoked timecourses of 2 channels across our task-conditioned models. Stimulus onset is at 0 seconds and offset is at 500ms. The peak occurring after 500 ms indicates a visual response to the offset of the stimulus (removal of the image). Shading indicates variability across trials.

The models in this section were trained on a single sample subject. As mentioned before, we used the task label timeseries from the training data when generating data with our models. If the models properly incorporate this conditioning, the generated data should reflect aligned task-related activity similar to real data.

By simple epoching of the generated timeseries based on the known task labels, we can

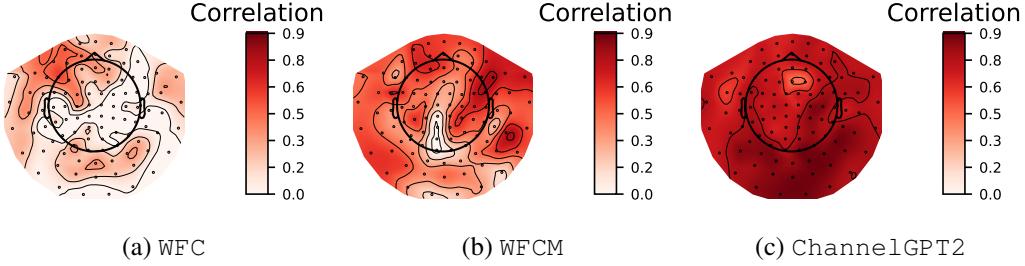


Figure 6: Correlation between the time-courses of the mean (over individual epochs) evoked responses from the real data and mean evoked responses generated by each model. The correlation values are visualised across sensors. WFC refers to WavenetFullChannel and WFMC refers to WavenetFullChannelMix. Darker reds indicate higher correlation.

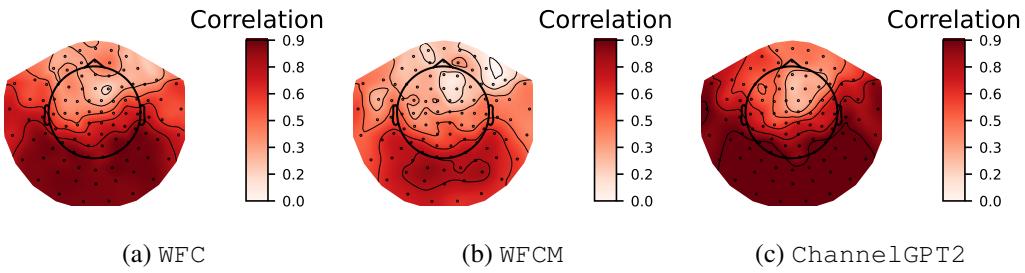


Figure 7: Correlation between the time-courses of the variance (over individual epochs) of the mean evoked responses from the real data and the variance of the mean evoked responses generated by each model. The correlation values are visualised across sensors. WFC refers to WavenetFullChannel and WFMC refers to WavenetFullChannelMix. Darker reds indicate higher correlation.

compute evoked responses in the data generated by our models. We do this for all models except AR(255) as it did not include task labels in its model. To compare the shape of average evoked responses, we averaged over all epochs in both real data and the generated timeseries. This results in data of shape  $\bar{\mathbf{X}} \in \mathbb{R}^{C \times T}$  where  $C = 306$  is the number of channels and  $T = 1000$  ms is the trial/epoch length.

The evoked responses across our models and the real data in a frontal and a visual channel are shown in Figure 5. While Wavenet models completely fail to capture the evoked time-course, ChannelGPT2 does a remarkably good job, especially in the visual channel. This is not surprising as the dataset is collected from a visual experiment, so most activity is visual. ChannelGPT2 closely matches both the amplitude and the timing of the evoked response peaks across the whole 1-second epoch. Variability across trials is also well matched.

To quantify the similarity between real and model generated evoked activity, we computed the correlation of the mean (across individual epochs) time-courses of the evoked response for each channel separately. Note that we averaged over the different MEG sensors (the magnetometers and gradiometers) found at the same location. The result of this is plotted in Figure 6, allowing insights into the spatial pattern of similarity.

As expected, ChannelGPT2 generates data with evoked responses that have much higher correlation with evoked responses from real data, and slightly higher correlation in visual ar-

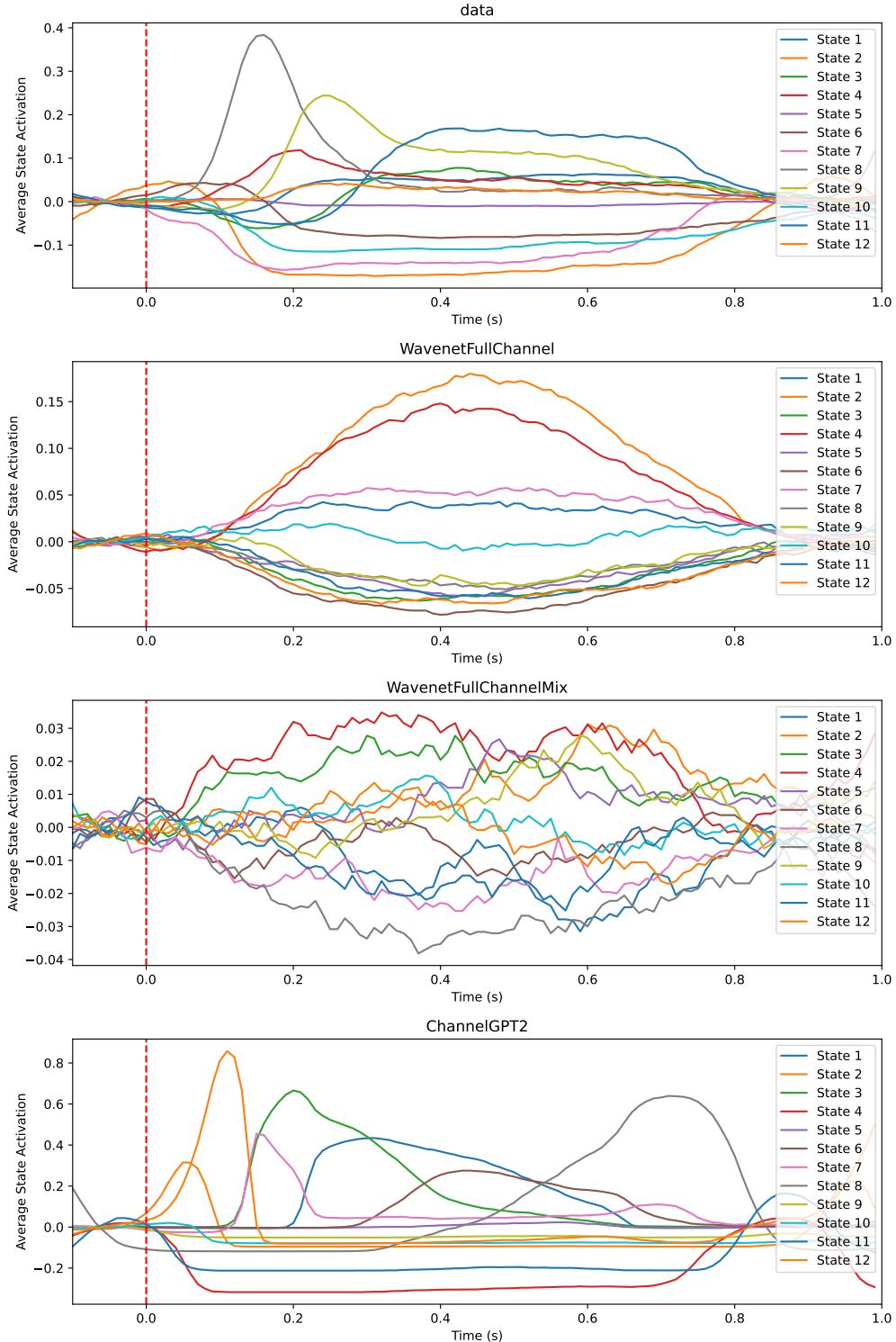


Figure 8: Evoked response state timecourses of HMMs trained on the real MEG data and generated data from each of our task-conditioned models. Note that states are not matched between models. Image presentation starts at 0 seconds and ends at 0.5 seconds.

eas compared to other channels, matching the known topography of visual evoked responses. In other models the correlation is low, and spatially better in frontal areas, likely because the evoked responses here are noisier providing an easier fit.

Figure 7 shows the correlation between the variance (over individual epochs) time-courses of the mean evoked response obtained from the actual data and the evoked responses obtained from data generated by each model. This captures a measure of the ability of the models to represent the trial-to-trial variability found in the real data. Again, ChannelGPT2 generates data that has the highest correlations with the real data, with higher values in channels in the back of the head, appropriately capturing the topography of response variability. Other models have similar spatial distribution, and notably WavenetFullChannel also produces evoked responses with variance partially matching the real data.

Finally, a different way to assess task-related activity is to examine the evoked state time-courses from the HMMs fitted on the real and model generated timeseries. Rather than looking at individual channels, this provides an overall view of which HMM state gets activated when, during individual trials. This is computed by simply epoching the state timecourse, and averaging over all trials. We plot these for the real data and each generated timeseries in Figure 8. As expected, the HMM trained on models other than ChannelGPT2 shows poor evoked state time-courses. ChannelGPT2 generated data produces states with similar evoked dynamics and variability as the real data.

## 2.5 Group-level ChannelGPT2 adapts generated data to individual subjects

Up to this point, all trainings and analyses were done on MEG data from a single subject. We next looked at whether adding data from multiple subjects improves modelling and generation performance. This is in line with the overall goal of training such foundational forecasting models on multiple large datasets. Here we took a first step in exploring this by scaling ChannelGPT2 to the 15 subjects in the Cichy et al. (2016) data, which we refer to as ChannelGPT2-group. For adapting to multiple subjects and to capture variability over subjects, we used subject embeddings as described in Section 3. The main reason for only evaluating ChannelGPT2 on group data is the comparatively much poorer performance of Wavenet-based models in evoked timeseries generation.

We were interested in whether the model generated evoked responses improved their similarity with the evoked responses from the real data, when using data from more subjects. To compare with the single-subject training we generated data using the subject embedding of that subject. The comparison of the evoked response of single-subject and group models for one 1 visual channel is shown in Figure 9. We found that generally ChannelGPT2-group produces evoked responses that are more smoothed than the single-subject model. This is possibly because the model learns to generate data that is closer to the average statistics over subjects, and while it can adapt its generation based on the subject label — this ability is not perfect.

To test our hypothesis regarding ChannelGPT2-group generating more of an average across subjects, we generated data for all subjects (using appropriate subject embeddings) and compared the grand average evoked responses with those extracted from the MEG data of all subjects. Two channels are plotted in Figure 10. The evoked response averaged over all subjects is much noisier because of the high between-subject variability. However, we can see that indeed ChannelGPT2-group can generate this well, perhaps slightly smoother

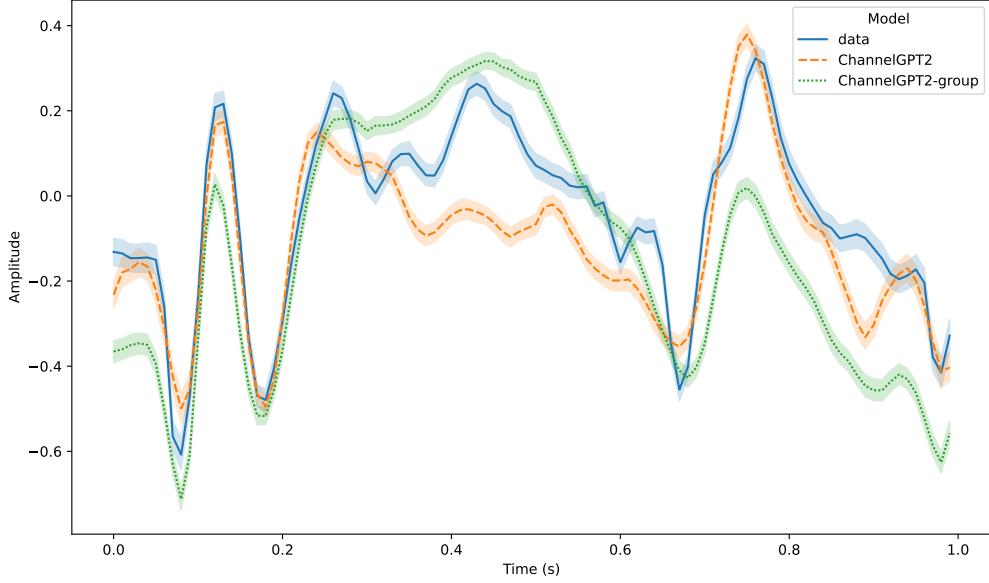


Figure 9: Comparison of evoked responses in a visual channel (MEG2332) across single subject and group models. The stimulus onset is at 0 s and the stimulus offset is at 500 ms. Shading indicates 95% confidence interval across trials.

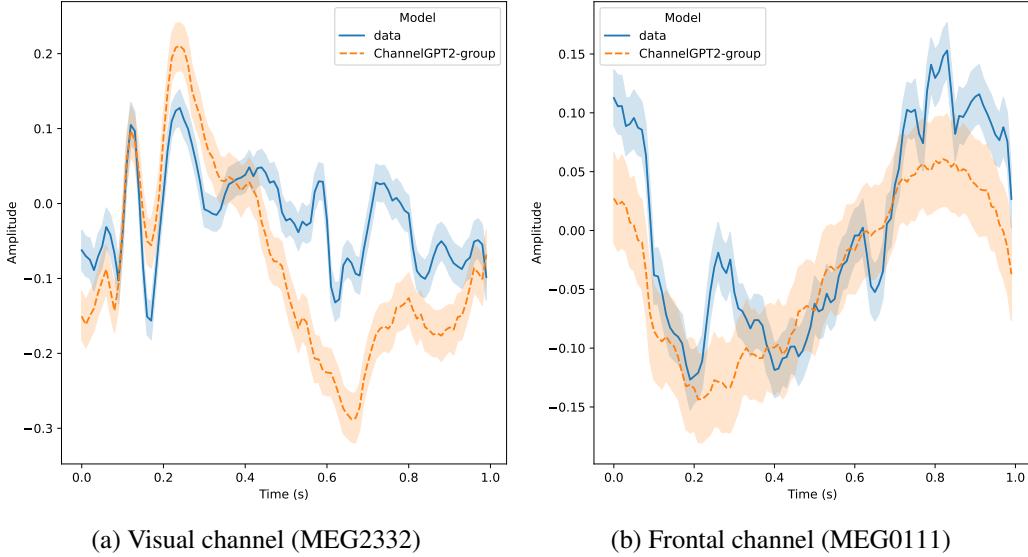


Figure 10: Comparison of evoked responses averaged across all subjects in the data (blue line) and the generated data from ChannelGPT2-group (orange line). The horizontal axis encompasses 1 second, stimulus onset is at 0 seconds and stimulus offset is at 0.5 seconds. Shading indicates 95% confidence interval across trials.

than the real data. Comparing these plots with Figure 9, it is also clear that it adapts its generation well to a specific subject compared to the group average.

A further way to test alignment between group-level evoked responses is to fit an HMM on the data of all subjects, and then infer state timecourses with this model on the generated

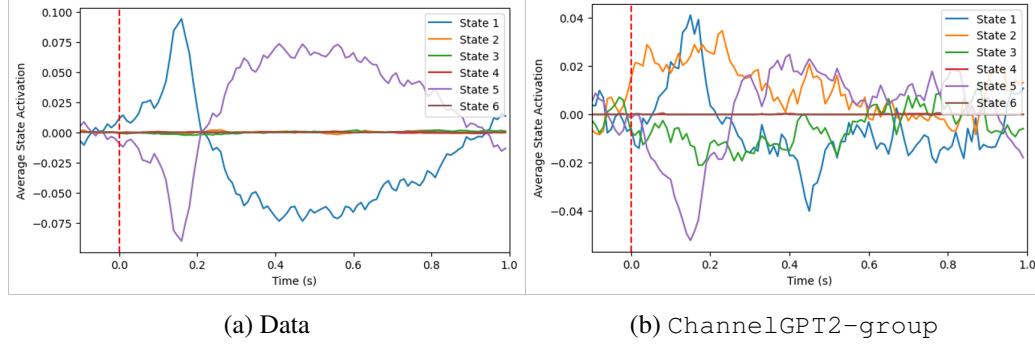


Figure 11: Comparison of evoked state timecourses inferred from the data of all subjects and from the generated data of ChannelGPT2-group for all subjects. State indices are matched between the two plots, as the same fitted HMM model was used.

data of all subjects from ChannelGPT2-group. By taking this approach we can directly match the evoked state timecourses between the real and generated timeseries. We trained an amplitude-envelope HMM (AE-HMM) with 6 states (Quinn et al., 2019) and show results in Figure 11. Two states that show strong activation during real task data show similar temporal signatures and amplitude changes in the generated data, albeit slightly noisier. In the generated data there are two additional states which seem to get activated during the trial. This indicates that while ChannelGPT2-group can capture some of the state-level dynamics, there is room for improvement.

Finally, we examine the variability in state time courses over individual trials. For this we trained an 8-state HMM on the real data of a single subject, and inferred the state timecourses on both the single-subject ChannelGPT2 and ChannelGPT2-group generated data (with the appropriate subject embedding), obtaining matched states. We hypothesised that even if the average evoked responses are similar to the real data, GPT2 may not be able to generate trials with variability in the temporal activation of states. Figure 12 shows that this is indeed true for the single-subject ChannelGPT2 generated data. ChannelGPT2-group responses seem to include much higher temporal variability in state activations, though still falling short of the real data. This indicates that the model can capture some trial-to-trial variability through its exposure to multiple subjects, but has difficulty fully matching the complexity of real neural data. More data may be needed to improve this aspect of generation.

## 2.6 Group-level ChannelGPT2 generates classifiable evoked responses

We have shown that deep learning models, and particularly the channel-independent Transformer-based model (ChannelGPT2), can generate data with spatial, temporal, and spectral signatures similar to real data. We were next interested whether such a foundational model can aid in a downstream task. Specifically, we look at the ability of ChannelGPT2 to aid in the decoding of experimental task conditions in Cichy et al. (2016).

We first investigated whether the task responses generated by the ChannelGPT2 model can be classified with performance comparable to trials of real data. This also further tests how well the model captures spatiotemporal task-related activity and information. The benefit of this approach is that if similar performance is obtained, then ChannelGPT2 could simulate

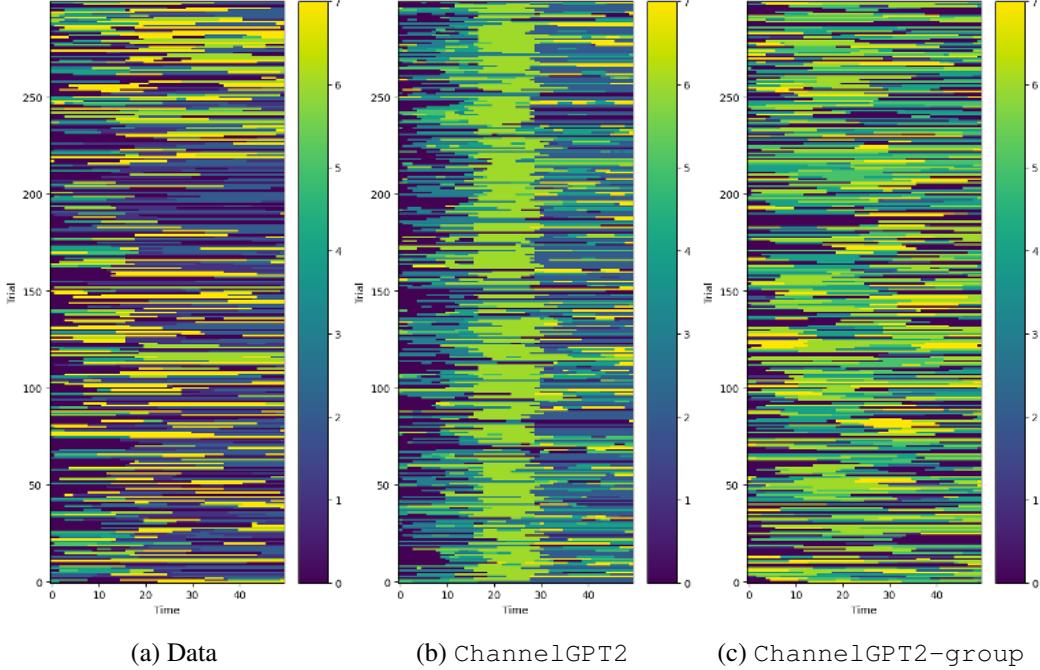


Figure 12: Comparison of trial-level variability in the evoked state timecourses of an HMM trained on real data and applied to the generated data of ChannelGPT2 and ChannelGPT2-group. Different colours represent different states (matched across models). Individual trials however are not matched and we cannot compare the plots at the trial-level, only as an aggregate visualisation of variability across trials.

an arbitrarily large number of trials to potentially improve decoding of real data through pretraining on this simulated data. This is a form of transfer learning where the decoding model, not the forecasting model, is transferred.

We generated 20 trials for all 118 conditions for 1 subject with both ChannelGPT2 and ChannelGPT2-group. We trained separate linear neural network models on the real data (20 trials/condition) and the generated datasets, with an appropriate 4:1 train and validation set ratio. This achieved 17.6%, 1.9%, and 7.2% validation accuracy for the real data, ChannelGPT2, and ChannelGPT2-group, respectively. Thus the group model generates more classifiable subject-specific task-responses, but still does not reach the classifiability of real data. This and previous analyses indicate the group model successfully leverages larger datasets to produce more accurate task-related activity.

We also tried obtaining a decoding model directly from the ChannelGPT2-group forecasting model using Bayes’ theorem. We found limited 5% accuracy over 1 subject’s validation set (versus 40-50% with a discriminative decoder). This generative decoding approach may require larger datasets or more sophisticated architectures.

## 2.7 Transfer learning

A key advantage of generated data is the ability to generate huge amounts of surrogate data. We generated additional datasets with 40 and 60 trials/condition using ChannelGPT2-group. Training a decoder on these achieved 21.7% and 44.2% ac-

curacy, respectively, exhibiting linear scaling of classification performance with simulated data amount. Critically, we assessed whether this simulated data can pretrain classifiers for transfer learning. We first pre-trained the neural network decoder on the 20-, 40-, and 60-trial generated datasets, then finetuned it (trained it further) on the real MEG dataset (20 trials/condition). As the simulated data used for pre-training increased, accuracy of the finetuned model improved rapidly. Zeroshot (no finetuning) performance on real data was above chance with 2%, 3%, and 4% accuracy, for increasing pretraining data quantities. Final accuracies after finetuning were 19.5%, 21.5%, and 23%, respectively. Thus, each additional 20 simulated trials/condition improved final decoding by ~2%. These results are summarised in Table 3.

## 2.8 Ablation experiments

Ablation studies are a common approach in machine learning to understand model behaviour by selectively removing or altering components of the model (Meyers et al., 2019). We performed ablation experiments with ChannelGPT2 to investigate how well it can generate task-related brain activity under varied conditions without further training.

First, we evaluated the model’s ability to adapt to different trial durations. The results reported thus far are for a ChannelGPT2 trained on trials lasting 0.5 seconds. We generated data using the same fitted ChannelGPT2 model but with trial durations of 0.2 s and 0.8 s. As shown in Figure 13, ChannelGPT2 accurately adapted to the shorter and longer trials. The evoked responses matched the expected time-courses, with appropriate truncation or lack of second peaks due to stimulus offset. This demonstrates the model’s ability to generalise to varied trial durations despite being trained on a fixed duration.

Next, we performed two experiments to determine whether ChannelGPT2 relies solely on timing information or also utilises the semantic content of the condition labels. First, we trained a model (ChannelGPT2-randomlabel) where the condition labels were shuffled randomly during training, breaking the semantic alignment between labels and evoked responses. Second, we trained a model (ChannelGPT2-1label) using a single

Trained on (no. trials)	Tested on MEG (20)	Tested on GPT2 (same no. trial data)
MEG (20)	17.6	-
GPT2 (20)	2	7.2
GPT2 (40)	3	21.7
GPT2 (60)	4	44.2
GPT2 (20) + MEG (20)	19.5	-
GPT2 (40) + MEG (20)	21.5	-
GPT2 (60) + MEG (20)	23	-

Table 3: Summary of transfer learning results. The first column shows the data used for training the decoder, with the number of trials per condition shown inside the parenthesis. GPT2 refers to the ChannelGPT2-group generated data, while GPT(.) + MEG (20) is the fine-tuned decoder on the MEG data. The other two columns represent the validation data on which the decoder performance is shown. Accuracy values are provided in percentages. Chance level is 100/118.

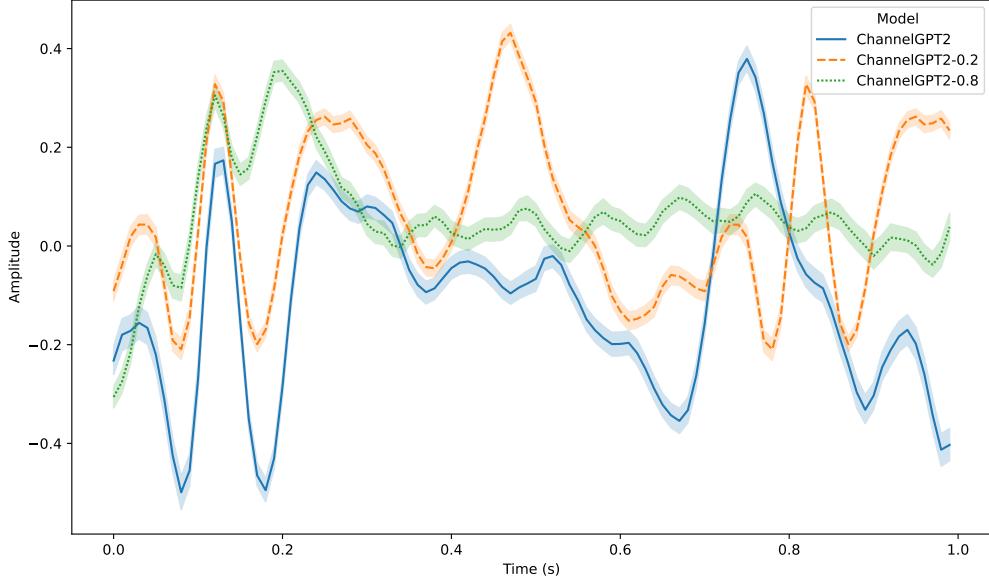


Figure 13: Evoked responses generated by ChannelGPT2 for trials of 0.2 s (orange), 0.5 s (blue), and 0.8 s (green). In all cases, the model was trained only on data containing trials of 0.5 s, but can then adapt appropriately to the different durations. The plotted channel is in the visual area (MEG2332).

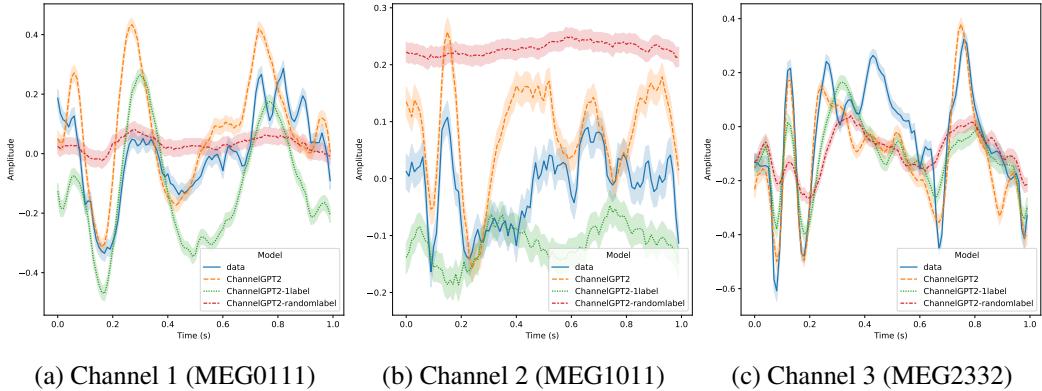


Figure 14: Evoked responses for ChannelGPT2 models trained with shuffled or single condition labels, indicating reliance on semantic content. Three representative channels are presented. MEG0111 is anterior-left, MEG1011 is anterior-central, and MEG2332 is posterior-central. See main text for an explanation of model types. Stimulus onset is at 0 seconds, with stimulus offset at 0.5 seconds.

condition label for all trials. This tests whether the model cheats by learning an average evoked response instead of adapting to each condition.

As evident in Figure 14, both models failed to generate distinct evoked responses for different semantic conditions. This demonstrates that ChannelGPT2 leverages both timing and semantic information in the conditioning labels, rather than simply learning a stereotyped temporal template. Quantitatively, evoked response correlation with real data dropped to 44% and 56% for ChannelGPT2-randomlabel and ChannelGPT2-1label, re-

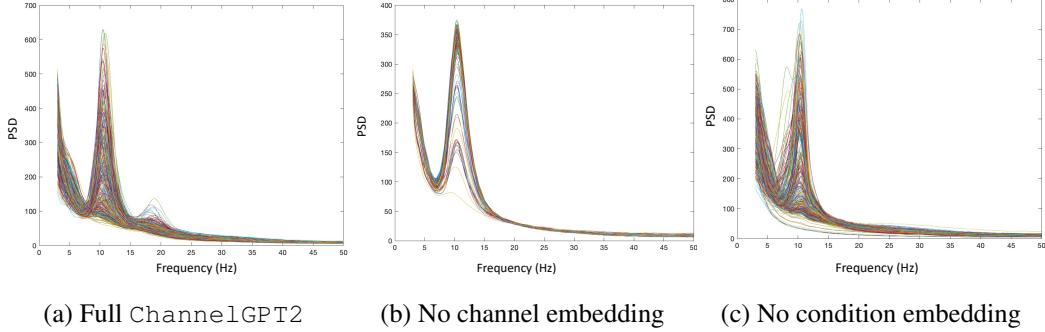


Figure 15: Generated power spectra for full ChannelGPT2 model (left) versus ablations. Both channel (middle) and condition embeddings (right) are critical for accurate spectral content.

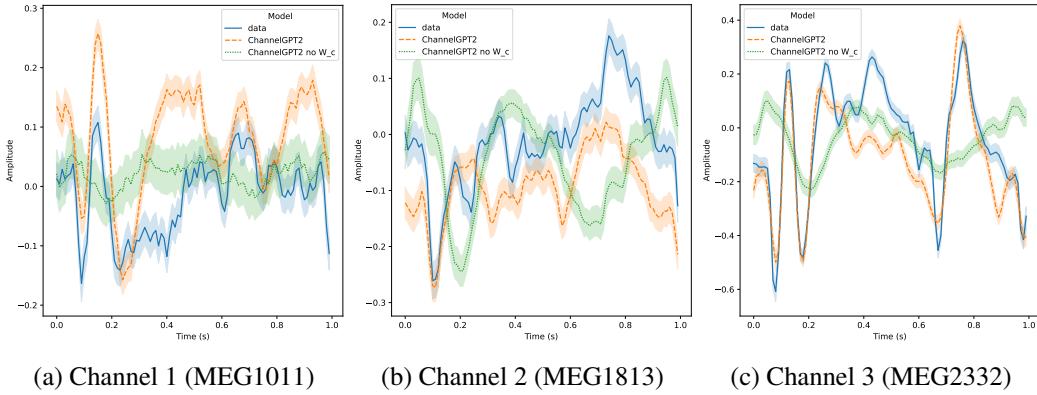


Figure 16: Comparison of generated evoked responses from ChannelGPT2 and the model with ablated channel embeddings (ChannelGPT2 no  $W_c$ ) across 3 representative channels. Without channel embeddings the model fails to adapt evoked responses to different channels. The stimulus onset is at 0 seconds and the offset is at 0.5 seconds.

spectively, compared to 74% for the full ChannelGPT2. Both the qualitative analysis and the correlation numbers indicate that ChannelGPT2-1label was somewhat closer to matching ChannelGPT2.

We also investigated the contributions of the channel and condition embeddings, by training two separate ablated models. As shown in Figure 15, removing the channel embeddings resulted in very similar PSD across channels in the generated data, indicating the model relies heavily on these embeddings to adapt generation per channel. The evoked responses in Figure 16 confirm that without channel embeddings, variability between channels is reduced. Removing the condition embeddings resulted in noisier power spectra of the generated data and no 20 Hz peak.

Finally, we found that the channel embeddings encode spatial relationships, as sensors that are near to each other in the real sensor montage tend to have more similar embeddings. This is shown through a t-SNE and PCA projection of the embedding space in Figure 17. Correlation between pairwise Euclidean distances of channels in physical space and embedding space was 0.45 (Figure 17c).

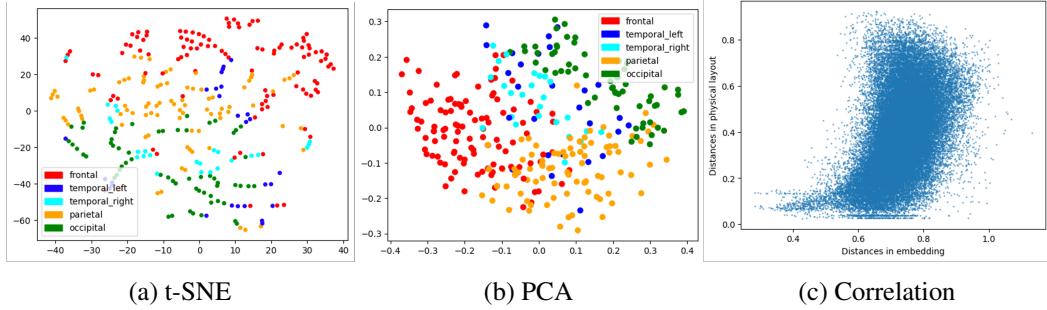


Figure 17: (a) (b) 2D projection of the channel embeddings from ChannelGPT2-group with t-SNE and PCA. Channels are coloured by their location on the scalp grouped into 5 major brain areas. (c) Plotting pairwise Euclidean distances of channels in real, physical space versus embedding space. Sensors that are near to each other in the real sensor montage tend to have more similar embeddings. Each point represents a different pair of channels. Correlation is 0.45.

### 3 Methods

#### 3.1 Multi-channel Wavenet

Here we describe how we adapted the Wavenet architecture (van den Oord et al., 2016) for electrophysiological data. Wavenet models the conditional probability of each time sample given all preceding samples autoregressively:

$$p(\mathbf{X}) = \prod_{t=1}^T p(\mathbf{x}_t | \mathbf{x}_1, \dots, \mathbf{x}_{t-1}) \quad (1)$$

where  $\mathbf{x}_t$  is the sample at time  $t$  and  $T$  is the total sequence length. The network predicts a categorical distribution over tokenised samples using a softmax output layer. Throughout this paper we use tokenisation and quantisation interchangeably. Both have the aim of discretising a continuous quantity into a finite set of distinct bins/levels/tokens.

In the original paper, the audio waveform is tokenised using a quantisation to 8 bits following a  $\mu$ -law companding transform (Lewis and MTSA, 1997):

$$f(\mathbf{x}_t) = \text{sign}(\mathbf{x}_t) \frac{\ln(1 + \mu|\mathbf{x}_t|)}{\ln(1 + \mu)} \quad (2)$$

where  $\mu$  controls the number of quantisation levels, set to 255 as in the original Wavenet.  $f(\cdot)$  is applied to each value of  $\mathbf{x}_t$  independently. This nonlinear transformation improves reconstruction versus uniform quantisation of the raw input, as it skews the distribution such that more levels are allocated to smaller magnitudes. For MEG data, we observe similar benefits when applying this transform prior to quantisation. Note that the input must be scaled to  $(-1, 1)$  first, and clipping outliers above some threshold helps ensure a more uniform mapping.

Critically, tokenisation, in this case through quantisation, enables modelling of probability distributions over data and sampling, instead of just point estimates from MSE-based training.

Cross-entropy loss also avoids the mean-prediction bias induced by MSE (Banville et al., 2021).

When adapting Wavenet to M/EEG, a key challenge is the multi-channel nature of the data. We devise two versions: `WavenetFullChannel` as univariate, and `WavenetFullChannelMix` as multivariate. In both, each channel is transformed and tokenised independently to form the input to the models.

In `WavenetFullChannel`, we first apply an embedding layer to the tokenised data, learned separately per channel. The embedding layer represents each discrete bin as a high-dimensional continuous vector, enabling powerful representations in the convolutional layers whose input channels match the embedding size. To be clear in this univariate approach the same model is applied to each channel. However, a different embedding layer is learned for each channel, meaning that for example the quantised value of 0.42 in channel x will have a different vector representation than in channel y. This helps the model differentiate between channels.

The embedding operation is given below:

$$\forall c \in 1, 2, \dots, C : \mathbf{X}_e^{(c)} = \mathbf{W}^{(c)} \mathbf{X}^{(c)} \quad (3)$$

$$\mathbf{H}_0 = \text{Concatenate}(\mathbf{X}_e^{(1)}, \mathbf{X}_e^{(2)}, \dots, \mathbf{X}_e^{(C)}) \quad (4)$$

Here,  $\mathbf{X}^{(c)} \in \mathbb{R}^{Q \times T}$  is the tokenised one-hot input and  $\mathbf{W}^{(c)} \in \mathbb{R}^{E \times Q}$  is the embedding layer of channel  $c$  mapping tokens  $Q$  to embeddings of size  $E$ . Concatenate concatenates along the channel dimension.

$\mathbf{H}_0 \in \mathbb{R}^{C \times E \times T}$  is the resulting input to Wavenet with  $C$  as the batch dimension. Thus, the same model is applied independently to each channel in parallel. At output, a distribution is predicted simultaneously for each channel at  $T + 1$ . The model is optimised to accurately predict all channels.

`WavenetFullChannelMix` includes an extra linear layer after summing the skip representations to mix information across the channel dimension:

$$\mathbf{S} = \sum_{l=1}^L \mathbf{S}^{(l)} \quad (5)$$

$$\mathbf{S} = \mathbf{S}.\text{permute}(1, 2, 0) \quad (6)$$

$$\mathbf{S}_{out} = \mathbf{S}\mathbf{W}_m \quad (7)$$

where  $\mathbf{W}_m \in \mathbb{R}^{C \times C}$  is the mixing weight matrix, and  $\mathbf{S}^{(l)}$  is the output of the skip connection at layer  $l$ . The permutation is needed to apply the projection to the appropriate channel dimension. After this  $\mathbf{S}_{out}$  is permuted back to the original dimension order and the rest proceeds identically to `WavenetFullChannel`.

In the original Wavenet, audio generation can be conditioned on additional inputs through embedding-based global conditioning or time-aligned local conditioning. For some experiments, we augment the model with local features of task stimuli or subject labels, first embedded into continuous vectors:

$$\mathbf{H}_y = \mathbf{Y}\mathbf{W}_y \quad (8)$$

$$\mathbf{H}_o = \mathbf{O}\mathbf{W}_o \quad (9)$$

$$\mathbf{H}_c = \text{Concatenate}(\mathbf{H}_y, \mathbf{H}_o) \quad (10)$$

where  $\mathbf{Y} \in \mathbb{R}^{T \times N}$  contains the condition index  $n \in (1, \dots, N)$  at each time point, and  $\mathbf{O} \in \mathbb{R}^{T \times S}$  contains the subject index  $s \in (1, \dots, S)$  at each time point  $t \in (1, \dots, T)$ .  $\mathbf{W}_y \in \mathbb{R}^{N \times E_n}$  and  $\mathbf{W}_o \in \mathbb{R}^{S \times E_s}$  are embedding matrices mapping the labels to learned continuous vectors of size  $E_n$  and  $E_s$ , respectively. The subject index is the same across time points of the recording from the same subject. The condition index is set to the (visual) stimuli presented (e.g., one of the 118 images in Cichy et al. (2016)), for exactly those time points when the stimulus is on. At any other time, the task condition embedding  $\mathbf{H}_y$  is set to 0.

$\mathbf{H}_c$  is the conditioning vector fed into Wavenet at each layer:

$$\mathbf{Z}^{(l)} = \tanh \left( \mathbf{W}_f^{(l)} * \mathbf{H}^{(l)} + \mathbf{W}_c^{(l)} * \mathbf{H}_c \right) \odot \sigma \left( \mathbf{W}_g^{(l)} * \mathbf{H}^{(l)} + \mathbf{W}_c^{(l)} * \mathbf{H}_c \right) \quad (11)$$

where  $\mathbf{W}_c^{(l)}$  (1x1 convolution) projects  $\mathbf{H}_c$  before adding it to the input representation ( $\mathbf{H}^{(l)}$ ).  $\mathbf{W}_f^{(l)}$  is the filter convolution weight,  $\mathbf{W}_g^{(l)}$  is the gate convolution weight, and  $\mathbf{Z}^{(l)}$  is the output representation at layer  $l$ .  $\odot$  is element-wise multiplication. This formulation conditions the prediction on both past brain activity and stimuli:

$$p(\mathbf{X}|\mathbf{Y}, \mathbf{O}) = \prod_{t=1}^T p(\mathbf{x}_t | \mathbf{x}_1, \dots, \mathbf{x}_{t-1}, \mathbf{y}_1, \dots, \mathbf{y}_{t-1}, \mathbf{o}_1, \dots, \mathbf{o}_{t-1}) \quad (12)$$

In single-subject models we only use the task labels  $\mathbf{Y}$ .

The full Wavenet architecture can either be interpreted as forecasting with extra conditioning or as a generative encoder augmented with past brain activity. In addition, the probabilistic formulation allows converting the model into a decoder using Bayes' rule, enabling both forecasting and decoding within the same framework:

$$p(Y|X) = \frac{p(X=x|Y)p(Y)}{p(X=x)} \quad (13)$$

where  $X$  is the random variable representing the data,  $Y$  is the random variable representing the task labels, and  $x$  is a particular sample of  $X$ .  $p(Y)$  is the task label prior distribution which in the 118-image dataset is uniform.  $p(X=x|Y)$  is the likelihood of the data given the label which we get from the above formulation of Wavenet. The only tricky part is  $p(X=x)$  as this requires marginalisation over  $Y$ . In the case of the 118-image dataset this means that we have to run the trained model with all of the possible task labels to obtain  $p(X=x)$ :

$$p(X = x) = \sum_{i=1}^N p(X = x|Y = i)p(Y = i) \quad (14)$$

Thus, in a single self-supervised deep learning model we have flexibly encapsulated forecasting, encoding, and decoding, all three of the main modelling methods of M/EEG data. This unification of modelling approaches was inspired by a GitHub repository applying similar ideas to images<sup>2</sup>. The inverted decoder formulation also allows for iterative estimation of  $p(Y|X)$  at each timestep. The author of the GitHub repository has applied this method to estimating the probability of image labels (digits 0 to 9) from pixel images, as more and more of the image was fed into the model.

### 3.2 Multi-channel GPT2

We set out to design a Transformer model suited for M/EEG data, while keeping the key elements that made it successful in language modelling. Specifically, we use GPT-2, a popular autoregressive Transformer variant. When adapting GPT-2 to continuous multivariate time series, the main challenges are at the input and output layers interfacing the model with the data. A particularly detailed visual description of GPT2 is given in Alammar (2019).

Altogether a GPT2 layer is the combination of the self-attention and feedforward layers:

$$\mathbf{H}^{(0)} = \mathbf{X}\mathbf{W}_e + \mathbf{W}_p \quad (15)$$

$$\mathbf{Z}^{(l)} = \text{LN}(\mathbf{H}^{(l)} + \text{MultiHeadAttention}(\mathbf{H}^{(l)})) \quad (16)$$

$$\mathbf{H}^{(l+1)} = \text{Dropout}(\text{LN}(\mathbf{Z}^{(l)} + \text{FFN}(\mathbf{Z}^{(l)}))) \quad (17)$$

$$\hat{\mathbf{Y}} = \text{softmax}(\mathbf{H}^{(L)}\mathbf{W}_e^T) \quad (18)$$

where  $\mathbf{W}_e \in \mathbb{R}^{Q \times E}$  embeds the discrete tokens  $\mathbf{X} \in \mathbb{R}^{T \times Q}$  into  $E$  dimensions. LN is Layer Normalisation, a regularisation technique which normalises all activations within a layer to zero mean and unit variance.  $\mathbf{W}_p \in \mathbb{R}^{T \times E}$  contains positional encodings, providing the model with sequential order information. This is needed as GPT2 lacks recurrent or convolutional elements. Each vector in  $\mathbf{W}_p$  indexed by  $t \in (1, \dots, T)$  contains a distinct  $E$ -dimensional representation of position  $t$ . The output  $\mathbf{H}^{(L)}$  is projected back to the vocabulary via the transpose embedding matrix (weight tying). Alternatively, a separate output projection can be learned. The softmax output gives a token probability distribution.

GPT-2 is trained via supervised learning to predict the next token given previous context, minimising cross-entropy loss between model outputs  $\hat{\mathbf{Y}}$  and ground truth targets  $\mathbf{Y}$ . To enable autoregressive training,  $\mathbf{Y}$  is set to  $\mathbf{X}$  shifted one timestep ahead. Crucially, to prevent information leakage from future timesteps  $t + 1, \dots, T$ , causal masking is applied in each self-attention layer, setting outputs that would reveal future information at position  $t$  to zero.

To apply GPT2 to our continuous multichannel time series data, we take a similar approach as with Wavenet by tokenising each channel independently using the same method as before. This serves as our equivalent of the discrete set of tokens in language modelling. The same

---

<sup>2</sup><https://github.com/cheind/autoregressive>

GPT2 model is applied to each channel in parallel by setting the channel dimension as the batch dimension. We call this ChannelGPT2.

The input to the model includes the position embedding as well as subject and task-stimulus embeddings. We also add a label/embedding telling GPT2 which channel the current time series is coming from:

$$\mathbf{H}^{(0)} = \mathbf{X}\mathbf{W}_e + \mathbf{W}_p + \mathbf{Y}\mathbf{W}_y + \mathbf{O}\mathbf{W}_o + \mathbf{W}_c \quad (19)$$

where  $+$  denotes element-wise addition,  $\mathbf{X} \in \mathbb{R}^{C \times T \times Q}$  is the tokenised input,  $\mathbf{W}_c \in \mathbb{R}^{C \times T \times E}$  are the learned channel embeddings of size  $E$ , which are distinct for each channel  $c \in 1, \dots, C$  but constant across time  $t$ .  $\mathbf{Y}$  and  $\mathbf{O}$  are the task and subject index matrices, mapped to their respective embeddings. As with the positional encoding  $\mathbf{W}_p$ , we simply add all embeddings (task, subject, channel) into a single representation. Note that instead of having channel-specific embeddings of the tokenised input  $\mathbf{X}$  we learn the same mapping  $\mathbf{W}_e \in \mathbb{R}^{Q \times E}$  across channels. Channel information is provided to the model through the channel embeddings.

A serious limitation of this channel-independent GPT2 model is that when predicting a single channel, it does not receive information from other channels. This is analogous to a univariate autoregressive model and ignores crucial cross-channel dependencies in the data. To be clear we often use the term univariate AR modelling in the sense that a separate AR model is trained on each channel. In the case of channel-independent Wavenet and GPT2 models, we train one and the same model on all channels.

### 3.3 Model interpretation

To evaluate whether Wavenet and GPT2 models accurately capture brain dynamics beyond just predictive performance, we develop several analysis techniques to interrogate what these models learn.

**Data generation** Different models have distinct generation procedures. Linear AR models take Gaussian noise as input and generate one timestep at a time. Gaussian noise is added to the output, which is appended to the input sequence. This recursive process is described by:

$$\mathbf{x}_t = \epsilon_t + f(\mathbf{X}_{t-K:t-1}) \quad (20)$$

This intuitively treats the model  $f$  as a black-box infinite impulse response (IIR) filter, where  $\epsilon_t \sim \mathcal{N}(0, 1)$ , and  $K$  is the receptive field size.

For tokenised models (Wavenet and GPT2), we generate data by sampling from the predicted probability distribution and recursively feeding the sample back as input. We use top-p sampling, which samples from the ordered outputs whose cumulative probability mass exceeds  $p\%$  (Holtzman et al., 2020). top-p sampling often works better in practice than using the full distribution by avoiding generation of low-probability tokens, and thus reducing noise.

**Comparison metrics** We compare generated timeseries to real data using power spectral density (PSD), and Hidden Markov Model (HMM) statistics. We fit a separate 12-state time-domain embedding HMM (TDE-HMM) to each multivariate generated timeseries (Vidaurre et al., 2018b). We used the `osl-dynamics` package (Gohil et al., 2023), and set the number of embeddings to 15, the PCA projection dimensionality of the channels to 80 and the sequence length to 2000. We trained the HMMs for 20 epochs with an initial learning rate of 0.02. Once fit we compare the state timecourses of generated timeseries and real data. We also compute four summary statistics (fractional occupancy, mean lifetime, mean interval, switching rate), and compare the distributions of these statistics across the 12 states. Finally we also compare the power spectra of the timeseries corresponding to each state.

**Evoked activity** For task-conditioned models, we assess reconstruction of task-dependent dynamics by feeding in task labels during generation and examining evoked responses. We compare visually average evoked responses from generated and real data, and also quantify the correlation between generated and real evoked responses. Furthermore, we epoch the state timecourses of the HMMs fit to generated and real data and compare evoked state timecourses. The variability in the evoked state timecourses is also compared.

**Classification** To evaluate how well models capture task activity, we apply standard decoding models (e.g., linear classification) to generated trials and compare performance to real data. We also evaluate the generalisability of decoders trained on generated data to real data. Strong similarity in these metrics would indicate accurate modelling of task responses. For decoding both generated and real data we use a simple four-layer linear neural network introduced by (Csaky et al., 2023a). The first layer of the model performs a dimensionality reduction across channels.

**Ablation** By removing certain model components and evaluating performance, ablation studies assess the contribution of different architectural factors. We perform ablations on conditioning and channel embeddings.

### 3.4 Model and training details

As our dataset of choice, we used the continuous 15-subject, 118-image data from Cichy et al. (2016). Here each subject views 118 distinct images, with 30 trials/image. For each subject, the data was bandpass filtered between 1 and 50 Hz, and a notch filter was applied to remove line noise. Subsequently, independent component analysis (ICA) artifact rejection was performed with a dimensionality of 64. Components were visually inspected for each subject, and those that exhibited clear artefactual features (e.g. eye or cardiac signals) were removed. The data was then downsampled to 100 Hz. The continuous data was split into non-overlapping validation, test, and training sets. The validation and test sets included 4 trials of each of the 118 conditions, while the training set contained the remaining 22 trials. This non-overlapping uniform splitting of the continuous data was possible due to the experimental setup during data recording.

For each model the data from each channel was tokenised independently to 256 bins using a quantisation via the mu-law algorithm discussed in Section 3.1. To achieve uniform quantisation, we first standardised each continuous-data channel, clipped values higher than

4 or lower than -4, applied per-channel maximum absolute scaling to map the data to the range (-1, 1), and finally applied the mu-law transform and 8-bit quantisation.

Our aim was to evaluate several models and methods on this dataset. Due to computational constraints and limited iteration speed over experiments and methods, all experiments were performed on a single representative subject, except in Section 2.5 where we explore our GPT2 model on all 15 subjects, and the classification and transfer learning analyses in Section 2.6 and Section 2.7.

We trained univariate AR(255) models on each channel. The AR order was chosen to match the receptive field of the other models. We implemented and trained the AR model using a single linear convolutional neural network layer. Note that we did also assess multivariate AR models (results not shown), but this did not improve performance compared to the univariate AR. We trained `WavenetFullChannel` with a matched receptive field of 255, two stacks of dilation blocks (7 layers per block, doubling dilation factors), 256 hidden channels, 1024 skip channels, no dropout, and a 20-dimensional task embedding. `WavenetFullChannelMix` had the same architecture but 128 hidden channels and 512 skip channels. We used early stopping on the validation set. This means that we ran training until overfitting was observed, and then analysed the model version with the lowest validation loss. All our analyses were performed on the distinct test set.

Our Channel-independent GPT2 (`ChannelGPT2`) had a variable receptive field between 128 and 256. This means that during training the model encountered examples that had a sequence length between 128 and 256, rather than all examples having the same length. GPT2 is normally trained to output all timesteps in a sequence of length  $T$ , given previous timesteps. However, this means that for the second timestep, the receptive field is only 1. Ideally, we wanted to match the training setup of our Wavenet models, where the receptive field is always 256. However, this would significantly slow down training as the whole forward and backward pass must be recomputed at each timestep. We opted for a trade-off, where we set the minimum receptive field to 128, ensuring efficient training and that the model is not trained to predict shorter sequence lengths.

The embedding size of all inputs (token vocabulary, position, task, channel) was set to 96, and we used 12 GPT2 layers, with 12 attention heads. We used Huggingface’s implementation<sup>3</sup>, so the rest of the parameters were the same as in their configuration. Dropout was set to 0 and we used early stopping on the validation set.

On average the mu-law quantisation achieved low reconstruction error. We tested the reconstructed data by performing evoked analysis, and classification of the task responses, and achieved comparable performance to the raw data (results not shown). Thus, both types of tokenisation add negligible quality loss to the data.

For group trainings we used the same hyperparameters as for the single-subject trainings, except for increasing the embedding size to 240. `ChannelGPT2-group` proved difficult to overfit, meaning that using more data acted as a regulariser, and we stopped training when validation losses did not improve for 5 consecutive epochs.

---

<sup>3</sup><https://github.com/huggingface>

## 4 Discussion

We have presented our initial efforts at developing a general forecasting model for M/EEG data. After carefully evaluating the trade-offs between various modelling approaches, we settled on two main architectures: one based on Wavenet (van den Oord et al., 2016), and one based on GPT-2 (Radford et al., 2019). These models have proven successful in the audio and natural language domains, which share similarities with the time series nature of brain signals. We systematically compared different variants of our proposed models on a visual-task MEG dataset.

We found that the forecasting performance was comparable between Wavenet and AR models according to next-timestep prediction metrics. This suggests such metrics may be limited in their ability to effectively evaluate model dynamics beyond one-step prediction. Perhaps looking at these metrics when recursively generating multiple timesteps in the future might be more informative. Generated data analysis provided more discerning model comparisons. While the channel-independent AR and Wavenet models accurately reproduced the overall power spectral density, only the Transformer-based models captured more abstract multivariate statistics like HMM state dynamics.

Critically, the ChannelGPT2 model-generated data closely matched real MEG recordings across both temporal and spectral domains. Analysis of the discovered latent brain states showed ChannelGPT2 reproduced variable oscillatory states similar to those inferred from human recordings (Vidaurre et al., 2018a). Each state captured distinct spectral content, while the linear and Wavenet-based models failed to achieve this degree of heterogeneity in their dynamics. It is possible that this does not indicate a failing of the Wavenet architecture, but rather that different conditioning methods may be needed. One such approach that we have not tested is using the same type of channel embeddings as for ChannelGPT2.

Multiple analyses consistently demonstrated ChannelGPT2’s strengths in realistic conditional timeseries generation. ChannelGPT2-generated evoked responses had high correlation to real MEG data. However, modelling single-trial variability and between-subject differences remain challenging areas needing further work. Scaling to multiple subjects showed promise. The model was able to adapt its generated data based on the input subject label and generate task trials with variability more similar to real recordings than a single-subject model.

Ablation studies quantified the importance of channel embeddings and task conditioning for accurate MEG modelling. Removing channel embeddings resulted in near identical generation across sensors, failing to capture spatial heterogeneity. Analysis of ChannelGPT2’s channel embeddings revealed spatial relationships between sensors were learned, with proximal channels having more similar embeddings. With incorrect or with no task labels, ChannelGPT2 produced noisy evoked responses, indicating the model leverages both timing and label semantics. Furthermore, the model trained on 0.5s trials only, was able to produce reasonable responses to longer or shorter trials, showcasing generalisation. These results demonstrates the value of multi-faceted conditioning for realistic brain data modelling.

A key investigation involved analysing the classification of generated data according to the condition labels. The trials generated by the group-level model were classified with much higher accuracy (closer to real data) than those of the single-subject model. We

also demonstrated that generated data can improve decoding of real trials via transfer learning (Torrey and Shavlik, 2010), with benefits scaling with generated data quantity. The classifiability of generated trials and transfer learning results highlight the utility of forecasting models like ChannelGPT2 for decoding real MEG data. Further analysis could involve permutation feature importance (Altmann et al., 2010) of the decoding model trained on generated data to gain insights into learned representations.

#### 4.1 Comparison with previous works

The Deep Recurrent Encoder (DRE) proposed by Chehab et al. (2022) is a highly relevant architecture to our approaches, as it demonstrates the advantages of modelling spatiotemporal dynamics for encoding neural data. DRE aims to predict MEG brain responses to visual word stimuli. While motivated as an encoding model, DRE can also be viewed through the lens of forecasting, with the addition of auxiliary task features. Forecasting holds inherent advantages over pure encoding, as it enables reconstructing real data and modelling complex spatiotemporal relationships, beyond just learning abstract representations.

Banville et al. (2021) investigate three SSL tasks for learning from unlabelled EEG recordings. Each task is trained via a contrastive loss function, where the model learns to pull positive pair examples closer in a representation space while pushing negative pairs apart. They demonstrate that the representations learned via SSL on unlabelled EEG data transfer well to supervised downstream tasks, consistently improving over limited label training and matching full supervision performance.

Building on this Kostas et al. (2021) propose combining self-supervised contrastive learning with Transformer networks to enable pre-training on large amounts of unlabelled EEG data. Their approach, BErt-inspired Neural Data Representations (BENDR), uses a Transformer encoder architecture applied to learned representations of raw EEG segments. First, a temporal convolutional network extracts initial representations of the EEG time series, referred to as BENDR features. Next, a Transformer encoder module takes the BENDR features as input. Contiguous segments of the BENDR representations are randomly masked, and the model is trained via a contrastive loss to predict the original features. Fine-tuning the pretrained model significantly improves performance on supervised EEG analysis tasks compared to training just on the downstream datasets.

In a similar vein, Wang et al. (2023) propose BrainBERT, a Transformer model tailored for self-supervised pretraining on unlabelled intracranial field potential recordings. BrainBERT is trained to reconstruct randomly masked patches of time-frequency spectrograms computed from the raw voltage measurements. Compared to BENDR, BrainBERT works directly on spectrogram representations rather than learned features. Similar to BENDR, the pretrained BrainBERT model serves as a general purpose feature extractor. By training linear classifiers on top of BrainBERT embeddings, Wang et al. (2023) achieve large performance gains on neural decoding tasks compared to raw or hand-engineered input features. Critically, BrainBERT generalises very well to new subjects and electrode locations.

In the image domain, tokenisation is often abandoned, and a linear projection directly maps image patches to continuous vector representations (Dosovitskiy et al., 2020). Similarly, Nie et al. (2022) have designed a channel-independent Transformer architecture applied to overlapping patches of continuous time series for forecasting. While this facilitates the input, without tokens categorical outputs cannot be generated. We posit that maintaining operations

over tokens and categorical outputs are desirable GPT2 features for M/EEG data. This is because we would like to output probability distributions and train using the cross-entropy loss. Previous applications of Transformer models to M/EEG data do not enforce a discrete tokenisation of the input, nor are they capable to generate new data auto-regressively. The latter allows for better interpretability of learned dynamics and artificial data simulation.

The tokenisation can happen either before or after mixing information across channels. The latter matches GPT2’s original design. One example of this is vector quantisation, which is used to tokenise multiple channels in Jukebox, a successful autoregressive Transformer model used on audio data (Dhariwal et al., 2020). Before training the Transformer, a hierarchical VQ-VAE (vector quantized variational autoencoder (Van Den Oord et al., 2017)) learns discrete codes (tokens) from raw audio. Once trained, VQ-VAE can map a continuous time series to a discrete token sequence  $\mathbf{z}$ . In the second step of Jukebox, the VQ-VAE is kept fixed, and the discrete tokens are used to learn an autoregressive Transformer.

Importantly, VQ-VAE is applied to single-channel audio to compress the temporal dimension into discrete codes. As our aim is to mix information across the channels we would primarily want to apply vector quantisation to the channel dimension, to have a discrete token at each timestep, or perhaps across a few timesteps. While an adaptation of this could work on MEG data, we experimented with a simpler non-deep learning method. We tried applying vector tokenisation on small groups of channels using the Residual Quantiser algorithm (Babenko and Lempitsky, 2014). Unfortunately, this approach to include information from all channels in the input sequence resulted in worse generation capabilities. See Supplementary Section A.1 for details on this method, and Supplementary Section A.3 and Section A.4 for results.

## 4.2 Limitations and future work

The full potential of self-supervised learning is only realised with large-scale data. This remains challenging for brain imaging compared to vision and language. Lowering barriers to data access and promoting data sharing is critical to realise the promise of foundation models in neuroimaging (Poldrack and Gorgolewski, 2014).

A core limitation of the channel-independent GPT2 model is no direct leveraging of cross-channel information for each sensor prediction. Our approach leveraging vector quantisation performed worse. We think that maintaining the innate inductive biases of Transformers, which emphasise 1D sequence modelling on embeddings of discrete tokens, is paramount. Different architectures or more data may enable proper utilisation of cross-channel dependencies. We tried various other approaches to mixing channel information beyond those reported, without success. For the Wavenet model, we incorporated all channels in the input by concatenating embeddings, and for the GPT2 models, we tried mixing channels with convolutions. We tried concatenating the output of each channel and then predicting from this shared output using a different projection for each channel. We also attempted to increase receptive field, dropout, and model size. One limitation in our approaches is the use of a next-timestep prediction loss. Future work should continue exploring architectures and different self-supervised or multi-timestep losses to leverage cross-channel information and improve modelling capabilities.

Some of our findings substantiated that predicting the next timestep may not serve as a robust measure of model performance. Future research should contemplate adopting multi-timestep

or contrastive loss frameworks. A plausible strategy could involve deploying the VQ-VAE model across both channel and temporal dimensions, aiming to distill a more coarse sequence of discrete tokens. Nevertheless, any quantisation-centric approach must carefully consider reconstruction error. We posit that a significant portion of the signal dynamics should be entrusted to the Transformer, given its adeptness in capturing complex dynamics.

A constraint in our modelling approach is its reliance on categorical task stimuli labels. Such an approach, while effective in our context, does not readily lend itself to scalability across diverse tasks and datasets. However, it is conceivable to construct robust representations tailored for various stimulus modalities—ranging from images to audio. These representations can then serve as conditioning embeddings. As shown by Défossez et al. (2022), tools such as wav2vec (Baevski et al., 2020) can be leveraged to derive informative representations of auditory stimuli.

Transfer learning also requires more thorough evaluation across diverse decoding tasks. It would be important to also investigate other more direct finetuning or transfer learning approaches of the forecasting model akin to vision or language domains. These could involve additional output layers and losses for finetuning on downstream tasks.

In conclusion, this work demonstrates that deep forecasting models can accurately reproduce complex neural dynamics of both ongoing and task-related activity and provides an extensive analysis methodology. Future work should explore more flexible conditioning, study different self-supervised and transfer learning frameworks, and critically, apply similar analyses when scaling up across diverse, large electrophysiology datasets. This has the potential to enable powerful transfer learning and advance foundational brain modelling and decoding.

## Acknowledgments

This research was supported by the NIHR Oxford Health Biomedical Research Centre. The views expressed are those of the author(s) and not necessarily those of the NIHR or the Department of Health and Social Care. RC is supported by a Wellcome Centre Integrative Neuroimaging Studentship. MVE’s research is supported by the Wellcome Trust (215573/Z/19/Z). OPJ is supported by the UK MRC (MR/X00757X/1). MWW’s research is supported by the Wellcome Trust (106183/Z/14/Z, 215573/Z/19/Z), the New Therapeutics in Alzheimer’s Diseases (NTAD) study supported by UK MRC and the Dementia Platform UK (RG94383/RG89702) and the EU-project euSNN (MSCA-ITN H2020-860563). The Wellcome Centre for Integrative Neuroimaging is supported by core funding from the Wellcome Trust (203139/Z/16/Z).

## References

- Alammar, J. (2019). The illustrated gpt-2 (visualizing transformer language models). [Jalammar. github. io. https://jalammar.github.io/illustrated-gpt2](https://jalammar.github.io/illustrated-gpt2).
- Altmann, A., Tološi, L., Sander, O., and Lengauer, T. (2010). Permutation importance: a corrected feature importance measure. *Bioinformatics*, 26(10):1340–1347.
- Babenko, A. and Lempitsky, V. (2014). Additive quantization for extreme vector compression. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 931–938.

- Baevski, A., Zhou, Y., Mohamed, A., and Auli, M. (2020). wav2vec 2.0: A framework for self-supervised learning of speech representations. *Advances in neural information processing systems*, 33:12449–12460.
- Baker, A. P., Brookes, M. J., Rezek, I. A., Smith, S. M., Behrens, T., Probert Smith, P. J., and Woolrich, M. (2014). Fast transient networks in spontaneous human brain activity. *elife*, 3:e01867.
- Banville, H., Chehab, O., Hyvärinen, A., Engemann, D.-A., and Gramfort, A. (2021). Uncovering the structure of clinical eeg signals with self-supervised learning. *Journal of Neural Engineering*, 18(4):046020.
- Beltagy, I., Peters, M. E., and Cohan, A. (2020). Longformer: The long-document transformer. [arXiv preprint arXiv:2004.05150](#).
- Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al. (2020). Language models are few-shot learners. [arXiv preprint arXiv:2005.14165](#).
- Chehab, O., Défossez, A., Loiseau, J.-C., Gramfort, A., and King, J.-R. (2022). Deep Recurrent Encoder: A scalable end-to-end network to model brain signals. *Neurons, Behavior, Data Analysis and Theory*, 1.
- Cichy, R. M., Khosla, A., Pantazis, D., Torralba, A., and Oliva, A. (2016). Comparison of deep neural networks to spatio-temporal cortical dynamics of human visual object recognition reveals hierarchical correspondence. *Scientific Reports*, 6(1):1–13.
- Csaky, R., van Es, M. W., Jones, O. P., and Woolrich, M. (2023a). Interpretable many-class decoding for meg. *NeuroImage*, 282:120396.
- Csaky, R., van Es, M. W. J., Jones, O. P., and Woolrich, M. (2023b). Group-level brain decoding with deep learning. *Human Brain Mapping*, 44(17):6105–6119.
- Cui, W., Jeong, W., Thölke, P., Medani, T., Jerbi, K., Joshi, A. A., and Leahy, R. M. (2023). Neuro-gpt: Developing a foundation model for eeg. [arXiv preprint arXiv:2311.03764](#).
- Défossez, A., Caucheteux, C., Rapin, J., Kabeli, O., and King, J.-R. (2022). Decoding speech from non-invasive brain recordings. [arXiv preprint arXiv:2208.12266](#).
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proc. of NAACL*.
- Dhariwal, P., Jun, H., Payne, C., Kim, J. W., Radford, A., and Sutskever, I. (2020). Jukebox: A generative model for music. [arXiv preprint arXiv:2005.00341](#).
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al. (2020). An image is worth 16x16 words: Transformers for image recognition at scale. [arXiv preprint arXiv:2010.11929](#).
- Fedus, W., Zoph, B., and Shazeer, N. (2022). Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *The Journal of Machine Learning Research*, 23(1):5232–5270.

- Gohil, C., Huang, R., Roberts, E., van Es, M. W., Quinn, A. J., Vidaurre, D., and Woolrich, M. W. (2023). osl-dynamics: A toolbox for modelling fast dynamic brain activity. [bioRxiv](#), pages 2023–08.
- Gohil, C., Roberts, E., Timms, R., Skates, A., Higgins, C., Quinn, A., Pervaiz, U., van Amersfoort, J., Notin, P., Gal, Y., et al. (2022). Mixtures of large-scale dynamic functional brain network modes. [NeuroImage](#), 263:119595.
- Hartigan, J. A. and Wong, M. A. (1979). Algorithm as 136: A k-means clustering algorithm. [Journal of the royal statistical society. series c \(applied statistics\)](#), 28(1):100–108.
- Holtzman, A., Buys, J., Du, L., Forbes, M., and Choi, Y. (2020). The curious case of neural text degeneration. In [International Conference on Learning Representations](#).
- Kaplan, J., McCandlish, S., Henighan, T., Brown, T. B., Chess, B., Child, R., Gray, S., Radford, A., Wu, J., and Amodei, D. (2020). Scaling laws for neural language models. [arXiv preprint arXiv:2001.08361](#).
- Kitaev, N., Kaiser, Ł., and Levskaya, A. (2020). Reformer: The efficient transformer. [arXiv preprint arXiv:2001.04451](#).
- Kostas, D., Aroca-Ouellette, S., and Rudzicz, F. (2021). Bendr: using transformers and a contrastive self-supervised learning task to learn from massive amounts of eeg data. [Frontiers in Human Neuroscience](#), page 253.
- Lewis, M. and MTSA, S. (1997). A-law and mu-law companding implementations using the tms320c54x. [Application Note SPRA163A](#), Texas Instrum., Dallas, TX, USA.
- Liu, S., Lu, H., and Shao, J. (2015). Improved residual vector quantization for high-dimensional approximate nearest neighbor search. [arXiv preprint arXiv:1509.05195](#).
- Meyes, R., Lu, M., de Puiseau, C. W., and Meisen, T. (2019). Ablation studies in artificial neural networks. [arXiv preprint arXiv:1901.08644](#).
- Nie, Y., Nguyen, N. H., Sinthong, P., and Kalagnanam, J. (2022). A time series is worth 64 words: Long-term forecasting with transformers. [arXiv preprint arXiv:2211.14730](#).
- Poldrack, R. A. and Gorgolewski, K. J. (2014). Making big data open: data sharing in neuroimaging. [Nature neuroscience](#), 17(11):1510–1517.
- Quinn, A. J., van Ede, F., Brookes, M. J., Heideman, S. G., Nowak, M., Seedat, Z. A., Vidaurre, D., Zich, C., Nobre, A. C., and Woolrich, M. W. (2019). Unpacking transient event dynamics in electrophysiological power spectra. [Brain topography](#), 32(6):1020–1034.
- Rabiner, L. R. (1989). A tutorial on hidden markov models and selected applications in speech recognition. [Proceedings of the IEEE](#), 77(2):257–286.
- Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., Sutskever, I., et al. (2019). Language models are unsupervised multitask learners. [OpenAI blog](#), 1(8):9.
- Sutton, R. (2019). The bitter lesson. [Incomplete Ideas \(blog\)](#), 13(1).

- Torrey, L. and Shavlik, J. (2010). Transfer learning. In Handbook of research on machine learning applications and trends: algorithms, methods, and techniques, pages 242–264. IGI global.
- van den Oord, A., Dieleman, S., Zen, H., Simonyan, K., Vinyals, O., Graves, A., Kalchbrenner, N., Senior, A., and Kavukcuoglu, K. (2016). WaveNet: A Generative Model for Raw Audio. In Proc. 9th ISCA Workshop on Speech Synthesis Workshop (SSW 9), page 125.
- Van Den Oord, A., Vinyals, O., et al. (2017). Neural discrete representation learning. Advances in neural information processing systems, 30.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L. u., and Polosukhin, I. (2017). Attention is all you need. In Guyon, I., Luxburg, U. V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., and Garnett, R., editors, Advances in Neural Information Processing Systems 30, pages 5998–6008. Curran Associates, Inc.
- Vidaurre, D., Abeysuriya, R., Becker, R., Quinn, A. J., Alfaro-Almagro, F., Smith, S. M., and Woolrich, M. W. (2018a). Discovering dynamic brain networks from big data in rest and task. NeuroImage, 180:646–656.
- Vidaurre, D., Hunt, L. T., Quinn, A. J., Hunt, B. A., Brookes, M. J., Nobre, A. C., and Woolrich, M. W. (2018b). Spontaneous cortical activity transiently organises into frequency specific phase-coupling networks. Nature Communications, 9(1):1–13.
- Wang, C., Subramaniam, V., Yaari, A. U., Kreiman, G., Katz, B., Cases, I., and Barbu, A. (2023). Brainbert: Self-supervised representation learning for intracranial recordings. arXiv preprint arXiv:2302.14367.
- Wang, S., Li, B. Z., Khabsa, M., Fang, H., and Ma, H. (2020). Linformer: Self-attention with linear complexity. arXiv preprint arXiv:2006.04768.
- Wen, Q., Zhou, T., Zhang, C., Chen, W., Ma, Z., Yan, J., and Sun, L. (2022). Transformers in time series: A survey. arXiv preprint arXiv:2202.07125.
- Yuan, Z., Zhang, D., Chen, J., Gu, G., and Yang, Y. (2024). Brant-2: Foundation model for brain signals. arXiv preprint arXiv:2402.10251.
- Zhou, H., Zhang, S., Peng, J., Zhang, S., Li, J., Xiong, H., and Zhang, W. (2021). Informer: Beyond efficient transformer for long sequence time-series forecasting. In Proceedings of the AAAI conference on artificial intelligence, volume 35, pages 11106–11115.

## A Supplementary Material

### A.1 FlatGPT2

FlatGPT2 is motivated by the observation that language models include extra information such as context within the sequence, instead of the feature space. We provide the detailed procedure below.

Directly vector quantising 300 channels to any vocabulary size would result in poor reconstruction. In FlatGPT2 we perform the tokenisation on small groups of channels instead. First, we compute the covariance over channels in the training data. Then, we apply K-means clustering (Hartigan and Wong, 1979) on the covariance matrix to group channels into buckets. This ensures that each bucket contains channels with high covariance. This is important because tokenising a feature space (group of channels) with high covariance can be done with fewer tokens while maintaining low reconstruction error. We set the number of clusters ( $B = 30$ ) based on manual tuning on the training data. Each cluster/bucket can contain a variable number of channels, usually between 5 and 20.

After assigning channels to buckets we apply the Residual Quantiser algorithm (Babenko and Lempitsky, 2014) from the faiss library<sup>4</sup> to each bucket  $b$  separately. This is a powerful additive quantiser (Liu et al., 2015) that achieves good reconstruction error with a relatively small vocabulary size  $V$ . Note that the total number of tokens, i.e. the vocabulary size will be  $BV$ , since we have  $B$  quantisers. Once fit to the training data the quantiser is fixed and can be applied to new data.

Mathematically, the covariance is obtained by:

$$\forall i, j \in 1, \dots, C \quad \mathbf{C}_{ij} = \frac{1}{T} \sum_{t=1}^T (x_{t,i} - \mu_i)(x_{t,j} - \mu_j) \quad (21)$$

Where  $x_{t,i}$  is the  $i^{th}$  channel at timestep  $t$ ,  $\mu_i$  is the mean of channel  $i$  over all timesteps, and  $C$  is the total number of channels.  $\mathbf{C}$  is a symmetric matrix, and thus the feature and variable dimensions of K-means are the same. K-means computes buckets  $\mathcal{C}_1, \dots, \mathcal{C}_B$  which partition channels  $C$  into distinct sets with high within-bucket covariance.

The residual quantiser  $Q_b$  learns a codebook  $\mathbf{C}_b \in \mathbb{R}^{V \times |\mathcal{C}_b|}$  for each bucket  $\mathcal{C}_b$ :

$$\forall t \in 1, \dots, T \quad z_{t,b} = Q_b(\mathbf{x}_{t,b}; \mathbf{C}_b) \quad (22)$$

Where  $z_{t,b}$  is the quantised representation (token/code) at timestep  $t$  of the channels  $\mathbf{x}_{t,b} \in \mathbb{R}^{|\mathcal{C}_b|}$  in  $\mathcal{C}_b$ . The encoding in the quantiser is sequential, thus at stage  $m$  of the encoding of  $\mathbf{x}_{t,b}$ , the quantiser picks the entry  $i_m$  that best reconstructs the residual of  $\mathbf{x}_{t,b}$  w.r.t. the previous encoding steps:

$$i_m = \operatorname{argmin}_j \|\mathbf{T}_m(j) - (\mathbf{x}_{t,b} - \mathbf{T}_1[i_1] + \dots + \mathbf{T}_{m-1}[i_{m-1}])\|^2 \quad (23)$$

---

<sup>4</sup><https://github.com/facebookresearch/faiss/wiki/Additive-quantizers>

where  $\mathbf{T}_m$  is a table of size  $K_m$  containing  $|\mathcal{C}_b|$  dimensional vectors. For notational simplicity we omit the index  $b$  from  $i_m$  and  $\mathbf{T}_m$  in the above. The quantisation provides a vector  $[i_1, \dots, i_M]$ , where each element  $i_m$  comes from a set of size  $\lceil \log_2(K_m) \rceil$  bits. This bit vector representation can be easily transformed to token indices ranging from 1 to  $V = \sum_{m=1}^M \lceil \log_2(K_m) \rceil$ . Note that this table description of the discrete code is just a different representation of the overall codebook  $\mathbf{C}_b$ . A code  $[i_1, \dots, i_M]$  can be reconstructed to obtain  $\hat{\mathbf{x}}_{t,b}$  by retrieving the corresponding vectors from  $(\mathbf{T}_1, \dots, \mathbf{T}_m)$  and adding them up. Reconstruction error is computed by comparing  $\hat{\mathbf{x}}_t$  and  $\mathbf{x}_t$ .

By using 30 buckets we obtain 30 tokens per timestep, which is already a 10-fold reduction of the original dimension space, but we have not reached our initial goal of having 1 token per timestep. To achieve this, we flatten the feature dimension (buckets) when feeding tokens to GPT2, hence the name FlatGPT2. Our total sequence length then becomes  $B \cdot T$ , where  $B$  is the number of buckets and  $T$  is the number of timesteps. This approach is also motivated by the observation that language models include extra information such as context within the sequence, instead of the feature space. Thus, when predicting the token of bucket  $b$ , we treat the previous timesteps of the other buckets as contextual information.

We also add an extra separator token  $z_{sep}$  between sequences of buckets corresponding to the same timestep to facilitate distinction between the bucket and time dimensions. An input sequence to FlatGPT2 consists of tokens  $z_{t,b}$  following a fixed order:

$$\mathbf{z} = (z_{sep}, z_{t=1,b=1}, z_{t=1,b=2}, \dots, z_{t=1,b=B}, \quad (24)$$

$$z_{sep}, z_{t=2,b=1}, z_{t=2,b=2}, \dots, z_{t=2,b=B}, \quad (25)$$

$$z_{sep}, \dots, z_{t=T,b=B}) \quad (26)$$

For each codebook  $\mathbf{C}_b$  a separate embedding  $\mathbf{W}_{e,b} \in \mathbb{R}^{V \times E}$  is learned. As in ChannelGPT2 we add the appropriate conditioning embeddings to the input embedding with appropriate flattening across the channel/bucket dimension:

$$\mathbf{H}^{(0)} = \mathbf{Z}\mathbf{W}_e + \mathbf{W}_p + \mathbf{Y}\mathbf{W}_y + \mathbf{O}\mathbf{W}_o + \mathbf{W}_c + \mathbf{W}_t \quad (27)$$

where  $+$  denotes element-wise addition and  $\mathbf{Z} \in \mathbb{R}^{(B+1)T \times V}$  is the one-hot version of  $\mathbf{z}$ . The task labels  $\mathbf{Y}$  can vary across time, but are the same across the buckets of one timepoint.  $\mathbf{W}_c$  now contains distinct embeddings of buckets  $b \in (1, \dots, B)$ , which are the same across timesteps. We also augment the input with  $\mathbf{W}_t$ , containing distinct embeddings for timesteps  $t \in (1, \dots, T)$ , which are the same across buckets. This is the timestep version of  $\mathbf{W}_c$ .

As usual, the model is trained to autoregressively predict the next token in the sequence given all previous inputs. At timestep  $t$  and bucket  $b$  the model has access to the tokens  $\mathbf{z}_{1:t-1}$  from all buckets (and thus information from all channels), and the tokens  $\mathbf{z}_{t,1:b}$ , and has to predict token  $z_{t,b+1}$ . The buckets of the same timestep are predicted sequentially, thus, bucket ordering could influence results. We use an arbitrary bucket ordering and do not experiment with different orderings of the input sequence.

Note that at the last bucket  $B$  in each timestep the prediction should be token  $z_{sep}$ , however, we simply discard this prediction during loss computation, as we do not require the model to

Description	Parameter	Typical value
Number of buckets	$B$	30
Number of code tables	$M$	2
Number of bits per code table	$\lceil \log_2(K_m) \rceil$	7
Vocabulary size per bucket	$V = \sum_{m=1}^M K_m$	16384

Table 4: Hyperparameters of the vector quantisation part of FlatGPT2.

predict separator tokens. The structure of the sequence already constrains the predictions such that a new timestep begins after every  $B$  tokens. Conversely, when computing the prediction at input token  $z_{sep}$ , the target is the token with bucket  $b = 1$  of the next timestep. This is useful as in theory we could start the recursive generation of data with a single  $z_{sep}$  token.

At the output, the transpose of  $\mathbf{W}_e$  can be used to predict probabilities over the vocabulary, or a separate linear projection can be learned. Note that because each codebook  $\mathbf{C}_b$  has a separate vocabulary of size  $V$  assigned to it, we can speed up the output softmax by only computing probabilities over codes/tokens in  $\mathbf{C}_b$  instead of the total vocabulary of size  $BV$ .

FlatGPT2 contains important hyperparameters that affect design choices and performance (Table 4). Increasing the number of buckets  $B$  improves reconstruction error, as the vector quantiser has to quantise less channels, but increases the length of the input sequence to FlatGPT2, and the total size of the vocabulary  $BV$ . The number of code tables  $M$  and the number of bits per code table define the size of the vocabulary  $V = \sum_{m=1}^M K_m$ . These were manually tuned, but generally, we observed that using fewer code tables with a higher number of bits achieves lower reconstruction error. For example, a vocabulary size of 16 bits can be achieved with both two 8-bit code tables and four 4-bit code tables. The trade-off is that using fewer code tables (with more bits) significantly increases computation time. Increasing the vocabulary  $V$  (through the number of code tables and bits per table) improves reconstruction error, as more codes are available for quantising a bucket of channels. However, this increases the total vocabulary  $BV$  of the model, resulting in a larger model.

In summary, key modifications compared to ChannelGPT2 include vector quantisation (tokenisation) of channel groups, and flattening the channel dimension into the sequence. While in theory we could have flattened the full channel dimension without bucketing, this would have resulted in a 10x longer sequence length. However, we are limited by memory constraints since a standard GPT2 model scales quadratically with the sequence length. Memory-efficient Transformer variants are an active research area (Kitaev et al., 2020; Beltagy et al., 2020; Wang et al., 2020), but they have other drawbacks, and we leave their application to M/EEG data to future work.

For FlatGPT2 we set the (temporal) receptive field to be between 120 and 240 because of memory constraints. Note that the total (actual) receptive field of the model is the temporal receptive field multiplied by the number of buckets + 1. All embedding sizes were set to 96, and we used 8 GPT2 layers, with 8 attention heads. Dropout was set to 0 and we used early stopping on the validation set. The quantisation parameters are given in Table 4.

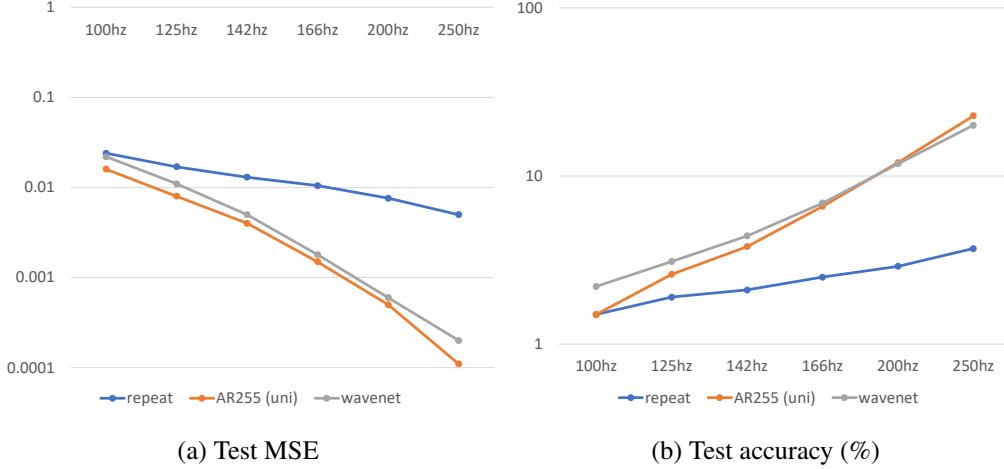


Figure 18: Comparing AR(255) and WavenetFullChannelMix (wavenet) across increasing sampling rates of the data. *repeat* refers to the repeat baseline. Accuracy is given in percentages.

## A.2 Effect of sampling rate on forecasting performance

We further analysed sampling rate effects on forecasting performance in Figure 18. We trained the AR(255) and WavenetFullChannelMix models on increasing sampling rates of the data from 100 Hz to 250 Hz. The lowpass filter was kept the same at 50 Hz. The receptive fields were kept the same in terms of timesteps, thus they decreased accordingly in terms of actual time in seconds. As expected, both AR and Wavenet models improved markedly with higher sampling rates, as the prediction task became easier when timesteps were closer together. The performance gap between models and the repeating baseline also grew with sampling rate. However, these trends are likely influenced by both the changing prediction interval and filtering artefacts. It is unlikely that such marked improvement would be caused by better modelling of higher-frequency content. Varying the low-pass cut-off with sampling rate reduced performance, suggesting filtering effects dominate. Removal of noise with lower lowpass filters is also a possible explanation.

## A.3 Generated covariance

On most metrics presented in the main paper FlatGPT2 performs considerably worse than ChannelGPT2. As the PSD is a channel-independent measure, we also looked at generated data covariance which captures the interactions between different channels (Figure 19). This reveals that the only model capable of closely matching the data covariance is FlatGPT2. All other models produce data with covariances much closer to 0. This is perhaps expected for channel-independent models which generate data independently for each channel, but somewhat surprising for WavenetFullChannelMix. Even though FlatGPT2 may not produce accurate spectral data, by having information about other channels in the input it does an excellent job at capturing covariance. This highlights the trade-offs between different model architectures.

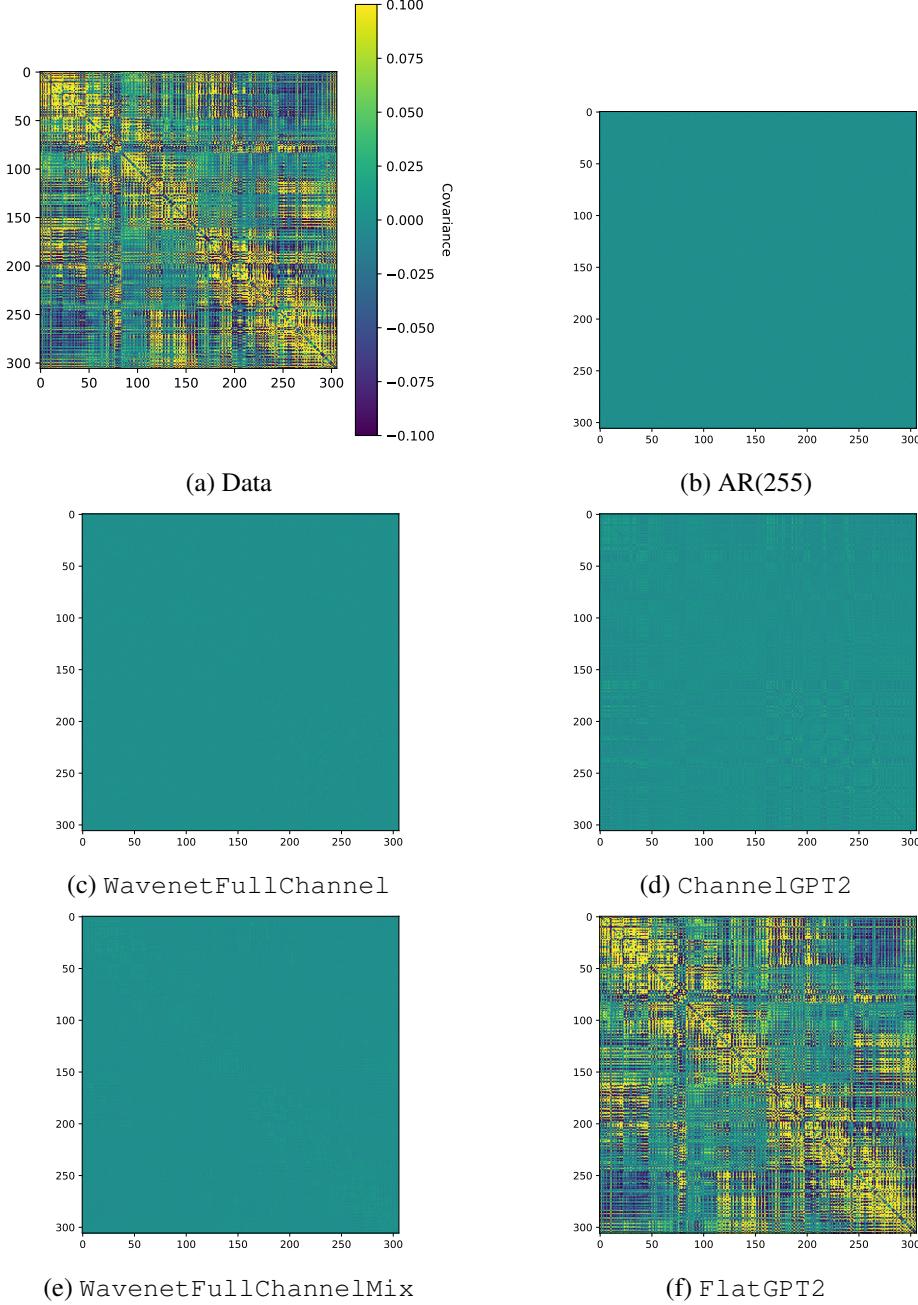


Figure 19: Covariance of generated data between channels (vertical and horizontal axes). All plots have the same scaling as (a).

#### A.4 FlatGPT2 on group data

Unfortunately, even scaling FlatGPT2 did not improve evoked generation. However, we did find that the spectral content of the generated data matched the real data much better than the single-subject version of FlatGPT2 (Figure 20). FlatGPT2-group seemed to scale particularly well with model size as larger models achieved lower and lower loss, improving test accuracy by multiple folds (16.1% top-1 and 40.1% top-5 accuracy) over single-subject

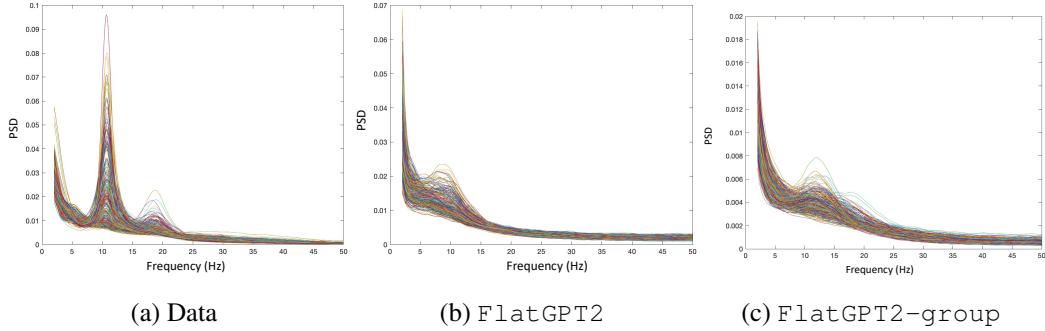


Figure 20: Comparison of generated data PSD across data single-subject FlatGPT2 and FlatGPT2-group. Each line represents a different MEG channel.

FlatGPT2 (3% top-1 and 10.8% top-5 accuracy). This is interesting behaviour compared to ChannelGPT2-group which did not improve much on our forecasting metrics. It remains to be seen whether even more data and larger models are needed to make this type of architecture viable.