

TLDR: Synthetic EEG/ECOG data using DDPMs (instead of the newer LDM approach). They also don't do conditional generation, and retrain everthing for each evaluation setting.

Generating realistic neurophysiological time series with denoising diffusion probabilistic models

Julius Vetter¹, Jakob H. Macke^{1, 2, *}, and Richard Gao^{1, *}

¹Machine Learning in Science, University of Tübingen and Tübingen AI Center

²Max Planck Institute for Intelligent Systems, Tübingen

*equal supervision

Abstract

In recent years, deep generative models have had a profound impact in engineering and sciences, revolutionizing domains such as image and audio generation, as well as advancing our ability to model scientific data. In particular, Denoising Diffusion Probabilistic Models (DDPMs) have been shown to accurately model time series as complex high-dimensional probability distributions. Experimental and clinical neuroscience also stand to benefit from this progress, since accurate modeling of neurophysiological time series, such as electroencephalography (EEG), electrocorticography (ECOG), and local field potential (LFP) recordings, and their synthetic generation can enable or improve a variety of neuroscientific applications.

Here, we present a method for modeling multi-channel and densely sampled neurophysiological recordings using DDPMs, which can be flexibly applied to different recording modalities and experimental configurations. First, we show that DDPMs can generate realistic synthetic data for a variety of datasets including different recording techniques (LFP, ECoG, EEG) and species (rat, macaque, human). DDPM-generated time series accurately capture single- and multi-channel statistics such as frequency spectra and phase-amplitude coupling, as well as fine-grained and dataset-specific features such as sharp wave-ripples. In addition, synthetic time series can be generated based on additional information like experimental conditions or brain states. We demonstrate the utility and flexibility of DDPMs in several neuroscience-specific analyses, such as brain-state classification and imputation of missing channels to improve neural decoding. In summary, DDPMs can serve as accurate generative models of neurophysiological recordings, and have a broad utility in the probabilistic generation of synthetic time series for neuroscientific applications.

1 Introduction

Statistical models that can accurately reconstruct complex datasets have long been useful in many areas of science. These so-called generative models allow scientists to interpret and analyze statistical dependencies between data features, as well as between data and mechanistic latent variables. Such models further grant the ability to generate realistic synthetic data, which is valuable in its own right, especially when it is possible to incorporate (or “condition on”) additional information or observations into the model. Examples of applications include imputing missing data given observed data, forecasting the future given past data, as well as creating synthetic datasets as input to a simulator and to replace or augment scarce training data for subsequent machine learning models, all of which rely on a (conditional) generative model that can produce realistic synthetic data.

The generation of synthetic neurophysiological recordings has received a great deal of attention over the years, predominantly with a focus on generating spike trains, i.e., binary sequences of action potential times: Krumin and Shoham (2009), Macke et al. (2009), Gutnisky and Josić (2010) use thresholding of Gaussian processes and similar techniques to generate spike trains with a pre-specified correlation structure. More generally, time series surrogate methods (Schreiber and Schmitz, 2000, Venema et al., 2006) can be used to generate synthetic neurophysiological recordings that mimic the real data in a well-defined set of features, but lose complex, nonlinear features by design. Other works have employed deep neural networks as non-linear

38 encoding models of brain response, e.g., predicting EEG (Gifford et al., 2022) response from viewed images.
39 However, they are deterministic by design and cannot be used to probabilistically simulate neural time series,
40 nor can they be used without conditioning information such as visual stimuli.

41 Recently, generative models that leverage deep neural networks for probabilistic modeling have made a
42 significant impact across various domains, transforming the way we approach tasks such as image and audio
43 generation (Kong et al., 2021, Dhariwal and Nichol, 2021, Rombach et al., 2022), time series imputation
44 (Fortuin et al., 2020), as well as in scientific applications ranging from molecular design (Sanchez-Lengeling
45 and Aspuru-Guzik, 2018) to black hole imaging (Sun and Bouman, 2021). Neuroscience research, in par-
46 ticular, has benefited from various types of deep generative models. For example, variational autoencoders
47 (VAEs) have been applied to infer low-dimensional representations of single-trial neural population dynamics
48 (Pandarinath et al., 2018), while generative adversarial networks (GANs) have been used for the task of
49 spike-train generation (Molano-Mazon et al., 2018, Ramesh et al., 2019), as well as to decode images from
50 single neuron and fMRI data (Ponce et al., 2019, Lin et al., 2022). The recently proposed denoising diffusion
51 probabilistic models (DDPMs) have also been applied to improve neural decoding performance, in particular
52 leveraging latent diffusion models (Rombach et al., 2022) to predict viewed images from fMRI data (Takagi
53 and Nishimoto, 2022, Chen et al., 2023).

54 Despite their utility for encoding and decoding spiking and imaging data, deep generative models have
55 not yet found empirical success in realistic reproduction of a type of data that is ubiquitous in neuroscience:
56 multivariate and densely sampled electrophysiological time series. Such recordings include non-invasive scalp
57 electroencephalography (EEG), as well as intracranial electrocorticography (ECOG) and local field potential
58 (LFP), which are routinely recorded during scientific research and clinical brain monitoring.

59 A general-purpose model that can conditionally synthesize such neurophysiological data would be valua-
60 able in many scenarios. For example, brain recordings often contain missing values due to sensor noise or
61 misplacement (Silva et al., 2012), which limits their use in neuroscience-specific applications and can lead
62 to them being discarded altogether. Such practical problems are exacerbated by the fact that individual
63 neurophysiological datasets are often limited in availability due to restrictions in clinical contexts. Deep
64 generative models for time series can partially alleviate these issues by providing synthetic data, for example,
65 to facilitate the training of a brain-computer interface (BCI) or as input data for scientific simulations with
66 reduced privacy concerns (Walonuski et al., 2018, Yoon et al., 2019, Abbasi et al., 2020). However, while deep
67 generative models and diffusion models in particular have been used for modeling complex high-dimensional
68 time series in a number of settings (Lin et al., 2023), including probabilistic imputation (Tashiro et al., 2021,
69 Alcaraz and Strodthoff, 2022) and forecasting (Rasul et al., 2021, Biloš et al., 2022), their general utility and
70 performance in modeling neurophysiological recordings have thus far been unexplored.

71 In this work, we demonstrate that denoising diffusion probabilistic models (DDPMs) (Ho et al., 2020)—or
72 diffusion models for short—can accurately generate multivariate and densely sampled continuous neurophys-
73 iological recordings, and further showcase their utility in a variety of neuroscience-specific applications. In
74 particular, we use deep neural networks with structured convolutions (Li et al., 2023) and Ornstein-Uhlenbeck
75 (OU) processes as diffusion processes (Biloš et al., 2022) to improve the modeling of such time series with
76 DDPMs, which allows trained models to produce synthetic recordings with realistic power spectra and even
77 complex features such as sharp wave-ripples and cross-frequency coupling. By conditioning on observed values
78 and other behavioral or experimental variables, synthetic recordings automatically capture cross-area interac-
79 tion between many channels, which is used to impute missing channels and improve performance in a human
80 BCI application at test time. Finally, DDPMs provide additional benefits accessible to generative models,
81 enabling behavioral state classification and outlier detection without additional network training. Code for
82 all our experiments is available at https://github.com/mackelab/neural_timeseries_diffusion.

83 2 Results

84 2.1 Denoising Diffusion Probabilistic Models for neurophysiological time series

85 We use DDPMs (Ho et al., 2020) to learn generative models of electrophysiological recordings. DDPMs are
86 defined by a forward (noising) process and a reverse (denoising) process. The noising process successively
87 adds random noise from a pre-specified distribution (usually independent Gaussian) to the real data through

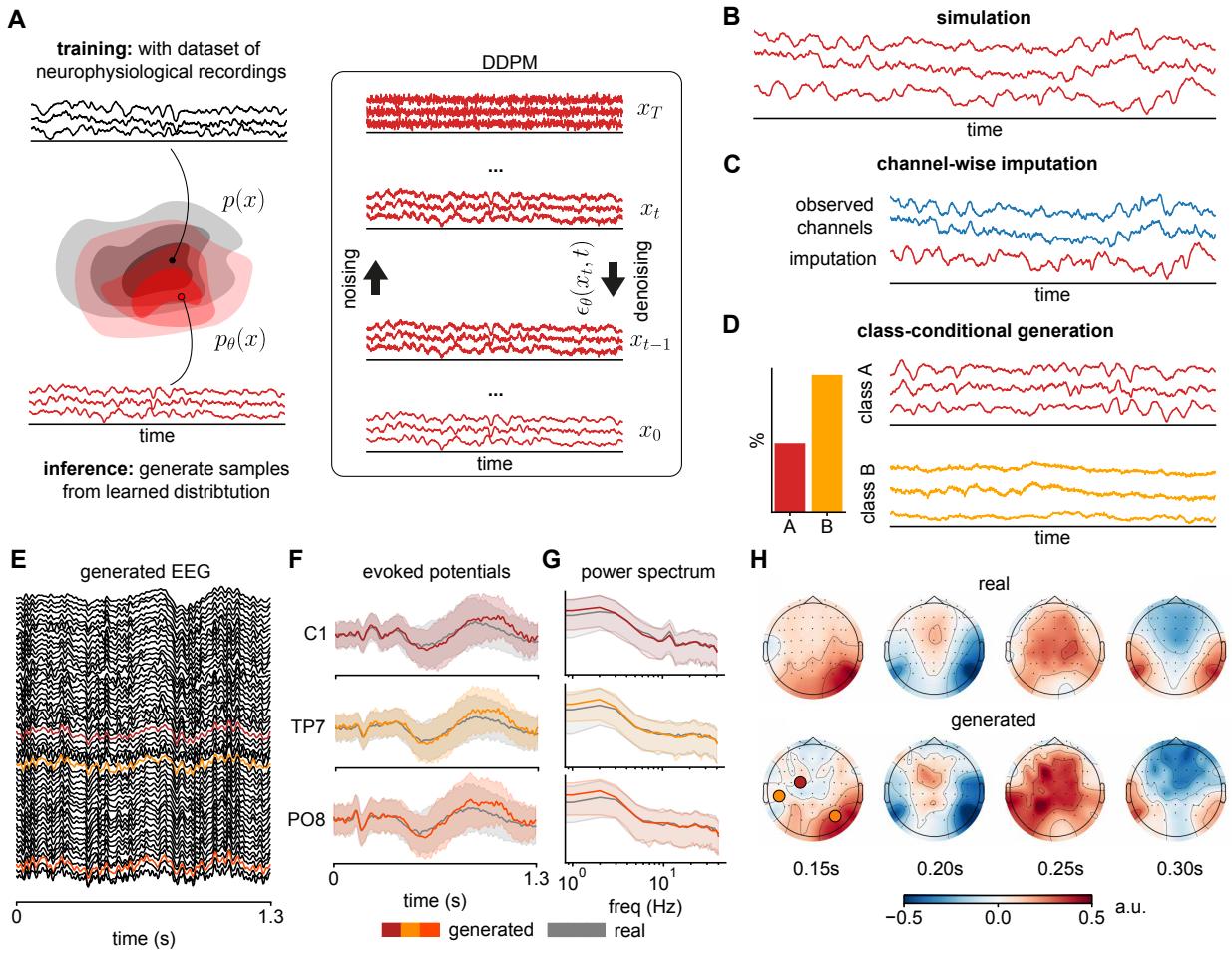


Figure 1: **Overview of diffusion models for neurophysiological recordings and subsequent applications, and an example of DDPM-generated EEG signal.** (A) A denoising diffusion probabilistic model (DDPM) $p_\theta(x)$ is trained on a dataset of neurophysiological recordings. It attempts to generate samples from the data distribution $p(x)$, underlying the training data, by successively denoising samples from a pre-specified Gaussian distribution, using a neural network $\epsilon_\theta(x_t, t)$ as the denoiser. In the context of neurophysiological recordings, DDPMs can be used for various different tasks. Examples include (B) simulation of neurophysiological recordings, (C) imputation of missing values in these recordings, and (D) class-conditional generation of recordings from different experimental conditions or brain states. Since DDPMs allow the computation of likelihoods, the class-conditional model can also be used to perform tasks like classification or outlier detection. (E) An example DDPM-generated trial of 56-channel EEG. (F) Trial-average of three channels show close overlap between real (gray) and generated (colored) evoked potentials (mean and standard deviation across trials), (G) power spectra (median and 10% / 90% percentiles), and (H) spatio-temporal relationships reflected in scalp topography.

88 a series of diffusion steps. The reverse process uses a deep neural network as a “denoiser”, which is trained to
89 reverse the forward process at each diffusion step (Fig. 1A). To generate synthetic data, samples of random
90 Gaussian noise are drawn and gradually denoised by the trained network (Fig. 1A box, top to bottom) over
91 the same number of diffusion steps, resulting in synthetic data that follows the statistical distribution of the
92 real data (i.e., unconditional generation or simulation, Fig. 1B). The training and inference procedure can also
93 incorporate additional information through conditioning. For example, when real data from other channels is
94 given alongside the initial noise sample for the generated channel, the denoising process is guided to generate
95 synthetic data that is likely given the other channels, which can be used to perform imputation of missing
96 channels or interpolation (Fig. 1C). Similarly, when behavioral or experimental state variables are included
97 during training, they can be used to guide the denoising process during time series generation (Fig. 1D). In
98 addition, such a conditionally trained model can also act as a classifier or outlier detector, since the model
99 can be queried for the likelihood of observing a data sample given either conditions.

100 We follow the standard DDPM training and inference procedure (Ho et al., 2020), while incorporating two
101 recent innovations in time series modeling: First, since neurophysiological time series are densely sampled (for
102 example, sampling rate of 500 Hz or more), learning long-range dependencies over hundreds of time points in
103 just one second of data is a non-trivial challenge. We address this problem by using structured convolutions in
104 the denoising network, which use a collection of multi-scale convolution kernels that have shown performance
105 and efficiency improvements on long-range time series tasks (Li et al., 2023). Second, neurophysiological
106 recordings of all modalities (EEG, ECoG, LFP) follow a $1/f$ -like power law in the frequency domain under
107 a range of behavioral and brain states (Pritchard, 1992, He et al., 2010, Gao et al., 2017, Donoghue et al.,
108 2020), in contrast to the flat spectrum of the Gaussian white noise commonly used in DDPMs. Therefore,
109 we sometimes use the Ornstein-Uhlenbeck (OU) process (i.e., colored noise) in the forward noising process
110 (Biloš et al., 2022) to incorporate this prior knowledge, and empirically demonstrate improvements over the
111 standard white noise in a number of experiments. Full details on DDPMs, as well as our denoising network
112 architecture and parameterization of the OU process can be found in Methods (Section 4).

113 As a first demonstration, we train a diffusion model on a single participant’s EEG recorded during a trial-
114 structured task (BCI Challenge @ NER 2015, Margaux et al. (2012)). The trained DDPM is able to simulate
115 (i.e., unconditionally generate) realistic single-trial EEG data (Fig. 1E) and capture features in the real trial-
116 averaged time series and power spectra (Fig. 1F, G), in particular the evolution of the evoked potentials, and
117 a 12 Hz oscillation in the C1-channel. In addition, spatio-temporal relationships across the entire scalp are
118 preserved, as shown in the topographic plots over time (Fig. 1H). This introductory demonstration illustrates
119 how DDPMs can realistically generate high-dimensional neurophysiological time series, and in the following
120 sections we further evaluate and exploit this capability on a variety of datasets and applications.

121 2.2 Overview of DDPM applications: experiments and datasets

122 Beyond the first example, we apply our model to three different datasets to test whether diffusion models
123 can generate realistic synthetic brain recordings in different settings, and to demonstrate their utility in a
124 variety of tasks relevant to neuroscience. On all three datasets, DDPMs accurately capture the statistics of
125 real recordings—in particular, the distribution of channel-wise power spectral densities (PSDs), as well as
126 the multivariate cross-spectrum—across species, recording modalities, and brain areas. Critically, the trained
127 models do not simply overfit to or memorize samples from the training set, which we confirm by measuring the
128 pairwise distance between generated and training data (see Appendix). We show that DDPMs are a general-
129 purpose generative model that accurately captures idiosyncratic features in each of the three datasets, and
130 can be applied in a variety of tasks with minor modifications to network architecture or training procedure:

131 The first dataset consists of 3-channel local field potential (LFP) recordings from rat prefrontal cortex,
132 thalamus, and hippocampus (Varela and Wilson, 2020). These recordings contain sharp wave-ripples that
133 exhibit complex cross-frequency and cross-channel dependencies, which are correctly captured by our model
134 in both simulation and imputation settings. The second dataset consists of human electrocorticography
135 (ECoG) recorded during naturalistic behavior from 12 participants undergoing epilepsy monitoring (AJILE12,
136 Peterson et al. (2022)). In addition to accurately generating high channel-count ECoG data, we show that
137 DDPMs can sensibly impute corrupted or missing channels in an artificially induced missing-data setting
138 (similar to Talukder et al. (2022)), which substantially improves performance in a neural decoding task
139 compared to a baseline (zero-)imputation. The third dataset consists of whole-brain ECoG recordings from

140 a macaque monkey under two different brain states (awake and anesthetized). DDPMs can conditionally
141 generate recordings given one of these two states, as well as evaluate class-specific likelihoods, which in turn
142 can be used to perform classification and outlier detection. The rest of this section presents detailed results
143 from each of the three experiments.

144 2.3 DDPM captures cross-frequency and cross-region dependencies of sharp 145 wave-ripples

146 In the first experiment, we apply our DDPM to model LFP recorded from freely behaving rats (sampled at
147 600 Hz) in three locations: medial prefrontal cortex (mPFC), hippocampal CA1 region, and the thalamic
148 nucleus reuniens (RE) (Varela and Wilson, 2019). These recordings contain neurophysiological features of
149 interest, such as slow oscillations, spindles, and sharp wave-ripples (SWR), which are region-specific and
150 occur with a specific phase relationship relative to one another (Varela and Wilson, 2020). Thus, we aim to
151 generate synthetic data and assess whether our model is able to correctly capture the single-channel voltage
152 distribution and power spectrum, but also complex features such as phase-amplitude coupling between slow
153 oscillations and ripples across channels.

154 After training the model on 3480 2-second time series, we generate the same number of synthetic samples
155 for comparison. Visually, generated traces look indistinguishable from real data, and contain clearly visible
156 SWRs (example in Fig. 2B). Across all generated samples, the distribution of voltage values and power
157 spectral density between real and synthetic data match closely in all three regions (Fig. 2C), and components
158 such as 60 Hz line-noise are also reproduced correctly.

159 To assess whether our DDPM can reproduce more complex cross-regional and ripple-specific features in
160 the real data, we measure the degree of phase-amplitude coupling (PAC) between the amplitude of CA1
161 ripple band (100–275 Hz) and the phase of delta (1–4 Hz), theta (4–9 Hz), and spindle (6–14 Hz) frequencies
162 in the mPFC. We additionally isolate this analysis to the actual occurrences of ripples by thresholding the
163 ripple band power and computing a “phase count coupling” (PCC) histogram between detected CA1 ripples
164 and the instantaneous phase of the aforementioned driver frequencies in the mPFC (similar to Varela and
165 Wilson (2020)). Overall, we see that ripples in both generated and real data have remarkably similar phase
166 preferences, and for all three frequencies in the mPFC, peaking between -45 to 45 degrees (Fig. 2D). As a
167 baseline comparison, we also created channel-wise surrogates (Venema et al., 2006) that preserve the single-
168 channel PSD, but do not show any preferred phase coupling.

169 While the above demonstrates that generated samples reproduce the correct temporal relationships be-
170 tween CA1 ripples and slower fluctuations in the mPFC, we further investigate whether SWRs are correctly
171 captured by our model by conditionally generating (i.e., imputing) CA1 traces given real mPFC and RE
172 traces for a set of held-out evaluation time series. In other words, when provided with the other channels
173 as context, can our model correctly infer when a ripple would have occurred in CA1? We use a technique
174 known as “inpainting” to impute the CA1 channel (Lugmayr et al., 2022). This technique allows DDPMs
175 to fill in missing values by guiding the diffusion process along observed values in the other dimensions (i.e.,
176 channels). To evaluate performance for a given (held-out) evaluation sample, we check if both the real and
177 the imputed CA1 traces contain a ripple by thresholding ripple band power, and compute the F1 score for
178 this classification task across the whole evaluation set. To test whether our model is better than chance at
179 generating CA1 ripples given the appropriate context, we perform a permutation test by shuffling the binary
180 labels of the real CA1 data (i.e., “was there a ripple or not?”) and recomputing the F1 score (Ojala and
181 Garriga, 2010).

182 We provide an example of successfully imputed ripple in Fig. 2E, where the generated CA1 trace contains
183 a ripple at exactly the same time as when the real ripple occurred around the 800ms mark. Repeating the
184 imputation procedure 100 times on the same mPFC and RE traces, we see that the average ripple band power
185 in the imputed CA1 traces is highest at the time that coincides with when the real ripple occurred (Fig. 2F).
186 Across all evaluation samples, our model conditionally generates ripples significantly better than chance
187 ($F1 = 0.45$, $p_{perm} < 0.001$, Fig. 2F inset). Thus, our model successfully captures within-channel temporal
188 characteristics, as well as cross-frequency and cross-channel relationships observed in rat LFP during SWR
189 occurrences, and can therefore be used for probabilistic generation or simulation of highly complex brain
190 recordings. This ability to produce high-quality simulations also extends to more high-dimensional examples
191 such as the 56-channel EEG data or the 128-channel macaque ECoG data (see Appendix).

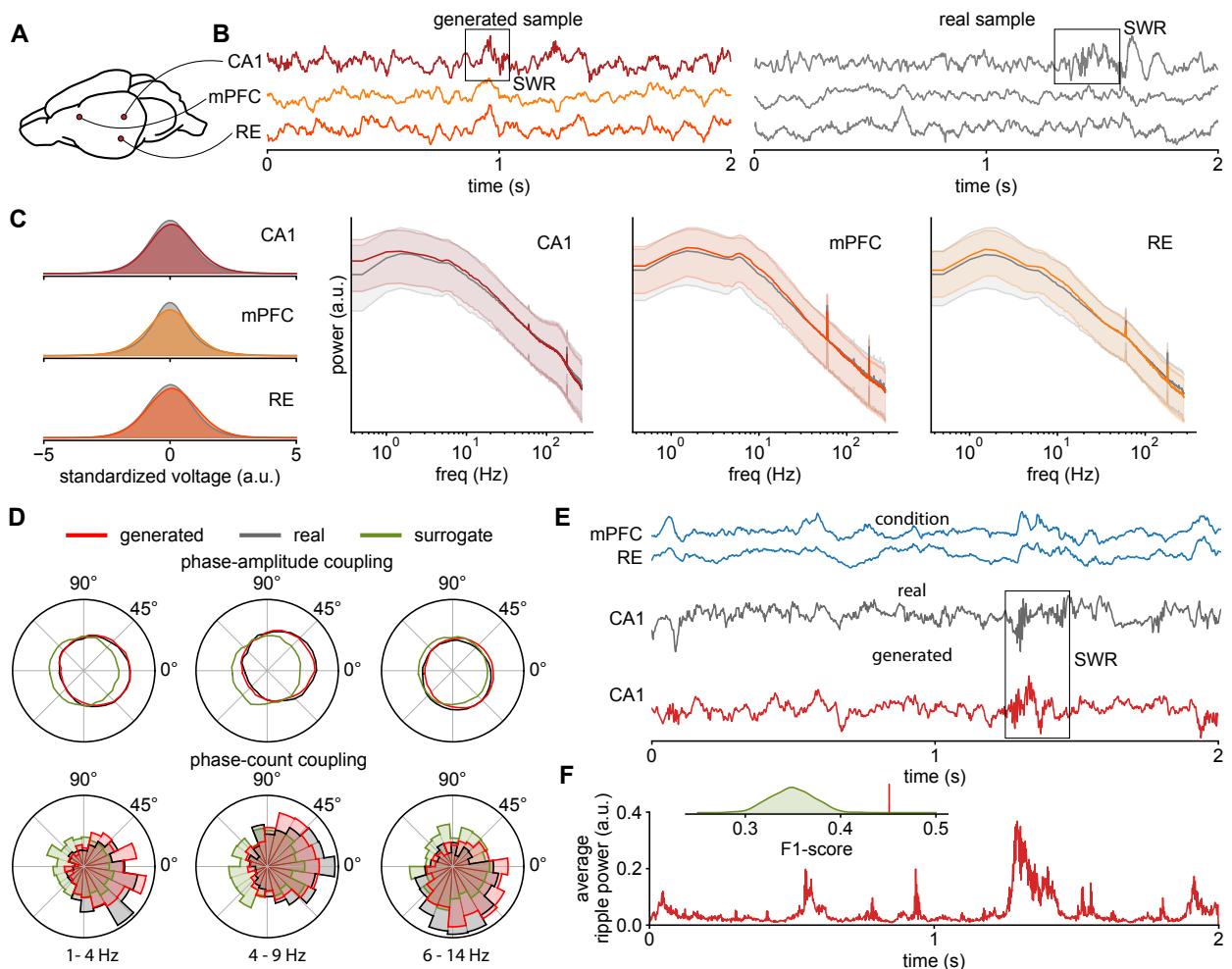


Figure 2: DDPM captures rat cortical-hippocampal LFP dynamics. (A) Schematic of rat brain showing recording locations: CA1, mPFC, and RE. (B) Example of time series generated by our model (left) and a real sample (right). Both examples contain sharp wave-ripples. (C) Marginal distributions over the standardized voltage of each channel for both real and generated data (left), as well as median and 10% / 90% percentiles of real and generated power spectra for each channel (right). (D) Phase-amplitude and phase-count coupling of CA1 ripples for different driver frequencies in the mPFC for real, surrogate, and DDPM-generated data. (E) Example of conditionally generated CA1 trace given real mPFC and RE traces. The real and conditionally generated traces contain a ripple in the same position. (F) Average power in the SWR frequency band (100–275 Hz) over 100 samples from the model (conditioned on the same mPFC and RE traces as in (E)) has a large peak at the same position, indicating that many SWR are generated at this location. Inset: Across the evaluation dataset, the prediction performance (“is an SWR present in both the real and predicted time series?”) measured in terms of F1 score is significantly different from the distribution of scores for randomly permuted predictions.

192 2.4 DDPM-generated imputations improve neural decoding with missing data

193 In our second experiment, we use the AJILE12 dataset (Peterson et al., 2022), which consists of electrocorticography (ECoG) recordings from twelve human participants during naturalistic behavior while undergoing 194 epilepsy monitoring (electrode coverage of one example participant in Fig. 3A). Using an artificially induced 195 missing data scenario (following a similar setup as in Talukder et al. (2022)), we demonstrate how diffusion- 196 based imputation can subsequently improve classifier performance in a neural decoding task. 197

198 For each participant, we first train a diffusion model on all channels from their first six days of recording, 199 with the number of 4-second long training time windows varying between 148 to 1512 per participant. Unlike 200 the model used for rat LFP data, the diffusion model for an AJILE participant is conditionally trained 201 by randomly masking out channels and conditioning on the remaining channels to perform imputation. 202 This conditional mask-based training can further improve imputation performance over the previously used 203 unconditional training without masks (Tashiro et al., 2021). Details on the dataset, network architecture, 204 and conditional versus unconditional training in Methods (Section 4).

205 After training the model, we test its ability to impute “missing” channels of held-out evaluation samples 206 from the seventh day of recording by generating synthetic data for half the channels while conditioning on 207 real data from the other half (i.e., 50% missingness, Fig. 3B). Conditionally generated synthetic traces look 208 visually similar to real data in the corresponding channels (Fig. 3C, 4 highlighted channels), while synthetic 209 PSDs closely match that of the real data for the imputed channels across all evaluation time series (Fig. 3D), 210 demonstrating that DDPMs can be used for participant-specific imputation even with a substantial amount 211 of channels missing. For the AJILE12 dataset, the quality of imputed traces is significantly improved by 212 using OU processes instead of white noise (experiments on the effect of using the OU process versus white 213 noise in the Appendix).

214 To demonstrate the utility of diffusion-based imputation, we use a binary neural decoding task of classifying 215 the 4-second long time windows that were recorded under rest from those recorded during movement. 216 We use per-participant random forest classifiers that were also trained on the first six days of each participant’s 217 recording as decoders (similar to Peterson et al. (2021)), and evaluate classification performance on 218 fully-observed evaluation windows to establish a performance upper bound (Fig. 3F, dashed). To evaluate 219 imputation quality, for each participant, 50%, 70%, and 90% of the channels are dropped out at random from 220 all windows of the evaluation set, and then probabilistically imputed using the trained diffusion model (as in 221 Fig. 3B) before classifying the imputed time series using the same random forest classifier. We repeat this 222 random channel dropping procedure 5 times, and report the average and standard deviation in Fig. 3F (red). 223 We also compare diffusion-based imputation to a zero-imputation baseline (Fig. 3F, blue)—a strategy often 224 employed in such neural decoding paradigms when faced with missing data (Talukder et al., 2022).

225 Across 12 participants and 3 missingness levels, we observe that diffusion-based imputation leads to better 226 classification performance than zero-imputation in 27 out of 36 cases. Remarkably, for some participants 227 (e.g., P05, P07, P09, Fig. 3F), diffusion-based imputations almost recover the original decoding performance 228 using all channels. Overall, classification performance after imputation appears highly participant-dependent. 229 Nevertheless, diffusion-based imputation significantly improves neural decoding accuracy compared to zero- 230 imputation, even with a significant proportion of channels missing, while achieving similar or better performance 231 compared to more tailored imputation methods for neurophysiological recordings (e.g., Talukder et al. 232 (2022)).

233 2.5 Class-conditional DDPMs can generate and evaluate brain state-dependent 234 recordings

235 In our third experiment, we use DDPM to model whole-brain surface electrocorticography (ECoG) from non- 236 human primate (macaque monkey) sampled at 1000 Hz (Yanagawa et al., 2013). In this experiment, we model 237 12 randomly selected channels (out of 128, Fig. 4A) during two distinct behavioral states: awake, a control 238 condition where the animal is head- and arm-restrained but otherwise freely behaving, and anesthetized, 239 where the animal is fully unconscious following the induction of general anesthesia. Our class-conditional 240 diffusion model can accurately generate synthetic recordings under both states, which can be used to train 241 a classifier without giving it access to real data. In addition, we can evaluate likelihoods under the learned 242 generative model, which can be used for state classification and outlier detection of real recordings.

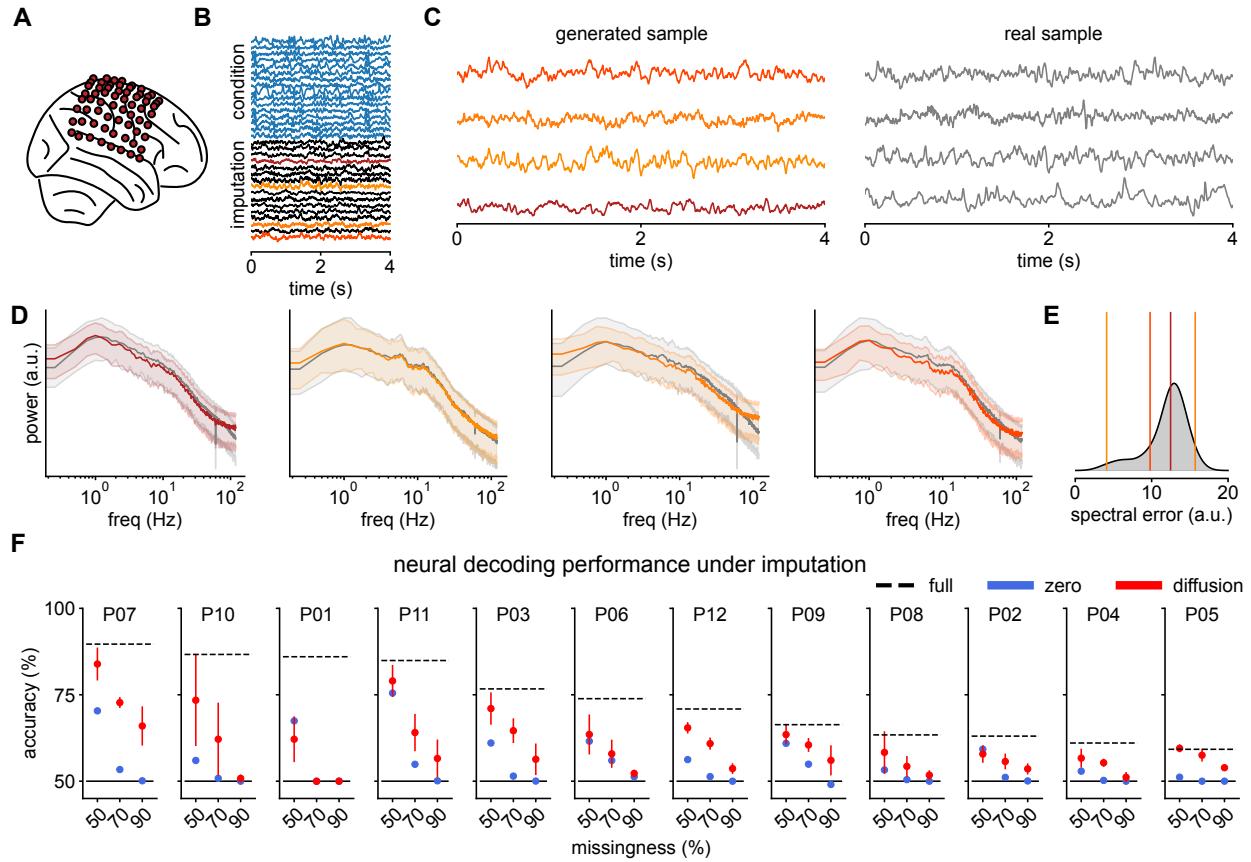


Figure 3: DDPM-generated imputations of human ECoG recordings improve BCI decoding. (A) Electrode layout of participant P07 from the AJILE12 dataset. (B) Example imputation of a 4-second window from P07. The first 32 channels are used as conditioning information (blue), while the remaining 32 channels are imputed (conditionally generated); four are highlighted and shown in (C). (C) Four randomly selected channels from the imputation (left) together with the corresponding real sample (right). (D) Median and 10% / 90% percentiles of real and generated power spectra for the four channels together, and (E) the distribution of spectral errors over all imputed channels. (F) Test performance of the neural decoding model for each of the 12 participants, with fully observed data (dashed), as well as with zero-imputation (blue) and DDPM-based imputation (red) under different amounts of missingness. Participants are sorted by the decoder performance on fully observed data.

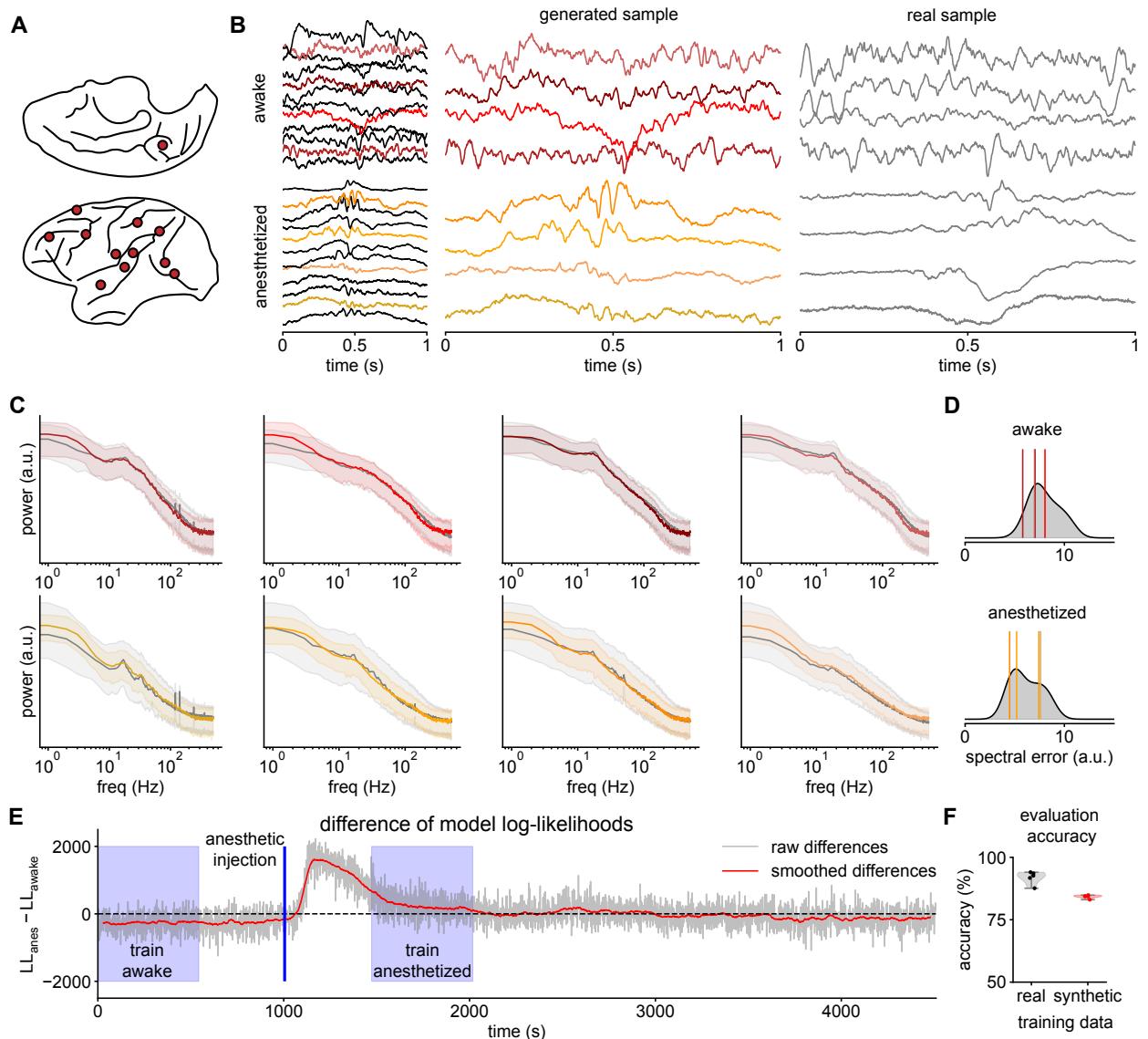


Figure 4: Conditional DDPM learns brain state-dependent macaque ECoG recordings. (A) Electrode locations in the macaque brain of the 12 randomly selected channels used for modeling. (B) Examples of conditionally generated time series for both the awake and anesthetized condition. 4 randomly selected channels are shown in more detail (left) together with real samples (right). (C) Median and 10% / 90% percentiles of real and generated PSDs for the four channels (left) stratified by condition, together with (D) the corresponding distribution spectral errors of all 12 channels (right). (E) Difference of log-likelihoods over 1.5 hours of recording. A positive difference indicates a classification towards “anesthetized”, a negative difference towards “awake”. Periods used for model training are shaded in blue. (F) Real-data test accuracy of a brain-state classifier trained on real (black) vs. DDPM-generated synthetic data (red).

243 After training our model on a total of 870 1-second long windows from both the awake and anesthetized
244 conditions, we conditionally generate the same number of synthetic time series for evaluation. For both the
245 awake and anesthetized condition, the real and generated example time series are visually indistinguishable
246 from each other (Fig. 4B, 4 out of 12 channels are highlighted). Notably, synthetic data during the awake
247 state contains more higher frequency fluctuations and oscillations lasting several cycles, while abrupt lower
248 frequencies fluctuations dominate during the anesthetized state. Across all 12 modeled channels, the PSDs
249 of real and synthetic data closely match in both conditions (Fig. 4C, D).

250 In addition to generating synthetic recordings in a state-dependent manner, class-conditional DDPMs
251 can also evaluate the likelihoods of recordings under each condition, i.e., how likely is it to observe a given
252 recording under awake vs. anesthetized states? Using the difference of log-likelihoods from our trained
253 models to classify previously unseen evaluation time windows from either state, we achieve an evaluation
254 accuracy of 83% on the balanced classification task. Additionally, when we evaluate likelihoods over a 1.5-
255 hour segment of the recording in a sliding-window manner, we observe that the difference in log-likelihood
256 between anesthetized and awake conditions is negative to start, and sharply increases after the induction
257 of general anesthesia—even though the diffusion model does not have access to this information, nor was
258 it trained on nearby segments of the recording (Fig. 4D, training segments in blue). As time elapses, the
259 difference of log-likelihoods slowly decreases towards the “awake level” (negative) as the anesthesia wears off,
260 demonstrating how such generative models can be applied for continuous state prediction.

261 Additionally, unlike with purely discriminative classifiers, the ability to evaluate likelihoods allows us to
262 use our generative model to perform outlier detection, which we demonstrate through two experiments. In
263 the first experiment, for all time series in the evaluation set, we replace half of the channels with white noise
264 to simulate complete signal loss in some of the recording electrodes. In the second experiment, we instead
265 reverse half of the channels time-wise, which is unlikely to happen in the real world, but creates time series
266 that are visually and statistically similar to the real data, thus representing a difficult out-of-distribution test
267 scenario. Using the likelihood of the real and outlier data provided by the diffusion model as a score for
268 classification, the AUC for time-wise flipped outliers is 0.63, while the AUC for white noise outliers is 1.0.
269 When comparing pairs of time-wise flipped outliers with their corresponding real time series, the likelihood
270 of the real time series is larger in 98% of the cases, demonstrating the model’s ability to distinguish data that
271 have the same marginal distribution per channel, but are nevertheless out-of-distribution.

272 Finally, to further demonstrate the utility of using diffusion models to generate realistic neurophysiological
273 recordings, we perform brain state classification based on synthetic data. In many situations, the original data
274 is sensitive due to privacy concerns, e.g., when collected as a part of medical examination, and therefore cannot
275 be publicly shared, but may otherwise be useful for applications such as training of biometric algorithms. In
276 such situations, synthetic recordings that preserve specific aspects of the recording while anonymizing others
277 may be valuable. Here, we use the trained DDPM to generate synthetic data from awake and anesthetized
278 conditions, and train a neural network-based discriminative classifier (as described in Wang et al. (2017))
279 solely on this synthetic data. The trained classifier is then evaluated on a held-out evaluation set of real
280 recordings, while never having seen real data before. Since the diffusion model captures the conditional
281 distribution of the two classes well, the classifier trained on synthetic data achieves an accuracy of 84% when
282 classifying real data (Fig. 4E), which is close to the accuracy of a classifier trained on real data (92%). We
283 thus demonstrate the ability to create artificial datasets that can be shared to facilitate further research with
284 reduced privacy concerns.

285 3 Discussion

286 3.1 Summary

287 Our experiments show that denoising diffusion probabilistic models (DDPMs), together with a flexible con-
288 volutional denoising network and the use of the Ornstein-Uhlenbeck processes, are a powerful tool that
289 can accurately model highly multivariate and densely sampled neurophysiological data, providing features,
290 and more importantly, realistic synthetic data, useful for applications relevant to neuroscience. The power of
291 DDPMs lies in their generality, as they provide a general-purpose generative model that can be used for many
292 different applications ranging from unconditional (i.e., simulation) and conditional generation, to imputation

293 and likelihood evaluation.

294 We demonstrate these applications here on three different datasets of neurophysiological recordings with
295 channel counts ranging from 3 to 128 channels, from a variety of brain regions, recording modalities, and
296 species. Overall, our models accurately capture data distributions in time and frequency domains, as well as
297 detailed cross-channel and cross-frequency dependencies like sharp wave-ripples. On an imputation task, they
298 achieve a performance that is on par with or better than tailored imputation methods for neurophysiological
299 recordings. Finally, we show successful applications of class conditional models in classification, outlier
300 detection, and synthetic training data generation.

301 These features make DDPMs a viable tool for neuroscientists and clinicians studying the complex dynamics
302 embedded in neurophysiological recordings. After successful training, DDPMs can be used to generate
303 unlimited amounts of data as input to neuroscience simulators. When neurophysiological data are subject
304 to strict privacy requirements, as may be the case with human clinical data, neuroscientists and clinicians
305 can use DDPMs to generate artificial datasets that can be shared with collaborators with less concern. In
306 addition, the accurate imputations provided by DDPMs can eliminate the need to discard large amounts of
307 data or to modify an already trained neural decoder. Finally, DDPMs output log-likelihoods that provide
308 users with a model to perform simple classification tasks or outlier detection at no additional training cost.

309 3.2 Other potential applications

310 Trained DDPMs can be used for a variety of other applications that we do not show in this work, such as
311 training data augmentation, as well as imputation of more complex missingness patterns, which includes the
312 task of time series forecasting.

313 Instead of completely replacing training datasets with artificial ones, as described in the data privacy
314 use case (Fig. 4F), it is possible to augment existing real training data with synthetic data generated by a
315 DDPM to improve the performance of another model (e.g., a classifier). This additional training data could
316 lead to performance improvements, especially if the relevant model is difficult to regularize and prone to
317 overfitting. DDPMs can provide tailored “noise”-augmented data samples that capture complex relationships
318 other surrogate methods may not. Throughout our experiments, the DDPMs were robust to overfitting and
319 were able to achieve good sample quality with only a few hundred training time windows of one to four seconds
320 long. In addition, DDPMs can be trained on a separate, larger dataset to generate synthetic “pre-training”
321 data for a classifier model that is later fine-tuned on task-specific but smaller datasets, i.e., to facilitate
322 transfer learning.

323 DDPMs can also handle more complex missingness patterns. In this work, we have demonstrated the
324 ability of DDPMs to perform cross-channel imputation, although completely missing channels represent only
325 one “pattern” of missingness that can occur in real neurophysiological recordings (i.e., dead electrode). Many
326 other patterns of missingness are possible, such as when some time points within a recorded channel are
327 missing due to momentary artifacts. Imputation in this context would mean filling in these missing measurements
328 given information from other channels *and* observed values in the channel itself. These more complex
329 missingness patterns can be realized with diffusion models under both training strategies: For unconditional
330 training without masks, all possible missingness patterns work out of the box via “inpainting” (i.e., Fig. 2E).
331 For conditional mask-based training, the training strategy needs to be generalized from dropping out whole
332 channels to more complex missingness patterns. Tashiro et al. (2021) discuss several strategies for this purpose.
333 A special case is the common task of time series forecasting, where the “missing” values are simply
334 unobserved time points in the future. Like imputation, time series forecasting can then be realized by diffusion
335 models trained unconditionally without masks or conditionally with masks.

336 An additional benefit of our fully convolutional denoising network is that it allows users to generate
337 neurophysiological recordings of arbitrary length, despite training on fixed-size time series. To generate long
338 synthetic recordings, the denoising network is simply applied to the entire time series without any further
339 technical modifications (though dependencies longer than the timescale of training segment may not be
340 correctly modeled). We have not used this capability here, but it will be important when generating input
341 to a simulation where long synthetic recordings are required.

342 3.3 Limitations

343 Although diffusion models are a very popular type of deep generative models today, they are not without
344 limitations. Because DDPMs generate samples by successively denoising them over (typically) hundreds of
345 denoising steps, inference is computationally intensive, especially when compared to other families of deep
346 generative models such as VAEs or GANs. Much recent work has focused on speeding up the inference time
347 of diffusion models, often by tolerating a small degradation in sample quality (Song et al., 2021a, Karras
348 et al., 2022, Lu et al., 2022). While in this work we applied the standard formulation of DDPMs without
349 further speedup, all progress in this direction is directly applicable to our setting and provides a potential
350 remedy for the comparatively long inference times.

351 Furthermore, diffusion models typically operate on the dimensionality of the original data space. In
352 neuroscience, however, there has been much interest in lower-dimensional state-space representations of high-
353 dimensional neurophysiological recordings, especially of neuronal population spiking data (Vyas et al., 2020),
354 but more recently also of continuous time series such as LFPs (Gallego-Carracedo et al., 2022). One way to
355 obtain such low-dimensional neural representations with generative models is to use VAEs (e.g., like in Pan-
356 darinath et al. (2018)), which model the original data distribution via a mapping to a lower-dimensional space
357 of fixed dimensionality. This encoding property is not automatically fulfilled by DDPMs, but combinations
358 of encoder-decoder architectures with a diffusion model in the latent space have been successfully applied to
359 image generation and fMRI decoding (Rombach et al., 2022, Chen et al., 2023). Similar approaches are possi-
360 ble for neurophysiological recordings and may provide low-dimensional representations, along with increased
361 computational efficiency. In particular, a shared latent space could be beneficial in applications similar to
362 our imputation experiment on the AJILE12 dataset (Fig. 3). Unlike the autoencoder-based model described
363 in Talukder et al. (2022), we trained an individual DDPM for each of the twelve participants. Therefore, a
364 potential extension of our current model is to jointly train a latent diffusion model using participant-specific
365 encoder and decoder layers, which could increase performance and robustness due to its ability to generalize
366 over different participants.

367 As discussed earlier, synthetic training data can alleviate privacy concerns. However, our DDPM-generated
368 synthetic data is private in the sense that the minimum difference between time series in the real training
369 data is similar to the minimum difference between time series in the real and generated data. While intuitive,
370 this measure does not provide theoretical guarantees against any form of reverse engineering or memorization
371 (van den Burg and Williams, 2021).

372 Lastly, a persisting challenge in modeling neurophysiological recordings is to capture very long-range
373 dependencies over timescales of minutes or even hours. In our training setup, the time series extracted
374 from a long neurophysiological recording are treated as independent samples by the model. Thus, temporal
375 dependencies beyond the length of the time series used for training are not captured. Naturally, for any given
376 dataset, long-range and cross-scale dependencies are much harder to capture due to data sparsity. Thus,
377 building deep generative models that are more data efficient and can capture long-range dependencies by
378 training on only a few samples is a challenging but exciting avenue for further research.

379 3.4 Conclusion

380 Together with a flexible convolutional denoising network and the use of the Ornstein-Uhlenbeck processes,
381 DDPMs provide a powerful and flexible type of generative model for neurophysiological recordings that can
382 generate high-quality samples, and can be used in a variety of applications relevant to neuroscience like
383 simulation, imputation, or brain-state classification.

384 4 Methods

385 4.1 Denoising Diffusion Probabilistic Models

386 To train our diffusion models, we consider a training dataset of N time windows extracted from neurophysio-
387 logical recordings $\{\mathbf{x}^{(i)}\}_{i=1}^N$ where $\mathbf{x}^{(i)} \in \mathbb{R}^{C \times L}$ with C channels of fixed length L time steps. During sampling
388 or inference, arbitrary time series lengths are admissible due to our fully convolutional network architecture.

389 Background on Denoising Diffusion Probabilistic Models Our primary goal is to learn a parametrized
390 probability distribution $p_\theta(\mathbf{x})$ on $\mathbb{R}^{C \times L}$ that is close to the original data distribution $p(\mathbf{x})$. $p_\theta(\mathbf{x})$ can then be
391 used for many different tasks like imputation, generation, or likelihood evaluation. We use denoising diffusion
392 probabilistic models (DDPM) (Ho et al., 2020):

393 DDPMs work by modeling two processes, the reverse (denoising) and forward (noising) diffusion process,
394 over a sequence of latent variables \mathbf{x}_t , $t = 1, \dots, T$, where T is a pre-specified process length. Note that T
395 is the number of diffusion steps, a hyperparameter in the generative model, and not the length of the time
396 series being modeled (L above). The forward process is given by a Markov chain

$$q(\mathbf{x}_{1:T} \mid \mathbf{x}_0) := \prod_{t=1}^T q(\mathbf{x}_t \mid \mathbf{x}_{t-1}), \quad (4.1)$$

397 where the noising distribution

$$q(\mathbf{x}_t \mid \mathbf{x}_{t-1}) := \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}) \quad (4.2)$$

398 iteratively adds Gaussian noise controlled by the noise level $\beta_t > 0$, resulting in an evolution from the data
399 distribution, $p(\mathbf{x}_0)$, to the noise distribution, $p(\mathbf{x}_T)$, over the course of T diffusion steps. This forward process
400 is chosen by the user beforehand and thus not learned.

401 As its name suggests, the reverse process reverses the forward process and is given by

$$\begin{aligned} p_\theta(\mathbf{x}_{0:T}) &:= p(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1} \mid \mathbf{x}_t), \quad \mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \\ p_\theta(\mathbf{x}_{t-1} \mid \mathbf{x}_t) &:= \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t, t), \sigma_\theta(\mathbf{x}_t, t)\mathbf{I}), \end{aligned} \quad (4.3)$$

402 i.e., the reverse conditional distribution is also Gaussian and its mean and variance are parametrized functions
403 of the current diffusion step t and the noised data point \mathbf{x}_t , and therefore need to be learned. DDPMs
404 parameterize $p_\theta(\mathbf{x}_{t-1} \mid \mathbf{x}_t)$ as a denoising distribution,

$$\begin{aligned} \mu_\theta(\mathbf{x}_t, t) &= \frac{1}{\alpha_t} \left(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(\mathbf{x}_t, t) \right), \\ \sigma_\theta(\mathbf{x}_t, t) &= \left(\frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \sqrt{\beta_t} \right)^{\frac{1}{2}}, \end{aligned} \quad (4.4)$$

405 where ϵ_θ is a neural network with learnable weights θ .

406 In other words, the neural network $\epsilon_\theta(\mathbf{x}_t, t, \text{cond})$ acts as “denoiser” that, together with the diffusion time
407 step t and, optionally, additional conditioning information cond , produces a denoised version of the data point
408 (in this case a time series) as the basis for the next denoising step.

409 The denoising network can be optimized by computing a variational lower bound on $p_\theta(\mathbf{x}_t \mid \mathbf{x}_{t-1})$, which
410 results in the full (simplified) training objective:

$$\min_{\theta} \mathcal{L}(\theta) := \min \mathbb{E}_{\mathbf{x}_0, \epsilon} \left[\|\epsilon - \epsilon_\theta(\mathbf{x}_t, t)\|_2^2 \right] \quad (4.5)$$

411 with $\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon$.

412 **Denoising Diffusion Probabilistic Models for neurophysiological time series** For our application
413 of modeling neurophysiological time series, we extend the standard DDPM setting with two modifications:

- 414** 1. First, we use a flexible convolutional neural network architecture based on structured convolutions as
415 our denoising network (Li et al., 2023). This architecture is able to accurately model long (thousands
416 of time points) and highly multivariate time series (over 100 channels).
- 417** 2. We replace the white noise process used in the standard DDPM by general Gaussian processes. This
418 allows us to include “prior information” about the modeled time series into the generative process
419 (Biloš et al., 2022). In this paper, we specifically focus on the Ornstein-Uhlenbeck process due to its
420 advantageous numerical properties and its relevance for neurophysiological recordings, in particular the
421 inverse power law scaling over frequencies.

422 Network architecture with structured convolutions Our denoising network consists of interleaved
423 layers of structured convolutions (Li et al., 2023) and linear layers that “mix” information between channels.

424 Structured long convolutions are created by concatenating several linearly interpolated and exponentially
425 downscaled convolutional kernels. The resulting convolutional kernel can be as long as the time series itself.
426 This construction is inspired by recent advances in deep state space models (Alcaraz and Strodthoff, 2022,
427 Gu et al., 2022) and allows for long yet parameter-efficient convolutional kernels that can model complex
428 and long-range dependencies within a given time series. To make the use of these large convolutional kernels
429 tractable, the convolution operation is computed by element-wise multiplication in the Fourier domain after
430 applying a Fast Fourier Transform. This approach drastically reduces the computation time for large kernels
431 and long time series.

432 The input layer is a standard convolutional layer (LeCun et al., 1995) that maps the input as well as the
433 diffusion time embedding and other conditional information into a high-dimensional latent space. We param-
434 eterize the size of the latent space per channel, i.e., we specify the number of latent dimensions per channel.
435 The full dimensionality of the latent space is then given by the number of channels times the latent dimension
436 per channel. The output layer is also a standard convolutional layer that maps the hidden activations to
437 the desired output dimensionality. The hidden layers consist of blocks of structured convolutions followed by
438 linear layers. The linear layers “mix” the output from the structured convolutions, which operate separately
439 on each hidden channel. Unlike the original implementation of structured convolutions, we introduce the
440 number of scales as a new hyperparameter. This means that the overall size of the structured convolutional
441 kernel is given by $\text{kernel_size} \cdot 2^{\text{num_scales}-1}$.

442 In our experiments, especially when training on neurophysiological recordings with high channel count,
443 we noticed that fully parametrized linear mixing layers can cause convergence issues. The number of weights
444 in the linear layers grows quadratically with the full latent dimensions and can quickly dwarf the number
445 weights in convolutional layers. To avoid this overparameterization and alleviate the convergence issues, we
446 sparsify the weight matrices of the linear layers. To do so, we decrease the number of active weights in off
447 diagonal blocks (see Appendix Fig. 5). The size of the off diagonal blocks is a hyperparameter.

448 Throughout our architecture, we use GELU activation functions (Hendrycks and Gimpel, 2016) and batch
449 normalization (Ioffe and Szegedy, 2015). Importantly, our architecture does not perform any compression in
450 time. This allows the denoising architecture to operate on time series of arbitrary length.

451 Alternative noise processes as diffusion process It is straightforward to replace the white noise process
452 commonly used in DDPMs with other Gaussian processes (Biloš et al., 2022). Instead of using white noise
453 from $\mathcal{N}(\mathbf{0}, I)$, the (forward) diffusion process is performed with noise sampled from a general Gaussian process
454 $\mathcal{N}(\mathbf{0}, \Sigma)$ with prespecified covariance Σ . The objective from Eq. (4.5) is also changed to include Σ :

$$\min_{\theta} \mathcal{L}(\theta) := \min \mathbb{E}_{\mathbf{x}_0, \epsilon} [(\epsilon - \epsilon_{\theta}(\mathbf{x}_n, n))^T \Sigma^{-1} (\epsilon - \epsilon_{\theta}(\mathbf{x}_n, n))], \quad (4.6)$$

455 where $\mathbf{x}_n = \sqrt{\bar{\alpha}_0} + (1 - \bar{\alpha}_n)\epsilon$ and $\epsilon \sim \mathcal{N}(\mathbf{0}, \Sigma)$. This replacement is useful for providing the model with
456 domain knowledge about the recordings and can lead to better empirical performance.

457 One such example of domain knowledge in the context of neurophysiological recordings is continuity.
458 Neurophysiological recordings are typically continuous and do not exhibit large and abrupt jumps. In addition,
459 for many neurophysiological recordings, the frequency spectrum exhibits inverse power law scaling (Pritchard,
460 1992).

461 In this paper, we use the Ornstein-Uhlenbeck (OU) process, which is a continuous noise process with
462 power-law frequency spectrum to encode this domain knowledge. The OU process is a stationary Gauss-
463 Markov process with many applications in physics and financial mathematics (Chan et al., 1992). As the
464 “roughest” member of the family of Matern kernels, it is continuous, but not differentiable.

465 Using the OU process is computationally feasible, even for long time series. While in general sampling
466 from Gaussian Processes scales cubically with the length of the time series, for the OU process it is possible
467 in linear time. This is because the OU process can be described equivalently by the stochastic differential
468 equation,

$$dx_{\tau} = -\rho x_{\tau} d\tau + \sigma dW_{\tau}, \quad (4.7)$$

469 where W_{τ} is the Wiener process. Furthermore, since the OU process is a stationary Gauss-Markov process,
470 its precision matrix is banded and the modified DDPM training objective can also be computed in linear

471 time. For both sampling and computing the training objective, we apply the OU process independently for
472 each time series channel. However, we find that this modification does not improve sample quality in all
473 datasets. See experiments on the effect of using the OU process versus white noise in the Appendix.

474 4.2 Model training and inference

475 Our convolutional denoising networks are trained with stochastic gradient descent on the objective given in
476 Eq. (4.6). We train all our models using the AdamW optimizer (Loshchilov and Hutter, 2019) and choose
477 hyperparameters by inspecting the quality of the generated data with respect to the training set. See the
478 Appendix for all network and training hyperparameters. To sample and impute using the trained denoising
479 network, a sample from the Gaussian distribution (white noise or OU) is gradually denoised following the
480 standard DDPM denoising procedure (Ho et al., 2020).

481 DDPMs naturally lend themselves to imputation by guiding the diffusion process based on observed data
482 points (Ho et al., 2020). For image-based diffusion models this process is usually known as “inpainting”
483 (Lugmayr et al., 2022). It works by replacing the observed part of the image—or time series—with the
484 analytical latent state of the forward diffusion process after each denoising step. This way, the reverse
485 diffusion process is “guided” along the observed part of the time series. However, this approach does not
486 provide the correct conditional distributions of the form $p(\mathbf{x}_{\text{observed}} | \mathbf{x}_{\text{missing}})$, which can hurt imputation
487 performance (Tashiro et al., 2021).

488 To alleviate this issue, it is possible to train conditional diffusion models directly: Here, the condition is
489 given as an input mask to the denoising network to compute the correct conditional distribution. To account
490 for different patterns of missingness at imputation time, many condition input masks need to be sampled
491 during training. In this work, we use the random sampling strategy (Tashiro et al., 2021) and randomly drop
492 out channels for a given time series following a two step process: First, we uniformly sample the number of
493 channels to drop, and then uniformly drop out channels according to this number. The denoising network
494 is then conditioned on the remaining channels and the training objective is computed. We use this masked-
495 based conditional training only for the AJILE12 dataset. For the other two datasets, the models are trained
496 without masks. All models are implemented and trained with the PyTorch library (Paszke et al., 2019).

497 4.3 Datasets

498 We perform our three experiments on three different datasets of publicly available neurophysiological recordings.

500 The dataset for our first experiment (Fig. 2) is a three-electrode LFP recording from rats recorded at
501 600 Hz at the medial prefrontal cortex (mPFC), hippocampal CA1 region, and the thalamic nucleus reunions
502 (RE) (Varela and Wilson, 2020)¹. It consists of a total of 3.5 hours of recordings from three different rats.
503 During the recording, rats were freely behaving with the possibility to rest and underwent cycles of sleep
504 and wakefulness. For our experiment, we used the LFP data from the first and third recording (52 and 93
505 minutes).

506 The AJILE12 dataset used for our second experiment (Fig. 3) consists of human ECoG data recorded
507 during naturalistic behavior from 12 different participants undergoing epilepsy monitoring (Peterson et al.,
508 2022)². It was recorded at 500 Hz over several days, and a total of 1280 hours. The length of recording per
509 participant ranges from 70 to 120 hours. The number of electrodes per participant ranges from 64 to 126. In
510 addition to the neurophysiological recordings, a variety of behavioral and movement event related metadata
511 are provided.

512 The dataset for our third experiment (Fig. 4) consists of ECoG recordings from non-human primates
513 (macaque monkey) recorded at 1000 Hz via surface 128 electrodes placed inside the cranium, directly on top
514 of the cortex (Yanagawa et al., 2013)³. We use data from only one macaque (Chibi). During the recording, the
515 macaques were injected with an anesthetic (propofol) and different neurophysiological states were annotated.
516 We use the data from the anesthetized and awake-eyes-closed condition (73 min of recording).

¹<https://crcns.org/data-sets/hc/hc-24>

²<https://dandiarchive.org/dandiset/000055/0.220127.0436>

³<http://www.neurotycho.org/anesthesia-task>

517 For our initial example (Fig. 1E-H), we use data from participant P02 of the BCI Challenge @ NER 2015
 518 ([Margaux et al., 2012](#))⁴. The data was recorded at 200Hz with 56 passive EEG sensors placed with the
 519 extended 10-20 system. The participants were tested under the “P-300 Speller” paradigm, where words are
 520 spelled out letter-by-letter by flashing screen items and measuring the evoked response ([Farwell and Donchin,](#)
 521 [1988](#)).

522 4.4 Data preprocessing

523 For our datasets of rat LFP and macaque ECoG recordings, we perform channel-by-channel normalization by
 524 standardizing the recording over its entire duration. After normalization, we create the set of time windows
 525 that are used to train and evaluate the model. For the rat LFP recordings, we took 2-second time windows
 526 consisting of 1200 time points resulting in a total of 4350 time series. Similarly, for the macaque ECoG
 527 dataset, we took 1-second time windows consisting of 1000 time points resulting in a total of 1088 time series.
 528 After creating the time series, we then randomly divide them into a training and an evaluation set (80%
 529 training, 20% evaluation).

530 The data of each participant in the AJILE12 dataset consists of seven days of recording. We follow
 531 the setup and preprocessing described in [Peterson et al. \(2021\)](#): The ECoG traces are downsampled before
 532 extracting time windows corresponding to either movement or rest. For each participant, the final dataset
 533 consists of 4-second time windows with a sampling frequency of 250 Hz with an equal amount of time
 534 series recorded under rest and movement (i.e., class-balanced). Again, we normalize the time series channel-
 535 by-channel by aggregating over all time series in the training set. The normalization of the training set,
 536 which consists of the first six days of recording, is applied to the evaluation set, i.e., the seventh day of
 537 recording. Furthermore, extreme outlier channels, where the difference between the minimum and maximum
 538 value is five or more times the interquartile range, are detected and masked during training and evaluation.
 539 The number of recorded channels varies widely between participants, ranging from 64 to 126. Similarly, the
 540 number of training and evaluation time series varies considerably between participants. However, the network
 541 hyperparameters are shared between participants and no attempt was made to fine-tune the imputation for
 542 an individual participant.

543 Unlike the DDPM, the random forest neural decoder is trained and evaluated on the unnormalized
 544 AJILE12 data (see [Peterson et al. \(2021\)](#) for details). Imputations are obtained by imputing normalized
 545 time series and then reversing the normalization for both the observed and imputed channels using the
 546 normalization constants from the training data.

547 For our EEG dataset, we extract 1.3 seconds of the recording immediately following a presented stimulus.
 548 Since we only use data from participant P02, this results in 340 time series (i.e., trials) consisting of 260
 549 time points each. After extraction, we band-pass filter the time series with a frequency window between 1
 550 and 40 Hz using a 5th order Butterworth filter. We then again perform channel-by-channel normalization by
 551 aggregating over the time series.

552 4.5 Analyses and Metrics

553 **Power spectral densities and spectral error** Throughout the experiments, we compare the distributions
 554 of (cross-) power spectral densities (PSDs) between channels of real and generated data. PSDs are computed
 555 with a Fourier transform over the whole the time series. For a given one-dimensional time series x and
 556 frequency f , the PSD is denoted by $S_{xx}(f)$. Since we are always interested in capturing the full distribution
 557 of PSDs over a set of time series, we compute the point-wise median and percentiles: Here, point-wise means
 558 that, for a set of real or generated time series $\{x_i\}_{i=1}^N$, the median $\text{median}[S_{xx}(f)]$ over all PSDs is computed
 559 at each frequency f , and similarly for the 10% and 90% percentiles.

560 To quantify how well the distribution of PSDs of the real and generated data match in a single metric,
 561 we compute the spectral error in log space. The spectral error is defined as

$$ES \left(\{x_i\}_{i=1}^N, \{y_i\}_{i=1}^N \right) := \sqrt{\sum_{f \in F} \left(\text{median}_{\{x_i\}_{i=1}^N} [\log S_{x_i x_i}(f)] - \text{median}_{\{y_i\}_{i=1}^N} [\log S_{y_i y_i}(f)] \right)^2}, \quad (4.8)$$

⁴<https://www.kaggle.com/competitions/inria-bci-challenge/>

562 where $f \in F$ are frequency bins and $\text{median}[\log S_{xx}(f)]$ is the median log-power at a fixed frequency f over
563 all generated (y_i) or real (x_i) time series. All computations related to the power spectrum are performed
564 with the numpy (Harris et al., 2020) and Scipy (Virtanen et al., 2020) libraries.

565 **Sharp wave-ripple detection and cross-channel couplings** For our experiments with the rat LFP
566 data, we detect ripples in the CA1 channel. Following the heuristic approach of Varela and Wilson (2020), we
567 perform this detection by band-pass filtering both real and generated CA1 traces within the ripple frequency
568 band from 100 to 275 Hz using a finite impulse response (FIR) filter. Ripples are detected when the power
569 of the time series is more than three standard deviations greater than the mean. The mean and standard
570 deviation of the filtered power are calculated over the entire set of real time series. Crucially, to ensure a fair
571 comparison between a set of real and a set of generated CA1 traces, the mean and standard deviation from
572 the real data are used to detect SWRs in the generated data. SWRs with less than 10 time points (17ms) are
573 fused. After fusing, a SWR must have a minimum length of 20 time points (33ms). Detected ripples below
574 this threshold are discarded.

575 In our SWR prediction experiment, we impute the CA1 channel given the other two channels and detect
576 SWRs in the real and imputed channel for all time windows in the evaluation set. If there are no SWRs or at
577 least one SWR in both the real and generated channel, the prediction is counted as correct. We then compute
578 the F1 score based on the number of correct and incorrect predictions. To test if the prediction performance
579 is significantly better than random, we shuffle the set of generated time series and evaluate the F1 score. This
580 procedure is repeated a 1000 times to compute the permutation test statistic (Ojala and Garriga, 2010).

581 We also use detected SWRs to compute a “phase-occurrence coupling” against the phase of delta (1–4 Hz),
582 theta (4–9 Hz), and spindle (6–14 Hz) frequencies in the mPFC. The mPFC phase is calculated using the
583 Hilbert transform after band-pass filtering (FIR, Hamming) with the corresponding frequency band. Given
584 the mPFC phase, the time point with maximal power within a detected SWR is used to extract the phase of
585 the SWR occurrence, and the counts are binned into 21 equally spaced phase-bins (from -180 to 180 degrees)
586 to create the histograms (Fig. 2D).

587 Finally, we compute the phase-amplitude coupling (PAC) between the CA1 amplitude in the ripple band
588 (100–275 Hz) and the mPFC phase of the previously mentioned driver frequencies. Phase and amplitude
589 in the respective frequency bands are calculated using the Hilbert transform, and the average amplitudes
590 occurring at 31 equally spaced phase-bins (from -180 to 180 degrees) are reported (Fig. 2D).

591 **Details on the imputation experiment and neural decoder** For our imputation experiment, we follow
592 the setup described in Talukder et al. (2022): We impute artificially dropped channels and investigate the
593 effect of this imputation on the performance of a neural decoder. As a baseline imputation method, we use
594 zero imputation. The neural decoding task is to discriminate between time series recorded at rest and time
595 series recorded during movement. The neural decoder used is a random forest, which works by “flattening” a
596 time series into a large feature vector. We train three decoders over three training folds as in Peterson et al.
597 (2021) using the same hyperparameters (maximum tree depth and number of estimators). No attempt was
598 made to further optimize the random forest models. Finally, the average test performance over all training
599 folds is calculated for all settings, i.e., fully observed, zero imputation, and DDPM-based imputation.

600 To evaluate the performance of DDPM-based imputation, we randomly select 50%, 70%, or 90% of the
601 channels from a given participant. For each time series in the evaluation set, this set of channels is dropped
602 and then imputed with zeroes or imputed using the trained DDPM. We then apply the previously trained
603 decoder to both imputed evaluation sets and compute the evaluation accuracy. This process is repeated
604 5 times for each missingness level to evaluate the imputation performance over different sets of dropped
605 channels. We report the mean and standard deviation over the 5 repetitions (Fig. 3).

606 **Likelihood computation and brain-state classification** Likelihood computation in DDPMs can be
607 realized by solving the probability flow equation (Song et al., 2021b). In our class-conditional model for the
608 macaque ECoG data, we compute log-likelihoods for both the awake and anesthetized states. That is, given
609 a time series \mathbf{x} , we get $\log p_\theta(\mathbf{x}|c = \text{awake})$ and $\log p_\theta(\mathbf{x}|c = \text{anesthetized})$. These values can be used to
610 perform classification by selecting the class with the higher log-likelihood for a given sample. We perform
611 this classification for all time series in our class-balanced evaluation set and report the overall classification

accuracy. For the time-resolved brain state classification in Fig. 4E, we simply compute this log-likelihood difference over 1-second sliding windows.

For outlier detection, we are interested in the total log-likelihood $\log p_\theta(\mathbf{x})$ of sample \mathbf{x} . We compute this total log-likelihood for each evaluation time series and correspondingly generated outliers (white noise or time-wise flipped) according to the law of total probability. Then, the area under the receiver operator characteristic (AUC) of the time series against both types of outliers is computed.

For the brain state classification with synthetic data, we use the convolutional neural network classifier described in Wang et al. (2017). The classifier is trained twice, once on real data and once on an equal amount of DDPM-generated data, and then evaluated on a held-out evaluation set of real data. Both the training and evaluation datasets consist of an equal number of awake and anesthetized time windows, and we compute the evaluation accuracy to compare the performance between the real and synthetic settings. The classifier is trained using the AdamW optimizer (Loshchilov and Hutter, 2019) and binary cross entropy loss with identical network and training hyperparameters in both settings.

5 Acknowledgments

JV is supported by the International Max Planck Research School for Intelligent Systems (IMPRS-IS) and the AI4Med-BW graduate program. RG is supported by the European Union's Horizon 2020 research and innovation program under the Marie Skłodowska-Curie grant agreement No. 101030918 (AutoMIND). JHM, RG are members of the Machine Learning Cluster of Excellence, EXC number 2064/1-390727645. This work was supported the German Federal Ministry of Education and Research (BMBF): Tübingen AI Center, FKZ: 01IS18039A. We would like to thank Auguste Schulz and Matthijs Pals for feedback on the manuscript and Mackelab members for discussions throughout the project.

References

- Samira Abbasi, Selva Maran, and Dieter Jaeger. A general method to generate artificial spike train populations matching recorded neurons. *Journal of Computational Neuroscience*, 2020.
- Juan Lopez Alcaraz and Nils Strodthoff. Diffusion-based time series imputation and forecasting with structured state space models. *Transactions on Machine Learning Research*, 2022.
- Marin Biloš, Kashif Rasul, Anderson Schneider, Yuriy Nevmyvaka, and Stephan Günnemann. Modeling temporal data as continuous functions with process diffusion. *arXiv preprint arXiv:2211.02590*, 2022.
- Kalok C Chan, G Andrew Karolyi, Francis A Longstaff, and Anthony B Sanders. An empirical comparison of alternative models of the short-term interest rate. *The Journal of Finance*, 1992.
- Zijiao Chen, Jiaxin Qing, Tiange Xiang, Wan Lin Yue, and Juan Helen Zhou. Seeing beyond the brain: Conditional diffusion model with sparse masked modeling for vision decoding. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023.
- Prafulla Dhariwal and Alexander Nichol. Diffusion models beat GANs on image synthesis. *Advances in Neural Information Processing Systems*, 2021.
- Thomas Donoghue, Matar Haller, Erik J Peterson, Paroma Varma, Priyadarshini Sebastian, Richard Gao, Torben Noto, Antonio H Lara, Joni D Wallis, Robert T Knight, Avgusta Shestyuk, and Bradley Voytek. Parameterizing neural power spectra into periodic and aperiodic components. *Nature Neuroscience*, 2020.
- Lawrence Ashley Farwell and Emanuel Donchin. Talking off the top of your head: Toward a mental prosthesis utilizing event-related brain potentials. *Electroencephalography and clinical Neurophysiology*, 1988.
- Vincent Fortuin, Dmitry Baranchuk, Gunnar Rätsch, and Stephan Mandt. GP-VAE: Deep probabilistic time series imputation. *International Conference on Artificial Intelligence and Statistics*, 2020.

- 654 Cecilia Gallego-Carracedo, Matthew G Perich, Raeed H Chowdhury, Lee E Miller, and Juan Álvaro Gallego.
655 Local field potentials reflect cortical population dynamics in a region-specific and frequency-dependent
656 manner. *eLife*, 2022.
- 657 Richard Gao, Erik J Peterson, and Bradley Voytek. Inferring synaptic excitation/inhibition balance from
658 field potentials. *NeuroImage*, 2017.
- 659 Alessandro T Gifford, Kshitij Dwivedi, Gemma Roig, and Radoslaw M Cichy. A large and rich EEG dataset
660 for modeling human visual object recognition. *NeuroImage*, 2022.
- 661 Albert Gu, Karan Goel, and Christopher Ré. Efficiently modeling long sequences with structured state spaces.
662 *International Conference on Learning Representations*, 2022.
- 663 Diego A Gutnisky and Krešimir Josić. Generation of spatiotemporally correlated spike trains and local field
664 potentials using a multivariate autoregressive process. *Journal of Neurophysiology*, 2010.
- 665 Charles R Harris, K Jarrod Millman, Stéfan J Van Der Walt, Ralf Gommers, Pauli Virtanen, David Cour-
666 napeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J Smith, et al. Array programming with
667 NumPy. *Nature*, 2020.
- 668 Biyu J He, John M Zempel, Abraham Z Snyder, and Marcus E Raichle. The temporal structures and
669 functional significance of scale-free brain activity. *Neuron*, 2010.
- 670 Dan Hendrycks and Kevin Gimpel. Gaussian error linear units. *arXiv preprint arXiv:1606.08415*, 2016.
- 671 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural*
672 *Information Processing Systems*, 2020.
- 673 Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing
674 internal covariate shift. *International Conference on Machine Learning*, pages 448–456, 2015.
- 675 Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based
676 generative models. *Advances in Neural Information Processing Systems*, 2022.
- 677 Zhifeng Kong, Wei Ping, Jiaji Huang, Kexin Zhao, and Bryan Catanzaro. Diffwave: A versatile diffusion
678 model for audio synthesis. *International Conference on Learning Representations*, 2021.
- 679 Michael Krumin and Shy Shoham. Generation of spike trains with controlled auto-and cross-correlation
680 functions. *Neural Computation*, 2009.
- 681 Yann LeCun, Yoshua Bengio, et al. Convolutional networks for images, speech, and time series. *The handbook*
682 *of brain theory and neural networks*, 1995.
- 683 Yuhong Li, Tianle Cai, Yi Zhang, Deming Chen, and Debadeepta Dey. What makes convolutional models
684 great on long sequence modeling? *International Conference on Learning Representations*, 2023.
- 685 Lequan Lin, Zhengkun Li, Ruikun Li, Xuliang Li, and Junbin Gao. Diffusion models for time series applica-
686 tions: A survey. *arXiv preprint arXiv:2305.00624*, 2023.
- 687 Sikun Lin, Thomas Sprague, and Ambuj K Singh. Mind reader: Reconstructing complex images from brain
688 activities. *Advances in Neural Information Processing Systems*, 2022.
- 689 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *International Conference on*
690 *Learning Representations*, 2019.
- 691 Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. DPM-Solver: A fast ODE solver
692 for diffusion probabilistic model sampling in around 10 steps. *Advances in Neural Information Processing*
693 *Systems*, 2022.

- 694 Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. Repaint:
695 Inpainting using denoising diffusion probabilistic models. *Proceedings of the IEEE/CVF Conference on*
696 *Computer Vision and Pattern Recognition*, 2022.
- 697 Jakob H Macke, Philipp Berens, Alexander S Ecker, Andreas S Tolias, and Matthias Bethge. Generating
698 spike trains with specified correlation coefficients. *Neural Computation*, 2009.
- 699 Perrin Margaux, Maby Emmanuel, Daligault Sébastien, Bertrand Olivier, and Mattout Jérémie. Objective
700 and subjective evaluation of online error correction during P300-based spelling. *Advances in Human-*
701 *Computer Interaction*, 2012.
- 702 Manuel Molano-Mazon, Arno Onken, Eugenio Piasini, and Stefano Panzeri. Synthesizing realistic neural
703 population activity patterns using generative adversarial networks. *International Conference on Learning*
704 *Representations*, 2018.
- 705 Markus Ojala and Gemma C. Garriga. Permutation tests for studying classifier performance. *Journal of*
706 *Machine Learning Research*, 2010.
- 707 Chethan Pandarinath, Daniel J O'Shea, Jasmine Collins, Rafal Jozefowicz, Sergey D Stavisky, Jonathan C
708 Kao, Eric M Trautmann, Matthew T Kaufman, Stephen I Ryu, Leigh R Hochberg, et al. Inferring single-
709 trial neural population dynamics using sequential auto-encoders. *Nature Methods*, 2018.
- 710 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen,
711 Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep
712 learning library. *Advances in neural information processing systems*, 2019.
- 713 Steven M Peterson, Zoe Steine-Hanson, Nathan Davis, Rajesh PN Rao, and Bingni W Brunton. General-
714 ized neural decoders for transfer learning across participants and recording modalities. *Journal of Neural*
715 *Engineering*, 2021.
- 716 Steven M Peterson, Satpreet H Singh, Benjamin Dichter, Michael Scheid, Rajesh PN Rao, and Bingni W
717 Brunton. AJILE12: Long-term naturalistic human intracranial neural recordings and pose. *Scientific data*,
718 2022.
- 719 Carlos R Ponce, Will Xiao, Peter F Schade, Till S Hartmann, Gabriel Kreiman, and Margaret S Livingstone.
720 Evolving images for visual neurons using a deep generative network reveals coding principles and neuronal
721 preferences. *Cell*, 2019.
- 722 W S Pritchard. The brain in fractal time: 1/f-like power spectrum scaling of the human electroencephalogram.
723 *The International journal of neuroscience*, 1992.
- 724 Poornima Ramesh, Mohamad Atayi, and Jakob H Macke. Adversarial training of neural encoding models on
725 population spike trains. *Real Neurons & Hidden Units: Future directions at the intersection of neuroscience*
726 *and artificial intelligence @ NeurIPS 2019*, 2019.
- 727 Kashif Rasul, Calvin Seward, Ingmar Schuster, and Roland Vollgraf. Autoregressive denoising diffusion models
728 for multivariate probabilistic time series forecasting. *International Conference on Machine Learning*, 2021.
- 729 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution
730 image synthesis with latent diffusion models. *Proceedings of the IEEE/CVF Conference on Computer*
731 *Vision and Pattern Recognition*, 2022.
- 732 Benjamin Sanchez-Lengeling and Alán Aspuru-Guzik. Inverse molecular design using machine learning:
733 Generative models for matter engineering. *Science*, 2018.
- 734 Thomas Schreiber and Andreas Schmitz. Surrogate time series. *Physica D: Nonlinear Phenomena*, 2000.
- 735 Ikaro Silva, George Moody, Daniel J Scott, Leo A Celi, and Roger G Mark. Predicting in-hospital mortality
736 of icu patients: The physionet/computing in cardiology challenge 2012. *Computing in Cardiology*, 2012.

- 737 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *International Conference on Learning Representations*, 2021a.
- 739 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole.
740 Score-based generative modeling through stochastic differential equations. *International Conference on Learning Representations*, 2021b.
- 742 He Sun and Katherine L Bouman. Deep probabilistic imaging: Uncertainty quantification and multi-modal
743 solution characterization for computational imaging. *Proceedings of the AAAI Conference on Artificial Intelligence*, 2021.
- 745 Yu Takagi and Shinji Nishimoto. High-resolution image reconstruction with latent diffusion models from
746 human brain activity. *bioRxiv*, 2022.
- 747 Sabera Talukder, Jennifer J Sun, Matthew Leonard, Bingni W Brunton, and Yisong Yue. Deep neural
748 imputation: A framework for recovering incomplete brain recordings. *NeurIPS 2022 Workshop on Learning from Time Series for Health*, 2022.
- 750 Yusuke Tashiro, Jiaming Song, Yang Song, and Stefano Ermon. CSDI: Conditional score-based diffusion
751 models for probabilistic time series imputation. *Advances in Neural Information Processing Systems*, 2021.
- 752 Gerrit van den Burg and Chris Williams. On memorization in probabilistic deep generative models. *Advances in Neural Information Processing Systems*, 2021.
- 754 Carmen Varela and Matthew A Wilson. Simultaneous extracellular recordings from midline thalamic nuclei,
755 medial prefrontal cortex and CA1 from rats cycling through bouts of sleep and wakefulness. *CRCNS.org*,
756 2019.
- 757 Carmen Varela and Matthew A Wilson. mPFC spindle cycles organize sparse thalamic activation and recently
758 active CA1 cells during non-REM sleep. *eLife*, 2020.
- 759 Victor Venema, Felix Ament, and Clemens Simmer. A stochastic iterative amplitude adjusted Fourier trans-
760 form algorithm with improved accuracy. *Nonlinear Processes in Geophysics*, 2006.
- 761 Pauli Virtanen, Ralf Gommers, Travis E Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau, Evgeni
762 Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, et al. SciPy 1.0: Fundamental algorithms
763 for scientific computing in python. *Nature Methods*, 2020.
- 764 Saurabh Vyas, Matthew D Golub, David Sussillo, and Krishna V Shenoy. Computation through neural
765 population dynamics. *Annual Review of Neuroscience*, 2020.
- 766 Jason Walonoski, Mark Kramer, Joseph Nichols, Andre Quina, Chris Moesel, Dylan Hall, Carlton Duffett,
767 Kudakwashe Dube, Thomas Gallagher, and Scott McLachlan. Synthea: An approach, method, and software
768 mechanism for generating synthetic patients and the synthetic electronic health care record. *Journal of the American Medical Informatics Association*, 2018.
- 770 Zhiguang Wang, Weizhong Yan, and Tim Oates. Time series classification from scratch with deep neural
771 networks: A strong baseline. *2017 International joint conference on neural networks (IJCNN)*, 2017.
- 772 Toru Yanagawa, Zenas C Chao, Naomi Hasegawa, and Naotaka Fujii. Large-scale information flow in conscious
773 and unconscious states: An ECoG study in monkeys. *PloS one*, 2013.
- 774 Jinsung Yoon, Daniel Jarrett, and Mihaela Van der Schaar. Time-series generative adversarial networks.
775 *Advances in Neural Information Processing Systems*, 2019.

776 6 Appendix

777 6.1 Ornstein-Uhlenbeck process as diffusion process

778 In our experiments, we sometimes use the OU process instead of the classically used white noise as our
 779 diffusion process. This substitution is possible because the OU process belongs to the family of Gaussian
 780 processes. See [Biloš et al. \(2022\)](#) for more details.

781 The training and sampling schemes are shown in Algorithm 1 and Algorithm 2 and differ slightly from
 782 those given in [Biloš et al. \(2022\)](#). Note that the extension to general Gaussian processes has a computational
 783 overhead compared to the original training objective. During training, the computation of the Mahalanobis
 784 distance scales quadratically with the length of the input. Sampling from a Gaussian process also requires
 785 computing a Cholesky decomposition. However, this decomposition only needs to be computed once, and can
 786 also be used to compute the Mahalanobis efficiently by solving a lower triangular linear system based on the
 787 Cholesky decomposition.

Algorithm 1 Training

```

1:  $L = \text{Cholesky}(\Sigma)$ 
2: repeat
3:  $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ 
4:  $t \sim \text{Uniform}(\{1, \dots, T\})$ 
5:  $\epsilon \sim \mathcal{N}(0, \Sigma)$ 
6: Take gradient descent step on
    $\nabla_{\theta}(\epsilon - \epsilon_{\theta}(\mathbf{x}_t, t))^T \Sigma^{-1} (\epsilon - \epsilon_{\theta}(\mathbf{x}_t, t))$ 
7: until converged

```

Algorithm 2 Sampling

```

1:  $\mathbf{x}_T \sim \mathcal{N}(0, \Sigma)$ 
2: for  $t = T, \dots, 1$  do
3:    $z \sim \mathcal{N}(0, \Sigma)$  if  $t > 1$ , else  $z = 0$ 
4:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} (\mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\alpha_t}} \epsilon_{\theta}(\mathbf{x}_t, t)) + \sigma_t z$ 
5: end for
6: return  $\mathbf{x}_0$ 

```

788 As described in the main text, the computational requirements scale more favorably for the special case
 789 of the OU process. Its formulation as a stochastic differential equation allows sampling in linear time.
 790 Furthermore, the OU process has a tridiagonal precision matrix, which allows the the Mahalanobis distance
 791 in the training objective to be computed in linear time as well.

792 6.2 Architecture details and hyperparameter settings

793 The hyperparameters used for our experiments are given in Table 1, Table 2, Table 3 and Table 4. An
 794 illustration of the sparsified linear mixing layers is given in Fig. 5.

795 6.3 Additional results

796 Here we provide some additional results. In addition to PSDs, we can also compute cross-PSDs between
 797 different channels. For the LFP data recorded from rats, the median cross-PSDs match as well (Fig. 6).

798 Furthermore, for all datasets, the distribution of Euclidean distances between real and generated data is
 799 similar to the distribution of distances within. Crucially, the minimum distance between and within the real
 800 and generated time series is away from zero (Fig. 7). Therefore, the trained DDPMs do not overfit to the
 801 training data.

802 We also provide all distributions of PSDs of AJILE12 participant P07 in the main text (Fig. 10). Addi-
 803 tionally, all distributions of PSDs for two other participants, participant P01 with 94 channels (Fig. 11) and
 804 participant P12 with 126 channels (Fig. 12) are provided. In both cases, the median spectra as well as 10%
 805 and 90% percentiles match well. Recall that the neural decoding performance on data imputed by our model

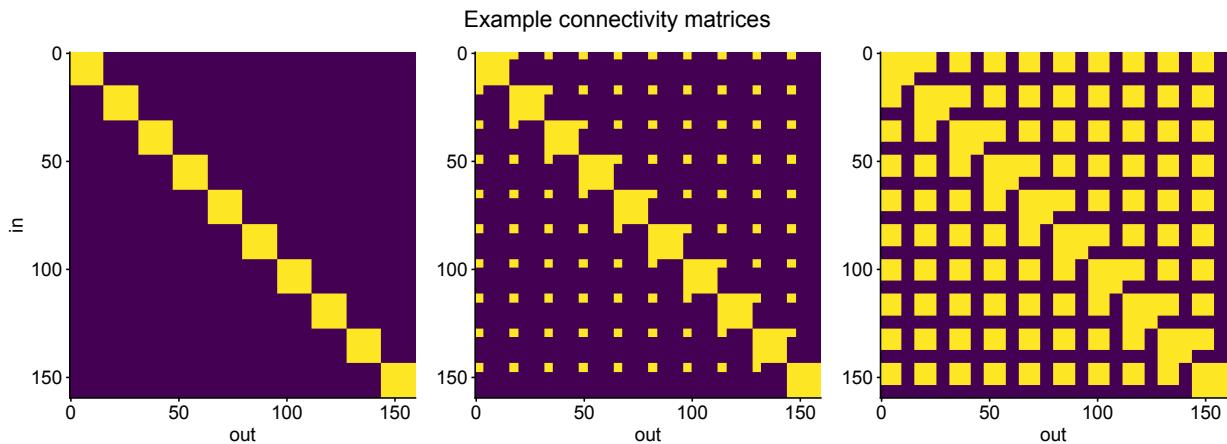


Figure 5: The linear “mixing” layers that follow the structured convolutions are parametrized as block diagonal matrices. To avoid overparameterization of the network for neurophysiological recordings with very large channel counts, we sparsify the corresponding matrices by reducing the size of the off-diagonal blocks. Three examples with increasing off-diagonal block size are shown. Yellow indicates that the weight is part of the network, blue indicates no connection.

806 did not improve over the zero imputation baseline for participant P01. Nevertheless, the overall spectra of the
 807 real and imputed data match. This means that while the model was able to capture the overall distribution
 808 of time series, it failed to learn the class-conditional mapping between observed and missing channels.

809 As an additional experiment with high channel count, we generated time series from the awake condition
 810 of the macaque ECoG data on all 128 recorded channels. The distribution of real and generated PSDs closely
 811 match across all channels, but our model sometimes overestimates the power in the low-frequency region
 812 (Fig. 9).

813 **OU process versus white noise** Our use of the OU process provides a modification that can further
 814 improve the quality of generated samples in some cases. Here, we provide an analysis in terms of the spectral
 815 error for the AJILE12 dataset and the unconditional generation of the macaque ECoG data (awake only) on
 816 the full 128 channels. We retrain the models using the standard independent Gaussian process (white noise).
 817 All other training and network hyperparameters are kept fixed.

818 For both the (awake-only) macaque ECoG data and all AJILE12 participants, the spectral errors are
 819 significantly larger when using the white noise instead of the OU process (Fig. 13). This difference in
 820 performance is largely due to the low power, high frequency noise that is overestimated when white noise is
 821 used. However, for the rat LFP, the awake/anesthetized macaque ECoG data, and the BCI EEG data, there
 822 was no advantage to using the OU process over white noise.

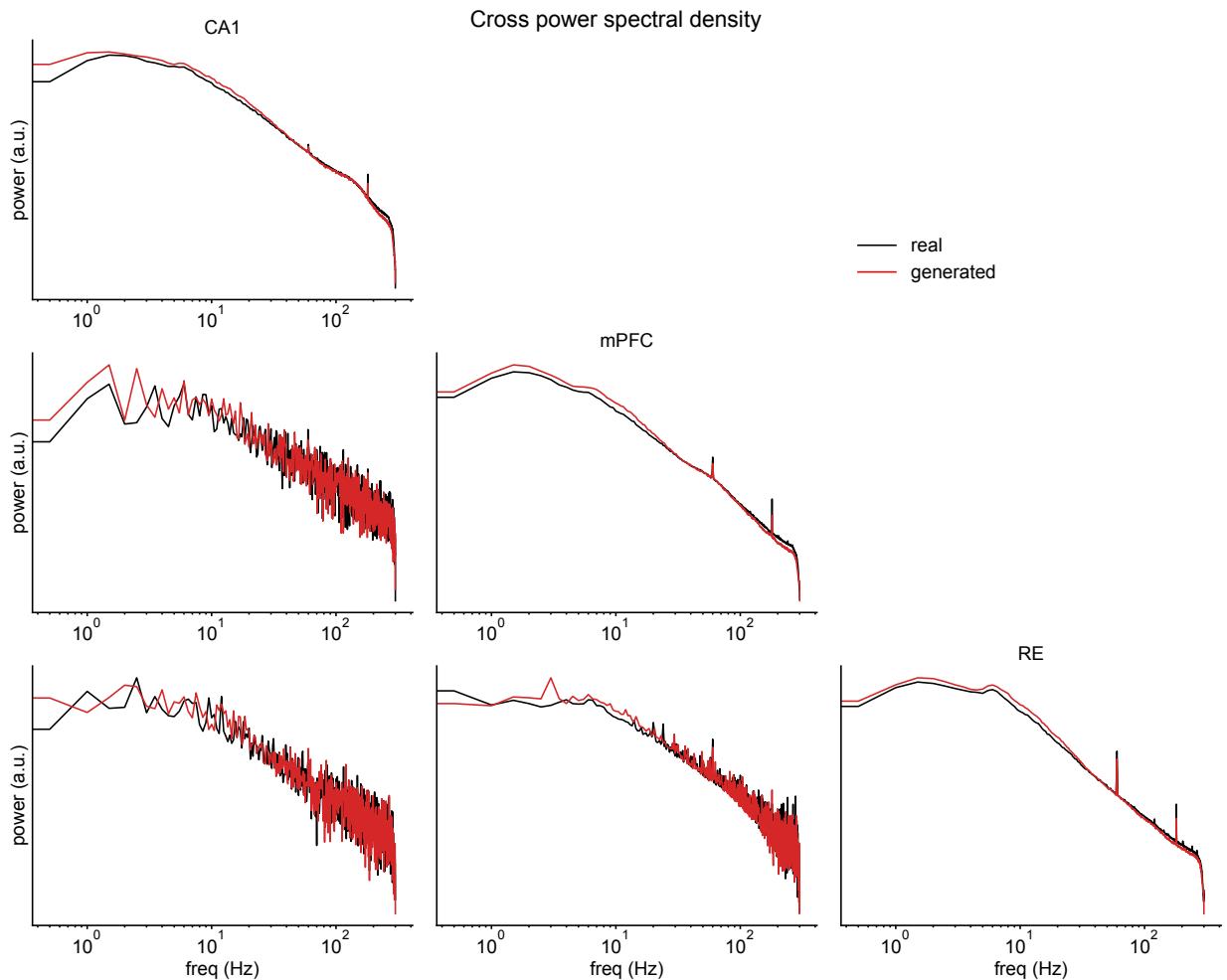


Figure 6: Median (cross-)PSDs for the rat LFP dataset for real and generated data.

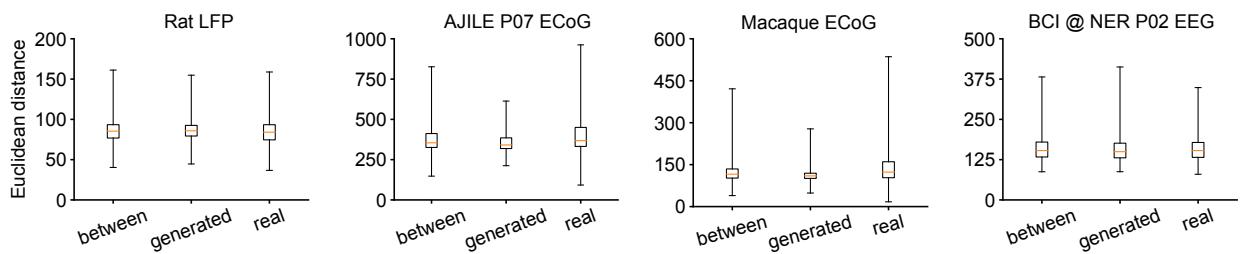


Figure 7: Distribution of distances between all real and generated time series, as well as within all real and generated time series. The box shows the median and 25% / 75% percentiles, and the whiskers show the minimum and maximum distance. The minimum distance between and within real and generated time series is away from zero.

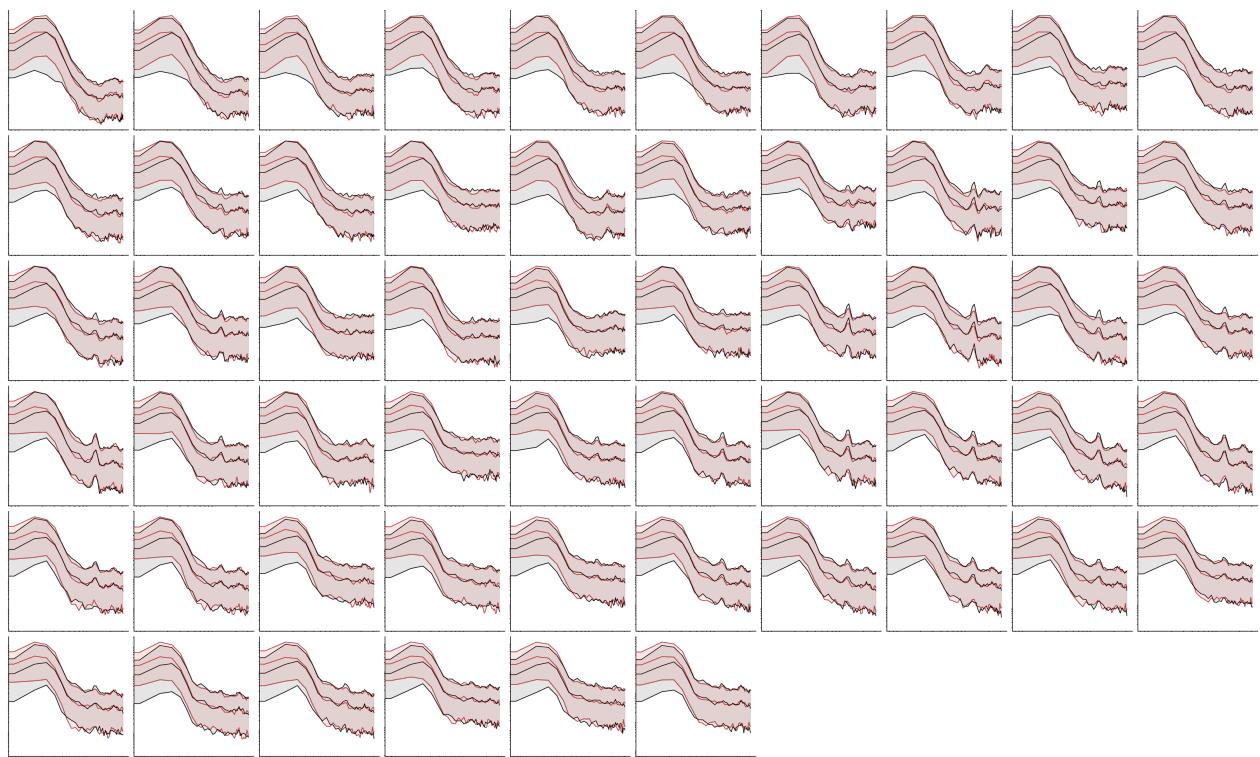


Figure 8: Full spectra of 56 channel EEG BCI data. Median power as well as 10% / 90% percentiles are shown.

Hyperparameter	Value
In kernel size	32
Out kernel size	32
SC layers	3
SC kernel size	32
SC scales	5
Decay min	1.0
Decay max	4.0
Heads	3
Latent dim. per channel	64
Off-diag. size	64
Noise process	White noise
Schedule	Quadratic
Diffusion steps T	500
B0	0.0001
B1	0.06
Time embedding dim.	16
Loss function	Mean-squared error (MSE)
Optimizer	AdamW
Learning rate	0.0001
Weight decay	0.01
Train batch size	32
Epochs	500

Table 1: Hyperparameters used to train models on the rat LFP dataset for simulation.

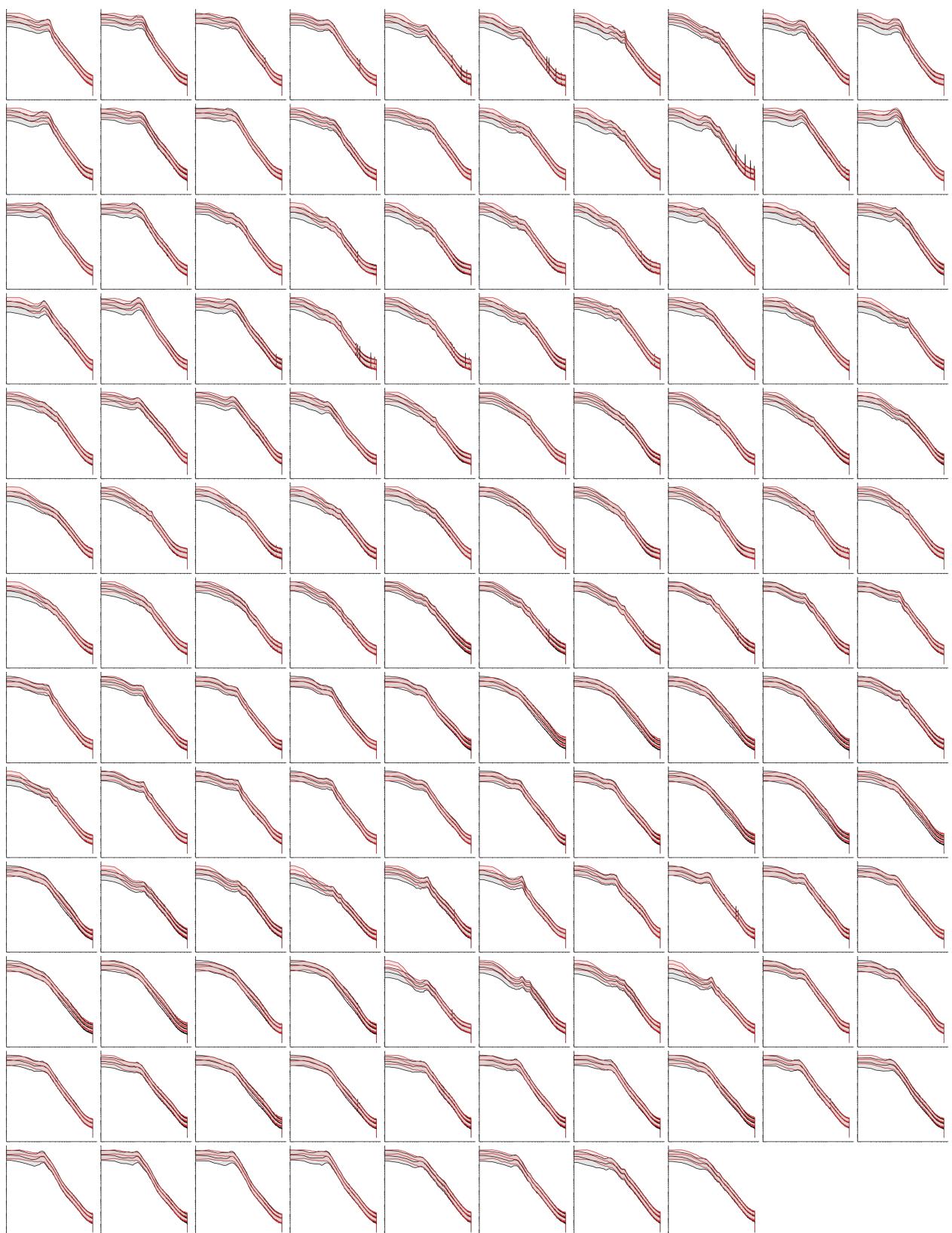


Figure 9: Full spectra of 128 channel macaque data in awake condition. Median power as well as 10% / 90% percentiles are shown.

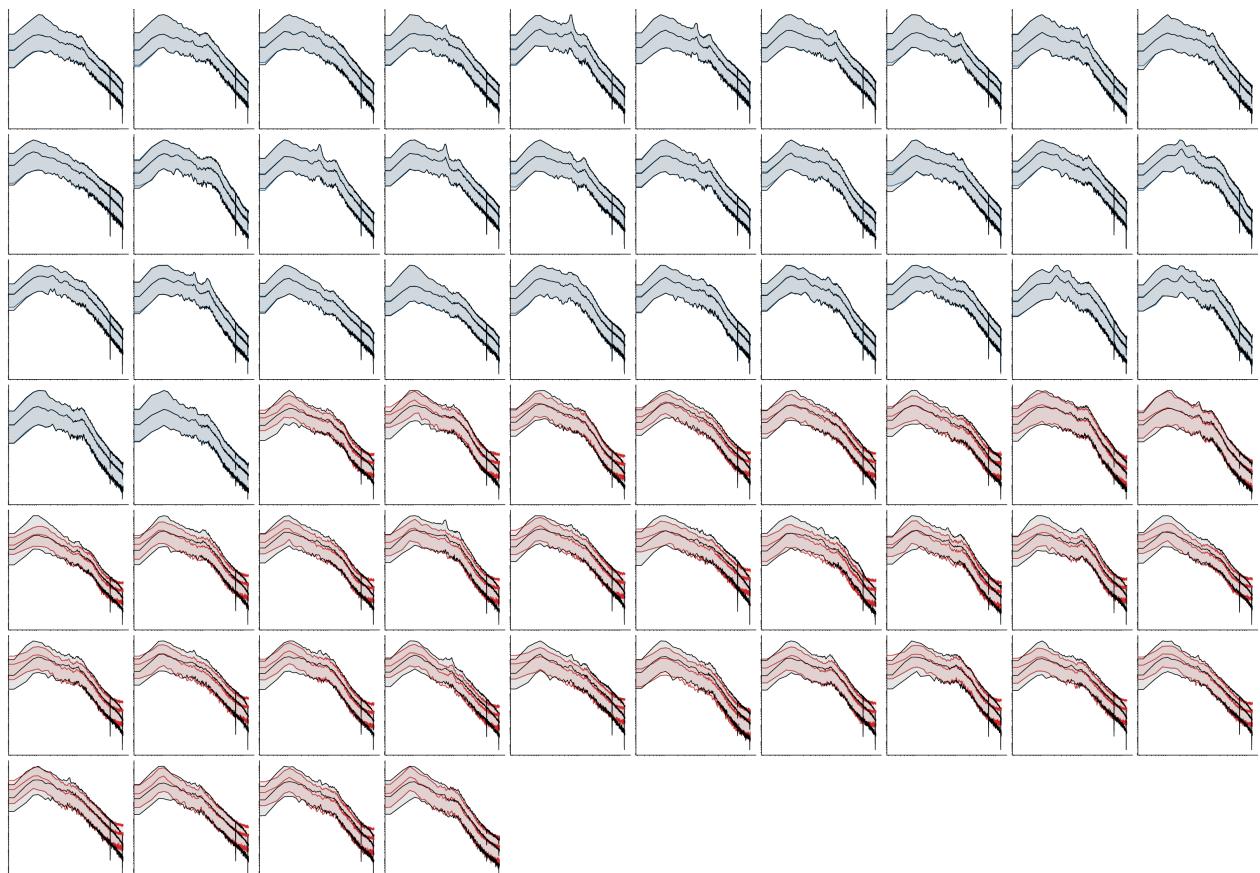


Figure 10: Full spectra of P07 from the AJILE12 dataset shown in the main text. The LFP traces consist of 64 channels. Here, we imputed the second half of the channels given the first half for all time series in the evaluation set. Median power as well as 10% / 90% percentiles are shown.

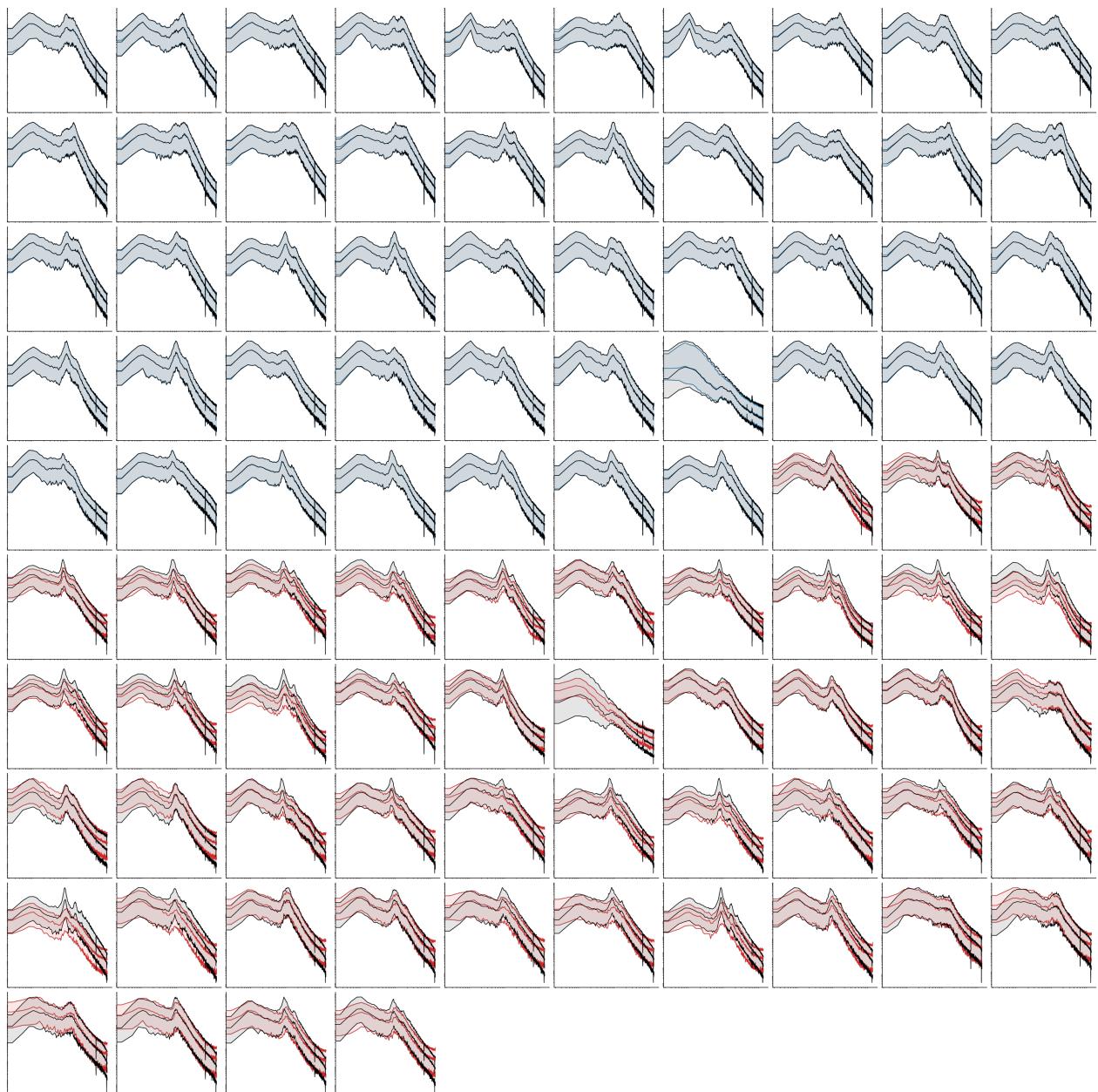


Figure 11: Full spectra of P01 from the AJILE12 dataset. The LFP traces consist of 94 channels. Here, we imputed the second half of the channels given the first half for all time series in the evaluation set. Median power as well as 10% / 90% percentiles are shown.

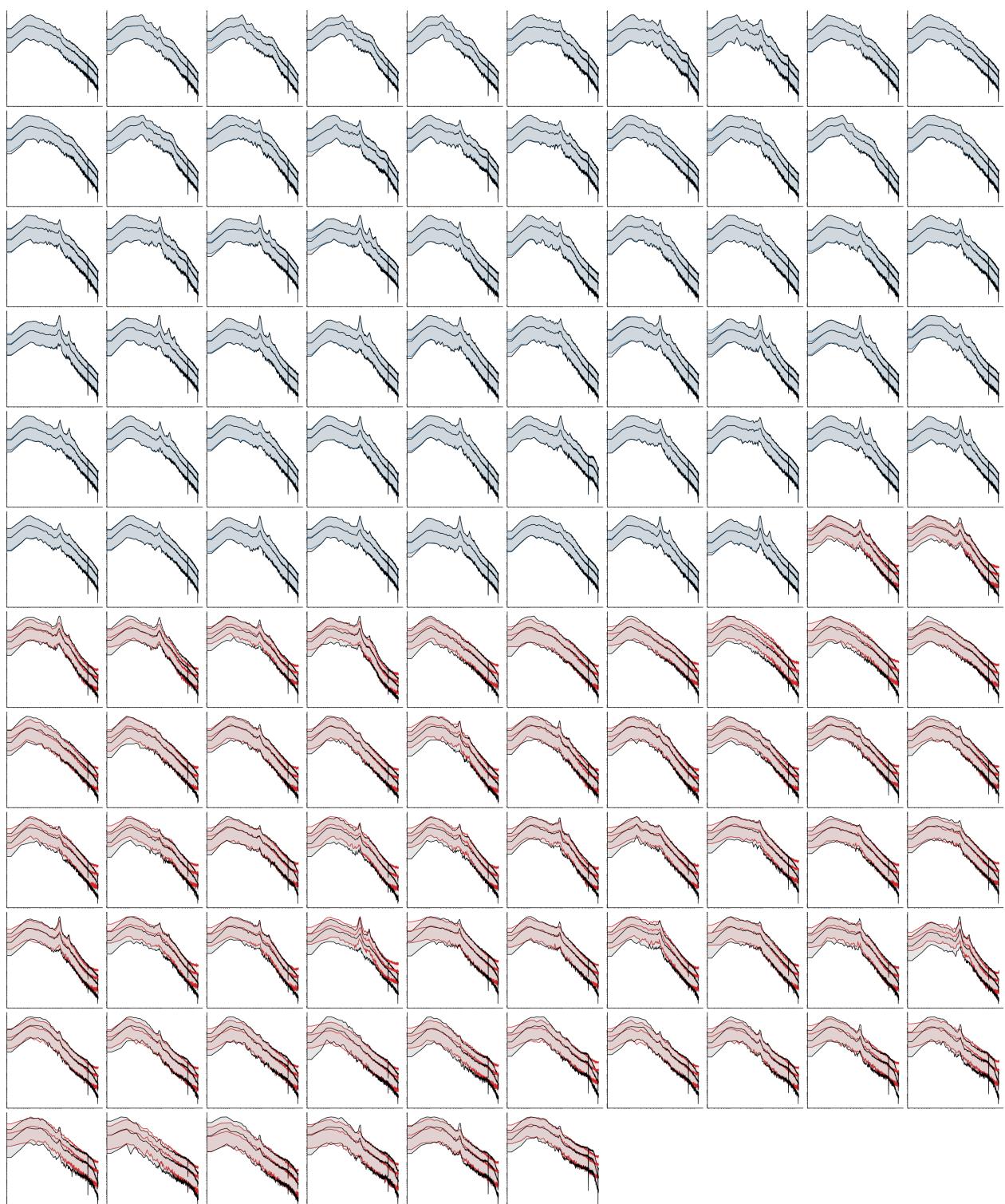


Figure 12: Full spectra of P12 from the AJILE12 dataset. The LFP traces consist of 126 channels. Here, we imputed the second half of the channels given the first half for all time series in the evaluation set. Median power as well as 10% / 90% percentiles are shown.

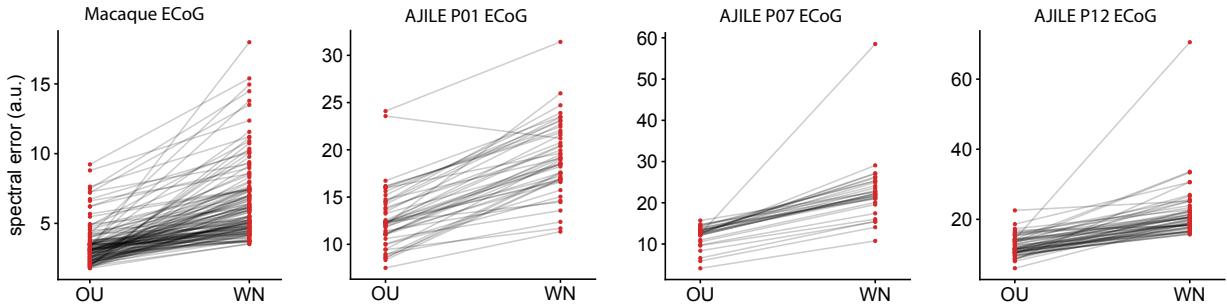


Figure 13: Comparison of OU versus white noise spectral errors of all generated / imputed channels for the macaque ECoG as well as AJILE12 participants P01, P07, and P12. All other training and network hyperparameters were fixed. The white noise process results in a higher spectral error.

Hyperparameter	Value
In kernel size	1
Out kernel size	1
SC layers	3
SC kernel size	53
SC scales	4
Decay min	2.0
Decay max	2.0
Heads	3
Latent dim. per channel	16
Off-diag. size	4
Noise process	OU with $\rho = 10$
Schedule	Quadratic
Diffusion steps T	50
B0	0.0001
B1	0.5
Time embedding dim.	16
Loss function	MSE
Optimizer	AdamW
Learning rate	0.0005
Weight decay	0.01
Train batch size	32
Epochs	250/500/1000

Table 2: Hyperparameters used to train models on the 12 different AJILE12 datasets for our imputation experiment. The same hyperparameters were used for all 12 participants, with the exception of the number of training epochs, to account for the large differences in training data size between participants.

Hyperparameter	Value
In kernel size	32
Out kernel size	32
SC layers	3
SC kernel size	101
SC scales	1
Decay min	2.0
Decay max	2.0
Latent dim. per channel	32
Off-diag. size	4
Noise process	White noise
Schedule	Linear
Diffusion steps T	500
B0	0.0001
B1	0.03
Time embedding dim.	16
Loss function	MSE
Optimizer	AdamW
Learning rate	0.0004
Weight decay	0.01
Train batch size	32
Epochs	500

Table 3: Hyperparameters used to train models on the macaque ECoG dataset for conditional generation of awake and anesthetized state.

Hyperparameter	Value
In kernel size	65
Out kernel size	65
SC layers	3
SC kernel size	65
SC scales	1
Decay min	2.0
Decay max	2.0
Latent dim. per channel	16
Off-diag. size	8
Noise process	White noise
Schedule	Linear
Diffusion steps T	1000
B0	0.0001
B1	0.02
Time embedding dim.	16
Loss function	MSE
Optimizer	AdamW
Learning rate	0.0004
Weight decay	0.01
Train batch size	32
Epochs	1000

Table 4: Hyperparameters used to train models on the EEG BCI dataset.