001 002 003

004

009 010

011

012

013

014

015

016

017

018

019

021

024

025

026 027 028

029

031

034

040

041

042

043

044

045

046 047

A GENERALIST INTRACORTICAL MOTOR DECODER

Anonymous authors Paper under double-blind review

ABSTRACT

Mapping the relationship between neural activity and motor behavior is a central aim of sensorimotor neuroscience and neurotechnology. Most progress to this end has relied on restricting complexity: studying specific simple behaviors, in limited subjects, with interpretable computational models. However, current trends in deep learning suggest that modeling a breadth of neural and behavioral data all at once is not only possible, but that such a model would also benefit downstream analysis of related data. We accordingly developed Neural Data Transformer 3 (NDT3) as a foundation model for motor decoding of neural data from intracortical microelectrodes. We pretrained NDT3 with 2000 hours of neural population spiking activity paired with diverse motor covariates from over 30 monkeys and humans from 10 labs. Pretrained NDT3 is broadly useful, benefiting decoding on 8 downstream decoding tasks and generalizing to a variety of neural distribution shifts. However, we find signs that scaling over diverse neural datasets may be challenging, as scaling from 200 to 2000 hours already requires increasing model size to 350M parameters to avoid model saturation, and several downstream datasets scarcely benefit from scale. We provide two demonstrations that this scaling is at least partially limited by variability in input and output spaces across neural datasets, which pretraining alone may not resolve.

1 INTRODUCTION



Figure 1. A. NDT3 is a deep network for decoding intracortical spiking activity into low-dimensional time series for various motor effectors¹. **B.** We aggregate decoding performance on downstream tasks with variable amounts of data (from Fig. 11). While from-scratch models only reliably outperforms a linear baseline after 15 minutes of data, tuning a 350M param. NDT3 pretrained with 2000 hours of data is consistently better at all scales.

Intracortical neural data collection is growing rapidly. This growth comprises not only larger individual datasets with more neurons and higher behavioral complexity (Urai et al., 2022; Stevenson, 2023), but also an increase in the collective number of datasets. This wealth of data presents an opportunity to develop insights and applications that span multiple datasets at once, provided we can reconcile their inherent diversity. Large deep networks appear very suitable for this task, so much

¹Photo courtesy of REDACT and The Chicago Tribune.

so that the creation of deep networks operating on broad domain data has been termed foundation modeling (Bommasani et al., 2022). Efforts to create foundation models are now proliferating beyond their origins in natural language processing (NLP) and computer vision (CV) into many domains of engineering and science (Wang et al., 2023c). Here, we propose a foundation model for motor decoding from intracortical spiking activity, which we call Neural Data Transformer 3 (NDT3).

Motor decoding is a valuable initial domain for characterizing neural data foundation models. 060 Academic, clinical, and industrial efforts to create iBCIs for neuroprosthetics provide a path for 061 scaling data collection from hundreds to thousands of subject-hours, and also fuel a need for pretrained 062 models that generalize quickly and perform robustly for new users and settings. Behavior prediction 063 metrics for BCI performance are also more intuitive for benchmarking progress than neural data 064 prediction or the abstract goal of providing scientific insight (e.g. with latent variable models or in silico models) (Pei et al., 2021; Wang et al., 2023b). Finally, recent work has shown that deep 065 networks are able to transfer learn across motor cortical datasets collected at different timepoints, 066 subjects, or tasks (Azabou et al., 2024; Ye et al., 2023; Schneider et al., 2023). These ingredients 067 provide the motivation and means for scaling neural data modeling. 068

069 However, scaling may be constrained by the design and heterogeneity of modern neural datasets. By design, we refer to the fact that many neural datasets restrict behavioral complexity to probe specific 071 hypotheses. These restrictions impoverish not only the behavioral signals but also the observed neural data (Gao and Ganguli, 2015), providing us a narrow window through which to understand 072 the general relationship between neural activity and behavior. Beyond the limitations of individual 073 datasets, each neural dataset inherently contains unique variability distinguishing them from others. 074 This is most salient when comparing across the datasets aggregated in pretraining, where different 075 neurons are recorded in each subject and distinct output dimensions are required for each effector. To 076 illustrate why these factors together challenge scaling, consider a 2-neuron toy setting, where both 077 neurons fire noisily. One neuron fires on leftward motion and the other fires on rightward motion. No amount of scaling on other datasets could reduce the data needed from this setting to determine 079 which each neuron's preferred direction, but neither is the problem trivial due to stochastic firing.

To assess the value of scaled pretraining on heterogeneous spiking activity, we developed Neural Data 081 Transformer 3 (NDT3). NDT3 uses simple tokenization strategies to enable pretraining over diverse datasets and fine-tuning to new tasks without introducing any new parameters (Fig. 1A). We pretrained 083 NDT3 using up to 2000 hours of neural and behavioral data from motor neuroscience experiments 084 with monkeys and clinical iBCI trials with humans. We then evaluated NDT3's decoding performance 085 on eight diverse motor tasks (Section 3.1) and find that tuning NDT3 yields models that either improve or match task-specific models trained from scratch, with prominent gains when task data is 087 under 1.5 hours (Fig. 1B). Further, these gains persist under a number of distribution shifts (Section 088 3.3). These benefits may enable both more complex experimental design and potentially decrease the burden of decoder training for people using iBCIs. However, our results also suggest that neural 089 data heterogeneity may be limiting scaling. Scaling pretraining data to 2K hours required raising 090 model capacity to 350M parameters to mitigate performance saturation, and some tasks accrue no 091 benefits from scale at all. We identify NDT3's sensitivity to the specific inputs and outputs seen 092 during fine-tuning as two limits to be overcome for more productive neural data foundation models. 093

093 094

2 Approach

096

098

2.1 Data

NDT3 models datasets of paired neural spiking activity and behavior (Fig. 2). Given our focus on motor decoding, most of the data comes from devices implanted in motor cortex of various monkeys and humans (Fig. 2A). These devices are intracortical multielectrode arrays or probes that record 30 kHz extracellular potentials. Spikes are extracted from these potentials, typically by bandpass-filtering the data between 300 and 3000 Hz, and marking a spike when the voltage signal crosses a preset threshold value. The neural data in our pretraining are diverse (Fig. 2B top). Data can have markedly different profiles across electrodes due to being from different electrode arrays in the same subject (left), have many silent channels (middle), or be densely active due to noise (right).

107 The typical behaviors in the pretraining data are different types of upper-limb reaching and grasping, nearly all from experimental paradigms that consist of short, repeated trials. While neural data were



124 Figure 2. NDT3 Data and Model Design: A. NDT3 models paired neural spiking activity and behavioral 125 covariate timeseries. We plot the distribution of 2000 hours of pretraining data volume by subjects (top) and covariate dimensionality (bottom). B. Examples of the neural and behavioral data for each of the three types 126 of behavioral covariates in pretraining: kinematics, EMG (electromyography), or forces. Not all modeled 127 dimensions in data are meaningfully task-related (right, grey behavior). C. Neural spiking activity is tokenized 128 in time by binning the number of spikes every 20 ms, and in "space" using patches of channels (usually 32), as 129 in NDT2 (Ye et al., 2023). Behavior is low-dimensional in our data, so we use 1 token per behavior dimension, 130 also per 20 ms timestep. NDT3 also pretrains on data from BCI control, which we annotate with two additional tokens. The phase token indicates whether the user is controlling or observing the behavior and the reward token 131 indicates if the BCI task was completed. D. NDT3 models tokens in a single flat stream with linear readins and 132 readouts. Every real-world timestep (shown by the blue cutout) yields several tokens, which we order to allow 133 causal decoding in evaluation. In evaluation, we omit return and phase tokens and zero-mask behavior tokens. 134

- 135
- 136

138

always recorded from microelectrodes, motor covariate signals came from various sensors. In monkey 139 datasets, these sensors measure actual limb activity (e.g., Fig. 2B, left: limb kinematics from optical 140 tracking, middle: electromyography (EMG)). In human datasets, physical movements are typically 141 not possible, so the data's behavior signals are programmatically generated. These signals are "paired" 142 with the neural data in that they are cued or otherwise instructed to the person, who will attempt or 143 imagine the corresponding behavior, such as grasping at a specified force level (Fig. 2B right). This 144 force panel also shows that in pretraining, we cannot always automatically discern the primary task 145 covariates (e.g., blue line, force, in the panel) from other recorded behavioral variables (grey). Thus, 146 some behavior variables may unpredictable. Finally, we include closed loop iBCI data, where some behavior is generated by an iBCI decoder (not NDT3, see modeling strategy in Section 2.2). 147

The pretraining datasets are composed of archives from several experimental labs and some public datasets, and contain data from non-human primate neuroscience experiments and human clinical trials for neuroprosthetics. The grassroots nature of this aggregate dataset presents a heterogeneity in neural data processing, motor effectors, and experimental setup, most comparable to aggregate robotics datasets (OpenX et al., 2024). We detail the pretraining data composition in Section C.4.

153 We minimize preprocessing of these data to maximize the applicability of our generalist model. 154 Kinematic signals are typically converted to velocities, and all behavior (kinematics, EMG, force) 155 is normalized per dataset such that the maximum absolute value of each variable is one. Data are 156 cut or concatenated into fixed length sequences, without additional annotation of data discontinuity. 157 This strategy, common in language modeling (Geiping and Goldstein, 2022), homogenizes the data 158 for improved GPU utilization while maintaining throughput of real data. We used a length of two 159 seconds as it is roughly the timescale of the behavior in our data (Fig. 2A). Sequences with no spikes or covariate variability are discarded. In total, this yielded about 3 million sequences, 1 billion 160 neural tokens, or 1750 hours, which we sample uniformly for pretraining. We round this to 2 khrs in 161 subsequent text for simplicity.

162 2.2 MODEL

164 NDT3 is a causal Transformer with linear readin and readout layers for its various modalities, similar to GATO or TDMs (Reed et al., 2022; Schubert et al., 2023; Chameleon, 2024). For use with 165 a Transformer, the data must be tokenized (Fig. 2C). We tokenize neural data by patching spike 166 counts (Ye et al., 2023); each token is a flattened vector of the binned spikes in a chosen temporal 167 resolution (20 ms) and spatial dimension (32 channels). For example, neural activity sampled from 168 an electrode array with 100 channels would patch into $4 = \lfloor 100/32 \rfloor$ 32D neural tokens per 20 ms timestep. As the behavioral variables are already low-dimensional, we simply assign 1 token per 170 dimension at the same temporal resolution as neural data. Finally, we add tokens marking whether the 171 behavior are generated by a BCI system or by physical limb movement. While measured kinematics, 172 EMG, or force will reflect a natural relationship with neural activity, behavioral data from BCI tasks 173 are controlled by a program or learned decoder. We frame BCI-driven behavior as a suboptimal 174 demonstration (Merel et al., 2016), and adopt a scheme inspired by Decision Transformers (Chen 175 et al., 2021; Lee et al., 2022). In this scheme, we use a Phase token to track the timesteps where 176 behavior is at least driven by neural activity and under decoder control, or only under programmatic, open loop control. We also use a Return token reflecting controller quality based on task completion. 177 Note that these signals are only considered for pretraining, and are ablated entirely from the model 178 at evaluation. Similarly, input behavior tokens are masked out in inference, so that the model input 179 only indicates how many tokens must be predicted. NDT3 is trained with mean-squared error for 180 prediction of behavioral variables, and categorical cross-entropy losses for prediction of neural spike 181 count and reward. 182

All modalities are flattened into a single token stream, with the order of tokens in each real-world timestep respecting a canonical order required for control (Fig. 2D). As in GATO, individual tokens are annotated with learned position embeddings identifying token modality and sub-modality "position."
 We additionally use rotary embeddings (Su et al., 2023) to track real-world timesteps.

Pretraining and Fine-Tuning We pretrain NDT3 models over variable pretraining data and in sizes of 45M and 350M parameters to assess the impact of data and model scaling. Pretraining is early stopped according to validation loss or terminated at a maximum of 400 epochs. The 200 hour, 45M model trains for 480 A100-hours while the 2000 hour (2kh) 350M model takes 20K A100-hours. Fine-tuning maintains the pretraining objectives and updates all parameters.

192

193 2.3 EVALUATION STRATEGY

194 Evaluation datasets and tuning Our main evaluation (Section 3.1) uses four human and four 195 monkey datasets sampling varied upper limb movements, which we detail in Section C.4. Each 196 dataset contains multiple sessions of data, typically from a single monkey or human. We will refer to 197 each such setting as a "task," distinguished from the behavioral procedure performed in each dataset. Each session has unique variability, so fine-tuning procedure may greatly impact decoding results. 199 Prior work (Azabou et al., 2024; Ye et al., 2023; Zhang et al., 2024) ran focused evaluations by tuning 200 and evaluating separately for each evaluation session. To manage compute and storage demands 201 and to reflect that real world datasets are rarely collected or analyzed in isolation, we fine-tune NDT3 jointly over data combined from multiple evaluation sessions. Fig. 12 shows this joint tuning 202 outperforms focused tuning for multi-session data. 203

204 Baselines We compare against Wiener filters (WF) and NDT2. WFs are a conventional linear method 205 for both motor decoding and control in iBCI devices (Pandarinath and Bensmaia, 2022), and we 206 implement them as ridge regression with multi-timestep history. NDT2 is a Transformer that uses MAE-style (He et al., 2021) self-supervision to learn across multiple neural datasets. We detail the 207 differences between NDT2 and NDT3 in Section C.3. We compare to NDT2 prepared both from 208 scratch and tuned from the public checkpoint pretrained on 100 hours of human data. Note for 209 tractability we - Other Transformers have been proposed for scaling over spiking data (Azabou et al., 210 2024; Zhang et al., 2024), but the field yet lacks consensus benchmarks to evaluate these proposals. 211

Downstream Hyperparameters We tune all deep networks (NDT2 and NDT3) over 3 learning rates.
This sweep is limited to make computational demands tractable, but also demonstrates the versatility
of the base model. Importantly, the same search space is used for all tasks; we list the space and
show its sufficiency relative to wider sweeps in Section C.2. The best learning rate is chosen based
on average validation score over three random seeds, and we report their mean on the evaluation data.

²¹⁶ 3 RESULTS

226 227

228

229

230

231

232

233

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

NDT3's pretraining effort advances prior intracortical models an order of magnitude in data and
 model scale, from 200 to 2000 hours and 10M+ to 100M+ parameters. In Section 3.1, we show the
 increased data scale saturates aggregate downstream performance unless simultaneously increasing
 model scale. We propose that the performance drop from scaling data alone is due to high variability
 across intracortical motor decoding datasets. In Section 3.2, we show how this variability reflects in
 NDT3's sensitivity to shifts in data input or output. Section 3.3 concludes by showing that despite
 this challenge for further data scaling, NDT3's pretraining already provides gains that generalize to
 various novel settings, establishing NDT3 as a useful foundation for motor decoding.

Α. В. C. 4D Bimanual 2D Self-paced R² R² 2K hrs (inc. 200 hrs R^2 R 200 hrs from test monkey Cursor Reach from test m 2000 Pretraining Volume (hrs) Test dataset 0.30 0.70 0.70 0.6 Interference 0.25 0.65 0.65 200 0.20 0.5 Params 45M 0.15 350M 0.60 0.60 0.4 200 Epochs 100 100 200 Epochs 0.10 Session: 60 s Session: 5 min Pretrained 45M Wiener Filter ·--- 🗙 --Task: 42 min LEGEND Task: 12 hr NDT3 Scratch ·----* 0.05 Pretrained 350M + 0.3 10 50 100% 25 25 50100% 10 F. 8 task x 3-5 scales Avg. ^{45M param_350M param} D. Ε. R^2 R^2 hr Normalized Baseline 1.5 NDT3 Scratch 0.3 45M 1 5 h 1.0 45M 200 hi 0.2 45M 2 khi 0.1 350M 200 hi 0.5 200 } NDT2 100 MDT3 0.0 45M 1.5 Hr 2.5 11 2451 NDT2 200 11 200 11 2424 45M 200 m 45M2 KM Nr. 0.0 00 10¹ 10² Minutes of Task Data

3.1 MULTI-SCALE EVALUATION ACROSS MOTOR DECODING TASKS

Figure 3. Evaluation on diverse motor tasks: A single legend and color scheme is used throughout. A. 250 Test-split pretraining R^2 compared for 3 models. All model pretraining data includes 1.5 hours of calibration 251 data for the test dataset. We compare a model with just this data (Test dataset only) vs using 200 hours of 252 additional data either from the test monkey or from over 10 other monkeys. Only the additional test monkey data improves over the calibration model. Models terminate at different points due to early stopping. B. Pretraining 253 R^2 for models with up to 2000 hours (2 khrs) of pretraining data. The 2 khr model degrades in performance vs 254 the 200 hr model at 45M parameters and merely maintains performance at 350M parameters. C. Examples of 255 good and bad data-scaling in downstream multiscale evaluation on two datasets. The bottom right text shows 256 time in each evaluation session and total time in each dataset. The x-axis scales this full dataset down by 257 random subsampling. Shading shows standard deviation on 3 tuning seeds. Increasing pretraining data yields performance gains at all downstream scales in the 4D task, but effects are unclear in the self-paced reach task. ▼ 258 indicate outliers clipped for clarity. **D.** Downstream performance averaged for 31 settings comprised of different 259 downstream datasets and scales, for different NDT3s and baselines. 45M NDT3s improve with data from 1.5 hrs 260 to 200 hrs but saturate at 2 khrs. Increasing model size to 350M parameters enables further gain. E. p-values 261 computed from FDR-corrected pairwise t-tests for each pair of models. The 350M 2 khr NDT3 significantly 262 outperforms other pretrained NDT3s, except the 350M 200 hr NDT3, and is the only model to do so. NDT2s 263 omitted for brevity, see Fig. 13. F. Per-task performance, normalized by the 350M 200 hr NDT3 performance, is shown against task time for different NDT3 models. Each vertical band shows models trained on the same 264 evaluation setting, e.g. dashed lines show the evaluations from the self-paced reaching dataset. Model variability 265 vanishes by 1.5 hours. 266

To set expectations for how data scale and model size will impact model performance, we first examine pretraining curves computed on a test split. This test split contains data from multiple sessions of 2D reaching behavior mainly from one monkey. Models are able to learn this test task in pretraining as they are given 1.5 hours of data separately sampled from these sessions. Fig. 3A

5

270 compares the test performance of a model using just this 1.5 hours of calibration data with those of two 271 200 hour models. The "Other datasets" 200 hour model used data from 10 other monkeys performing 272 a variety of reaching tasks, but did not improve over the minimal 1.5 hour model on the test data. 273 In contrast, using 200 hours from the test monkey (from a separate set of experiments with similar 274 behavior) achieved a small improvement in performance. Thus, only closely related data appears to benefit a model that already uses sufficient task-specific data, in this case 1.5 hours. Fig. 3B further 275 shows that scaling to 2 khr can actually degrade test performance, suggesting dissimilar pretraining 276 data can interfere with learning of the test task. This interference is mitigated by increasing model 277 size to 350M parameters, consistent with prior recommendations to scale model size and dataset size 278 in tandem (Dosovitskiy et al., 2021; Kolesnikov et al., 2020; Aghajanyan et al., 2023). However, the 279 test task does not improve beyond what is achieved by providing 200 hours from the test monkey. 280

This upstream saturation motivated a downstream evaluation conducted at multiple data scales. We 281 illustrate this evaluation for two tasks in Fig. 3C. These two datasets are from a human performing 282 open loop iBCI calibration for bimanual cursor use (Deo et al., 2024), and a monkey performing 283 self-paced reach to random targets (O'Doherty et al., 2017). In both cases, the individuals are held-out 284 from pretraining entirely, so the task-specific data is only seen in tuning. In the bimanual task, NDT3 285 performance improves with increased pretraining data at all downstream data scales. Encouragingly, 286 and distinct from the saturation seen in pretraining, joint scaling of data and model size to 2 khr and 287 350M parameters improves performance in the bimanual task. The self-paced reach dataset shows a 288 much less clear result. For example, the from-scratch NDT3 is competitive at all data scales. 289

These tasks show just two of several different trends on the eight evaluation datasets we study. We 290 defer discussion of individual tasks and their variability to Section B.5, and next consider summary 291 performance. Fig. 3D was produced by tuning over 2000 models in 31 evaluation settings, and 292 identifies an overall benefit of pretraining scale. To begin, NDT3 from scratch outperforms the WF 293 and NDT2, whether pretrained or not. We discuss NDT2's poor performance in Section B.7. This 294 from-scratch NDT3 performance can be improved with minimal pretraining (1.5 hrs), consistent 295 with findings in computer vision (Entezari et al., 2023), and continues improving up to 200 hours of 296 pretraining data. Further scaling to 2000 hours is also helpful, but only when paired with increased 297 model size to 350M parameters, as in upstream evaluation. The gain of the 350M 2 khr model over other models is statistically significant for all except the 350M 200 hr model (Fig. 3E). Other pairs of 298 pretrained NDT3s are not significantly different at the current scale of evaluation. 299

300 While scaling shows an overall positive trend, benefits vary greatly across downstream settings. 301 While we have already seen that comparing benefits by experimental task is challenging, it is known 302 that pretraining is most beneficial at low downstream data scale. We confirm this in our setting 303 by plotting normalized task performance against downstream scale in Fig. 3F. As in unnormalized 304 performance (Fig. 1B) and upstream results (Fig. 3B), the distinction between pretrained models and from-scratch models vanishes by 1.5 hours, or 0.5-5K data points. Vanishing benefits by 5K data 305 points, while far from trivial, implies that NDT3's pretraining is relatively impotent by CV or NLP 306 standards (Kolesnikov et al., 2020; Wang et al., 2019), and in practical terms can be exceeded after 307 a few sessions of data collection. Importantly, since scaling downstream data past 1.5 hours still 308 reliably improves performance, we cannot attribute the saturation of pretraining to insufficient signal 309 in each dataset's neural activity. 310

- 311
- 312
- 313
- 314
- 315 316

3.2 PRETRAINING IS LIMITED BY INPUT AND OUTPUT VARIABILITY IN TUNING

317 318

In Section 3.1 we observe that the benefit of scaling pretraining appears to saturate at relatively low downstream data scales. It is possible that this limit reflects intrinsic variability across neural datasets, with a long tail of specialized features that are needed for the best performance in each dataset. Alternatively, our modeling decisions around architecture, hyperparameters, and post-training, may all have significantly limited NDT3's scaling. As a first step towards dissociating these factors, we next analyze NDT3's sensitivity to the specific neural inputs and covariate outputs seen in tuning.



342 Figure 4. NDT3 fails in certain novel input and output configurations. A: Cross-session transfer persists after pretraining, but cross-subject does not. We test NDT3 in a downstream task with one evaluation session 343 from a monkey self-paced reaching dataset. Training uses 1 minute from the evaluation session and additional 344 data from other sessions (Cross-Session) with the same monkey or from sessions from a different monkey for 345 the same behavior (Cross-Subject). B: Shuffling inputs ablates cross-session data to resemble cross-subject 346 transfer. uses the same cross-session neural data but permutes input dimensions (recording channels). Shuffle 347 channel randomly permutes inputs, half-token shift rolls channels so that each channel i uses data from i + 16, and shuffle token permutes data patchwise, keeping channels from the same patch together. Channel shuffling 348 and half-token shifts both are sufficient to reduce cross-session transfer to the same level as cross-subject transfer. 349 All panels show the baseline performance achieved by the model with just 1 minute of test-session data, x-axis 350 shows additional cross-context data provided. C-F. Pretraining does not improve angular extrapolation. C. We 351 study a new dataset where a monkey performs an isometric center out task, with exerted forces mapped to cursor 352 positions in 8 different angles. **D.** We split data into three held-in and five held-out angles. Behavior is cleanly 353 separated across conditions. The neural data for each condition can also be visualized separably by projecting them to 2D plane computed by combining PCA and LDA. E. Predictions derived either by fitting a Wiener Filter 354 to the projected neural data (Linear, PCA-IDA), or from NDT3 (Scratch, 350M 2khr). While the linear model 355 generalizes to held-out angles, NDT3 predictions are restricted to held-in ranges. F. Pretraining quantifiably 356 improves over from-scratch in all conditions, but far underperforms the generalization of PCA-LDA. 357

Input order sensitivity may limit cross-subject transfer. In neural data, the effectiveness of transfer 358 learning is greatly reduced when using cross-subject data compared to cross-session data (Ye et al., 359 2023). This suggests limited scaling may be caused in part by the vanishing utility of cross-subject 360 data, even when the data is collected in identical experimental setups and thus controlled for other 361 variables. We can illustrate this by comparing cross-session and cross-subject transfer after large-scale 362 pretraining. On a monkey 2D reaching dataset (O'Doherty et al., 2017) in Fig. 4A, we tune NDT3 with calibration data from one test session and additional cross-session or cross-subject data. As in 364 Section 3.1, cross-session data is still highly beneficial even after pretraining, but cross-subject data 365 is only helpful for from-scratch models. Pretrained models in the cross-subject setting instead begin 366 and plateau at a performance that is just slightly better than the best cross-subject performance in from-scratch models. These results suggest that NDT3's pretraining has already learned the features 367 that cross-subject transfer provides in this task, supporting the idea that scaling is limited by a low 368 (task-dependent) ceiling on cross-subject transfer. 369

370 This low transfer stems from cross-dataset variability, but may be exacerbated by NDT3's design 371 and inductive bias. To dissect NDT3's sensitivities, we observe that since cross-subject data contain 372 different neurons, cross-subject transfer must at least accommodate changes in the specific semantics of each data dimension, a problem which we term "sensor variability." We can isolate how effectively 373 NDT3 resolves sensor variability by transferring with cross-session data while permuting the test-374 session's neural dimensions. Fig. 4B Channel shuffle shows that input permutation cripples the ability 375 of NDT3, whether pretrained or not, to learn from cross-session data. We next apply structured 376 ablations of input order, as in Neyshabur et al. (2020). We find that even the small alteration of a 377 half-token shift in channel order is sufficient to reproduce the effect of full shuffling (Fig. 4B center). This shows NDT3's cross-session transfer depends greatly on the specific token dimension semantics, and may drive the observed limits in scaling. Finally, we consider a shuffle that only alters the test session's neural token order with respect to cross-session data, hypothesizing Transformers can more easily correct token ordering alterations. Indeed, pretrained and from-scratch models recover some gains when cross-session data is only token-shuffled (Fig. 4B right). This manipulation shows the influence of model design on transfer.

384 For reference, POYO (Azabou et al., 2024) adopts a graph-based Perceiver design motivated by 385 this sensor variability challenge, which makes it a promising candidate to scale. However, their 386 experiments only show nontrivial cross-subject transfer rather than robustness to input perturbations 387 as we conduct here. Since NDT3 also achieves cross-subject transfer, stronger evidence is needed to 388 support that Perceiver-style models may scale better. Finally, we note that to control for potential confounds from tuning over heterogeneous data, we evaluated both sequential and joint tuning 389 strategies in these experiments and reported the better approach in each panel. A full comparison is 390 given in Section B.6. 391

392 Pretraining does not enable angular extrapolation. The increased data efficiency of pretrained 393 models suggests that NDT3 could decode a new subject's behavior without sampling the full range 394 of neural and behavioral data. For simple behaviors, this expectation provides a validation of 395 whether pretraining has succeeded at all. Center-out reaching provides this simple litmus test as its underlying neural activity can be visualized in a planar subspace (Churchland et al., 2012). This 396 visualization doubles as an explicit prior for decoders aiming to generalize to held-out angles. We 397 expect that NDT3's pretraining could learn this prior, enabling generalization to reach directions 398 yet unseen in a new subject. Note that non-pretrained deep networks fail at such held-out angle 399 generalization (Rizzoglio et al., 2022). To evaluate this hypothesis, we analyze a single session of 400 an isometric monkey center-out dataset, where the monkey exerts forces in one of eight angles and 401 its force level is mapped to cursor position (Fig. 4C). We separate this data into 3 held-in and 5 402 held-out angles, as shown in Fig. 4D. We can visualize the corresponding separability of the neural 403 data by first applying principal components analysis (PCA) and then linear discriminant analysis 404 (LDA, see Section B.1 for methods). These held-out neural data are still organized radially, by reach 405 angle (Fig. 4D bottom right).

406 This reduced view of the neural data implies we can build a linear decoder that generalizes to held-out 407 angles. Indeed, a Wiener Filter on this 2D neural data predicts behaviors that, while generally low in 408 quality, are coarsely aligned with the correct reach angle (Fig. 4E). In contrast, NDT3 from scratch 409 predictions for held-out angles are constrained to their nearest held-in angles. Pretraining provides 410 mild improvements for the interpolated angles at $\pm 45^{\circ}$, but largely replicates the failure to predict 411 extrapolated angles. We quantify prediction accuracy in Fig. 4F, showing that while pretrained 412 NDT3 improves over from-scratch NDT3 overall, pretraining still far underperforms a simple prior in generalizing to held-out angles. We repeat this evaluation in two more settings in Section B.1. 413 Importantly, this demonstration leaves open the question of whether pretraining has failed to learn 414 the utility of dimensionality reduction, or whether aligning NDT3 to yield behavioral generalization 415 in tuning will require further model post-training. 416

These two studies on model input and output highlight the difficulty of pinpointing whether NDT3
scaling is limited by data or methodology. However, they do provide basic tests of model capability,
namely robustness to channel shifts and generalization to unseen behaviors, that we expect future
approaches will need to overcome to achieve better neural data foundation models.

421 422

423

3.3 WHERE DOES PRETRAINING HELP?

Despite challenges for scaling further, NDT3's pretraining has learned a prior from hundreds of hours
 of neural data. We next provide examples where this pretraining does usefully generalize.

Neural distribution shifts Neural data is nonstationary, with shifts rising from a mix of controlled
experimental variables to more speculative factors. For example, the firing rate of different channels
will evolve over the course of an hour, implying a distribution shift associated with change in time
(Fig. 5A top left). Shifts also occur between activity in two arm postures or whether finger motion
occurs under spring load or not (Fig. 5A top middle and right). Since these shifts are common in
neural data, pretraining gains should ideally be robust to their effect. We thus tune models on data
from one setting (in-distribution, ID) measure the performance of models in that same setting and the



Figure 5. Generalizability of pretraining gains. A. Models fine-tuned in one distribution of data are evaluated in-distribution (ID) and out-of-distribution (OOD). Top plots show the distribution across channels of neural firing rates from OOD and ID trials, normalized by average ID firing rates. Lower plots scatter OOD vs ID performance, with each point being a single model with different hyperparameters. . The time shift uses two human cursor datasets collected one hour apart. Models were tuned in each block and were evaluated in the second block. Pose shift uses a monkey center-out reach task which was performed with the hand starting in different locations in the workspace. Spring Load uses a dataset of monkey 1D finger motion with or without spring force feedback. B. Models are evaluated on a human open-loop cursor dataset prepared in two ways. Trialized training receives inputs according to trial boundaries, varying from 2-4 seconds in length. Continuous training receives random 1 second snippets (that can cross trial boundaries). Trialized evaluation matches trialized training, and continuous evaluation is done by streaming up to 1 second of history. ▼ indicates points below 0.0. Continuously trained models perform well in both evaluation settings, while models trained on trialized data fail in continuous evaluation. C. Multiscale fine-tuning performance of NDT3 on datasets recorded outside motor cortex, namely S1 (Somatosensory) and FEF/MT (Oculomotor).

432

433

434

435

437

439

440

441

445

446

447

448

449

450

451

452

453

454

455

458 459

shifted setting (out-of-distribution, OOD). Positive correlation of performance in all cases imply the 460 ID gains conferred by pretraining persist OOD. This ID-OOD correlation is consistent with Miller 461 et al. (2021), implying a potentially fruitful relationship between the distribution shifts characterized 462 in neural data and those studied in computer vision. More practically, these examples suggest that 463 pretraining benefits are not dependent on narrow features specific to the choice of tuning dataset.

464 Trial structure DNNs have been observed to overfit the temporal structure of experimental data, 465 challenging their use in iBCI control (Deo et al., 2024; Costello et al., 2023). For example, a DNN 466 might learn there is always no motion before the start of a behavioral trial, independent of the neural 467 activity. To date, these claims have been studied exclusively on un-pretrained deep networks. We next 468 assess how pretraining affects this overfitting in open loop human cursor control data by comparing a 469 continuous and trialized setting. In the continuous setting, we cut random one second intervals of 470 data in training and continuously stream up to one second of data in evaluation. The trialized setting formats data to respect trial boundaries, so the model always sees data aligned to the start of behavior. 471

472 In Fig. 5B, we show that models using both trialized training and evaluation outperform models with 473 both continuous training and evaluation. This implies NDT3 will learn to exploit clear trial structure 474 in data. However, while trialized from-scratch models become subtrivial under continuous evaluation 475 (solid blue line is off-panel), pretrained models degrade more gracefully. For example, the 350M 476 2 khr model evaluated continuously only performs slightly worse with trialized tuning than with 477 continuous tuning. Pretraining NDT3 thus reduces its dependence on trial structure, which should benefit both data analysis and iBCI control. Note however the contrast in these results with Fig. 4C, 478 which show that DNNs clearly do overfit to tuning data in some cases. These nuances underscore the 479 importance of rigorously evaluating model generalization in future work. 480

481 New brain areas In Fig. 5C, we return to multiscale fine-tuning to test how NDT3, pretrained on 482 motor cortex, performs in somatosensory cortex (S1) and oculomotor areas (FEF and MT). Pretraining 483 provides a large boost over from-scratch models is high in S1, but also nontrivial in the Oculomotor dataset. While the former can be attributed to the close interaction of sensorimotor areas, the latter 484 implies NDT3 has learned a broader prior. While encouraging, our results thus far suggest this prior 485 could reflect neurophysiology (e.g. declining subject focus over time (Steinmetz et al., 2019)), but might also reflect common experimental artifacts like trial structure. For example, this Oculomotor
 dataset contains 4 behavioral conditions, which may benefit from the tendency to learn classifiers
 shown in Fig. 4C rather than a prior on neural dynamics.

490 491

492

4 DISCUSSION

493 Many fields are now pursuing large scale deep learning as "a tide that lifts all boats" (Abnar et al., 494 2021), with the hope that improvements in effective pretraining will yield field-wide, downstream 495 improvements. Such a unifying abstraction may be timely for neuroscience, given the increasing 496 volume, diversity, and complexity of modern neural data. Joining other pretraining efforts on 497 varied modalities of neural data (Section A), we trained NDT3 on 2000 hours of paired neural 498 population activity from motor cortex and behavior, and then conducted a broad downstream decoding evaluation. Consistent with the broad foundation modeling narrative, we found the best aggregate 499 performance from increasing data scale and model size jointly (Section 3.1). However, these 500 benefits from pretraining vary with the downstream dataset, with several datasets having minimal 501 improvements from scale (Section B.5). This result may stem in various ways from our approach, 502 for example in insufficient breadth of hyperparameter sweeps, or too narrow of a focus on decoding 503 metrics. Alternatively, we have highlighted how improving downstream gains may require new 504 architectural innovations robust to input ordering shifts, and possible new training strategies to 505 promote generalization (Section 3.2). Overall, NDT3 establishes a strong baseline foundation model 506 for intracortical decoding from spiking activity, but highlights important directions for future scaling.

507 More broadly, we advocate for further consideration of how neural data can contribute to and gain 508 from the ongoing cross-disciplinary conversation on foundation modeling. For example, our input 509 and output sensitivity analyses were inspired by ML (Neyshabur et al., 2020; Pham et al., 2021) and 510 neuroscientific literature (Gallego et al., 2020; Sadtler et al., 2014), respectively. Improving scaling 511 in neural data may benefit from insights developed from characterizing multimodal interference more 512 broadly (Aghajanyan et al., 2023; Liu et al., 2024). Inversely, neural distribution shifts have the 513 advantage of being carefully studied, and so the appearance of correlated ID-OOD performance in 514 neural data, as also appears in CV, NLP, and other AI domains (Taori et al., 2020), may refine our 515 understanding of when such correlation will occur, and thus when foundation models will be effective. 516 Our hypothesized challenge of sensor variability should be particularly interesting to compare across the biosignals community, which must overcome analogous variability to achieve our shared goal of 517 achieving user-general models. 518

519 520

521

4.1 ETHICS STATEMENT

522 The animal datasets used in this work were collected for other studies that were approved by 523 Institutional Animal Care and Use Committees. Human datasets were also collected for other 524 studies, with Institutional Review Board approval and as part of clinical trials conducted under 525 FDA Investigational Device Exemptions. Informed consent was obtained prior to any experimental procedures. We discuss the potential for NDT3 to reduce user burden for iBCI-based neuroprosthetics, 526 though the dissemination of pretrained models on these data raise the risk that the original human 527 data may be recoverable from model weights. Since this seems technically challenging at this point, 528 and since the source data are restricted to binned spiking activity to begin with, we deem the risk low 529 enough to justify the potential scientific benefit of sharing our pretrained models. 530

531 532

4.2 REPRODUCIBILITY STATEMENT

Advancing neural data foundation modeling will require a flourishing open-source ecosystem, including data, models, and evaluations. While we will release our models and codebase, our work currently has limited reproducibility given our inability to release pretraining data. Similarly, we have tried to use open evaluations where possible, but several evaluation datasets remain private. We expect that field-wide trends toward open data releases, and larger scale academic (Koch et al., 2022) or academic-industrial collaborations, can alleviate this limitation in the near future.

540 REFERENCES

542

543 544

545

546

550

551 552

553

554

555

558

559

560

577

578

579

580

582

583

589

- S. Abnar, M. Dehghani, B. Neyshabur, and H. Sedghi. Exploring the limits of large scale pre-training, 2021. URL https://arxiv.org/abs/2110.02095.
- E. H. Adelson, J. R. Bergen, et al. *The plenoptic function and the elements of early vision*, volume 2. Vision and Modeling Group, Media Laboratory, Massachusetts Institute of ..., 1991.
- A. Aghajanyan, L. Yu, A. Conneau, W.-N. Hsu, K. Hambardzumyan, S. Zhang, S. Roller, N. Goyal, O. Levy, and L. Zettlemoyer. Scaling laws for generative mixed-modal language models. In *International Conference on Machine Learning*, pages 265–279. PMLR, 2023.
 - R. Antonello, A. Vaidya, and A. G. Huth. Scaling laws for language encoding models in fmri, 2024. URL https://arxiv.org/abs/2305.11863.
 - M. Azabou, V. Arora, V. Ganesh, X. Mao, S. Nachimuthu, M. Mendelson, B. Richards, M. Perich, G. Lajoie, and E. Dyer. A unified, scalable framework for neural population decoding. *Advances in Neural Information Processing Systems*, 36, 2024.
- G. Bachmann and V. Nagarajan. The pitfalls of next-token prediction, 2024. URL https://arxiv.org/abs/2403.
 06963.
 - T. Benster, G. Wilson, R. Elisha, F. R. Willett, and S. Druckmann. A cross-modal approach to silent speech with llm-enhanced recognition. *arXiv preprint arXiv:2403.05583*, 2024.
- R. Bommasani, D. A. Hudson, E. Adeli, R. Altman, S. Arora, S. von Arx, M. S. Bernstein, J. Bohg, A. Bosselut, 561 E. Brunskill, E. Brynjolfsson, S. Buch, D. Card, R. Castellon, N. Chatterji, A. Chen, K. Creel, J. Q. 562 Davis, D. Demszky, C. Donahue, M. Doumbouya, E. Durmus, S. Ermon, J. Etchemendy, K. Ethayarajh, 563 L. Fei-Fei, C. Finn, T. Gale, L. Gillespie, K. Goel, N. Goodman, S. Grossman, N. Guha, T. Hashimoto, 564 P. Henderson, J. Hewitt, D. E. Ho, J. Hong, K. Hsu, J. Huang, T. Icard, S. Jain, D. Jurafsky, P. Kalluri, S. Karamcheti, G. Keeling, F. Khani, O. Khattab, P. W. Koh, M. Krass, R. Krishna, R. Kuditipudi, A. Kumar, 565 F. Ladhak, M. Lee, T. Lee, J. Leskovec, I. Levent, X. L. Li, X. Li, T. Ma, A. Malik, C. D. Manning, 566 S. Mirchandani, E. Mitchell, Z. Munyikwa, S. Nair, A. Narayan, D. Narayanan, B. Newman, A. Nie, J. C. 567 Niebles, H. Nilforoshan, J. Nyarko, G. Ogut, L. Orr, I. Papadimitriou, J. S. Park, C. Piech, E. Portelance, C. Potts, A. Raghunathan, R. Reich, H. Ren, F. Rong, Y. Roohani, C. Ruiz, J. Ryan, C. Ré, D. Sadigh, 569 S. Sagawa, K. Santhanam, A. Shih, K. Srinivasan, A. Tamkin, R. Taori, A. W. Thomas, F. Tramèr, R. E. Wang, W. Wang, B. Wu, J. Wu, Y. Wu, S. M. Xie, M. Yasunaga, J. You, M. Zaharia, M. Zhang, T. Zhang, X. Zhang, 570 Y. Zhang, L. Zheng, K. Zhou, and P. Liang. On the opportunities and risks of foundation models, 2022. URL 571 https://arxiv.org/abs/2108.07258. 572
- 573 M. P. Branco, L. M. De Boer, N. F. Ramsey, and M. J. Vansteensel. Encoding of kinetic and kinematic 574 movement parameters in the sensorimotor cortex: A brain-computer interface perspective. *European Journal* 575 *of Neuroscience*, 50(5):2755–2772, Sept. 2019. ISSN 0953-816X, 1460-9568. doi: 10.1111/ejn.14342. URL 576 https://onlinelibrary.wiley.com/doi/10.1111/ejn.14342.
 - J. O. Caro, A. H. de Oliveira Fonseca, S. A. Rizvi, M. Rosati, C. Averill, J. L. Cross, P. Mittal, E. Zappala, R. M. Dhodapkar, C. Abdallah, and D. van Dijk. BrainLM: A foundation model for brain activity recordings. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=RwI7ZEfR27.
 - M. Caron, H. Touvron, I. Misra, H. Jégou, J. Mairal, P. Bojanowski, and A. Joulin. Emerging properties in self-supervised vision transformers, 2021. URL https://arxiv.org/abs/2104.14294.
- T. Chameleon. Chameleon: Mixed-modal early-fusion foundation models. *arXiv preprint arXiv:2405.09818*, 2024.
- G. Chau, C. Wang, S. Talukder, V. Subramaniam, S. Soedarmadji, Y. Yue, B. Katz, and A. Barbu. Population transformer: Learning population-level representations of intracranial activity, 2024. URL https://arxiv.org/abs/2406.03044.
- L. Chen, K. Lu, A. Rajeswaran, K. Lee, A. Grover, M. Laskin, P. Abbeel, A. Srinivas, and I. Mordatch. Decision transformer: Reinforcement learning via sequence modeling, 2021. URL https://arxiv.org/abs/2106. 01345.
- 593 M. M. Churchland, J. P. Cunningham, M. T. Kaufman, J. D. Foster, P. Nuyujukian, S. I. Ryu, and K. V. Shenoy. Neural population dynamics during reaching. *Nature*, 487(7405):51–56, 2012.

594 J. T. Costello, H. Temmar, L. H. Cubillos, M. J. Mender, D. M. Wallace, M. S. Willsey, P. G. Patil, and C. A. 595 Chestek. Balancing memorization and generalization in rnns for high performance brain-machine interfaces. 596 bioRxiv, pages 2023-05, 2023. 597 T. Dao. Flashattention-2: Faster attention with better parallelism and work partitioning, 2023. URL https: 598 //arxiv.org/abs/2307.08691. 600 A. Das, W. Kong, R. Sen, and Y. Zhou. A decoder-only foundation model for time-series forecasting, 2024. URL https://arxiv.org/abs/2310.10688. 601 602 M. Dehghani, J. Djolonga, B. Mustafa, P. Padlewski, J. Heek, J. Gilmer, A. Steiner, M. Caron, R. Geirhos, 603 I. Alabdulmohsin, R. Jenatton, L. Beyer, M. Tschannen, A. Arnab, X. Wang, C. Riquelme, M. Minderer, 604 J. Puigcerver, U. Evci, M. Kumar, S. van Steenkiste, G. F. Elsayed, A. Mahendran, F. Yu, A. Oliver, F. Huot, J. Bastings, M. P. Collier, A. Gritsenko, V. Birodkar, C. Vasconcelos, Y. Tay, T. Mensink, A. Kolesnikov, 605 F. Pavetić, D. Tran, T. Kipf, M. Lučić, X. Zhai, D. Keysers, J. Harmsen, and N. Houlsby. Scaling vision 606 transformers to 22 billion parameters, 2023. URL https://arxiv.org/abs/2302.05442. 607 608 D. R. Deo, F. R. Willett, D. T. Avansino, L. R. Hochberg, J. M. Henderson, and K. V. Shenoy. Brain control of 609 bimanual movement enabled by recurrent neural networks. Scientific Reports, 14(1):1598, 2024. 610 A. Doerig, R. Sommers, K. Seeliger, B. Richards, J. Ismael, G. Lindsay, K. Kording, T. Konkle, M. A. J. V. 611 Gerven, N. Kriegeskorte, and T. C. Kietzmann. The neuroconnectionist research programme, 2022. URL 612 https://arxiv.org/abs/2209.03718. 613 A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, 614 G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby. An image is worth 16x16 words: Transformers 615 for image recognition at scale. In International Conference on Learning Representations, 2021. URL 616 https://openreview.net/forum?id=YicbFdNTTy. 617 R. Entezari, M. Wortsman, O. Saukh, M. M. Shariatnia, H. Sedghi, and L. Schmidt. The role of pre-training data 618 in transfer learning. arXiv preprint arXiv:2302.13602, 2023. 619 620 C. Fan, N. Hahn, F. Kamdar, D. Avansino, G. H. Wilson, L. Hochberg, K. V. Shenoy, J. M. Henderson, and 621 F. R. Willett. Plug-and-play stability for intracortical brain-computer interfaces: A one-year demonstration 622 of seamless brain-to-text communication. In Thirty-seventh Conference on Neural Information Processing Systems, 2023. URL https://openreview.net/forum?id=STqaMqhtDi. 623 624 J. Farebrother, J. Orbay, Q. Vuong, A. A. Taïga, Y. Chebotar, T. Xiao, A. Irpan, S. Levine, P. S. Castro, A. Faust, 625 A. Kumar, and R. Agarwal. Stop regressing: Training value functions via classification for scalable deep rl, 626 2024. URL https://arxiv.org/abs/2403.03950. 627 R. D. Flint, E. W. Lindberg, L. R. Jordan, L. E. Miller, and M. W. Slutzky. Accurate decoding of reaching 628 movements from field potentials in the absence of spikes. Journal of neural engineering, 9(4):046006, 2012. 629 J. A. Gallego, M. G. Perich, R. H. Chowdhury, S. A. Solla, and L. E. Miller. Long-term stability of cortical 630 population dynamics underlying consistent behavior. Nature Neuroscience, 23(2):260-270, Feb 2020. ISSN 631 1546-1726. doi: 10.1038/s41593-019-0555-4. URL https://www.nature.com/articles/s41593-019-0555-4. 632 633 P. Gao and S. Ganguli. On simplicity and complexity in the brave new world of large-scale neuroscience. Current 634 opinion in neurobiology, 32:148-155, 2015. 635 J. Geiping and T. Goldstein. Cramming: Training a language model on a single gpu in one day, 2022. URL 636 https://arxiv.org/abs/2212.14034. 637 638 M. Geva, R. Schuster, J. Berant, and O. Levy. Transformer feed-forward layers are key-value memories, 2021. URL https://arxiv.org/abs/2012.14913. 639 640 A. Graves, S. Fernández, F. Gomez, and J. Schmidhuber. Connectionist temporal classification: labelling 641 unsegmented sequence data with recurrent neural networks. In Proceedings of the 23rd international 642 conference on Machine learning, pages 369-376, 2006. 643 K. He, X. Chen, S. Xie, Y. Li, P. Dollár, and R. Girshick. Masked autoencoders are scalable vision learners, 644 2021. URL https://arxiv.org/abs/2111.06377. 645 646 B. Jarosiewicz, A. A. Sarma, D. Bacher, N. Y. Masse, J. D. Simeral, B. Sorice, E. M. Oakley, C. Blabe, 647 C. Pandarinath, V. Gilja, et al. Virtual typing by people with tetraplegia using a self-calibrating intracortical

brain-computer interface. Science translational medicine, 7(313):313ra179–313ra179, 2015.

648 W. Jiang, L. Zhao, and B. liang Lu. Large brain model for learning generic representations with tremendous 649 EEG data in BCI. In The Twelfth International Conference on Learning Representations, 2024. URL 650 https://openreview.net/forum?id=QzTpTRVtrP. 651 J. Kaplan, S. McCandlish, T. Henighan, T. B. Brown, B. Chess, R. Child, S. Gray, A. Radford, J. Wu, and 652 D. Amodei. Scaling laws for neural language models, 2020. 653 B. M. Karpowicz, J. Ye, C. Fan, P. Tostado-Marcos, F. Rizzoglio, C. Washington, T. Scodeler, D. de Lucena, S. R. 654 Nason-Tomaszewski, M. J. Mender, X. Ma, E. M. Arneodo, L. R. Hochberg, C. A. Chestek, J. M. Henderson, 655 T. Q. Gentner, V. Gilja, L. E. Miller, A. G. Rouse, R. A. Gaunt, J. L. Collinger, and C. Pandarinath. Few-shot 656 algorithms for consistent neural decoding (falcon) benchmark. bioRxiv, 2024. doi: 10.1101/2024.09.15. 657 613126. URL https://www.biorxiv.org/content/early/2024/09/16/2024.09.15.613126. 658 C. Koch, K. Svoboda, A. Bernard, M. A. Basso, A. K. Churchland, A. L. Fairhall, P. A. Groblewski, J. A. Lecoq, 659 Z. F. Mainen, M. W. Mathis, et al. Next-generation brain observatories. *Neuron*, 110(22):3661–3666, 2022. 660 661 A. Kolesnikov, L. Beyer, X. Zhai, J. Puigcerver, J. Yung, S. Gelly, and N. Houlsby. Big transfer (bit): General visual representation learning, 2020. URL https://arxiv.org/abs/1912.11370. 662 663 D. Kostas, S. Aroca-Ouellette, and F. Rudzicz. BENDR: Using transformers and a contrastive self-supervised 664 learning task to learn from massive amounts of EEG data. Frontiers in Human Neuroscience, 15, 2021. 665 ISSN 1662-5161. doi: 10.3389/fnhum.2021.653659. URL https://www.frontiersin.org/articles/10. 666 3389/fnhum.2021.653659. 667 C. Lea, R. Vidal, A. Reiter, and G. D. Hager. Temporal convolutional networks: A unified approach to action 668 segmentation, 2016. URL https://arxiv.org/abs/1608.08242. 669 K.-H. Lee, O. Nachum, M. Yang, L. Lee, D. Freeman, W. Xu, S. Guadarrama, I. Fischer, E. Jang, H. Michalewski, 670 and I. Mordatch. Multi-game decision transformers, 2022. URL https://arxiv.org/abs/2205.15241. 671 672 J. Liu, T. Wang, P. Cui, and H. Namkoong. On the need of a modeling language for distribution shifts: 673 Illustrations on tabular datasets, 2024. URL https://arxiv.org/abs/2307.05284. 674 X. Ma, F. Rizzoglio, E. J. Perreault, L. E. Miller, and A. Kennedy. Using adversarial networks to extend 675 brain computer interface decoding accuracy over time. Aug 2022. doi: 10.1101/2022.08.26.504777. URL 676 https://www.biorxiv.org/content/10.1101/2022.08.26.504777v1. 677 P. J. Marino, L. Bahureksa, C. Fernández Fisac, E. R. Oby, A. L. Smoulder, A. Motiwala, A. D. Degenhart, 678 E. M. Grigsby, W. M. Joiner, S. M. Chase, et al. A posture subspace in primary motor cortex. *bioRxiv*, pages 679 2024-08, 2024. 680 M. J. Mender, S. R. Nason-Tomaszewski, H. Temmar, J. T. Costello, D. M. Wallace, M. S. Willsey, N. Ganesh Ku-681 mar, T. A. Kung, P. Patil, and C. A. Chestek. The impact of task context on predicting finger movements in a 682 brain-machine interface. eLife, 12:e82598, jun 2023. ISSN 2050-084X. doi: 10.7554/eLife.82598. URL 683 https://doi.org/10.7554/eLife.82598. 684 685 J. Merel, D. Carlson, L. Paninski, and J. P. Cunningham. Neuroprosthetic decoder training as imitation learning. PLoS computational biology, 12(5):e1004948, 2016. 686 J. Miller, R. Taori, A. Raghunathan, S. Sagawa, P. W. Koh, V. Shankar, P. Liang, Y. Carmon, and L. Schmidt. 688 Accuracy on the line: On the strong correlation between out-of-distribution and in-distribution generalization, 689 2021. URL https://arxiv.org/abs/2107.04649. 690 M. Mitchell, S. Wu, A. Zaldivar, P. Barnes, L. Vasserman, B. Hutchinson, E. Spitzer, I. D. Raji, and T. Gebru. 691 Model cards for model reporting. In Proceedings of the Conference on Fairness, Accountability, and 692 Transparency, FAT* '19. ACM, Jan. 2019. doi: 10.1145/3287560.3287596. URL http://dx.doi.org/10. 693 1145/3287560.3287596. 694 B. Neyshabur, H. Sedghi, and C. Zhang. What is being transferred in transfer learning? Advances in neural 695 information processing systems, 33:512-523, 2020. 696 K. K. Noneman and J. Patrick Mayo. Decoding continuous tracking eye movements from cortical spiking 697 activity. International Journal of Neural Systems, page S0129065724500709, Oct. 2024. ISSN 0129-0657, 698 1793-6462. doi: 10.1142/S0129065724500709. URL https://www.worldscientific.com/doi/10.1142/ 699 S0129065724500709. 700 701 J. E. O'Doherty, M. M. B. Cardoso, J. G. Makin, and P. N. Sabes. Nonhuman primate reaching with multichannel sensorimotor cortex electrophysiology, May 2017. URL https://doi.org/10.5281/zenodo.788569.

702 C. OpenX, A. O'Neill, A. Rehman, A. Gupta, A. Maddukuri, A. Gupta, A. Padalkar, A. Lee, A. Pooley, A. Gupta, 703 A. Mandlekar, A. Jain, A. Tung, A. Bewley, A. Herzog, A. Irpan, A. Khazatsky, A. Rai, A. Gupta, A. Wang, 704 A. Kolobov, A. Singh, A. Garg, A. Kembhavi, A. Xie, A. Brohan, A. Raffin, A. Sharma, A. Yavary, A. Jain, A. Balakrishna, A. Wahid, B. Burgess-Limerick, B. Kim, B. Schölkopf, B. Wulfe, B. Ichter, C. Lu, C. Xu, 705 C. Le, C. Finn, C. Wang, C. Xu, C. Chi, C. Huang, C. Chan, C. Agia, C. Pan, C. Fu, C. Devin, D. Xu, 706 D. Morton, D. Driess, D. Chen, D. Pathak, D. Shah, D. Büchler, D. Jayaraman, D. Kalashnikov, D. Sadigh, 707 E. Johns, E. Foster, F. Liu, F. Ceola, F. Xia, F. Zhao, F. V. Frujeri, F. Stulp, G. Zhou, G. S. Sukhatme, 708 G. Salhotra, G. Yan, G. Feng, G. Schiavi, G. Berseth, G. Kahn, G. Yang, G. Wang, H. Su, H.-S. Fang, H. Shi, H. Bao, H. B. Amor, H. I. Christensen, H. Furuta, H. Bharadhwaj, H. Walke, H. Fang, H. Ha, I. Mordatch, 710 I. Radosavovic, I. Leal, J. Liang, J. Abou-Chakra, J. Kim, J. Drake, J. Peters, J. Schneider, J. Hsu, J. Vakil, J. Bohg, J. Bingham, J. Wu, J. Gao, J. Hu, J. Wu, J. Wu, J. Sun, J. Luo, J. Gu, J. Tan, J. Oh, J. Wu, J. Lu, J. Yang, 711 J. Malik, J. Silvério, J. Hejna, J. Booher, J. Tompson, J. Yang, J. Salvador, J. J. Lim, J. Han, K. Wang, K. Rao, 712 K. Pertsch, K. Hausman, K. Go, K. Gopalakrishnan, K. Goldberg, K. Byrne, K. Oslund, K. Kawaharazuka, 713 K. Black, K. Lin, K. Zhang, K. Ehsani, K. Lekkala, K. Ellis, K. Rana, K. Srinivasan, K. Fang, K. P. Singh, 714 K.-H. Zeng, K. Hatch, K. Hsu, L. Itti, L. Y. Chen, L. Pinto, L. Fei-Fei, L. Tan, L. J. Fan, L. Ott, L. Lee, L. Weihs, M. Chen, M. Lepert, M. Memmel, M. Tomizuka, M. Itkina, M. G. Castro, M. Spero, M. Du, M. Ahn, 715 M. C. Yip, M. Zhang, M. Ding, M. Heo, M. K. Srirama, M. Sharma, M. J. Kim, N. Kanazawa, N. Hansen, 716 N. Heess, N. J. Joshi, N. Suenderhauf, N. Liu, N. D. Palo, N. M. M. Shafiullah, O. Mees, O. Kroemer, 717 O. Bastani, P. R. Sanketi, P. T. Miller, P. Yin, P. Wohlhart, P. Xu, P. D. Fagan, P. Mitrano, P. Sermanet, 718 P. Abbeel, P. Sundaresan, Q. Chen, Q. Vuong, R. Rafailov, R. Tian, R. Doshi, R. Mart'in-Mart'in, R. Baijal, 719 R. Scalise, R. Hendrix, R. Lin, R. Qian, R. Zhang, R. Mendonca, R. Shah, R. Hoque, R. Julian, S. Bustamante, S. Kirmani, S. Levine, S. Lin, S. Moore, S. Bahl, S. Dass, S. Sonawani, S. Tulsiani, S. Song, S. Xu, S. Haldar, 720 S. Karamcheti, S. Adebola, S. Guist, S. Nasiriany, S. Schaal, S. Welker, S. Tian, S. Ramamoorthy, S. Dasari, 721 S. Belkhale, S. Park, S. Nair, S. Mirchandani, T. Osa, T. Gupta, T. Harada, T. Matsushima, T. Xiao, T. Kollar, 722 T. Yu, T. Ding, T. Davchev, T. Z. Zhao, T. Armstrong, T. Darrell, T. Chung, V. Jain, V. Kumar, V. Vanhoucke, 723 W. Zhan, W. Zhou, W. Burgard, X. Chen, X. Chen, X. Wang, X. Zhu, X. Geng, X. Liu, X. Liangwei, X. Li, 724 Y. Pang, Y. Lu, Y. J. Ma, Y. Kim, Y. Chebotar, Y. Zhou, Y. Zhu, Y. Wu, Y. Xu, Y. Wang, Y. Bisk, Y. Dou, Y. Cho, Y. Lee, Y. Cui, Y. Cao, Y.-H. Wu, Y. Tang, Y. Zhu, Y. Zhang, Y. Jiang, Y. Li, Y. Li, Y. Iwasawa, 725 Y. Matsuo, Z. Ma, Z. Xu, Z. J. Cui, Z. Zhang, Z. Fu, and Z. Lin. Open x-embodiment: Robotic learning 726 datasets and rt-x models, 2024. URL https://arxiv.org/abs/2310.08864. 727 728 C. Pandarinath and S. J. Bensmaia. The science and engineering behind sensitized brain-controlled bionic hands.

Physiological Reviews, 102(2):551–604, 2022.
 C Pandarinath P Nuvujukian C H Blabe B L Sorice I Saab F R Willett I R Hochberg K V Shenov.

- C. Pandarinath, P. Nuyujukian, C. H. Blabe, B. L. Sorice, J. Saab, F. R. Willett, L. R. Hochberg, K. V. Shenoy, and J. M. Henderson. High performance communication by people with paralysis using an intracortical brain-computer interface. *elife*, 6:e18554, 2017.
- F. C. Pei, J. Ye, D. M. Zoltowski, A. Wu, R. H. Chowdhury, H. Sohn, J. E. O'Doherty, K. V. Shenoy, M. Kaufman, M. M. Churchland, M. Jazayeri, L. E. Miller, J. W. Pillow, I. M. Park, E. L. Dyer, and C. Pandarinath. Neural latents benchmark '21: Evaluating latent variable models of neural population activity. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*, 2021. URL https://openreview.net/forum?id=KVMS3f14Rsv.
- S. M. Peterson, S. H. Singh, B. Dichter, K. Tan, C. DiBartolomeo, D. Theogarajan, P. Fisher, and J. Parvizi.
 Ajile12: Long-term naturalistic human intracranial neural recordings and pose. *Scientific Data*, 9(1):184, 2022.
 ISSN 2052-4463. doi: 10.1038/s41597-022-01280-y. URL https://doi.org/10.1038/s41597-022-01280-y.
- T. M. Pham, T. Bui, L. Mai, and A. Nguyen. Out of order: How important is the sequential order of words in a sentence in natural language understanding tasks?, 2021. URL https://arxiv.org/abs/2012.15180.
- K. M. Quick, J. L. Mischel, P. J. Loughlin, and A. P. Batista. The critical stability task: quantifying sensory-motor control during ongoing movement in nonhuman primates. *Journal of Neurophysiology*, 120(5):2164–2181, 2018.
- S. Reed, K. Zolna, E. Parisotto, S. G. Colmenarejo, A. Novikov, G. Barth-Maron, M. Gimenez, Y. Sulsky, J. Kay, J. T. Springenberg, T. Eccles, J. Bruce, A. Razavi, A. Edwards, N. Heess, Y. Chen, R. Hadsell, O. Vinyals, M. Bordbar, and N. de Freitas. A generalist agent, 2022. URL https://arxiv.org/abs/2205.06175.
- F. Rizzoglio, E. Altan, X. Ma, K. L. Bodkin, B. M. Dekleva, S. A. Solla, A. Kennedy, and L. E. Miller. Monkey-to-human transfer of brain-computer interface decoders. *bioRxiv*, 2022. doi: 10.1101/2022.11.12.515040. URL https://www.biorxiv.org/content/early/2022/11/13/2022.11.12.515040.
- P. T. Sadtler, K. M. Quick, M. D. Golub, S. M. Chase, S. I. Ryu, E. C. Tyler-Kabara, B. M. Yu, and A. P. Batista. Neural constraints on learning. *Nature*, 512(7515):423–426, 2014.
- 755 M. Sato, K. Tomeoka, I. Horiguchi, K. Arulkumaran, R. Kanai, and S. Sasai. Scaling law in neural data: Non-invasive speech decoding with 175 hours of eeg data, 2024. URL https://arxiv.org/abs/2407.07595.

756 757 758	S. Schneider, J. H. Lee, and M. W. Mathis. Learnable latent embeddings for joint behavioural and neural analysis. <i>Nature</i> , May 2023. ISSN 1476-4687. doi: 10.1038/s41586-023-06031-6. URL https://doi.org/10.1038/s41586-023-06031-6.
759 760 761	I. Schubert, J. Zhang, J. Bruce, S. Bechtle, E. Parisotto, M. Riedmiller, J. T. Springenberg, A. Byravan, L. Hasenclever, and N. Heess. A generalist dynamics model for control, 2023. URL https://arxiv.org/abs/2305_10912
762 763 764	N. M. M. Shafiullah, Z. J. Cui, A. Altanzaya, and L. Pinto. Behavior transformers: Cloning k modes with one stone, 2022. URL https://arxiv.org/abs/2206.11251.
765 766 767	Q. Simeon, L. Venâncio, M. A. Skuhersky, A. Nayebi, E. S. Boyden, and G. R. Yang. Scaling properties for artificial neural network models of a small nervous system. In <i>SoutheastCon 2024</i> , pages 516–524. IEEE, 2024.
768 769	N. A. Steinmetz, P. Zatka-Haas, M. Carandini, and K. D. Harris. Distributed coding of choice, action and engagement across the mouse brain. <i>Nature</i> , 576(7786):266–273, 2019.
770 771 772	I. H. Stevenson. Tracking advances in neural recording. Statistical Neuroscience Lab, University of Connecticut, 2023. URL https://stevenson.lab.uconn.edu/scaling/. Accessed September 6, 2024.
773 774	J. Su, Y. Lu, S. Pan, A. Murtadha, B. Wen, and Y. Liu. Roformer: Enhanced transformer with rotary position embedding, 2023. URL https://arxiv.org/abs/2104.09864.
775 776 777	S. J. Talukder, J. J. Sun, M. K. Leonard, B. W. Brunton, and Y. Yue. Deep neural imputation: A framework for recovering incomplete brain recordings. In <i>NeurIPS 2022 Workshop on Learning from Time Series for Health</i> , 2022. URL https://openreview.net/forum?id=c9qFg8UrIcn.
778 779	R. Taori, A. Dave, V. Shankar, N. Carlini, B. Recht, and L. Schmidt. Measuring robustness to natural distribution shifts in image classification, 2020. URL https://arxiv.org/abs/2007.00644.
780 781 782	R. Thapa, B. He, M. R. Kjaer, H. Moore, G. Ganjoo, E. Mignot, and J. Zou. Sleepfm: Multi-modal representation learning for sleep across brain activity, ecg and respiratory signals, 2024. URL https://arxiv.org/abs/2405. 17766.
783 784 785	A. W. Thomas, C. Ré, and R. A. Poldrack. Self-supervised learning of brain dynamics from broad neuroimaging data, 2023.
786 787	A. E. Urai, B. Doiron, A. M. Leifer, and A. K. Churchland. Large-scale neural recordings call for new insights to link brain and behavior. <i>Nature Neuroscience</i> , 25:11–19, 2022. doi: 10.1038/s41593-021-00980-9.
788 789 790	A. Wang, A. Singh, J. Michael, F. Hill, O. Levy, and S. R. Bowman. Glue: A multi-task benchmark and analysis platform for natural language understanding, 2019. URL https://arxiv.org/abs/1804.07461.
791 792 793	C. Wang, V. Subramaniam, A. U. Yaari, G. Kreiman, B. Katz, I. Cases, and A. Barbu. BrainBERT: Self- supervised representation learning for intracranial recordings. In <i>The Eleventh International Conference on</i> <i>Learning Representations</i> , 2023a. URL https://openreview.net/forum?id=xmcYx_reUn6.
794 795	E. Y. Wang, P. G. Fahey, K. Ponder, Z. Ding, A. Chang, T. Muhammad, S. Patel, Z. Ding, D. Tran, J. Fu, et al. Towards a foundation model of the mouse visual cortex. <i>bioRxiv</i> , 2023b.
796 797	H. Wang, T. Fu, Y. Du, W. Gao, K. Huang, Z. Liu, P. Chandak, S. Liu, P. Van Katwyk, A. Deac, et al. Scientific discovery in the age of artificial intelligence. <i>Nature</i> , 620(7972):47–60, 2023c.
798 799 800 801	T. Wang, A. Roberts, D. Hesslow, T. L. Scao, H. W. Chung, I. Beltagy, J. Launay, and C. Raffel. What language model architecture and pretraining objective work best for zero-shot generalization?, 2022. URL https://arxiv.org/abs/2204.05832.
802 803 804	F. R. Willett, D. T. Avansino, L. R. Hochberg, J. M. Henderson, and K. V. Shenoy. High-performance brain- to-text communication via handwriting. <i>Nature</i> , 593(7858):249–254, May 2021. ISSN 1476-4687. doi: 10.1038/s41586-021-03506-2. URL https://www.nature.com/articles/s41586-021-03506-2.
805 806	J. R. Wolpaw, N. Birbaumer, D. J. McFarland, G. Pfurtscheller, and T. M. Vaughan. Brain–computer interfaces for communication and control. <i>Clinical neurophysiology</i> , 113(6):767–791, 2002.
807 808	M. Wortsman, P. J. Liu, L. Xiao, K. E. Everett, A. A. Alemi, B. Adlam, J. D. Co-Reyes, I. Gur, A. Kumar, R. Novak, J. Pennington, J. Sohl-Dickstein, K. Xu, J. Lee, J. Gilmer, and S. Kornblith. Small-scale proxies

R. Novak, J. Pennington, J. Sohl-Dickstein, K. Xu, J. Lee, J. Gilmer, and S. Kornblith. Small-scale proxies
 for large-scale transformer training instabilities. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=d8w0pmvXbZ.

810 811	W. Xia, R. de Charette, C. Öztireli, and JH. Xue. Umbrae: Unified multimodal brain decoding. In <i>European Conference on Computer Vision (ECCV)</i> , 2024.
812 813	C. Yang, M. B. Westover, and J. Sun. Biot: Cross-data biosignal learning in the wild, 2023. URL https://arxiv.org/abs/2305_10351
814	// di x11.01 g/ db0/ 2000.10001.
815 816	C. Yang, J. Li, X. Niu, X. Du, S. Gao, H. Zhang, Z. Chen, X. Qu, R. Yuan, Y. Li, J. Liu, S. W. Huang, S. Yue, and G. Zhang. The fine line: Navigating large language model pretraining with down-streaming capability
817	analysis, 2024. URL https://arxiv.org/abs/2404.01204.
818	I.V. I.I. Collingon I. Wahha and D. Count. Naural data transforman 2. Multi contact protocining for noural
819	5. Fe, J. L. Conniger, L. wende, and K. Gaunt. Neural data transformer 2: Multi-context pretraining for neural spiking activity. In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> , 2023. URL
820	https://openreview.net/torum/id=CBBtMnilg.
821 822	Z. Yuan, F. Shen, M. Li, Y. Yu, C. Tan, and Y. Yang. Brainwave: A brain signal foundation model for clinical applications, 2024. URL https://arxiv.org/abs/2402.10251.
823	
824 825	Y. Zhang, Y. Wang, D. Jimenez-Beneto, Z. Wang, M. Azabou, B. Richards, O. Winter, T. I. B. Laboratory, E. Dyer, L. Paninski, et al. Towards a" universal translator" for neural dynamics at single-cell, single-spike
826	resolution. arXiv preprint arXiv:2407.14668, 2024.
827	
828	
820	
920	
030	
001	
032	
033	
034	
835	
030	
837	
000	
039	
040	
041	
042	
844	
8/5	
8/6	
847	
848	
849	
850	
851	
852	
853	
854	
855	
856	
857	
858	
850	
860	
861	
260	
002	
003	

A RELATED WORKS AND A PROPOSED TAXONOMY

866 Neural data is sufficiently diverse so as to support many distinct efforts to train large neural data 867 models. The scale of pretraining is somewhat larger in the non-implanted modalities, where data 868 is more abundant. The largest EEG models have reached a scale of 2.5K (Jiang et al., 2024) to 40K hours of data (Yuan et al., 2024), or higher volumes if also considering non-brain biosignals (EKG) (Yang et al., 2023; Thapa et al., 2024). Current fMRI models operate in the 1K (Thomas et al., 870 2023) to 7K (Caro et al., 2024) hour range. The largest models in these studies are in the 0.1B-1B 871 parameter range. Intracranial modalities, including sEEG (Wang et al., 2023a; Chau et al., 2024), 872 ECoG (Talukder et al., 2022; Peterson et al., 2022), and spiking activity (Wang et al., 2023b), have 873 thus far been studied at an order of magnitude smaller scales of data and model size (20-1000 hours, 874 <0.1B parameters). 875

Direct scaling on neural data modeling should be distinguished from NeuroAI efforts (Doerig et al., 2022) to measure how models of the human sensorimotor experience (e.g. language, vision, audio models) predict neural data (Antonello et al., 2024). However, as multimodal efforts begin to blur this distinction (Benster et al., 2024; Xia et al., 2024), care will be required to distinguish advances in modeling neural data, embodied data, or their interaction.

Comparing neural data models Current efforts to understand scaling in neural data Simeon et al. (2024); Sato et al. (2024) will have their reach limited by the specificity of every neural dataset. A
 meta-challenge for the field is understanding how different parameters (species, brain area, modality, task) impact scaling properties. This would be greatly aided by development of reporting practices for
 different neural data models. To facilitate comparison, we create a model card (Mitchell et al., 2019)
 for NDT3 in Section D. In addition to the standard model card, we propose reporting an additional taxonomy to aid comparisons across neural data models, using two concepts.

First: neural data models can be conceptualized as modeling slices of the *plenneural function*, inspired
by the plenoptic function in vision (Adelson et al., 1991). The plenoptic function is a model of an
idealized eye which parameterizes all possible images with 7 dimensions: 4D to describe the global
spacetime of the view, 2D to describe viewing angle (spherical) or coordinate (Cartesian) of the
image, and 1D for wavelength. Since neural data models are primarily interested in circumscribed
systems rather than the physical world, a similar global coordinate system (e.g. 4D for all possible
electric potentials) would be uninformative. We thus propose reporting more qualitative coordinates:

- 895
- 896
- 897
- 1. Identity: The network or individual being recorded.
- 2. Task: The behavior, stimuli, or other activity the network is reflecting.
- 3. Spacetime: Coordinates specified in a network-local coordinate frame (e.g. brain area).

Second: The modeled extent of this plenneural function is conveniently discretized in three resolutions in a Transformer-like sequence modeling framework: the token, the sequence, and the full training data. The token is the most granular unit of data being modeled; NDT3 models neural populations 32 neurons at a time, in 20ms bins. At the sequence input level, NDT3 models inputs from single humans or monkeys, across 128-256 neurons in 2 second snippets, while performing effectively one "movement." Finally, NDT3's pretraining spans dozens of individuals, records motor and premotor areas over 2.5K hours, over a variety of arm and hand movements.

906 907

908

909

B SUPPLEMENTARY RESULTS

B.1 FURTHER TESTS OF HELD-OUT ANGLE GENERALIZATION.

910 To further support our single-dataset illustration of attractor structure in pretrained NDT3, we evaluate 911 reach angle generalization across sessions in the isometric, force-based dataset (Monkey J) setting and 912 a second, manipulandum-based (Monkey J) setting with monkey movement. To begin, we visualize 913 the separability of these neural data with LDA (Fig. 6A Neural Data). We next train decoders on 914 every pair of angles separated by 90 degrees (one shown) and plot predictions on held-out trials 915 from all angles. NDT3 and WFs, here directly fit to high-D data instead of after PCA-LDA, both fail to extrapolate to held-out angles, consistent with Rizzoglio et al. (2022). We quantify prediction 916 performance in Fig. 4D. These plots again show that held-out generalization is subtrivial, while being 917 entirely consistent with pretraining's overall narrative of improved performance in all conditions.



Figure 6. A. Two monkey 2D center-out datasets with 8 angular conditions. LDA projections show monkey J's data is distinctly more separable than monkey V's. We then test generalization of behavioral decoding to held-out angles after training on 2 of 8 angles (boxed in red). All NDT3s and the WF produce predictions constrained between the held-out angles. **B.** We quantify performance for decoding on each angle with respect to distance from the central angle between the held-in angles. We average performance over all 8 central angles.



Figure 7. Ablation of covariate masking on an open 2D Cursor + Click dataset. Covariate inputs are completely masked in inference for the default NDT3, and autoregressively generated in the ablation.

Methods for PCA-LDA. Our PCA-LDA preparation used in Fig. 4E-F adheres closely to our standard data preparation and is directly comparable to the inputs received by the deep networks. To begin, we smooth all spike counts with an exponential kernel, as in our Wiener Filter baseline (Section C.3). We then fit PCA to the high-dimensional (96D) neural activity at each timestep and extract the top 10 PCs. Finally, we fit linear discriminant analysis (LDA) to reduce this 10D PC space to neural activity to 2D. The training data for PCA and LDA are both restricted to only the held-in angles. Note that LDA uses categorical labels to separate the 3 held-in reach directions but is applied without class labels after fitting. While this train-time labeling is technically not given to NDT3, this does not influence our argument that NDT3 pretraining should be capable of angular generalization, as NDT3 does not confuse the reach angle of held-in predictions. The final Wiener Filter is fit to the 2D data and uses 1 second of history.

B.2 Ablations

We ablate the major design decisions made to enable NDT3's large scale pretraining. These ablations
give us confidence that NDT3 overcomes the basic challenges we encountered in development, but
compute restrictions prevent more exhaustive comparisons or exploration of model design space.
We encourage further work exploring the influence of different hyperparameters. In these plots, we
distinguish validation split performance and evaluation split performance, which is computed by
batch-mode prediction (not the costly streaming evaluation used throughout main experiments).

969 Covariate dropout We find the default next-step prediction objective fails for learning decoding
 970 of highly autocorrelated covariate timeseries, perhaps because simply relying on teacher-forced
 971 behavioral inputs provides a severe shortcut that prevents learning of a proper neural to behavior
 decoding map (Bachmann and Nagarajan, 2024). Different time-series models have addressed

989 990 991



Figure 8. Ablations evaluated based on upstream evaluation split. **A., B.** Ablation of BCI control tokens **C.** Ablation of neural objective and covariate MSE objective in favor of classification over quantized covariates.

992 this by adopting convolutional input-output layers (Lea et al., 2016), tokenizing along temporal 993 dimensions (Das et al., 2024), or learning with contrastive objectives (Chau et al., 2024; Kostas 994 et al., 2021). We avoid introducing architectural modifications and instead adopt a simple dropout 995 procedure that masks a portion of covariate inputs some fraction of the time. Specifically, on every training batch, two random numbers are drawn. The first, $M \sim U[0,1]$, determines what fraction of 996 covariate inputs should be masked. On 90% of batches, we also sample $T \sim U[0,2]$ seconds, such 997 that the mask is only applied after timestep T. That is, on 90% of batches, the model is provided a 998 prefix-prompt. We do not block losses on this prefix as in prefix-LMs (Wang et al., 2022). Pretraining 999 metrics for validation and evaluation are always computed with a prefix and full masking of non-prefix 1000 timesteps. In Fig. 7, we ablate covariate masking (which also removes the prefix logic), and tune on a 1001 2D Cursor + Click task. The ablated model performs subtrivially with student-forced predictions 1002 provided as input at test time. Note that the ablated model performs trivially with masked inputs (not 1003 shown).

1004 BCI-phase and return conditioning NDT3's pretraining includes several hundred hours of BCI 1005 control data, where the covariates were set by another decoder. We introduced phase and return conditioning tokens to differentiate the several types of BCI control data from recorded behavior. Specifically, in BCI data, NDT3 receives input tokens specifying what fraction of the behavior reflects 1008 neural input (BCI control is on) vs programmatic input (BCI control is off, as in open loop BCI 1009 calibration). Further, we provide inputs encoding reward (trial success) when trials change, and return 1010 (future reward over a 10 second horizon, which crosses data boundaries). This design is intended to evaluate the potential for a Decision-Transformer like offline learning strategy for improved online 1011 control, but we do not discuss this in this work. In Fig. 8A, B, we focus on whether these inputs 1012 improves pretraining loss and R^2 in validation splits, which has contains BCI data, and the held-out 1013 evaluation split containing only monkey behavior. The figures show that the ablation significantly 1014 decreases validation split performance, and causes a slightly earlier stopping point leading to worse 1015 evaluation performance. Note both models early before the full training budget of 400 epochs. 1016

Neural reconstruction objective All main NDT3 models used a neural reconstruction objective inherited from the self-supervised learning pretraining from NDT2. We ablate this choice post-hoc and see it may actually minorly harm pretraining (validation split), though the neural objective doesn't harm evaluation split decoding (Fig. 8C). Note the scalar weighting of neural vs covariate objectives were set to be roughly balanced in pretraining. Section B.3 provides a downstream analysis on the standalone value of the neural reconstruction objective.

MSE over classification In robotics and certain generalist models (Schubert et al., 2023), continuous action spaces are sometimes better decoded and controlled when quantized (Shafiullah et al., 2022).
 This is because MSE is an insufficient objective when the output distribution is multimodal (e.g. one of two possible paths in robotics). While it seems unlikely that the close relationship between



Figure 9. Scaling of unsupervised pretraining After pretraining models up to 200 hours with only the neural reconstruction objective, we fine-tune models in a similar multiscale evaluation as in Section 3.1, with a newly initialized behavioral readout. We use the CST task here. A. left shows neural loss scales with neural-only pretraining. The right panel plots the neural loss, also present in the standard evaluation, also scales with joint pretraining. B. left shows that with neural pretraining, downstream decoding performance saturates at a flat performance after 25 hours. This is compared against the nonsaturated scaling from joint pretraining. Colorbar is common for all plots, and Xs are from-scratch NDT3 models.

movement behavior and motor cortex is multimodal, multimodal behavior may be appropriate when pretrained on heterogeneous data, i.e. when similar neural activity corresponds to different behavior in two datasets. We attempted such a quantization, including HL-Gauss smoothing (Farebrother et al., 2024) which we found to help; but this does not recover the performance of the default MSE objective (Fig. 8C) on the evaluation split. We found this performance gap persisted under fine-tuning (not shown). This suggests that NDT3 is differentiates neural data inputs from different datasets.

Patch size NDT2 and NDT3 both tokenize neural data by patching them into fixed size clusters. It is unclear whether transfer learning might occur for sub-token features, which motivates the use of smaller tokens in larger datasets that might afford it (Caron et al., 2021). We change patch size to 16 and show this performs slightly worse in the 45M 200h model (Fig. 8)C. Smaller patches (and subsequent increased neural tokens) may be more beneficial in the larger scale models, but their benefit must be weighed against their increased compute burden.

1057

1043 1044

1058 1059

1059 B.3 ISOLATED SCALING IN NEURAL DATA 1060

1061 Due to NDT3's joint modeling of behavior and neural data, it is difficult to dissociate whether scaling 1062 gains in behavior come from improved behavioral or neural priors. To assess whether NDT3 can 1063 scale solely from neural data modeling, we pretrain a new set of 45M parameter models up to 200 1064 hours with only causal neural data modeling objective. As before, we then tune to a downstream decoding task, in this case, a Critical Stability Task dataset (CST). From a representation probing perspective, improved downstream performance implies higher quality neural representations. We 1066 use the standard single-stream autoregressive modeling objective as in the rest of this work in the 1067 downstream setting, we find direct linear probing of neural representations perform worse. Fig. 9 1068 compares the scaling on downstream neural and behavioral metrics after the standard fine-tuning 1069 procedure. 1070

Fig. 9A shows that downstream neural reconstruction improves with increased pretraining data either
when using only the neural objective or both neural and behavioral objectives (as in the standard
setting). The joint pretraining achieves advances neural metrics in all settings, illustrating that
decoding behavior is a complementary objective to neural data reconstruction even for representation
learning.

Fig. 9B contrasts decoding curves in the two pretraining settings, in that neural pretraining has saturated decoding after just 25 hours of pretraining. This is consistent with the interpretation that the behavioral readout reflects only one aspect of the neural data. Together with the neural metric plots, this analysis shows scaling over solely neural data is possible, but also that decoding behavior is a complementary pretraining objective for improving neural representation learning and decoding.



Figure 10. Three regimes of NDT3 training for handwriting decoding. We show validation loss and character error rates for example runs of from-scratch and fine-tuned NDT3s.

1107 B.4 PRETRAINING DOES NOT BENEFIT FALCON H2 (HANDWRITING) 1108

We also evaluated NDT3 for decoding of letters in a human-open loop handwriting task (FALCON H2). Although this is also a motor cortical decoding task, we excluded H2 from NDT3's aggregate evaluation since it is a sequence-to-sequence as opposed to continuous task. To apply NDT3 to this task, we pool neural tokens at each timestep and add a linear projection and optimize with a CTC loss (Graves et al., 2006). We maintain the default neural reconstruction loss and causal attention mask, and do not apply data augmentation.

1115 Note that RNNs are the current standard architecture for communication tasks like H2 (Karpowicz 1116 et al., 2024; Willett et al., 2021). Training and tuning was less stable than for our continuous decoding tasks and required more extensive hyperparameter tuning, perhaps because the overall dataset size 1117 remains small (<1k samples), specific parameters are listed in the codebase. We observe three regimes 1118 in both training and fine-tuning. First, the model can fail to achieve an initial learning period. Second, 1119 the model can achieve reasonable nontrivial solutions, comparable to expected performance for 1120 unaugmented RNNs (though we do not quantify this). Third, some models will exhibit learning 1121 instabilities that resolve in significantly improved performance. We illustrate these regimes in example 1122 validation curves below. Overall, the third regime is rarely achieved. More relevant to the main 1123 narrative of this work, fine-tuning appears to degrade both final solution quality and reduces the range 1124 of nontrivial hyperparameters (not shown). Investigating a sequence to sequence objective over CTC 1125 loss would be valuable future work.

- 1126
- 1127 1128

B.5 MULTISCALE DECODING ON INDIVIDUAL MOTOR TASKS

Fig. 11A plots model performance for each of the 31 evaluation settings we study in the eight
primary evaluation datasets we use. Studying any individual dataset will yield variable conclusions
on whether pretraining structure is helpful, underscoring the need for proposed foundation models
to be evaluated across many different datasets. Here specifically we see the most clear scaling with
pretraining data (color gradient with red on top) in the Critical Stability Task and Bimanual Task.
FALCON tasks and Self-paced Reach appear minimally affected by scaling pretraining data, in that



Figure 11. A. Fine-tuning evaluations for individual datasets. Performance on both held-out (left) and held-in (right) splits are shown side by side by FALCON datasets. We shade the standard deviation of 3 model seeds in fine-tuning. Different tasks show substantial variability in benefit from pretraining. B. We show example predictions of a pretrained (45M 200h) and from-scratch NDT3 for the 2D + Click Cursor task to give a sense of what different prediction performances mean in terms of open loop data prediction. Numbers in legend are the R^2 for that model's predictions in the shown snippet.

1134

- 1181
- 1182
- 1183
- 1184 1185
- 1186
- 1187



Figure 12. A replication of Fig. 4A/B, but showing results for both sequential and joint fine-tuning in each setting. As before, each model here tunes with some data from other settings (e.g. cross-subject data for the cross-subject panel), and a fixed amount of data on the test session. Sequential tuning first tunes on other-setting data, and then tunes on test session data. Joint tuning uses all data at once. In Fig. 4A/B, we only showed the better choice for each panel, i.e. we showed joint tuning for cross-session data, and sequential tuning for the other settings. In addition, we overlay how the cross-subject tuned models perform when applied to half-token shifted test data, confirming that subject and session data transfer similarly to shifted test data.

1205 either pretrained models are generally slightly above a from scratch model at all data scales with no 1206 particular best pretrained model. The 2D + Click and Grasp datasets uniquely show high variability 1207 in model performance and strong degradation of the 350M 2 khr model at low data scales. Grasp 1208 instability was so high that we trained 9 seeds instead of the standard 3 to better estimate model 1209 performance. We propose this degradation is due to the instability of full fine-tuning of large models 1210 at the extremely low data scales these datasets present (e.g. 2.5 minutes at the 25% scaling). Finally, 1211 we remind that the 2D + Click, FALCON H1, and 1D Grasp Force tasks are datasets from human 1212 participants that are included in the 2 khr pretraining. Surprisingly, we see no particular benefit to the 2 khr model. 1213

These scaling plots also provide more precise context for baseline performance. NDT2 performs particularly poorly in the low data regime, while Wiener Filters perform poorly in the high data regimes.

In Fig. 11B, we illustrate qualitative predictions on private datasets. These visualizations show a diversity in covariate timescales and structure. They also illustrate that the summary R^2 obscure several features of model predictions. For example, pretrained models in Cursor Y tend have false positive deflections in movement. R^2 also is not easily comparable in tasks with continuous dynamics (CST) vs. transient dynamics (Cursor G1).

1222 1223

1224

B.6 SEQUENTIAL TUNING IS SIMILAR TO JOINT TUNING

1225 In Section 3.2, we showed that channel shuffling and half-token shifts were sufficient to reduce 1226 cross-session transfer to the extent of cross-subject transfer. Here we add a methodological subtlety on how the tuning is done. Given cross-context data and a test session, we can either jointly tune 1227 on all data (as we do in our primary evaluation), or sequentially tune on the cross-context data 1228 and test session. We find sequential tuning is particularly necessary for successful subject transfer 1229 of from-scratch NDT3 models, but that it slightly underperformed joint tuning on cross-session 1230 models (Fig. 12A). Seeing that sequential tuning is mainly advantageous for from-scratch models, 1231 we speculate that sequential tuning is particularly helpful for filtering learning signals in highly 1232 heterogeneous data. In Fig. 4A/B, for clarity, we reported jointly tuned results for cross-session data 1233 and sequentially tuned results otherwise.

1234 1235

1236

B.7 AGGREGATE PERFORMANCE ON ALL NDT3 MODELS WITH SIGNIFICANCE TESTS

1237 We additionally report the average performance of an NDT2 model pretrained with 100 hours of 1238 human data and two NDT3 models pretrained with 25 hours and 70 hours. These models are placed 1239 in context with the models from Fig. 3D, in Fig. 13A. Note that for NDT3 models each successively 1240 larger data scale uses a strict superset of data from smaller scales. We also provide the p-values 1241 computed for the significance of the difference between each pair of models in Fig. 13B. P-values are 1242 computed as FDR-corrected pairwise t-tests. The 350M 2 khr model has p < 0.06 improvements



Figure 13. A. A replication of Fig. 3D including an additional 25 hr and 70 hr model. The two additional models show the precise performance we measure may be noisily related to to pretraining data scale. B. Heatmap of differences between performances of pairs of models with significance tested with FDR-corrected pairwise t-tests. Note that coloring is used here to indicate differences, not significance. Positive numbers with significance indicate the row model outperforms the column model.



Figure 14. Pretraining curves shown for 45M parameter NDT3s at 1.5 hour, 25 hour, and 70 hours, in addition to the curves for the 2 khr 350M parameter NDT3. We separately show metrics on neural and behavioral objectives through training. Since early stopping is used in model selection, we verify here that neither objective is overfits significantly, except for the 1.5 hour model's neural objective.

over all but the 350M 200 hr model. Interestingly, all other pretrained models, except the 25 hr model, appear equivalent, at least statistically. We presume this is due to the fact that our evaluation of 31 task settings may be insufficiently large.

NDT2 performs relatively poorly in our evaluation. This is true even when tuning from the public checkpoint trained on 100 hours of neural data from humans, though tuning does in general improve over the from-scratch NDT2 training. We believe pretrained NDT2's performance gap with NDT3 is partially due to NDT2's need to newly initialize decoding layers in each downstream task, which increases NDT2's dependence on thorough hyperparameter tuning. This makes NDT2 a poor candidate for a foundation model. Section C.3 and Section C.6 describe a number of methodological innovations that likely each contribute to the remaining performance differences between NDT2 and NDT3.

1292

1294

1293 B.8 NEURAL VS BEHAVIORAL OBJECTIVES

1295 In this work, the neural objective is present mainly as an auxiliary objective to improve downstream decoding. We see that neural and behavioral objectives are complementary in Section B.3.



Figure 15. Downstream evaluation conducted through pretraining of a 45M 25 hour NDT3, for the 25% scaling of the bimanual task, provided in context of the final evaluation achieved by the 10% and 50% bimanual task scale settings for the 45M 25 hour models, and the 25% scalings for the from-scratch NDT3 and best NDT3 (350M 2 khr). The benefits from pretraining are attained by epoch 20, and plateaus without overfitting for the rest of pretraining.

1320

¹³¹⁸ In Fig. 14, we provide some additional context, showing that neural objectives also improve through pretraining, i.e. that we are not overfitting the neural objective and thus degrading decoding.

1321 B.9 DOWNSTREAM PROBE THROUGH PRETRAINING

1323 Upstream and downstream performance are generally assumed to be correlated in pretraining studies. 1324 This assumption supports the use of a single evaluation at the end of pretraining, rather than throughout 1325 pretraining, particularly if pretraining is not overfit. It is possible, however, that the narrower tasks 1326 used in evaluation may be learned or even overfit earlier than the general pretraining task. To assess this, we conduct a small downstream probe of checkpoints every 20 epochs of pretraining for a 1327 45M 25 hour NDT3. This probe measures the performance specifically at the 25% scaling of the 1328 bimanual task. In Fig. 15, we show this progression and compare the variability of checkpoints in 1329 pretraining with the differences across different task scales and pretraining models. The downstream 1330 probe for this task shows quick benefit from pretraining, providing the full gain over from-scratch 1331 models by the first checkpoint we evaluate at epoch 20, and then plateaus. This contrasts with the 1332 pretraining plots in Fig. 14. The mismatch between upstream and downstream progress here differs 1333 from correlated upstream-downstream progress in single-epoch language model studies (Yang et al., 1334 2024), which may be due to differences in domain or pretraining scale.

1335

1336 B.10 EVALUATION ON THE NEURAL LATENTS BENCHMARK

The Neural Latents Benchmark (Pei et al., 2021) is a benchmark for evaluating latent variable modeling of neural activity. This evaluation differs from NDT3's direct decoding evaluations in that all of NLB's metrics are derived from acausally inferred neural firing rates, akin to a pixel-level objective in computer vision. Note specifically that the decoding metric officially reported in the NLB is derived from submitted firing rates on the evaluation server; report decoding performance on NLB datasets without modeling neural activity is not an expected use of the benchmark.

Despite its distinct setting, the NLB provides two well-established datasets on which to evaluate
models of motor cortical activity, providing context for NDT3's performance outside of decoding.
To apply NDT3 to the NLB, we tuned NDT3 in a new downstream task where insert new a token at
each timestep, from which we linearly decode firing rates of held out neurons. We performed this
tuning separately for each of the maze and random target tasks (RTT), reporting the resulting co-bps
scores in Table 1. NDT3 performs poorly. From-scratch training on the single-session benchmark
datasets underperform NDT1 in all tasks. Tuning from pretrained models improved performance,



Figure 16. For 3 monkeys datasets at 10% scale, we extend a HP sweep to 5 LRs and dropout in [0.0, 0.1, 0.3]1359 (vs default 0.1). For fine-tuning, we also sweep weight decay in [0.0001, 0.01, 0.1] (vs default 0.1), while for 1360 from-scratch models we also sweep Transformer width ([256, 512, 1024]) vs default 512. This yields a 45-model sweep on 1 seed. We compare the range of scores achieved by this larger sweep against the standard 3 LR x 3 1361 seed sweep. 1362

sample efficiency, and robustness to hyperparameters (only performance is reported here), but did not 1365 dramatically change model competitiveness. Zhang et al. (2024) reported the value of more diverse neuron-level objectives for neural activity prediction, though their submission also dramatically 1367 underperforms NLB SoTA, warranting critical examination of whether pretraining benefits low-level 1368 modeling of neural activity. 1369

1370	Dataset	RTT	Maze	Maze Large	Maze Medium	Maze Small
1371	NDT1	0.1643	0.3597	0.3739	0.3081	0.2788
1372	SoTA	0.2010	0.3650	0.3831	0.3329	0.3458
1373	NDT3 Scratch	0.1533	0.3093	0.2892	0.1859	0.1853
1374	NDT3 350M 2kh	0.1695	0.2775	0.2781	0.1937	0.2116
1075						

1375 Table 1. NDT1 vs NDT3 co-bps (higher is better) on NLB's MC Maze and MC RTT datasets in the 20ms split. 1376 The NDT3 results were produced were the standard sweeps used in this work, e.g. three random learning rates 1377 and three random seeds. SoTA results come from LangevinFlow for RTT and Maze and MINT for the Maze scaling datasets. 1378

1379 1380

1363 1364

С METHODS

1381 1382

C.1 METRICS AND EVALUATION 1383

1384 Throughout this work we evaluate offline decoding of continuous covariates timeseries. The metric 1385 we specifically use is the coefficient of determination, R^2 , as computed by scikit-learn's $r2_score$ 1386 function. R^2 is a useful metric over MSE as 1 represents perfect prediction and 0 is the score achieved by best-guess baseline, the mean of the data. In pretraining, R^2 is computed over the flat average of all 1387 covariate dimensions, since each datapoint has differing covariate dimensionalities. In evaluation, R^2 1388 is computed as a variance-weighted average of R^2 s in each covariate dimension. Another difference 1389 between training and evaluation metrics is that training predictions are made over batched data, while 1390 evaluation predictions are mostly computed in a *streaming* fashion. Streaming requires continuous 1391 neural data across different behavioral epochs, and so cannot be performed for the Oculomotor and 1392 CST datasets. We also omit it for the motor cortex self-paced reach dataset, which has a very large 1393 evaluation split. Streaming allows timesteps at the beginning of each sequence to leverage neural 1394 context from the preceding sequence, which raises performance slightly, as shown in the continuous 1395 vs trialized analysis (Section 3.3). We limit history in streaming evaluations to the max history seen 1396 in tuning (1 second).

1398 C.2 TRAINING 1399

1400 Pretraining hyperparameters were manually tuned in preliminary experiments at the 45M parameter 1401 models on small datasets. 350M models diverged at the chosen 4e - 4 peak LR, so we lowered peak LR to 1e - 4. For tuning, the explored LRs are 1e - 4, 3e - 4, 5e - 4 for training from scratch 1402 and 3e - 5, 1e - 4, 4e - 4 for fine-tuning. While this is far from an exhaustive search, we show 1403 in Fig. 16 that other regularization hyperparameters are set to reasonable defaults such that this

1404 sweep finds near optimal results for both a from scratch model and fine-tuning the 45M model. 1405 Fine-tuning, like pretraining, is early stopped with a patience of 100 epochs. Batch size is uniformly 1406 set to 16K in pretraining, and scaled to be roughly 10-20% of dataset size in fine-tuning. NDT3 1407 from-scratch models were trained at the 11M parameter range. Exact model configurations for 1408 different experiments are documented in the codebase.

1409 NDT3's simple architectural design allows us to train on batches from different tasks and dimen-1410 sionalities. To avoid excess padding in training, we concatenate pretraining data that is otherwise 1411 discontinuous (trialized) into 2 second data. We do not add any separator tokens, as this does not 1412 appear to have a performance impact for language models (Geiping and Goldstein, 2022). With 1413 mixed-precision training, the 350M parameter NDT3 can fit the 4-8K tokens in each input context in 1414 the memory of 40G NVIDIA A100 GPUs. Thus we can restrict NDT3's pretraining parallelism to data-parallelism. 1415

1416 Using Kaplan et al. (2020)'s equation for FLOP computation, $C_{\text{forward}} = 2N + 2n_{\text{laver}}n_{\text{ctx}}d_{\text{attn}}$, we 1417 compute the footprint of the 350M 2kh model. We use about 0.9B FLOPs per token in the forward 1418 pass, and about 0.9T neural tokens processed over training, which yields a pretraining footprint of 1419 about 2.4e21 FLOPs.

1420

1421

1422 C.3 BASELINES

1423

1424 **Wiener Filter** The Wiener Filter baseline was cross-validated over regularization strength. We also 1425 swept history of neural input up to the max length provided to NDT, and reported the R^2 of the 1426 best WF according to test data in primary evaluation (slightly advantaging the WF). Generalization 1427 plots in Section 3.3 report the performance of WF models at these different histories. For evaluating angular generalization, WFs were only swept up to 1s history due to memory limits; performance 1428 was not varying substantially with history so we do not expect this to have impacted conclusions. The 1429 WF was for simplicity directly fit on the concatenated trial data, which may have slightly negatively 1430 impacted its performance in trialized datasets (Oculomotor, CST, Generalization analyses). 1431

1432 In the primary evaluations in Section 3.1, we considered WFs fit either independently per session in a dataset or jointly on all sessions, which is helpful for sessions in very low data regimes. We report 1433 the better of the 2. In generalization analyses, for simplicity, we only report joint fits, which may 1434 cause a slight downward bias in performance. 1435

1436			- 0	0
1437	Dataset	Patience	Held-In R ²	Held-Out R^2
1438	H1	100	$0.567_{\pm 0.034}$	$0.453_{\pm 0.030}$
1439	H1 (reproduction)	250	$0.628_{\pm 0.011}$	$0.517_{\pm 0.016}$
1440	H1 ((Karpowicz et al., 2024))	250	0.62	0.52
1441	M2	100	$0.563_{\pm 0.015}$	$0.352_{\pm 0.028}$
1442	M2 (reproduction)	250 250	$0.582_{\pm 0.002}$	$0.391_{\pm 0.009}$
1443	M12 ((Karpowicz et al., 2024))	250	0.63	0.43

Table 2. NDT2 H1 and M2 results when trained with 100 epochs of patience (this work) in fine-tuning vs 250 as 1444 in Karpowicz et al. (2024). We report mean and standard deviation of 3 model seeds on the FALCON evaluation 1445 (which is in turn a cross-session mean). 1446

1447 **NDT2** NDT2 baselines were prepared with its public codebase. We trained NDT2 models both 1448 from-scratch and from the public checkpoint pretrained on 100 hours of human data. Max context 1449 length and patience were held constant across the models. This restriction to a patience of 100 epochs accounts for some difference with the reported FALCON benchmark results in Karpowicz et al. 1450 (2024), as we note in Table 2. Other choices were left to NDT2 defaults. For example, NDT2 uses 1451 z-score normalization, which we kept. In from-scratch training, for simplicity, we jointly trained 1452 NDT2 with its neural reconstruction loss (masking of 25%) and a supervised decoding loss. This is 1453 true for all eight evaluation tasks except CST. In the CST task, we used only the supervised decoding 1454 loss, as the token dropout used in reconstruction can dropout all neural input. 1455

For hyperparameter tuning, we matched NDT3's tuning budget for the pretrained NDT2 checkpoint 1456 by only exploring 3 learning rates. Given mediocre NDT2 from-scratch performance, we swept 1457 NDT2 over 2 model sizes in addition to the standard 3 learning rates. We set the NDT2 from-scratch



Figure 17. NDT2 and NDT3 architectures differ mainly in the conversion of NDT2's representation learning backbone to NDT3's single multimodal stream. NDT3 directly intakes neural and behavioral tokens and predicts with a next-step objective. NDT2 employs explicit masking of input neural tokens and extracts neural and behavioral predictions at each timestep with cross-attention layers.

model sizes to 20M and 72M to be comparable with NDT3 45M, but note that the NDT2 pretrained
 checkpoint is only 6M parameters.

1478 NDT2 vs NDT3. NDT3 builds off of NDT2 but departs in several manners to enable more streamlined 1479 scaling and analysis over heterogeneous decoding tasks, which we overview in Fig. 17. Lower level 1480 technical changes are described in Section C.6. Both models relate neural data to behavior and train 1481 with both neural data and behavior prediction objectives. Both models use both these objectives in 1482 fine-tuning, but NDT2 only uses the neural objective in pretraining. NDT2's neural data objective is 1483 based on MAE (He et al., 2021), such that some fraction of input neural data tokens are masked and 1484 reconstructed in a decoder separate from the main backbone. However, this explicit masking was originally developed to study representation learning on images, not timeseries decoding. Causal 1485 domains like BCI control can also learn nontrivial representations simply through next-step prediction. 1486 NDT2 employed both explicit MAE masking and a causal attention mask in its backbone, which 1487 is redundant computationally and also reduces the context available to make predictions. NDT3 1488 thus dispenses with the masking mechanism and uses next-step prediction alone. NDT2 also differs 1489 from NDT3 in its readout of behavior. Again, since NDT2 studied representation learning, it used 1490 additional cross-attention readout layers to "probe" behavior predictions at each timestep. NDT3 1491 simplifies this two-part encoder-decoder design to a decoder-only design. In this flow, masked 1492 behavior tokens are provided at the input and filled in through the backbone, and varying the input 1493 tokens allows us to make predictions of the corresponding behavior dimensionality.

1494

1495 1496 C.4 Pretraining and Evaluation datasets

Pretraining datasets were comprised of historical data from several labs, the rough composition of which is shown in Fig. 2B. The evaluation behavior used during pretraining was reaching in 2 monkeys. The first monkey dataset came from a public release (Flint et al., 2012), and the second from a private dataset (REDACT lab). The latter had center-out reach in standard conditions and under visual feedback perturbations. The monkey in the second dataset is also present in the 1khr monkey and 2kh and up model dataset sizes, though performing in a different set of experiments.

1503 Inherent to the process of large-scale scraping is a loss of detail on what precise tasks were used, so 1504 we only have a qualitative description of tasks we believe are well represented. NDT3 trains on a wide 1505 variety of reaching behaviors from relatively constrained (2D center-out reaching to fixed number of 1506 targets) to relatively unconstrained (self-paced, more targets, potentially 3D) and under experimental 1507 manipulations (delayed onset, multiple targets, different error thresholds requiring more precision). These reaching behaviors are described in both endpoint kinematics and as EMG. A smaller fraction of pretraining data are isometric and force related (force exerted against manipulandums) for wrist 1509 and arm motion. Human datasets contain a variety of iBCI tasks, with closed loop datasets reflecting 1510 both high and low quality control. These tasks include reach and grasp behavior from 1-10 degrees of 1511 freedom, as well as some individuated finger tasks for clicking.

We detail evaluation datasets in Table 3. Three datasets come from the FALCON benchmark (Karpowicz et al., 2024), two are based on public datasets ((O'Doherty et al., 2017; Deo et al., 2024)), and three are private. Note we avoid the Neural Latents Benchmark (Pei et al., 2021) as it does not directly measure decoding performance. For each evaluation dataset, we specify a tuning split and an evaluation split. Only tuning split data is changed when varying data scale. Tuning and evaluation splits are block-contiguous, i.e. trials are not interleaved, for better downstream applicability.

1518 1519

C.5 GENERALIZATION ANALYSES AND FURTHER EVALUATIONS

Intra-session generalization Posture, spring, and angular generalization evaluate OOD performance in the standard setup of comparing in-distribution and out-of-distribution performance directly (with changes in the underlying evaluation dataset) The intra-session temporal shift analysis is evaluated in an inverted, slightly more rigorous setting. Specifically, we trained two sets of models on the two different temporal blocks, and evaluated on an evaluation split in the later block, rather than only training on the early block and evaluating on both blocks. This way, the OOD shift is measured with respect to the same evaluation dataset.

1528

1535 1536

1537

1538

1546

1547

1548

1529 C.6 ARCHITECTURAL DETAILS

NDT3 adopts several architectural innovations used in recent Transformer models. These were compared against baselines in preliminary experiments, but formal ablations in the final experimental setting were not conducted. We defer full description of the Transformer dimensions to the public codebase.

- FlashAttention 2 (Dao, 2023) is used to increase training and inference speeds. On the NERSC Perlmutter cluster, with FA2, 45M NDT3 trained at about 270M neural tokens per 40G A100 hour, 350M NDT3 trained at about 70M neural tokens per A100 hour. Note, FA2 also enables use of the 350M model for real-time (<20ms) inference latency.
- Positional Embeddings (Su et al., 2023): Rotary embeddings are applied to indicate the real-world timestep of every input token. Additionally, 48 categorical learned embeddings are reserved to distinguish token modality and position within a timestep (10 for neural, 16 for covariates, 16 for covariate constraints, 1 for reward/return, 1 for dummy tokens, remainder unused).
 - QK Normalization (Dehghani et al., 2023; Wortsman et al., 2024): An additional layer norm is applied to the query and key embeddings, before the rotary embeddings, which helped stabilize training of the 350M parameter models.
- No context embeddings (Ye et al., 2023): Differing from NDT2, no learned embeddings for disambiguating input datasets were prepended to each input. This was removed for simplicity. Per GATO (Reed et al., 2022) and language modeling practices, we instead leave task / dataset disambiguation to the modeling process: In pretraining, the covariate maskout strategy allows for many tasks to be specified in-context (as later behavior can be inferred on the basis of earlier neural-behavioral token relationships). In fine-tuning, the tuning dataset already uniquely specifies the function to be learned.
- Cross entropy loss for spiking data prediction: We used the standard cross entropy loss to classify spike count over the Poisson loss common in many neural data architectures. Since the overall ablation of neural objective shows no large impact in this work, it is likely that this decision should be evaluated with neural data related tasks rather than decoding.
- 1560

We document the Transformer model shapes considered in our work in Table 4. This shape is not systematically explored in our work, and is by historical artifact, slightly different than the shapes used in NLP/CV. Embedding parameters are negligible. One possible area of interest is that the feedforward expansion factor is 1 in our model, i.e. the MLP dimension is low. If MLPs do serve as memory stores in Transformers (Geva et al., 2021), increasing this shape may yield more performant model size scaling, given the heterogeneity of our datasets.

1566 D NDT3 MODEL CARD 1567

1568	The card is currently only provided in the codebase
1569	The card is currently only provided in the codebase.
1570	
1571	
1572	
1573	
1574	
1575	
1576	
1577	
1578	
1579	
1580	
1581	
1582	
1583	
1584	
1585	
1586	
1587	
1588	
1589	
1590	
1591	
1592	
1593	
1594	
1595	
1596	
1597	
1598	
1599	
1600	
1601	
1602	
1603	
1604	
1605	
1606	
1607	
1608	
1609	
1610	
1611	
1612	
1613	
1614	
1615	
1616	
1617	
1618	
1619	

Table 3. Evaluation datasets used for multiscale decoding and generalization analyses. The references provide extended description of the behavioral task. Dashed line separates datasets for Section 3.1 and for analysis. Datasets use unsorted multi-unit activity and are processed in 1s chops unless otherwise mentioned.

1624	Dataset	Description
1625	FALCON H1, M1,	3 separate single-subject multi-session datasets for different iBCI tasks.
1626	M2 (Karpowicz et al.,	Data comes in a high data split (held-in), and a low-data split (held-
1627	2024)	out), with the intention on identifying methods that can achieve parity in
1628		the two settings. H1 is an open loop human dataset for calibrating 7D
1629		reach-and-grasp in a robot arm. M1 is a monkey reach-and-grasp task
1630		to different objects with EMG recordings. M2 is a monkey 2D finger
1631		movement task with manipulandum-measured kinematics. Scaling scores
1632	Salf paced reach	are reported on the test set. Monkeys reach for random targets one at a time in a small planar.
1633	(RTT) (O'Doherty	workspace. We decode 2D arm velocity in monkey Indy. Has neu-
1634	$(\mathbf{R}\mathbf{I}\mathbf{I})(\mathbf{O}\mathbf{D}\mathbf{O}\mathbf{I}\mathbf{C}\mathbf{I}\mathbf{Y})$ et al. 2017)	ral data from M1 and S1 we use M1 in Section 3.1 and Section 3.2 and
1635	Di 10	S1 in Section 3.3.
1636	Bimanual Cursor	A human open loop dataset where the participant attempts movement of
1637	Control (Deo et al.,	one or both hands to control two cursors.
1638	2024) 2D Cursor + Click	Cursor control is a classic iBCI endpoint (Pandarinath et al. 2017; Wol
1639	(private)	naw et al. 2002: Jarosiewicz et al. 2015) Two human participants
1640	(ritine)	attempt movement according to visually cued cursor movement and au-
1041		diovisual click cues. We also use this dataset for trial structure analysis
1042		in Section 3.3.
1643	Grasp force (private)	A open-loop dataset with two human participants attempting isometric
1644		power grasps. Specifically, participants were asked to match force output
1646		according to visual cues in a Mujoco environment. Grasps cued were both
1647		static (instant onset, noid, and onset) or dynamic (gradually increasing force). This detect is valuable for human iBCI study because force
1648		modulation is required in many motor behaviors and grash force has
1649		primarily only been characterized in monkeys until now (Branco et al.,
1650		2019). Uses 2 second intervals due to long behavior timescale. We expect
1651		this dataset can be released by end of 2024.
1652	Critical Stability	A monkey dataset collected to study continuous control relative to ballistic
1653	Task (Quick et al.,	movement. The monkey balances a virtual cursor on a 1D workspace for
1654	2018) (private,	up to 6 seconds.
1655	Trialized, sorted)	A monkey contor out task but the monkey's hand is adjusted to one of 6
1656	Center-Out (Marino	different starting positions. We use the central position as center and the
1657	et al 2024) (private	rest as edge
1658	trialized, sorted)	
1659	Spring-load (Mender	A monkey moves fingers, clamped together in a manipulandum for effec-
1660	et al., 2023)	tive 1DoF, is neutral or under spring load.
1661	Center-out, Monkey	Used in Section 3.2. A monkey performs an isometric center out task.
1662	J (Ma et al., 2022)	Forces are measured by the manipulandum and converted to cursor veloc-
1663	(Irlalized) Center out Monkey	Ity signals.
1664	V (private trialized)	targets by moving a manipulandum (Kinarm)
1665	Oculomotor	A monkey visually tracks (via smooth pursuit) a target that moves from
1666	pursuit (Noneman	center of workspace to one of four directions. A few dozen neurons
1667	and Patrick Mayo,	are recorded on probes in each of frontal eye field (FEF) and area MT.
1668	2024) (private,	We decode pupil velocity. The small number of neurons in this dataset
1669	trialized, sorted)	required resetting NDT3 neural readin/readout layers.
1670	FALCON H2	Human open loop dataset where a participant attempts movement to
1671		write letters cued on a screen (Willett et al., 2021; Fan et al., 2023). The
1672		readin/readout layers (to use fewer neural tokens)
1673		reaching reactout layers (to use rewer neural toxens).

Model	Layers	Width	MLP Size	Heads	Parameters (M)
NDT2 PT (Ye et al., 2023)	4	256	256	4	6
NDT3 Base	6	1024	1024	8	45
NID TO DI					
NDT3 Big	12	2048	2048	16	350
NDT3 Big Table 4	12 Transform	2048 er Model S	2048 hapes used in the	16 his work.	350
NDT3 Big Table 4	12 Transform	2048 er Model S	2048 hapes used in th	16 his work.	350
NDT3 Big Table 4	12 Transform	2048 er Model S	2048 hapes used in th	16 his work.	350
NDT3 Big Table 4	12 Transform	2048 er Model S	2048 hapes used in t	16 his work.	350
NDT3 Big Table 4	12 • Transform	2048 er Model S	2048 hapes used in th	16 his work.	350
NDT3 Big Table 4	12 . Transform	2048 er Model S	2048 hapes used in t	16 his work.	350
NDT3 Big Table 4	12 . Transform	2048 er Model S	2048 hapes used in t	16 his work.	350
NDT3 Big Table 4	12 . Transform	2048 er Model S	2048 hapes used in t	16 his work.	350
NDT3 Big Table 4	12 . Transform	2048 er Model S	2048 hapes used in th	16 his work.	350
NDT3 Big Table 4	12 . Transform	2048 er Model S	2048 hapes used in th	16 his work.	350
NDT3 Big Table 4	12 . Transform	2048 er Model S	2048 hapes used in the	16 his work.	350
NDT3 Big Table 4	12 . Transform	2048 er Model S	2048 hapes used in the	16 his work.	350
NDT3 Big Table 4	12 . Transform	2048 er Model S	2048 hapes used in th	16 his work.	350
NDT3 Big Table 4	12 . Transform	2048 er Model S	2048 hapes used in t	16 his work.	350
NDT3 Big Table 4	12 . Transform	2048 er Model S	2048 hapes used in the	16 his work.	350
NDT3 Big Table 4	12 . Transform	2048 er Model S	2048 hapes used in the	16 his work.	350
NDT3 Big Table 4	12 . Transform	2048 er Model S	2048 hapes used in the	16 his work.	350
NDT3 Big Table 4	12 . Transform	2048 er Model S	2048 hapes used in the	16 his work.	350
NDT3 Big Table 4	12 . Transform	2048 er Model S	2048 hapes used in the	16 his work.	350
NDT3 Big Table 4	12 . Transform	2048 er Model S	2048 hapes used in the	16 his work.	350
NDT3 Big Table 4	12 . Transforme	2048 er Model S	2048 hapes used in the	16 his work.	350
NDT3 Big Table 4	12 . Transform	2048 er Model S	2048 hapes used in the	16 his work.	350
NDT3 Big Table 4	12 . Transform	2048 er Model S	2048 hapes used in the	16 his work.	350
NDT3 Big Table 4	12 . Transform	2048 er Model S	2048 hapes used in the	16 his work.	350