FORECASTING WHOLE-BRAIN NEURONAL ACTIVITY FROM VOLUMETRIC VIDEO

Anonymous authors

Paper under double-blind review

ABSTRACT

Large-scale neuronal activity recordings with fluorescent calcium indicators are increasingly common, yielding high-resolution 2D or 3D videos. Traditional analysis pipelines reduce this data to 1D traces by segmenting regions of interest, leading to inevitable information loss. Inspired by the success of deep learning on minimally processed data in other domains, we investigate the potential of forecasting neuronal activity directly from volumetric videos. To capture long-range dependencies in high-resolution volumetric whole-brain recordings, we design a model with large receptive fields, which allow it to integrate information from distant regions within the brain. We explore the effects of pre-training and perform extensive model selection, analyzing spatio-temporal trade-offs for generating accurate forecasts. Our model outperforms trace-based forecasting approaches on ZAPBench, a recently proposed benchmark on whole-brain activity prediction in zebrafish, demonstrating the advantages of preserving the spatial structure of neuronal activity.

024 025

026 027

004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

028 Recent advances in imaging techniques have enabled the recording of neuronal activity at unprecedented resolution and scale. Light-sheet imaging allows recording of whole-brain activity for small 029 animals, such as the larval zebrafish (Hillman et al., 2019). Raw recordings are in the form of volumetric videos, with hundreds of millions voxels per time step, recorded over hours. Typically, heavy 031 postprocessing is applied to reduce dimensionality of this data down to 1D time traces of activity for distinct regions of interest representing individual neurons or clusters of cells (Abbas & Masip, 033 2022). Inspired by the success of deep learning models in analyzing minimally processed data in 034 other fields, such as weather and climate forecasting (Rasp et al., 2020; Andrychowicz et al., 2023), we explore the potential of building predictive models directly on such volumetric videos, avoiding any information loss. 037

The ability to predict future behavior based on past observations is a cornerstone of scientific modeling across a diverse range of domains, ranging from physics to social sciences. Until now, it has not been applied in the context of whole-brain activity in a vertebrate. The recently introduced Zebrafish 040 Activity Prediction Benchmark (ZAPBench) (Anonymous, 2024) aims to change that, taking advan-041 tage of datasets that can now be acquired with modern microscopy techniques. ZAPBench provides 042 a rigorous evaluation enabled by the comparison of future brain activity predicted from past brain 043 activity to actual experimental recordings, thereby achieving an objective measure for evaluating 044 predictive models of brain function. The dataset used in ZAPBench is a whole-brain recording from a larval zebrafish, collected using a light-sheet microscopy setup (Vladimirov et al., 2014). The raw volume is made of ~ 1.5 trillion voxels, which is reduced in size by three orders of magnitude to a 046 trace matrix of time series by applying a neuron segmentation mask. ZAPBench is the first bench-047 mark that poses the forecasting problem for a significant fraction of neurons in a single brain and 048 provides the raw volumetric recordings for the experiments. 049

To test the viability of end-to-end forecasting on such data, we propose to use a video model based
on techniques that have not been applied to this domain previously. Since processing in the brain
is highly distributed (Urai et al., 2022; Naumann et al., 2016), we hypothesize that large receptive
fields are important. Furthermore, in comparison to models applied to activity traces, we expect the
following advantages. First, by utilizing the entire video as input, a video-based model is not reliant



Figure 1: We propose to model light-sheet microscopy recordings of neural activity directly as volumetric video for forecasting instead of extracting and modelling neuron traces. Specifically, we train a model directly on the video and mask the output to optimize the per-neuron mean absolute error (MAE). We find that a UNet performs particularly well for small temporal context and can more effectively utilize spatial contextual information than trace-based time series models.

on the precision of neuron segmentation masks (masks are still applied to predicted frames for direct comparison to trace-based methods). Second, the inherent grid structure of the video preserves the spatial relationships between neurons, information that is otherwise lost during trace extraction. Finally, such a model can leverage potentially relevant visual cues present in the voxels between segmented cells or within the voxels of individual cell masks to enhance forecasting accuracy.

Building a model for this problem poses fundamental engineering challenges. Standard video models operate on 2D frames, and the presence of the additional Z dimension naturally complicates scaling. In ZAPBench, a single XY slice of a 3D frame has a native size of 2048×1328 pixels, and is thus comparable to a frame of a natural 1080p video. Every 3D frame is composed of 72 such slices, increasing the volume of the input data by up to two orders of magnitude relative to such videos and resulting in several hundreds of megabytes per frame.

For our model, we choose a variant of the UNet (Ronneberger et al., 2015) and adapt it to 4D data. We also develop a data input pipeline where both input and model are spatially sharded across multiple hosts and accelerators. To maintain a manageable size of the intermediate activations, we represent temporal input frames as channels. This approach allows us to explore the impact of varying spatial context by manipulating the receptive field while keeping computational cost (FLOPS) roughly constant.

We perform extensive experiments to construct an effective video model for neuronal activity forecasting on ZAPBench. Despite the success of pre-training in other domains (Devlin, 2018; Bao et al., 2021), we find it not to be a useful technique for improving forecast accuracy, even when using an order of magnitude more data recorded from other specimens of the same species. Further, we investigate the effect of input resolution, spatial context, and temporal context of the model on the forecast accuracy. Surprisingly, we find that lowering input resolution by up to 4x can be beneficial for performance and observe a clear trade-off between spatial and temporal context.

Our models, which implicitly capture the spatial relationships within their field of view, can improve forecast accuracy beyond that achieved by trace-based models on ZAPBench, especially when only short temporal context is available. On ZAPBench, multivariate trace-based models, which can in principle learn functional relationships between cells, do not perform significantly better than univariate models that treat all cells independently and identically. Our proposed model is therefore the only multivariate model that can consistently outperform univariate models on this benchmark.¹

- 103 In summary, our contributions are as follows:
- 104 105

066

067

068

069

- 1. We propose to forecast zebrafish neuronal activity recorded using light-sheet microscopy directly in the native domain as volumetric video (3D + time).
- 106 107

¹Comparisons between video and trace-based models are also included in ZAPBench (Anonymous, 2024).

108

111 112

113

114

115 116

118

148 149 150 2. We empirically show that the input resolution and pre-training on additional volumetric videos from similar specimens have negligible impact on the results.

- 3. We perform exhaustive model selection to quantify the impact of spatial (XYZ) and temporal context size for activity forecasting accuracy and find a clear trade-off.
- 4. On ZAPBench, our proposed model is the only approach that consistently benefits from multivariate information, and therefore achieves leading performance for short temporal context.

117 2 FORECASTING NEURONAL ACTIVITY FROM VIDEO

119 We propose to forecast neuronal activity in the ZAPBench dataset (Anonymous, 2024) directly in 120 the volumetric video domain. Specifically, we utilize a temporal context of C video frames to 121 predict the subsequent H frames. Per-neuron forecasts and loss are then computed by applying the 122 segmentation mask to the predicted video frames. This contrasts with the traditional approach, which 123 applies the segmentation mask to the original video data to extract activity traces before performing 124 any forecasting. See Figure 1 for a comparison of these two approaches.

125 The ZAPBench dataset comprises high-resolution, whole-brain activity recordings of a larval zebrafish engaged in various behavioral tasks. Data was acquired using light-sheet fluorescence mi-126 croscopy, enabling real-time imaging of neuronal activity at cellular resolution. This was made 127 possible by using an animal genetically modified to express GCaMP (Dana et al., 2019), a flu-128 orescent calcium sensor, in the nuclei of its neurons. ZAPBench provides both preprocessed ac-129 tivity traces for approximately 70,000 neurons and the corresponding raw volumetric video data. 130 This raw data, denoted as Y, has dimensions of $2048 \times 1152 \times 72 \times 7879$ (XYZT) and a resolution of 131 $406 \text{ nm} \times 406 \text{ nm} \times 4 \mu \text{m} \times 914 \text{ ms}$. We use a center crop of 1328 voxels in Y due to negligible cell 132 activity in the border regions. Models forecasting H = 32 time steps are benchmarked using short 133 (C = 4) or long (C = 256) temporal context. 134

Anonymous (2024) preprocess the raw volumetric video by aligning each frame to a reference volume for stabilization so that the neuron segmentation masks can be statically applied throughout the experiment. Further, a standard " $\Delta F/F$ " normalization scheme is applied to the voxel intensities, with F denoting a baseline value (Mu et al., 2019; Zhang et al., 2023). The normalized signal is in the [-0.25, 1.5] range.

139 The neuron segmentation model is specifically trained 140 for the dataset and yields 71,721 neurons. Formally, 141 the segmentation mask can be considered as a map-142 ping seg: $\mathcal{N} \to 2^{\mathcal{S}}$ from integer identity of a neuron $\mathcal{N} = [71721]$ to a set of three-dimensional spatial indices, 143 which is an element of the power set of index locations 144 $\mathcal{S} = [2048] \times [1328] \times [72]$. The neuron activity at an ar-145 bitrary timestep t is then given by *averaging* the activity 146 over spatial locations associated with each cell, i.e.: 147

$$\mathbf{y}_{n}(t) = \frac{1}{|\operatorname{seg}(n)|} \sum_{s \in \operatorname{seg}(n)} \mathbf{Y}_{s}(t).$$
(1)

While this is a natural choice, it loses information related to cell size, position and spatial distribution of intensities within it, and completely discards voxels that are not part of any segmentation mask or incorrectly segmented. Figure 2 depicts these potential issues.

We instead apply a video model to the raw input frames
and directly forecast volumetric frames while optimizing
and measuring the mean absolute error (MAE) on the seg-



Figure 2: Illustration of potential loss of information when segmenting neurons. The colored objects are predicted segmentation masks. A fragment of a 2d slice of the activity video is shown in greyscale.

mented neurons. Prior work in neural response prediction (Schoppe et al., 2016; Cadena et al., 2019)
has proposed additional metrics that explicitly take trial-to-trial variability into account. The experimental setting used in ZAPBench did not allow for sufficiently numerous trial repetitions to make these metrics applicable, but we note them as an interesting direction to explore in future work if



Figure 3: Architecture and input sharding overview. A: We use a variation of the UNet architecture (Ronneberger et al., 2015) with 3D spatial input and treat the C input frames as channels. Further, we use a fixed number of features at every resolution to improve scalability. The network is conditioned on the time horizon H and outputs a single volumetric frame at a time, similar to MetNet-3 (Andrychowicz et al., 2023). To control for spatial context at constant FLOPS, four blocks at the lowest resolution can be replaced by one block of higher resolution. B: Data loading and the network are spatially sharded and allow for flexible scaling to full resolution inputs.

calcium recordings made with increased number of trials become available. In addition to MAE, we also report correlations between the predicted and actual activity in App. A.4.

185 One frame of the volumetric video can be described as $\mathbf{Y}(t) \in \mathcal{V}$ with $\mathcal{V} = \mathbb{R}^{2024} \times \mathbb{R}^{1152} \times \mathbb{R}^{72}$. A video model with a P-dimensional weight vector $\mathbf{w} \in \mathbb{R}^{P}$ can then be denoted as $\mathbf{f} : \mathcal{V}^{C} \times \mathbb{R}^{P} \to \mathbb{R}^{P}$ 187 \mathcal{V}^{H} . That means the video model f receives a 4D volumetric input with C frames, outputs H frames, 188 and is parameterized with weights w. We obtain the prediction of the h-th frame as 189

$$\hat{\mathbf{Y}}(t,h) = \mathbf{f}_h(\mathbf{Y}(t),\dots,\mathbf{Y}(t+C),\mathbf{w}),\tag{2}$$

191 and denote by $\hat{\mathbf{y}}(t,h)$ the corresponding 1D trace vector computed using Eq. 1. For a fair compar-192 ison with trace-based models in ZAPBench we optimize the trace-based MAE \mathcal{L} over all training 193 timesteps T_{train} with respect to the model parameters w 194

$$\mathcal{L}(\mathbf{w}) = \frac{1}{|T_{\text{train}}|H|\mathcal{N}|} \sum_{t \in T_{\text{train}}} \sum_{h \in [H]} \sum_{n \in \mathcal{N}} |\mathbf{y}_n(t+h) - \hat{\mathbf{y}}_n(t,h)|.$$
(3)

If we instead optimize the voxel-wise MAE, the models perform relatively worse when evaluated on the trace-based MAE because it corresponds to a different weighting of neurons by their size in number of voxels.

SCALABLE VOLUMETRIC VIDEO ARCHITECTURE 3

Efficiently training models consuming high-resolution volumetric video of varying input context sizes C requires a scalable architecture and data loading system. We achieve this by extending a standard UNet architecture (Ronneberger et al., 2015) to 4D by mapping temporal input context to features of the first convolutional layer, conditioning on lead-time to predict only single frames, and sharding both the model and the data loading process. Figure 3 shows the intermediate resolutions and representation sizes. The network comprises a series of pre-activation residual convolutional blocks (He et al., 2016) with fixed feature size F = 128, each with two group normalization 210 layers (Wu & He, 2018) using 16 groups, Swish activation (Ramachandran et al., 2017), and 3^3 convolutions for XYZ throughout.

211 212 213

214

175

176

177

178

179

181 182 183

190

195 196 197

199

200 201

202 203

204

205

206

207

208

- 3.1 TEMPORAL INPUT CONTEXT AS FEATURES
- Typically, video UNet variations use color channels as input features (Gao et al., 2022; Ho et al., 215 2022a) and convolve over frames using a temporal convolution (Ho et al., 2022b). This approach

is intractable in our case because of the additional z dimension. Instead, we treat the temporal input context of C frames as input features to the UNet. This confers the following advantages:
1) the temporal sizes of the input and output are decoupled, 2) the network parameter count is easily controlled, 3) representation sizes and computation requirements are reduced while using more features, and 4) early layers of the network have access to long-range temporal dependencies. Our model is similar to architectures used in standard time series models, which often treat temporal context as features (Zeng et al., 2023; Chen et al., 2023).

223 224

225

3.2 VARYING THE RECEPTIVE FIELD

226 We design a flexible UNet architecture that can adapt the receptive field while keeping the computa-227 tional cost (FLOPS) fixed. We find that full native resolution is not necessary for optimal prediction 228 accuracy (see Sec. 4.1), and thus downsample the input by a factor of 4 in XY using averaging. The 229 first resampling block then uses a factor 2 in XY to achieve roughly isotropic resolution in XYZ, while the following ones downsample equally in all dimensions. We always use four residual blocks 230 at the lowest resolution, and three at all other resolutions. This allows us to change the receptive 231 field while keeping the FLOPS roughly fixed by removing the four lowest resolution blocks and 232 instead adding one block to the respective next higher resolution. This is because one block at the 233 higher resolution requires as many FLOPS as four blocks after downsampling by a factor of two 234 in x and y. In an ablation, we show that controlling for FLOPS is sensible because increasing the 235 parameter count does not increase performance further (see Figure 7). 236

The receptive field along a dimension depends on the cumulative product of the downsampling factors and the number of convolutions at the lowest resolution,

receptive_field_dim = cum_downsampling_factor_dim × num_blocks × 4,

241 where the factor 4 is because every block has two convolutions, each of which increase the receptive 242 field by two. For a network that does not downsample at all, as for example used in Sec. 4.1.1, 243 to account for the input and output convolutions and the center voxel, we have to increase the re-244 ceptive field size by five. Therefore, the architecture depicted in Figure 3 has a receptive field of 245 (1024, 1024, 128) in XYZ comparable to the size of the complete frame. We tried to further enhance 246 the receptive field to cover the whole frame using a multi-axis vision transformer (Tu et al., 2022) at 247 the lowest resolution, but did not observe any accuracy gains. For the output, we upsample twice to obtain the original resolution, and use one residual block per resolution, but with a reduced feature 248 dimension of F' = 32 to keep hidden representations at a manageable size. 249

250 251

252

239

240

3.3 LEAD-TIME CONDITIONING

253 Instead of forecasting autoregressively or predicting the 254 complete horizon of H frames in a one-shot way, we 255 condition the network on an integer lead-time $h \in [H]$ and predict the corresponding single frame independently 256 as proposed by Andrychowicz et al. (2023) for weather 257 forecasting. During data loading, we sample a lead time 258 and the corresponding target frame uniformly at random 259 from [H]. This requires loading only a single target 260 frame per sample and, in line with previous results from 261 weather forecasting (Rasp et al., 2020), performs better 262 than frame-level autoregressive prediction on our prob-263 lem. Every convolutional block in the network is condi-264 tioned using a FiLM layer (Perez et al., 2018) on the lead 265 time encoded using a 32-dimensional sinusoidal embed-



Figure 4: Comparison of direct MAE and lead-time conditioned variants.

ding (Vaswani et al., 2017). Figure 4 shows that directly predicting all *H* frames tends to overfit
while lead-time conditioning performs equally well with both MAE and HL-Gauss (Farebrother
et al., 2024), a distributional regression objective that results in slightly faster model convergence.
However, in our experiments we use the conditioned MAE for its simplicity and because it does not
require binning, which might complicate pre-training on datasets of slightly different scale.

270 3.4 SHARDED DATA LOADING AND MODEL271

272 Despite the scalability features of the proposed UNet model for volumetric video, in practice apply-273 ing it requires distributing the input and hidden representations across accelerators and machines. We train all models in Sec. 4 using a single sample per batch, noting that this can already correspond 274 to several GBs of input data. We use spatial sharding in XY using the jax. Array API (Bradbury 275 et al., 2018) so that each box in Figure 3B is handled by an individual accelerator. We also imple-276 ment a custom data loader that distributes data loading across hosts so that each machine only loads the necessary boxes. To achieve this, we chunk our data in the zarr3 format (Miles et al., 2023) 278 and use the TensorStore API (TensorStore developers, 2024) to load and collate chunks. Our data 279 loader follows the jax sharding automatically. 280

281 282

283

4 EXPERIMENTAL RESULTS

We present experimental results evaluating the proposed volumetric video model on ZAPBench, a 284 benchmark for whole-brain neuronal activity prediction for a larval zebrafish (Anonymous, 2024). 285 Uniquely, ZAPBench provides the raw volumetric recordings for most of the neurons in the brain 286 enabling data-driven approaches like ours. First, in Sec. 4.1 we empirically select and validate 287 the final architecture variant used for the benchmark. In particular, we investigate the trade-off 288 between temporal context C and spatial context in the form of the receptive field to assess the need 289 for multivariate models. Further, we evaluate the feasibility of pre-training on additional zebrafish 290 specimens as well as the effect of input resolution. We identify the model depicted in Figure 3 as 291 a strong model for the short context size C = 4, where we achieve the best performance across 292 the benchmark, as presented in Sec. 4.2. For the long temporal context C = 256, we only see an 293 improvement of forecast accuracy in specific cases.

294

295 Hyperparameters. Unless stated otherwise, we train every model for 250k to 500k steps by optimizing the trace-based MAE with a batch size of 1 using the AdamW optimizer (Loshchilov et al., 296 2017) using an initial learning rate of 10^{-4} decayed using a cosine schedule (Loshchilov & Hut-297 ter, 2017) to 10^{-7} and a weight decay factor of 10^{-5} . Due to their tendency to overfit, we use a 298 dropout rate of 0.1 on the features for long-context models with C = 256. These hyperparameters 299 were optimized on the validation set during development. We choose checkpoints based on the val-300 idation performance monitored during training. We present experimental results in terms of mean 301 performance and report two standard errors over at least three random seeds that control data load-302 ing and parameter initialization. The only exception to this are the high resolution results presented 303 in Sec. 4.1.3, where we only report a single result because of their compute requirements. Most 304 individual training experiments use 16 A100 40GB GPUs.

305 306 307

4.1 MODEL SELECTION

We compare between different methods and models to improve performance on the ZAPBench benchmark. For Sec. 4.1.1 and 4.1.2, we downsample the volumetric frames by a factor of four in XY using averaging to 512×288×72. The segmentation mask is downsampled to the same shape using striding. We investigate the effect of spatial and temporal context and the potential for pre-training on related datasets. In Sec. 4.1.3, we use the full resolution ZAPBench targets and segmentation, and assess the importance of input resolution on performance.

314 315

4.1.1 SPATIAL VS. TEMPORAL CONTEXT

316 We use UNets with different numbers of downsampling blocks to vary the spatial context but keep 317 the FLOPS fixed (see Sec. 3), and find that there is a trade-off between the spatial (S) and tempo-318 ral (C) context. We compare models without any downsampling blocks, with two downsampling 319 blocks, and with four. The models have spatial contexts S of 21, 64, and 256 in XY, respectively. 320 Details on the architecture and computation of the receptive field can be found in App. A.1.1. Also 321 note that the spatial context at full resolution of these models would be $4 \times$ higher. Figure 5 shows that a short temporal context requires larger spatial context to obtain optimal performance. For long 322 temporal context, however, the models with large spatial context start to overfit and underperform. 323 The effect becomes apparent between a temporal context of 16 and 64. This result suggests that



Figure 5: Validation and test performance for varying temporal context sizes C as well as spatial context sizes S with networks having comparable FLOPS. We find that there is a trade-off between spatial and temporal context with a cross-over point between C = 16 and C = 64, where spatial context stops being useful and leads to overfitting. The periodicity of many conditions is roughly 64, which might explain spatial context becoming redundant. We report the mean and two standard errors.

video models are able to exploit multivariate information for short temporal context but provide little benefit for long context, where univariate models perform equally well (Anonymous, 2024).

4.1.2 PRE-TRAINING ON OTHER SPECIMENS

345 We attempt to pre-train a model on other specimens 346 recorded and preprocessed in a similar way to the zebrafish used for ZAPBench. We pre-train the model ei-347 ther on two additional specimens recorded in the same 348 experimental session, or on these two and six more 349 from two other sessions. Because there is no seg-350 mentation available for the other specimens, the model 351 is pre-trained using voxel-based MAE for 800k steps, 352 and then fine-tuned on the ZAPBench dataset for 200k 353 steps using the trace-based MAE. We use three different 354 learning rates, 10^{-4} , 10^{-5} , and 10^{-7} , for fine-tuning, 355 and select the best model by validation performance,

SETTING	TEST MAE
Train	0.02573 ± 0.00005
Pre-train +2	0.02590 ± 0.00005
Pre-train +8	0.02591 ± 0.00001
Train + Val	0.02534 ± 0.00010

Table 1: Training on more data from the same specimen ("Train + Val") improves performance more than pre-training and finetuning on others. Results shown for C = 4 and data $4 \times$ downsampled in XY.

which was obtained by fine-tuning with the lowest learning rate. Table 1 shows that pre-training with fine-tuning does not improve performance over standard training. However, training on $\sim 14\%$ more data from the same specimen does improve performance significantly. Confidence intervals shown are calculated as two standard errors.

360 361

333

334

335

336

337

338 339 340

341

342

4.1.3 EFFECT OF INPUT RESOLUTION

We assess the relevance of input resolution when forecasting neuronal activity, and find that, surprisingly, predicting from a lower resolution performs best. We compare three variants of our model: a model that predicts from data $4 \times$ downsampled in XY, as depicted in Figure 3, one that downsamples only by factor 2, and one that is parameterized at full native resolution. Therefore, the full-resolution model loads and processes $16 \times$ more data. We achieve almost perfectly linear scaling by using proportionally more compute resources, maintaining the same throughput thanks to the sharded data input pipeline and model (see Sec. 3). In all cases, we scale the field of view of the network so that its size in physical units remains constant between experiments (see App. A.1.3).

370 For C = 4, Table 2 shows that the model with the low-371 est input resolution obtains a trace-based test MAE that 372 is statistically identical with that of the model using in-373 termediate resolution inputs. However, the full resolu-374 tion model performs significantly worse. This suggests 375 that despite the short temporal context input resolution does not play a major role in improving performance, 376 and that the intracellular voxel-to-voxel variations in 377 the recorded images do not carry information useful for

Input	TEST MAE
Downsample $4 \times$	0.0267 ± 0.0002
Downsample $2 \times$	0.0268
Full resolution	0.0273

Table 2: Increasing input resolution does not improve performance, and decreases it slightly at full resolution. forecasting, which might have applications to the de-

sign of future zebrafish activity recording experiments. We suspect that the decreased performance
 of the full resolution model could be caused by the significantly increased input voxel-to-parameter
 ratio while keeping the number of training examples fixed.

382

384 385

386

387

388

389

390

391

392

393

394

407

408

409

410

4.2 PERFORMANCE ON ZAPBENCH

We evaluate the best-performing architectures on ZAPBench for both short and long context settings. In Figure 6, we report the trace-based MAE versus forecasting steps in comparison to the best-performing trace-based models (Anonymous, 2024). We average performance across test sets of different stimulus conditions. For trace-based models, TSMixer (Chen et al., 2023) achieves best performance for short context, C = 4, and a univariate MLP for long context, C = 256. We use the video-based architecture depicted in Figure 3 for the short context that has a spatial context of $1024 \times 1024 \times 72$ in XYZ, which is global except in X where it covers half the voxels. For the long temporal context, we use a model that does not downsample further than (4, 4, 1) at the input, which we found in Sec. 4.1.1 to work best for this case. This model has a spatial context of $64 \times 64 \times 21$, which corresponds to $26 \,\mu\text{m} \times 26 \,\mu\text{m} \times 84 \,\mu\text{m}$ in XYZ.



Figure 6: Comparison of volumetric video model with best-performing trace-based model for short (left) and long (right) context on the benchmark test set (averaged over eight conditions) and the experimental condition held out from the training data. We report the mean and two standard errors.

411 We find that the volumetric video models achieve the best performance in the short context C = 4412 setting. For C = 256, there is no significant difference between the univariate trace-based model 413 and the video model on the test set when evaluated with MAE, but the video model does improve 414 correlation metrics (App. A.4). This aligns with our observation in Sec. 4.1.1, where longer tem-415 poral context requires less spatial context for the same forecast accuracy. ZAPBench also holds out 416 one stimulus condition entirely from training. We find that video models generalize better on this holdout condition for one-step-ahead forecasts but not for longer horizons. In App. A.3, we further 417 show model performance separately for each experimental condition. For the short context, we find 418 that the video model performs better in six, equally well in one, and worse in two out of the nine 419 conditions. 420

421 More precisely, when evaluated with a context length C = 4 on both the test and holdout sets, 422 the video model demonstrates a significant improvement in one-step-ahead forecasting accuracy, 423 achieving about 10 percentage point reduction in error compared to the best performing trace-based 424 model. With C = 256, the video model exhibits marginally superior performance in the first few 425 forecasting steps, achieving up to a 2 percentage point reduction in MAE at the first step. Beyond 426 the initial steps, both models demonstrate comparable accuracy on the test set.

What explains the improved performance of the video model relative to the trace-based approaches?
In App. A.2 we report results of an experiment in which we masked out all unsegmented voxels,
which did not reduce the grand test-MAE. This suggests that segmentation quality is not a significant
limitation, that no significant information is contained in the unsegmented regions of the dataset,
and that he accuracy gains can be attributed to the better utilization of the spatial distribution of
the recorded fluorescence signal. Furthermore, the results in Table 2 and Figure 5 suggest that it

432 is specifically the correlations between cells in the recorded fluorescence signals, rather than the 433 distribution of signals within individual cells, that drives these improvements.

434 435

5 CONCLUSION

436 437

We presented an approach for forecasting of zebrafish neuronal activity based on utilizing the raw 438 neural recording data as a volumetric video to make predictions. We find that this method has several 439 advantages over traditional trace-based methods. In particular, video-based prediction leverages 440 spatial relationships between neurons that are hard to exploit when reducing the data to 1D traces. 441 This allows for more accurate predictions, especially when working with short temporal contexts. 442 This advance comes at the expense of a significant increase in computational cost (2-3 orders of 443 magnitude relative to trace-based models, see App. A.5 for details).

444 We report several findings that were contrary to our expectations. First, we surprisingly find that 445 using higher resolution input frames does not improve performance. Second, the commonly used 446 paradigm of pre-training on a larger data set and fine-tuning only leads to reduced forecast accuracy 447 in our experiments. In contrast, we observe that more data from the same specimen does improve 448 performance, so we hypothesize that pre-training may be complicated by distribution shifts between 449 specimens, such as differences in signal and noise levels. Finally, increasing model capacity does 450 not always translate to performance improvements but instead leads to overfitting for long temporal 451 context.

452 Future work could explore the use of probabilistic models, latent space representations, and more 453 sophisticated regularization methods and input augmentations to further improve the accuracy of 454 video forecasting for neuronal activity.

455 456

REPRODUCIBILIY STATEMENT

457 458 459

460 461

463

464

465

469

470

471

472

473

All relevant code and interactive visualizations of the predictions will be made publicly available following double blind review.

462 REFERENCES

- Waseem Abbas and David Masip. Computational methods for neuron segmentation in two-photon calcium imaging data: a survey. Applied Sciences, 12(14):6876, 2022.
- 466 Marcin Andrychowicz, Lasse Espeholt, Di Li, Samier Merchant, Alexander Merose, Fred Zyda, Shreya Agrawal, and Nal Kalchbrenner. Deep learning for day forecasts from sparse observations. 467 arXiv preprint arXiv:2306.06079, 2023. 468
 - Anonymous. Zapbench: A benchmark for whole-brain activity prediction in zebrafish. ICLR 2025 Submission, 2024.
 - Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. Beit: Bert pre-training of image transformers. arXiv preprint arXiv:2106.08254, 2021.
- 474 James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal 475 Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao 476 Zhang. JAX: composable transformations of Python+NumPy programs, 2018. URL http: 477 //github.com/jax-ml/jax. 478
- Santiago A Cadena, George H Denfield, Edgar Y Walker, Leon A Gatys, Andreas S Tolias, Matthias 479 Bethge, and Alexander S Ecker. Deep convolutional models improve predictions of macaque v1 480 responses to natural images. *PLoS computational biology*, 15(4):e1006897, 2019. 481
- 482 Si-An Chen, Chun-Liang Li, Nate Yoder, Sercan O Arik, and Tomas Pfister. TSMixer: An All-MLP 483 architecture for time series forecasting. Transactions on Machine Learning Research, 2023. 484
- Hod Dana, Yi Sun, Boaz Mohar, Brad K Hulse, Aaron M Kerlin, Jeremy P Hasseman, Getahun 485 Tsegaye, Arthur Tsang, Allan Wong, Ronak Patel, et al. High-performance calcium sensors for

517

imaging activity in neuronal populations and microcompartments. *Nature methods*, 16(7):649–657, 2019.

- Jacob Devlin. Bert: Pre-training of deep bidirectional transformers for language understanding.
 arXiv preprint arXiv:1810.04805, 2018.
- Jesse Farebrother, Jordi Orbay, Quan Vuong, Adrien Ali Taiga, Yevgen Chebotar, Ted Xiao, Alex
 Irpan, Sergey Levine, Pablo Samuel Castro, Aleksandra Faust, et al. Stop regressing: Training
 value functions via classification for scalable deep rl. In *Forty-first International Conference on Machine Learning*, 2024.
- Zhangyang Gao, Cheng Tan, Lirong Wu, and Stan Z Li. Simvp: Simpler yet better video prediction.
 In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 3170–3180, 2022.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV 14*, pp. 630–645. Springer, 2016.
- Elizabeth MC Hillman, Venkatakaushik Voleti, Wenze Li, and Hang Yu. Light-sheet microscopy in
 neuroscience. *Annual review of neuroscience*, 42(1):295–313, 2019.
- Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P
 Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High definition
 video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022a.
- Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J
 Fleet. Video diffusion models. *Advances in Neural Information Processing Systems*, 35:8633–
 8646, 2022b.
- 512
 513 Ilya Loshchilov and Frank Hutter. SGDR: Stochastic gradient descent with warm restarts. In *International Conference on Learning Representations*, 2017.
- Ilya Loshchilov, Frank Hutter, et al. Fixing weight decay regularization in adam. *arXiv preprint arXiv:1711.05101*, 5, 2017.
- Alistair Miles, Jonathan Striebel, and Jeremy Maitin-Shepard. Zarr v3. https://zarr.
 readthedocs.io/en/stable/spec/v3.html, 2023.
- Yu Mu, Davis V Bennett, Mikail Rubinov, Sujatha Narayan, Chao-Tsung Yang, Masashi Tanimoto, Brett D Mensh, Loren L Looger, and Misha B Ahrens. Glia accumulate evidence that actions are futile and suppress unsuccessful behavior. *Cell*, 178(1):27–43, 2019.
- Eva A Naumann, James E Fitzgerald, Timothy W Dunn, Jason Rihel, Haim Sompolinsky, and Flo rian Engert. From whole-brain data to functional circuit models: the zebrafish optomotor re sponse. *Cell*, 167(4):947–960, 2016.
- Ethan Perez, Florian Strub, Harm De Vries, Vincent Dumoulin, and Aaron Courville. Film: Visual reasoning with a general conditioning layer. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
- Prajit Ramachandran, Barret Zoph, and Quoc V Le. Searching for activation functions. arXiv preprint arXiv:1710.05941, 2017.
- Stephan Rasp, Peter D Dueben, Sebastian Scher, Jonathan A Weyn, Soukayna Mouatadid, and Nils
 Thuerey. Weatherbench: a benchmark data set for data-driven weather forecasting. *Journal of Advances in Modeling Earth Systems*, 12(11):e2020MS002203, 2020.
- Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomed ical image segmentation. In *Medical image computing and computer-assisted intervention– MICCAI 2015: 18th international conference, Munich, Germany, October 5-9, 2015, proceed- ings, part III 18*, pp. 234–241. Springer, 2015.

- Oliver Schoppe, Nicol S Harper, Ben DB Willmore, Andrew J King, and Jan WH Schnupp. Measuring the performance of neural models. *Frontiers in computational neuroscience*, 10:10, 2016.
- TensorStore developers. Tensorstore: Library for reading and writing large multi-dimensional ar rays., 2024. URL https://github.com/google/tensorstore.
- Zhengzhong Tu, Hossein Talebi, Han Zhang, Feng Yang, Peyman Milanfar, Alan Bovik, and Yinxiao
 Li. Maxvit: Multi-axis vision transformer. In *European conference on computer vision*, pp. 459–479. Springer, 2022.
- Anne E. Urai, Brent Doiron, Andrew M. Leifer, and Anne K. Churchland. Large-scale neural recordings call for new insights to link brain and behavior. *Nature Neuroscience*, 25(1):11–19, 2022. ISSN 1546-1726. doi: 10.1038/s41593-021-00980-9.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
 Ł ukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Infor- mation Processing Systems*, 2017.
- Nikita Vladimirov, Yu Mu, Takashi Kawashima, Davis V. Bennett, Chao-Tsung Yang, Loren L. Looger, Philipp J. Keller, Jeremy Freeman, and Misha B. Ahrens. Light-sheet functional imaging in fictively behaving zebrafish. *Nature Methods*, 11(9):883–884, 2014. ISSN 1548-7105. doi: 10.1038/nmeth.3040.
- Yuxin Wu and Kaiming He. Group normalization. In *Proceedings of the European conference on computer vision (ECCV)*, pp. 3–19, 2018.
- Ailing Zeng, Muxi Chen, Lei Zhang, and Qiang Xu. Are transformers effective for time series forecasting? In *Proceedings of the AAAI conference on artificial intelligence*, volume 37, pp. 11121–11128, 2023.
- Yan Zhang, Márton Rózsa, Yajie Liang, Daniel Bushey, Ziqiang Wei, Jihong Zheng, Daniel Reep,
 Gerard Joey Broussard, Arthur Tsang, Getahun Tsegaye, et al. Fast and sensitive gcamp calcium
 indicators for imaging neural populations. *Nature*, 615(7954):884–891, 2023.

594 A APPENDIX

A.1 ARCHITECTURAL DETAILS

598 Every network has an embedding 3^3 convolutional layer mapping from temporal context C to F features, and an output convolutional layer mapping from F' (when upsampling) or F to 1 feature, 600 giving the single lead-time conditioned forecast frame. In the downsampling pathway, we apply one 601 convolutional block at every resolution. During symmetric upsampling, we use three convolutional 602 blocks at the lowest resolution, and two for all higher resolutions. Before upsampling to superresolution (i.e., resolution that is higher than that of the input), we use a convolution to map from 603 F to F' features to reduce the size of intermediate representations. During super-resolution upsam-604 pling we use one convolutional block per resolution. Each convolutional block has a pre-activation 605 residual design with the following chained layers: group normalization, swish activation, 3^3 convo-606 lution, group normalization, conditioning on lead time using a FiLM layer, swish activation, optional 607 feature dropout (only used with rate 0.1 for C = 256), and lastly the second 3^3 convolution. The 608 UNet-structure is realized by adding the representations obtained during sequential downsampling 609 to the upsampled representation. The number of features at every resolution is fixed to F = 128, 610 except for the super-resolution upsampling, where it is F' = 32.

611 612

613

596

597

A.1.1 SPATIAL VS. TEMPORAL MODELS

⁶¹⁴ This study employs three distinct models based on the aforementioned design.

The first model, maintaining a consistent spatial dimension of $512 \times 288 \times 72$, forgoes downsampling and upsampling blocks. It incorporates four processing blocks at this resolution, along with two convolutional layers at the input and output stages. The receptive field, calculated as $S = 1 + (4 \times 2+2) \times 2 = 21$, is determined by considering the central voxel and adding 2 for each 3^3 convolution.

The second model, downsamples the input data to $64 \times 64 \times 32$ and has a receptive field of S = 64. This is derived from the cumulative downsampling factors of (4, 4, 2) in the X, Y, and Z dimensions, respectively, and applying Equation 4.

623 Similarly, the third model employs downsampling factors of (16, 16, 8), resulting in a 624 $256 \times 256 \times 128$ receptive field. This translates to a global receptive field along the Z-axis, a near-625 global receptive field along the Y-axis, and a receptive field encompassing half of the total extent 626 along the X-axis.

627

628 A.1.2 LEAD-TIME CONDITIONING

629 For the results shown in Figure 4, we use three different losses: direct MAE, conditioned MAE, 630 and conditioned HL-Gauss. Apart from the FiLM layers to condition on lead-time, the architecture 631 is the same in all cases, with the exception of the last layer which maps from F to the output 632 dimensionality. The output dimensionality for the direct MAE is the number of forecast timesteps H. 633 For the conditioned MAE, it is simply 1, as also described in Figure 3. For the conditioned HL-634 Gauss loss, it is 32, which is equal to H, and each output corresponds to a discretized bin of the data 635 range. The HL-Gauss loss transforms a real value by representing it as a weighted average of bin 636 mean-values, for details see (Farebrother et al., 2024).

637 638

639

A.1.3 MODELS FOR DIFFERENT INPUT RESOLUTIONS

To investigate the influence of input resolution on model performance, we conducted a comparative
analysis. We compared our primary model, which operates on data downsampled by a factor of 4 in
the XY plane, with two alternative configurations: one employing a downsampling factor of 2, and
another utilizing full-resolution input.

To ensure equitable comparison, the architectures of these models were kept broadly consistent, with necessary adjustments to accommodate the differing input resolutions, while maintaining a consistent full-resolution output frame. Specifically, for the model operating on $2 \times$ downsampled input, we augmented the architecture with three additional blocks at the input resolution and removed one upsampling block, relative to the architecture depicted in Figure 3. In contrast, the model utilizing

full-resolution input incorporated two initial downsampling blocks with factors (2, 2, 1) and omitted 649 any super-resolution components, resulting in a conventional UNet architecture. 650

A.2 IMPACT OF UNSEGMENTED VOXELS

653 In contrast to video models which analyze all voxels of the calcium movie, the trace extraction pro-654 cess ignores voxels that do not correspond to segmented cells. This potentially discards information that could be useful in forecasting. To test to which degree this is indeed the case, we trained the 655 C = 4 video forecasting model with the unsegmented voxels set to constant value (0). The grand 656 average test MAE for that model (0.02663 ± 0.00003) was not significantly different from that of the video model processing the complete volume (0.02672 ± 0.00010) . This indicates that the unseg-658 mented voxels are unlikely to contain information that could improve forecasts and that any gains relative to the trace-based models can be attributed to the utilization of the spatial distribution of the 660 underlying calcium signals within the segmented cells.

661 662 663

657

659

651

652

A.3 ADDITIONAL EXPERIMENTAL RESULTS

664 In Figure 8 and 9, we show a fine-grained version of the benchmark of the trace and video-based 665 models in the main text (Figure 6). The figures show the performance per experimental condition 666 the fish was exposed to. For more details on these conditions, we refer to ZAPBench (Anonymous, 667 2024). For short context, C = 4, we observe that the video-based model performs better on six experimental conditions, and worse for many steps ahead on the "dots", "taxis", and "open loop" 668 conditions. For long context, the video-based model performs almost identical. As in the main text, 669 we display two standard errors about the mean in with shaded regions. 670

671 Figure 10 reports performance relative to four trace-based models included in ZAPBench. Figure 11 672 illustrates MAE differences for a few example frames.

673 On the right in Figure 7, we further show an ablation to 674 confirm that the improvement of multivariate video mod-675 els is due to increased receptive field and not because of 676 using more parameters. In particular, in our experimen-677 tal setup in Sec. 4.1.1 we keep FLOPS fixed instead of 678 number of parameters. In the example on the right, we 679 instead increase the width by a factor of two leading to an incrase in FLOPS by a factor of 4 while keeping the 680 receptive field fixed. We observe that increasing FLOPS 681 at the same spatial context leads statistically to the same 682 performance. Therefore, the performance improvement 683 observed in Figure 5 is likely due to the increased recep-684 tive field, especially for short context. The example on 685 the right is for the case of C = 4. 686



Figure 7: Ablation on increasing parameter count instead of receptive field.

- 687
- 688
- 689 690
- 691
- 692
- 693



- 696
- 697
- 699
- 700



Figure 8: Comparison of volumetric video to best trace-based model on all conditions for short context, C = 4.



Figure 9: Comparison of volumetric video to best trace-based model on all conditions for long context, C = 256.



Figure 10: Comparison of volumetric video model with four trace-based models from ZAP-Bench (Anonymous, 2024) for short (left) and long (right) context on the benchmark test set (averaged over eight conditions) and the experimental condition held out from the training data. In remaining figures, we report performance relative to TSMixer and Time-Mix (a univariate MLP) for short and long context, respectively. Note that MAEs of TiDE on the holdout are higher than the axis limits, which is due to its reliance on stimulus covariates. We report the mean and two standard errors.



Figure 11: Illustration of MAE differences. Top row shows the MAE between target and predicted activity for a video model on five test set frames for the gain condition, C = 4, at 32 steps predicted ahead, with brighter colors indicating higher error. Bottom row shows corresponding MAEs on these frames for a trace-based model. When MAEs are averaged across all test set frames and neurons for this condition, the MAE difference between these models is approximately 0.005.

803

804

805



Figure 12: Illustration of the two types of correlation metrics (for a single neuron). Top: actually recorded activity (green) is aligned in experiment time t with predicted snippets (blue) of activity of length H = 32 starting from various offsets. Bottom: correlations are always computed over 32 time steps between predicted activity and corresponding real recording. In Corr_H, complete predicted snippets are correlated with the recordings, and then averaged over starting points. In Corr_W, snippets are assembled from predictions at a specific lead time h, and correlated with the corresponding recordings. Reported metrics are averaged over all neurons.

A.4 CORRELATION METRICS

832

833

834

835

836

837

838 839

849

850

851

852

856

Table 3: Test set $Corr_H \pm 2$ SE.		
Context	Video	Trace
C = 4	$\textbf{0.1511} \pm \textbf{0.0021}$	0.1080 ± 0.0090
C = 256	$\textbf{0.1874} \pm \textbf{0.0004}$	0.1650 ± 0.0022

To better measure the quality of the temporal structures predicted by the model, we also computed two types of correlation metrics Corr_W and Corr_H , which compare recorded and predicted activity over H = 32 steps, with the predictions assembled at constant lead time h or from a specific starting point t, respectively (see Figure 12).

The correlation metrics paint a picture broadly consistent with that shown by the MAE, except in the long-context regime C = 256 where the video model outperforms the trace-based models (see Table 3, Figure 13).

857 A.5 COMPUTATIONAL COST ESTIMATES

The loss ablation in Figure 4 required around 5k GPU hours, pre-training and fine-tuning as shown in Sec. 4.1.2 around 14k GPU hours, comparing spatial to temporal context in Sec. 4.1.1 around 50k GPU hours, and the final results including the ablation on input resolution another 30k GPU hours. This makes a total of roughly 100k GPU hours used for the experiments presented in the paper.

A single training run of the best performing video model for C = 4 required 36 h, whereas the model for C = 256 required 120 h, both using 16 A100 GPUs. This compute cost is two to three orders



Figure 13: Comparison of volumetric video model with best-performing trace-based model in terms of $Corr_W$ for short (left) and long (right) context on the benchmark test set (averaged over eight conditions), higher is better. We report the mean and two standard errors.

of magnitude higher than that incurred by training the baseline trace-based models, which require about 2 h on a single A100 GPU. However, video models require less raw data preprocessing relative to time series models, partially offsetting the increased cost.

