

Efficient coding of natural scenes improves neural system identification

Yongrong Qiu^{a,b,c}, David A. Klindt^d, Klaudia P. Szatko^{a,b,c,e}, Dominic Gonschorek^{a,b,f}, Larissa Hoefling^{a,b,e}, Timm Schubert^{a,b}, Laura Busse^{g,h}, Matthias Bethge^{b,e,i}, and Thomas Euler^{a,b,e,✉}

^aInstitute for Ophthalmic Research, U Tübingen, 72076 Tübingen, Germany

^bCentre for Integrative Neuroscience (CIN), U Tübingen, 72076 Tübingen, Germany

^cGraduate Training Centre of Neuroscience (GTC), International Max Planck Research School, U Tübingen, 72076 Tübingen, Germany

^dDepartment of Mathematical Sciences, Norwegian University of Science and Technology, 7491 Trondheim, Norway

^eBernstein Centre for Computational Neuroscience, 72076 Tübingen, Germany

^fResearch Training Group 2381, U Tübingen, 72076 Tübingen, Germany

^gDivision of Neurobiology, Faculty of Biology, LMU Munich, 82152 Planegg-Martinsried, Germany

^hBernstein Centre for Computational Neuroscience, 82152 Planegg-Martinsried, Germany

ⁱInstitute for Theoretical Physics, U Tübingen, 72076 Tübingen, Germany

1 *Neural system identification* aims at learning the response
2 function of neurons to arbitrary stimuli using experimen-
3 tally recorded data, but typically does not leverage normative
4 principles such as efficient coding of natural environments.
5 Visual systems, however, have evolved to efficiently process
6 input from the natural environment. Here, we present a nor-
7 mative network regularization for system identification mod-
8 els by incorporating, as a regularizer, the *efficient coding* hy-
9 pothesis, which states that neural response properties of sen-
10 sory representations are strongly shaped by the need to pre-
11 serve most of the stimulus information with limited resources.
12 Using this approach, we explored if a system identification
13 model can be improved by sharing its convolutional filters
14 with those of an autoencoder which aims to efficiently encode
15 natural stimuli. To this end, we built a hybrid model to pre-
16 dict the responses of retinal neurons to noise stimuli. This
17 approach did not only yield a higher performance than the
18 “stand-alone” system identification model, it also produced
19 more biologically-plausible filters. We found these results to
20 be consistent for retinal responses to different stimuli and
21 across model architectures. Moreover, our normatively reg-
22 ularized model performed particularly well in predicting re-
23 sponses of direction-of-motion sensitive retinal neurons. In
24 summary, our results support the hypothesis that efficiently
25 encoding environmental inputs can improve system identifi-
26 cation models of early visual processing.

27 Neural system identification | Efficient coding | Natural scene statistics
28 Correspondence: thomas.euler@cin.uni-tuebingen.de

29 Significance

30 Computational models use experimental data to learn
31 stimulus-response functions of neurons, but they are rarely
32 informed by normative coding principles, such as the idea
33 that sensory neural systems have evolved to efficiently pro-
34 cess natural stimuli. We here introduce a novel method
35 to incorporate natural scene statistics to predict responses
36 of retinal neurons to visual stimuli. We show that con-
37 sidering efficient representations of natural scenes im-
38 proves the model’s predictive performance and produces
39 biologically-plausible receptive fields. Generally, our ap-
40 proach provides a promising framework to test various
41 (normative) coding principles using experimental data for

42 understanding the computations of biological neural net-
43 works.

44 Introduction

45 In the past years, advances in experimental techniques
46 enabled detailed, large-scale measurements of activity at
47 many levels of sensory processing (1). As a consequence,
48 *neural system identification* (SI) approaches have flour-
49 ished (Fig. 1a top). They empirically fit the stimulus-
50 response (transfer) function of neurons based on experi-
51 mentally recorded data (2–4). A classic example is the
52 generalized linear model (GLM, (2, 5)), which consists of
53 a linear filter as a first order approximation of a neuron’s
54 response function (i.e., its receptive field; (6)), followed
55 by a point-wise nonlinear function for the neuron’s output.
56 To account for additional non-linearities (e.g., (7, 8)), sev-
57 eral extensions, such as linear-nonlinear cascades (9, 10),
58 have been proposed. More recently, deep neural network-
59 based SI approaches inspired by the hierarchical process-
60 ing along the visual pathway (11, 12) have been developed
61 (reviewed in (13–17)). While SI methods became particu-
62 larly successful in predicting responses of visual neurons
63 (18–22), they often require large amounts of training data
64 and, more critically, do rarely consider adaptations to the
65 natural environment.

66 However, like other senses, vision has evolved to promote
67 a species’ survival in its natural environment (23), driv-
68 ing visual circuits to efficiently represent information un-
69 der a number of constraints, including metabolic limits and
70 space restrictions (24, 25). As a consequence, the visual
71 system has adapted to natural statistics, as shown, for ex-
72 ample, by the fact that the distribution of orientation pref-
73 erences of visual neurons mirrors the dominance of cardin-
74 al orientations in natural scenes (26–28).

75 Such adaptations are at the heart of *efficient coding* (EC)
76 approaches (Fig. 1a bottom): They derive computational
77 principles underlying neural systems from the statistical
78 properties of natural stimuli and by incorporating biolog-
79 ical constraints (15, 24, 25, 29–31). Here, one popular strat-
80 egy starts from the assumption that early visual processing

81 serves to decorrelate the redundant signals in natural environments (32, 33). This theory can reproduce feature selectivity, e.g., difference-of-Gaussian (DoG) kernels that have similar receptive field (RF) properties as retinal ganglion cells (RGCs; (34)). Recently, deep neural networks-augmented EC approaches were proposed, such as convolutional autoencoders (35, 36), which are trained to optimally reconstruct inputs in the presence of an information "bottleneck" (i.e., from a constrained latent representation). Such convolutional autoencoders have been shown to yield center-surround spatial RFs with similar properties as those observed in RGCs when encoding either pink ($1/f$) noise or natural scenes (37, 38). Still, a downside of EC is that it is not always straight forward to experimentally measure coding efficiency and feature selectivity predicted by these approaches in neural systems (discussed in (39, 40)) and, hence, the interpretation of EC models with respect to the biological underpinnings remains challenging.

100 Notably, the intersection between EC and SI has long remained largely unexplored but lately shifted more into focus. For instance, Mlynarski and colleagues proposed a theoretical framework incorporating normative theories for statistical inference on neural data (41). Here, we tested whether the EC hypothesis can serve as a useful inductive bias for learning the response functions of neurons from high-dimensional data. To investigate this, we built a hybrid model combining a SI branch with an EC branch, forced the two branches to share filters (Fig. 1b) and asked, if knowledge about natural scene statistics could help predicting retinal responses. To this end, we experimentally recorded Ca^{2+} signals of neurons in the mouse retina while presenting it with visual stimuli and then used these responses to train the SI branch, which aims to predict retinal responses. We used natural movies that we recorded in mouse habitats outdoors to train the EC branch, which aims to represent natural scenes efficiently (38). We found a synergy between neural prediction and natural scene statistics: The hybrid approach did not only have a better predictive performance than a pure SI approach, it also produced more biologically-plausible filters. Our results demonstrate that predicting sensory responses benefits from considering adaptations to the natural environment.

125 Results

126 **Hybrid system identification and efficient coding models.** To test if learning an efficient representation of natural input could help predict neuronal responses in the early visual system, we employed *normative regularization*, i.e., statistical regularization that is informed by normative coding principles, such as the idea that sensory systems have evolved to efficiently process natural stimuli. Specifically, we used this strategy to incorporate EC as a regularizer and developed a hybrid model that combines SI-based neural prediction and EC in a single model. The two model branches are linked by shared convolutional fil-

137 ters (Fig. 1b).

138 The *SI branch* approximates the response functions of recorded neurons to a visual dense noise (see below), and was implemented using a convolutional neural network (CNN) (Fig. 2a). Here, we used an L2 regularization on the convolutional layers to encourage smooth filters (42) and an L1 regularization on the fully connected (FC) layer for sparse readouts ((19); for details, see Methods).

145 The *EC branch* was trained to efficiently reconstruct input stimuli (i.e., natural scenes) from a constrained latent representation. For this branch, we used a convolutional autoencoder network (CAE) that we published before (for details, see (38) and Methods). Also in the EC branch, we enforced smooth filters by using L2 regularization, and limited the bandwidth by adding Gaussian noise and imposing L1 regularization on the hidden activations. The latter regularization also encourages sparse representations.

154 In the *hybrid model*, we implemented interactions between the two branches by shared filters (symbolized by red circle in Fig. 1b). Both branches were trained in parallel, with a weighted sum of their respective losses (L_{SI} and L_{EC}) used as optimization objective. By changing the weighting of the two losses, we were able to control the relative contribution of two branches on shaping the shared filters, and test our hypothesis to which degree efficient representations of natural scenes improve neural predictions (Fig. 2a,b). Specifically, weight w was used to define the hybrid model's loss function as $L_{Hybrid} = w \cdot L_{SI} + (1 - w) \cdot L_{EC}$ (Methods). For $w = 1$, the EC branch had no influence on the shared filters and, hence, the hybrid model behaved like the pure SI model. Conversely, for $w = 0$, the SI branch had no influence on the shared filters and, hence, the hybrid model behaved like the pure EC model. Thus, the smaller the weight, the more the EC branch contributed to shaping the filters.

172 To evaluate the influence of stimulus statistics on neural response predictions, we fed not only natural stimuli to the EC branch, but also phase-scrambled natural stimuli as well as noise. We refer to these models as *hybrid-natural*, *hybrid-pha-scr* and *hybrid-noise* (Fig. 2c). Moreover, to examine whether the performance improvements could be attributed to simple low-pass filtering, we trained SI networks using spatial convolutional filters composed from different numbers of basis functions derived from principle component analysis (PCA) on natural images (Fig. 2d), or the discrete cosine transform (DCT). These models are referred to as *SI-PCA* and *SI-DCT* networks.

184 To train the SI branch of our hybrid framework, we recorded somatic Ca^{2+} responses from populations of cells in the ganglion cell layer (GCL) of the *ex-vivo* mouse retina to 9-minute long noise stimuli using two-photon imaging (Fig. 3a; Methods; (43, 44)). The GCL contains the RGCs, which represent the retina's output neurons and form in the mouse about 40 parallel feature channels to higher visual brain areas (reviewed in (23)). RGCs gain their specific response properties by integrat-

193 ing upstream input from distinct sets of bipolar cells and
194 amacrine cells. Note that the GCL also contains some "dis-
195 placed" amacrine cells (dACs; (43, 45)). If not indicated
196 otherwise, we did not distinguish between these two GCL
197 cell classes in our datasets. The noise stimulus contained
198 two chromatic components (UV, green) matching the spec-
199 tral sensitivities of mouse photoreceptors (46). We used
200 the data of $n=96$ GCL cells that passed our quality crite-
201 ria (Methods) to fit a pure SI model with factorized spatial
202 and temporal convolutional filters, whose predictive per-
203 formance served as our baseline (Fig. 3b left).

204 **Neural system identification benefits from natural**
205 **scene statistics.** First, we measured the predictive per-
206 formance of the hybrid-natural model on the validation
207 data (for hyperparameter tuning) by systematically vary-
208 ing the relative impact of the two branches by changing
209 the weight w . We found that the performance steadily in-
210 creased with increasing EC influence (i.e., decreasing w)
211 up to an optimum (peaking at $w = 0.2$; Fig. 3c, red), af-
212 ter which the SI had too little influence on the shared fil-
213 ters and the performance dropped. Note that the correla-
214 tion values for the validation data are relatively low be-
215 cause these predictions were calculated on a single-trial
216 basis (Methods).

217 Next, we replaced the natural input to the EC pathway by
218 phase-scrambled scenes (*hybrid-pha-scr*) and white noise
219 across space and chromatic channels (*hybrid-noise*). Like
220 for the hybrid-natural model, the performance of the two
221 control models also increased with increasing EC influ-
222 ence up to a certain point, peaking at $w = 0.3$ and $w =$
223 0.4 for hybrid-pha-scr and hybrid-noise, respectively (Fig.
224 3c). This indicates that when incorporating EC, all hybrid
225 model versions showed some improvement up to certain w
226 values, before performance sharply declined.

227 To test to what extent simple low-pass filtering contributes
228 to the performance improvement observed for the hybrid-
229 natural model, we quantified the performance of two addi-
230 tional SI models, one with PCA and the other one with
231 DCT bases. By varying the number of bases used, we
232 found a maximum in predictive performance at 16 and 4
233 bases for SI-PCA and SI-DCT (zig-zag ordering), respec-
234 tively (Suppl. Fig. S1b).

235 Finally, to compare the performance on the test data across
236 models, we picked for each model, the w or number of
237 bases with the best predictive performance for the vali-
238 dation data. We found that the hybrid model with natu-
239 ral inputs to the EC branch attained the best performance
240 among all tested models (Fig. 3d,e). The hybrid-natural
241 model's superior performance compared to the hybrid-
242 pha-scr model suggests that the benefit of learning natu-
243 ral scene statistics extends beyond second-order statistics
244 such as the $1/f$ power spectrum of natural images. Nev-
245 ertheless, the hybrid-pha-scr model performed better than
246 the hybrid-noise version, pointing at a general benefit of
247 learning second-order statistics in the EC branch. More-
248 over, the hybrid-natural model was consistently better than

249 low-pass filtering control models (*SI-PCA* and *SI-DCT*),
250 suggesting that simple low-pass filtering does not fully ex-
251 plain the benefits of sharing kernels with the EC branch
252 trained to efficiently represent natural stimuli.

253 Together, our results suggest that normative network reg-
254 ularization — in particular, based on natural statistics —
255 can improve the performance of neural SI models.

256 **Hybrid models with natural inputs learn the most bi-**
257 **ologically-plausible filters.** To confirm that our hybrid
258 models capture the properties of the recorded cells, we
259 estimated their RFs (Fig. 3b; Suppl. Fig. S1f; Meth-
260 ods). Indeed, we found that the models learned antago-
261 nistic center-surround RFs with biphasic temporal kernels,
262 reminiscent of RGC RFs found in other studies (2, 43). To
263 get insights to which degree our models resembled biolog-
264 ical vision systems, we next investigated the internal repre-
265 sentations by analyzing the filters of the models' subunits
266 (18, 47). To this end, we compared the shared spatial con-
267 volutional filters between our tested models. As neurons in
268 the retina and further upstream in the early visual system
269 often feature smooth, Gaussian or DoG shaped RFs (e.g.,
270 (43, 48, 49)), we considered models with such shared fil-
271 ters as more biological plausible than those with other filter
272 organizations.

273 Interestingly, while the learned neuronal RFs were quite
274 consistent between models (cf. Fig. 3b), their shared spa-
275 tial filters differed considerably (Fig. 3f,h). When us-
276 ing natural images in the EC branch (*hybrid-natural*), fil-
277 ters indeed became smoother and more Gaussian-shaped,
278 which may be a result of the regularization by the EC
279 branch on the SI branch and which may have contrib-
280 uted to the performance improvement of predicting re-
281 sponses. This effect persisted though reduced when phase-
282 scrambled images were used (*hybrid-pha-scr*). More-
283 over, for smaller w values (i.e., stronger EC influence),
284 Gaussian-shaped filters became more frequent in the
285 hybrid-natural but not in the hybrid-noise model (Fig. 3f,
286 upper vs. lower row). For the SI models with PCA or DCT
287 basis, we found all filters to be smooth as they profited
288 from low-pass filtering of the respective transformation.
289 However, compared to the hybrid-natural model, their fil-
290 ters were less frequently Gaussian-shaped (Fig. 3h).

291 To quantify these findings, we fitted 2D Gaussian func-
292 tions to the filters and measured the goodness of the fit
293 via the coefficient of determination (R-squared; Methods).
294 Notably, for all three hybrid models, the w with the best
295 Gaussian fit was the same w that also resulted in the best
296 response predictive performance ($w = 0.2$, $w = 0.3$, and
297 $w = 0.4$ for *hybrid-natural*, *hybrid-pha-scr*, and *hybrid-*
298 *noise*, respectively; Fig. 3g). The filters of the hybrid-
299 natural model resembled smooth 2D Gaussians more than
300 for any other model (Fig. 3i), including SI-PCA and SI-
301 DCT. The difference of fit quality between hybrid-natural
302 vs. hybrid-pha-scr and hybrid-pha-scr vs. hybrid-noise
303 may be related to higher-order statistics and second-order
304 statistics of natural scenes, respectively.

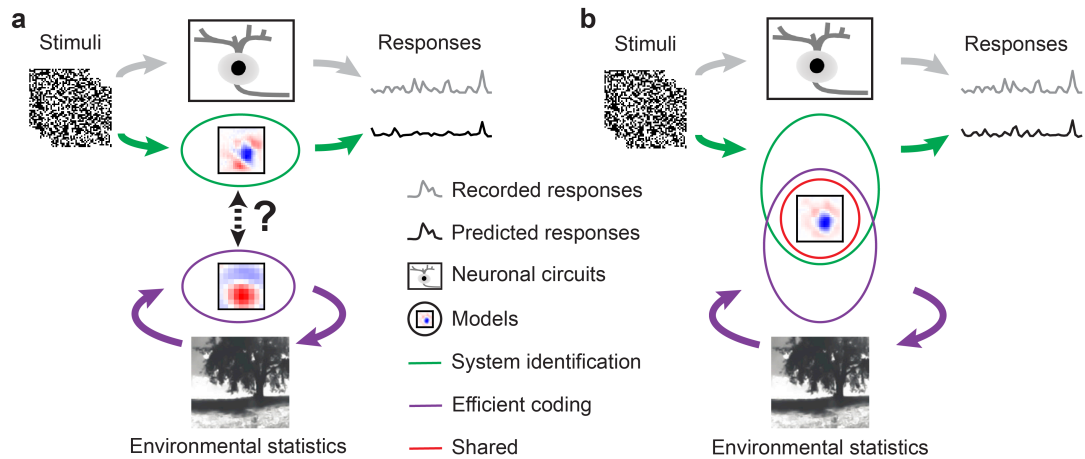


Fig. 1. Illustration of our hybrid model combining SI and EC. **a.** Illustration of two common approaches to studying visual systems: system identification (green branch) aims at predicting neural responses, whereas efficient coding (purple) seeks working out principles of the visual system based on environmental statistics. **b.** Our hybrid models combine system identification and efficient coding in a single model with shared filters (red circle) to predict neuronal responses.

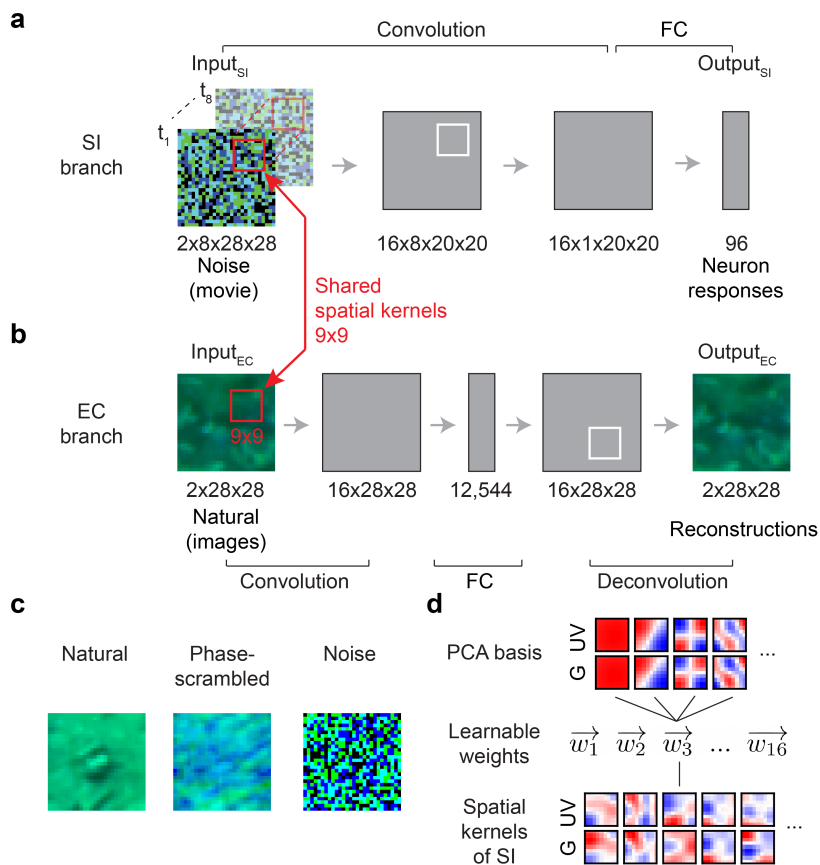


Fig. 2. Hybrid model with shared spatial filters. **a,b.** Schemata of SI model (a) and EC model (b) from Qiu et al. (38). The SI model branch consists of spatial and temporal convolutional layers, a fully connected (FC) layer and a nonlinear layer (see Methods). The EC model branch is a convolutional autoencoder (CAE), consisting of an encoder and a decoder network. In the hybrid model, the two branches were trained in parallel with shared spatial filters (red). Inputs_{SI}: 8-frame UV-green noise ($t_1 \dots t_8$); Output_{SI}: predicted GCL cell Ca²⁺ responses; Input_{EC}: UV-green natural images; Output_{EC}: reconstructed Input_{EC}. **c.** Example for the different inputs (natural images, phase-scrambled natural images, and noise) for the EC branch in hybrid models (*hybrid-natural*, *hybrid-pha-scr*, *hybrid-noise*). **d.** Using PCA filters as bases for spatial convolutional filters of the SI model; *SI-PCA* learned 16 weight vectors ($\vec{w}_1 \dots \vec{w}_{16}$) with same vector length as the number of PCA basis.

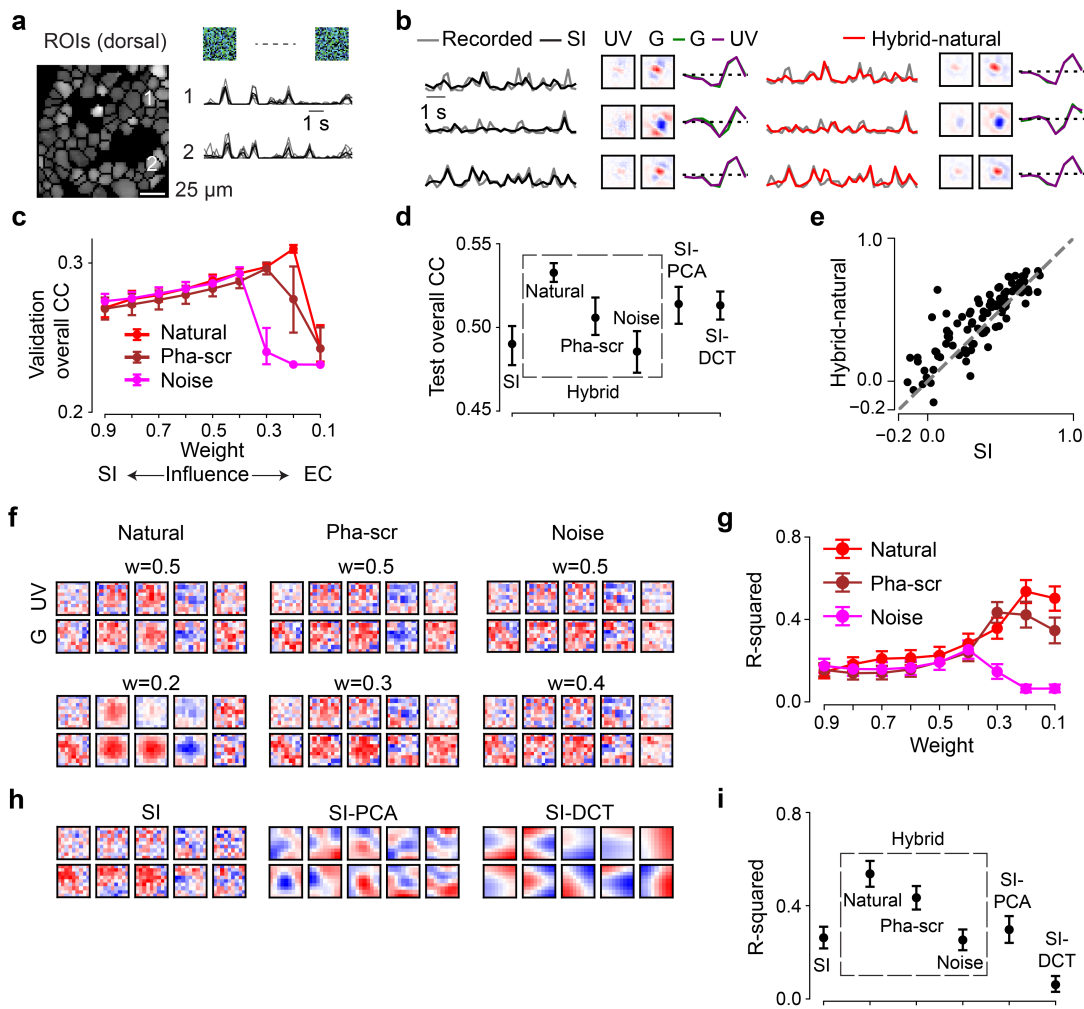


Fig. 3. Neural encoding tasks benefit from natural scene statistics. **a.** Region-of-interest (ROI) mask of one recording field in dorsal retina (left) and mean Ca^{2+} responses (black) of exemplary ROIs in response to 6 repeats of noise stimuli (single trials in gray). **b.** Three representative GCL cell responses (gray) to the noise stimulus (cf. Fig. 2a, left), together with predictions of best performing models on test data (black, SI; red, hybrid w/ natural scenes as input to the EC path, i.e., Input_{EC}), and learned spatio-temporal receptive fields (RFs) visualized by SVD. **c.** Model performance (linear correlation coefficient, CC; mean for $n=10$ random seeds per model) based on validation data for hybrid model with natural scenes (red), with phase-scrambled scenes (brown), or with noise (magenta) as Input_{EC} , and for different weights. **d.** Best performance (mean for $n=10$ random seeds per model) based on test data for SI, SI-PCA (16 bases), SI-DCT (4 bases), hybrid-natural ($w=0.2$), hybrid-pha-scr ($w=0.3$) and hybrid-noise ($w=0.4$); $p<0.0001$ for SI vs. hybrid-natural, $p=0.0085$ for SI-PCA vs. hybrid-natural, $p=0.0011$ for hybrid-natural vs. hybrid-pha-scr, two-sided permutation test, $n=10,000$ repeats). **e.** Scatter plot for model predictions based on test data for hybrid-natural ($w=0.2$) vs. SI at one random seed, with each dot representing one neuron. **f.** Representative spatial filters (shared convolutional filters) for hybrid models with different Input_{EC} and different weights. Upper: with optimal w ; lower: with different w . **g.** Mean R-squared of fitting a 2D Gaussian to spatial filters (cf. (f)), for hybrid model with natural scenes (red), with phase-scrambled scenes (brown), or with noise (magenta) as Input_{EC} , and for different w ($n=10$ random seeds per model). **h.** Representative spatial filters (shared convolutional filters) for SI, SI with PCA filters (16 bases) and SI with DCT filters (4 bases). **i.** Mean R-squared of fitting a 2D Gaussian to the spatial filters for one chromatic stimulus channel (green; $n=10$ random seeds per model; $p<0.0001$ for SI vs. hybrid-natural, $p<0.0001$ for SI-PCA vs. hybrid-natural, $p=0.0074$ for hybrid-natural vs. hybrid-pha-scr, two-sided permutation test, $n=10,000$ repeats). Error bars in (c),(d),(g),(i) represent 2.5 and 97.5 percentiles obtained from bootstrapping.

305 Taken together, our comparisons of the hidden spatial representations suggest that natural scene statistics promote
 306 latent feature representations akin to transformations in the early visual system.
 307
 308

309 **Efficient coding increases the data efficiency of system identification.** Next, we asked if the observed performance increase in the hybrid-natural vs. the baseline SI model was sensitive to the amount of training data, both with respect to their response predictions (Fig. 4a) and their learned spatial filters (Fig. 4b). To this end, we trained the SI and the hybrid-natural model ($w = 0.2$) with different amounts of data, ranging from 30% to 100%.
 310
 311
 312
 313
 314
 315
 316

317 Not unexpectedly, when more training data was used, predictive performance increased for both models (Fig. 4a top). However, we also found that the performance of the hybrid-natural model was consistently higher than that of the SI model, with the difference becoming significant for $\geq 60\%$ and peaking at around 90% training data (Fig. 4a bottom). Additionally, for both models the spatial filters became increasingly more Gaussian-like with more data (Fig. 4b). We also observed that the performance difference dropped for large dataset sizes — which, we expect, may be asymptotically near zero in the regime of infinite data.
 318
 319
 320
 321
 322
 323
 324
 325
 326
 327
 328

329 Together, these results suggest that a hybrid-natural model,

330 which has access to natural statistics, requires significantly
331 less training data than the baseline SI model.

332 **Hybrid models for testing temporal coding strate-**
333 **gies.** It has been suggested that early stages of visual pro-
334 cessing, rather than encoding a past stimulus, aim at pre-
335 dicting future stimuli in their temporal stream of inputs
336 (24, 50–52). Such a future prediction strategy is thought
337 to require a smaller dynamic range to be encoded than that
338 needed for representing past stimuli (past encoding), and
339 thus allows for lower energy consumption (53, 54). There-
340 fore, we next tested if the neural encoding task would profit
341 even more from natural statistics when spatio-temporal
342 (i.e., 3D) filters were shared between the hybrid model’s
343 two branches. We implemented both strategies — past
344 encoding and future prediction — in the EC branch, and
345 compared their influence on the SI task (55).

346 We modified the 2D SI model to use spatio-temporal (in-
347 stead of factorized spatial and temporal) convolutional fil-
348 ters to predict neural responses for 8-frame noise movies
349 (3D SI model; Suppl. Fig. S2a). Likewise, we employed
350 spatio-temporal convolutional filters for the EC branch. As
351 before, the two branches of the resulting hybrid model
352 were trained in parallel, but now sharing spatio-temporal
353 filters. In the past encoding case, the EC branch was
354 trained to reconstruct the 7th frame (at $t - 1$) of a contin-
355 uous 8-frame natural movie clip based on frames at $t - 7$
356 to t (*hybrid-natural-past*; Suppl. Fig. S2b,c). In the future
357 prediction case, the EC branch was trained to predict the
358 8th unseen frame based on the first 7 frames ($t - 7$ to $t - 1$)
359 of the clip (*hybrid-natural-future*; Suppl. Fig. S2d left).

360 Like for the 2D models, we varied w or the number of
361 bases and then selected the best model for each condition
362 (3D SI, *hybrid-natural-past*, *hybrid-natural-future*, and 3D
363 SI-PCA) based on validation performance. We next quan-
364 titatively compared the different models using the test data
365 (Fig. 5a,b; Suppl. Fig. S3c). We found that the 3D SI-
366 PCA model outperformed the 3D SI model, presumably
367 because the former profited from the low-pass filtering of
368 the PCA transformation. Importantly, both hybrid models
369 displayed a better performance than the 3D SI-PCA model.
370 While the *hybrid-natural-past* model performed slightly
371 better than its *hybrid-natural-future* counterpart, this dif-
372 ference was not statistically significant. In summary, both
373 the past encoding and future prediction strategy in the EC
374 branch turned out to be equally beneficial for the neural
375 encoding task and, as before, the benefit extended beyond
376 low-pass filtering effects. However, no performance in-
377 crease was achieved with respect to the 2D *hybrid-natural*
378 model (Fig. 5b vs. Fig. 3d).

379 We also analyzed the shared spatio-temporal filters using
380 the same metric as for the 2D case, which assesses the sim-
381 ilarity between spatial filters (after performing a low-rank
382 decomposition of 3D shared filters into spatial and tempo-
383 ral components; see Methods) and smooth 2D Gaussians
384 (Fig. 5c,d). Again, we found higher R-squared values for
385 the hybrid models and the 3D SI-PCA model compared

386 to the baseline SI case. Note that here, the 3D SI-PCA
387 model did not significantly differ from the two hybrid mod-
388 els, possibly due to a large number of bases ($n = 128$ vs.
389 $n = 16$ in the 2D case).

390 Next, we asked if the fact that we did not see a significant
391 advantage of 3D over 2D could be due to the relatively
392 slow (5 Hz) noise stimulus, which may drive insufficiently
393 temporal properties of the GCL cell responses. There-
394 fore, we recorded a new dataset ($n = 64$ cells) in which
395 we presented a 30-Hz dense noise stimulus and used it
396 with the 3D hybrid models. Like for 5-Hz noise, hybrid-
397 natural-past and hybrid-natural-future models performed
398 similarly on the validation data, with a peak in perfor-
399 mance at around $w = 0.7$ (Suppl. Fig. S4a), as well as on
400 the test data, where they were significantly better than the
401 3D SI model (Suppl. Fig. S4b). Moreover, both 3D hy-
402 brid models learned shared filters with similar R-squared
403 values, which were significantly higher than that of the 3D
404 SI model (Suppl. Fig. S4c). But again, the 3D models
405 performed only equally well compared to the 2D models.

406 In summary, the hybrid-natural models achieved a higher
407 performance for different noise stimuli (5-Hz vs. 30-Hz)
408 and different shared filter organizations (2D vs. 3D) than
409 all other tested models. Therefore, it is likely that their su-
410 perior predictive performance for neuronal responses and
411 their more biologically plausible filters resulted from the
412 EC branch having access to natural statistics.

413 **Direction-selective neurons benefit more than oth-**
414 **ers from hybrid models.** The retina encodes the visual
415 scene in a number of features that are represented by the
416 more than 40 different types of RGC whose outputs are
417 relayed in parallel to higher visual centers in the brain
418 (43, 56–59). Thus, we next asked, if access to natural
419 statistics allows our hybrid models to predict some cell
420 types better than others (Fig. 6). Earlier, it has been shown
421 that motion-relevant properties emerge in the efficient cod-
422 ing framework for both past encoding and future prediction
423 approaches (55). Therefore, we employed our 3D hybrid
424 models (cf. Fig. 5) and focused on direction-selective (DS)
425 cells (43, 60).

426 For this analysis, we used a set of $n=427$ GCL neurons,
427 whose responses were recorded not only to the 5-Hz noise
428 stimulus (for training the models) but also to full-field
429 chirp and moving bar stimuli. The latter two stimuli (Fig.
430 6a) enabled us to identify the functional type of each
431 recorded GCL neuron (43) using a cell type classifier (see
432 Methods; Suppl. Fig. S5).

433 To explore cell type-specific effects, we chose a dataset
434 size (30% of total recording time) for which the synergy
435 between neural SI and EC was particularly pronounced.
436 As expected, we found that both hybrid networks (*hybrid-*
437 *natural-past* and *hybrid-natural-future*) performed signifi-
438 cantly better than the SI model, with no significant differ-
439 ence between the two hybrid models (cf. Fig. 5b, Suppl.
440 Fig. S4b).

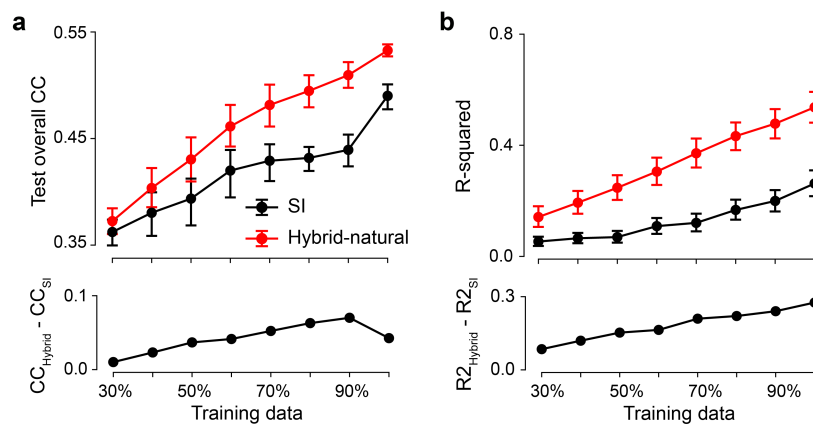


Fig. 4. Hybrid-natural models with better data efficiency for neural prediction. **a.** Mean model performance (top) based on test data for SI and hybrid-natural ($w=0.2$; $n=10$ random seeds) with different training data sizes and mean difference between SI and hybrid-natural (bottom). **b.** Mean R-squared (top) of fitting a 2D Gaussian to spatial filters for green stimulus channel for SI and hybrid-natural ($w=0.2$; $n=10$ random seeds) with different training data sizes, and the mean difference between R-squared for SI and hybrid-natural (bottom). Error bars represent 2.5 and 97.5 percentiles with bootstrapping.

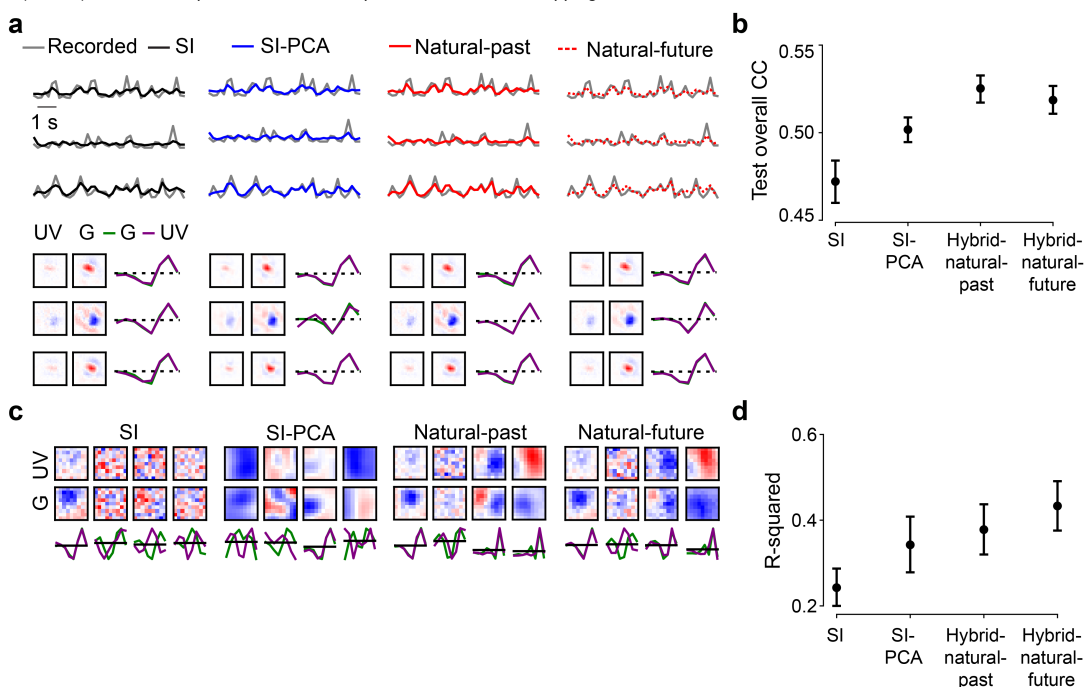


Fig. 5. Past encoding or future prediction strategies using 3D shared filters perform equally well. **a.** Top row: Responses of three exemplary GCL cells to 5-Hz noise stimulus (gray) and predictions of best performing models on test data (black, SI; blue, SI with PCA filters; red solid, hybrid for encoding the past; red dotted, hybrid for predicting the future). Bottom row: Respective learned RFs of the three cells (visualized by SVD). **b.** Mean model performance based on test data for SI, SI-PCA (128 bases), hybrid-natural-past, and hybrid-natural-future (both $w=0.4$; $n=10$ random seeds; $p<0.0001$ for SI vs. hybrid-natural-past, $p=0.0005$ for SI-PCA vs. hybrid-natural-past, $p=0.2563$ for hybrid-natural-past vs. hybrid-natural-future, two-sided permutation test, $n=10,000$ repeats). **c.** Representative shared spatial and temporal filters of 3D models ($n=1$ random seed, visualized by SVD; temporal kernels for UV and green stimulus channels indicated by purple and green, respectively). **d.** Mean R-squared of fitting a 2D Gaussian to shared spatial filters (for green stimulus channel; $n=10$ random seeds per model; $p=0.0003$ for SI vs. hybrid-natural-past, $p=0.4356$ for SI-PCA vs. hybrid-natural-past, $p=0.1895$ for hybrid-natural-past vs. hybrid-natural-future, two-sided permutation test, $n=10,000$ repeats). Error bars in (b),(d) represent 2.5 and 97.5 percentiles with bootstrapping.

441 First, we evaluated if any of the broader functional groups
 442 of GCL cells profited more from natural statistics than others.
 443 For this, we sorted the cells into 6 groups based on
 444 their response polarity (ON vs. OFF) and transience, and
 445 based on whether they were RGCs or dACs (for group
 446 sizes, see Fig. 6 legend). For all 6 groups, the hybrid
 447 models showed a better predictive performance than the SI
 448 model (Fig. 6b). However, no significant differences were
 449 observed between any pair of groups ($p>0.05$ for all pair-

450 wise comparisons, two-sided permutation test, $n=10,000$
 451 repeats; Fig. 6c) and the two hybrid models ($p>0.05$ for all
 452 pair-wise comparisons; Suppl. Fig. S6a).

453 Next, we grouped the cells into DS ($p<0.05$, direction
 454 tuning using a permutation test; $n=90$) and non-DS cells
 455 ($n=300$) based on their moving bar responses (Fig. 6a
 456 right). Note that $n=37$ neurons were excluded as they did
 457 not pass the quality check for chirp and moving-bar responses
 458 (Methods). We found that the predictive perfor-

459 mance for DS cells was significantly higher than that of
460 the non-DS cells for both hybrid-natural-past (Fig. 6d,e;
461 $p=0.0027$) and hybrid-natural-future (Suppl. Fig. S6b,c;
462 $p=0.0042$). To test whether this performance difference
463 was merely due to different signal-to-noise ratios in DS vs.
464 non-DS cells, we compared their response quality indices
465 (QI ; Methods). While DS cells had significantly higher
466 QI values for moving-bar responses (QI_{bar}) than non-DS
467 cells, we did not find any significant difference between the
468 two groups with respect to their noise (QI_{noise}) or chirp
469 responses (QI_{chirp} ; Suppl. Fig. S6e-g). These results sug-
470 gest that DS cells benefit more from the EC branch of the
471 hybrid models than non-DS cells, partially consistent with
472 earlier findings ((55); see also Discussion).

473 In summary, efficient coding of natural statistics served as
474 a beneficial normative regularization for all types of mouse
475 GCL cells and in particular DS cells, suggesting the poten-
476 tial role of motion statistics in the natural environment on
477 shaping neuronal response properties.

478 Discussion

479 In this study, we asked if access to natural scene statis-
480 tics can help predicting neural responses. To address this
481 question, we combined system identification (SI, (3)) and
482 efficient encoding (EC, (25)) methods into a normatively
483 regularized (hybrid) modeling framework. Specifically,
484 we used models that efficiently represent natural scenes
485 recorded in the mouse' habitat to regularize models that
486 predict retinal responses to visual stimuli. We analyzed
487 such hybrid models with shared spatial filters, and found
488 that natural images as input to the EC branch indeed im-
489 proved the performance in predicting retinal responses and
490 allowed the model to generate filters that resembled RFs
491 found in the early visual system. These improvements ex-
492 tend beyond those gained by simple low-pass filtering or
493 using second-order statistics of the natural scenes. Our hy-
494 brid models with shared spatio-temporal filters performed
495 similarly well as those with shared spatial filters, indepen-
496 dently of whether they used a past encoding or a future
497 prediction strategy. Notably, predictions for DS cells in
498 the mouse retina improved the most in the hybrid mod-
499 els with natural input. In summary, our results suggest
500 that sourcing information about an animal's environment
501 — e.g., through hybrid SI-EC models — helps building
502 more predictive and biologically-plausible models of neu-
503 ronal networks. More generally, our findings lend support
504 to the idea that knowledge of natural statistics is already
505 encoded in sensory circuits.

506 **Hybrid models improve data efficiency.** The differ-
507 ence in predictive performance between the hybrid and the
508 baseline SI model was significant and it depended on the
509 amount of available data, indicating that our hybrid model-
510 ing approach increased data efficiency. We note that both
511 the stimulus (dense noise) and the neural model system
512 (retinal neurons) present much easier SI problems than,
513 for instance, predicting more nonlinear neural responses to

514 natural stimuli (18, 61). For those more challenging prob-
515 lems at downstream visual areas, where neural response
516 functions and, hence, the neural prediction tasks, become
517 more complex (62), the data efficiency of a hybrid ap-
518 proach and the improvement from natural scene statistics
519 may be even higher.

520 **Biological plausibility and temporal coding princi-
521 ples in hybrid models.** The biological plausibility of
522 most learned models was positively correlated with their
523 predictive performance except some indeterminacy for SI-
524 DCT models, suggesting that more biologically plausible
525 filters increased performance. Note that we used the filters'
526 similarity to smooth 2D Gaussian functions as a measure
527 of biological plausibility, following the assumption that
528 RFs in the retina (and at early downstream stages of the vi-
529 sual system) often feature smooth, Gaussian-like structure
530 (43, 48, 49). However, a deep, systematic understanding
531 of artificial and neuronal networks and their hidden repre-
532 sentations likely calls for other methods besides of filter
533 inspection (discussed in (63)).

534 As the natural environment is not static, we also created
535 hybrid models that acknowledge the time domain by shar-
536 ing spatio-temporal filters. Surprisingly, both variants —
537 past encoding and future prediction — behaved quite sim-
538 ilar. However, in the stand-alone EC models (that is only
539 the respective EC branch), the temporal components of the
540 filters learned by the future prediction were much more
541 diverse than those of past encoding (Suppl. Fig. S2c,d
542 right). Interestingly, the differences between temporal fil-
543 ter of these stand-alone EC models decreased with the in-
544 corporation of the neural prediction task in the hybrid mod-
545 els.

546 The filter diversity in our 3D hybrid models is reminis-
547 cent of earlier findings by Chalk and colleagues (2018),
548 who reported the emergence of filters sensitive to motion
549 direction and motion speed in their past encoding and fu-
550 ture prediction EC models, respectively. However, in con-
551 trast to their results, we did not see a difference between
552 our hybrid-past and hybrid-future models with respect to
553 motion-sensitive filters: Both of them performed better in
554 predicting responses of DS vs. non-DS cells. Further work
555 is needed to understand that partial (mis)match between
556 our work and that by Chalk et al., and why specifically DS
557 cells profited from both our 3D hybrid models.

558 **Hybrid models of retinal signal processing.** It has
559 been suggested that natural stimuli drive more diverse neu-
560 ral responses, and more complex feature transformations
561 are required to determine the respective stimulus-response
562 functions ((18, 64), but also see (65)). Therefore, one fu-
563 ture direction may be to record retinal activity while pre-
564 senting natural movies (e.g., from (38)) and use it as input
565 for the SI branch of the hybrid model. Finding a more pro-
566 nounced performance improvement compared to the base-
567 line SI model would support the notion that the noise stim-
568 ulus we used in this study may have indeed limited the ben-

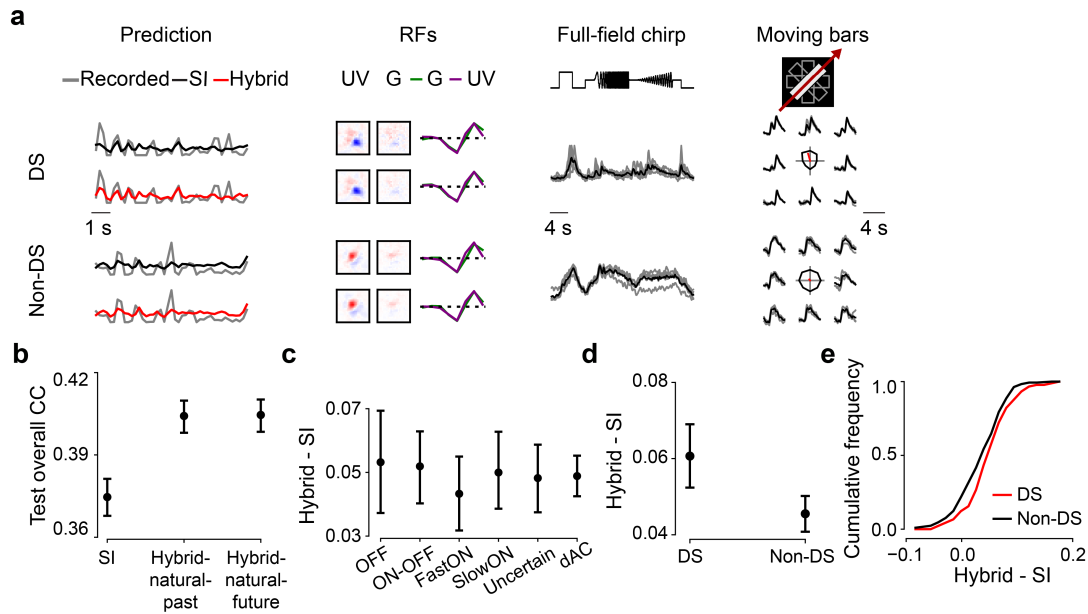


Fig. 6. Direction-selective (DS) neurons benefit more from hybrid models. **a.** Recorded (gray) and predicted (black, SI; red, hybrid-natural-past; response amplitude scaled with a constant 1.5 for better visualization) responses to noise, RFs, as well as full-field chirp responses and moving bar responses (gray, single trials; black, means) of representative DS and non-DS cells. Note that the RFs were dominated by UV stimulus channel because cells were recorded in ventral retina (see Methods). **b.** Mean model performance based on test data for SI, hybrid-natural-past and hybrid-natural-future (both $w = 0.7$; $n=10$ random seeds per model; trained with responses of $n=427$ GCL cells to 5-Hz noise stimulus; $p<0.0001$ for SI vs. hybrid-natural-past, $p=0.9307$ for hybrid-natural-past vs. hybrid-natural-future; two-sided permutation test, $n=10,000$ repeats). **c.** Difference in mean performance between hybrid-natural-past and SI based on test data for 6 broad functional groups of GCL cells (35 OFF, 59 ON-OFF, 49 fast-ON, 38 slow-ON, and 64 uncertain RGCs, as well as 145 dACs; see Methods and Results; $n=10$ random seeds per model). **d.** Like (b) but for $n=90$ DS and $n=300$ non-DS cells. **e.** Cumulative histogram of difference in mean prediction between hybrid-natural-past ($w = 0.7$) and SI on test data for DS (red) and non-DS cells (black), at one particular seed. Error bars in (b)–(d) represent 2.5 and 97.5 percentiles with bootstrapping.

569 efits from the EC branch (see above). Neural data to natu-
 570 ral stimuli would also allow us to revisit our hybrid mod-
 571 els with respect to the prediction of motion sensitive cells
 572 and the differences between our results and those from ear-
 573 lier work ((55); see above). Furthermore, such data may
 574 be useful for characterizing model generalization (domain
 575 transfer, see e.g., (61, 64)) by using responses to natural
 576 stimuli as unseen test data with a hybrid model trained with
 577 cell responses to noise stimuli.

578 For our current analysis, we used broad group assign-
 579 ments (e.g., FastON RGCs), which include several func-
 580 tional types of RGC (e.g., ON-step, ON-transient, ON-
 581 high-frequency etc; (43)) or dACs, but did not detect any
 582 differences in performance gain except for the DS neurons.
 583 Still, it is possible that distinct types of RGC profit more
 584 than others from the EC branch of our hybrid models. For
 585 example, the so-called W3 RGCs, for which the best stimu-
 586 lus found so far is a small dark moving spot (66), may
 587 not be “designed“ to efficiently represent natural stimuli
 588 but rather to extract survival-relevant features (i.e., detect-
 589 ing aerial predators). Here, we could build models with
 590 different normative regularization or tasks (i.e., detecting
 591 predators in images of the sky) and would expect that this
 592 RGC type profits little from efficiently encoding natural
 593 statistics in the hybrid model. Studying coding strategies
 594 across RGC types could contribute an important biological
 595 perspective to the perennial debate between efficient cod-
 596 ing (67) and feature detection (56) proponents.

597 Normative network regularization as a framework 598 for studying neural coding.

599 In this study, we regularized the filters of a SI model with a normative EC model to pre-
 600 dict visually-evoked responses of cells in the retina. Some
 601 forms of such normative regularization have also been dis-
 602 cussed and/or applied in earlier work. For example, Den-
 603 eve and Chalk (68) discussed the relations between SI (en-
 604 coding) models and EC, and argued that the latter may pro-
 605 mote shifting the focus in SI from the single-cell to to the
 606 population level, while Turner et al. (15) considered the in-
 607 tegration of stimulus-oriented approaches (such as EC) for
 608 discriminative tasks (such as object recognition). Later,
 609 Teti et al. (69) employed sparse coding with lateral inhibi-
 610 tion in simulations of neuronal activation in visual cortex.
 611 More recently, Młynarski et al. (41) proposed a probabilis-
 612 tic framework combining normative priors with statistical
 613 inference and demonstrated the usefulness of this approach
 614 for the analysis of diverse neuroscientific datasets. Our
 615 work not only provides further evidence to the feasibility
 616 of combining coding principles for identification of neu-
 617 ral response properties, it also demonstrates the benefits of
 618 leveraging natural scene statistics for neural prediction.

619 We expect that our hybrid modeling strategy may also
 620 work for different processing stages along the early visual
 621 pathway (and potentially other modalities, e.g., sound).
 622 This said, however, one needs to keep in mind that dif-
 623 ferent stages along the visual pathway have different tasks
 624 and constraints, and, thus, likely incorporate different ef-
 625 ficient coding principles: For instance, the retinal hard-

ware is space-limited and has to encode visual features in view of a bottleneck with limited bandwidth (optic nerve), whereas the primary visual cortex has comparably abundant resources which might serve for accurate probability estimation for behavioral tasks, such as novelty detection (discussed in (24, 70)). It is also worth to note that different visual processing stages (such as primary visual cortex vs. higher visual areas, or adaptation of visual coding to different behavioral states) may benefit from the hybrid modeling to a different degree, as efficient coding approaches learn filters that may be more relevant to stimulus-related features, but not high-level behavior goals (see discussion in (15)). Additionally, it would be interesting to compare our hybrid models with SI models regularized with other behavioral tasks such as object recognition (e.g., (11)) or predator detection (see above) for neural predictions along the ventral visual stream.

There is a long tradition of using SI models (reviewed in (3)) in predicting the responses of neurons to a great variety of stimuli (e.g., (2, 4, 18, 19, 71, 72)). Our results demonstrate how the EC hypothesis can be successfully leveraged as normative regularization for the identification of neural response properties. More generally, using EC as a flexible tool to impose regularization on modeling, the hybrid framework offers an opportunity to test different coding principles and unsupervised learning objectives with regards to experimental data for understanding neuronal processing.

Materials and Methods

Animal procedures and retinal activity recordings.

Animal procedures. All animal procedures were performed in accordance with the law governing animal protection issued by the German Federal Government (Tierschutzgesetz), approved by the governmental review board (Regierungspräsidium Tübingen, Baden-Württemberg, Konrad-Adenauer-Str. 20, 72072 Tübingen, Germany). We used $n=5$, 5-9 weeks old female C57BL/6 mice (wild-type; JAX 000664, Jackson Laboratory, USA). Due to the exploratory nature of our study, we did not use any statistical methods to predetermine sample size, nor did we perform blinding or randomization.

Animals were housed under a standard light-dark (12h:12h) cycle. All procedures were carried out under very dim red illumination (>650 nm). Prior to the start of the experiment, animals were dark-adapted for ≥ 1 h, then anesthetized with isoflurane (Baxter, Germany), and killed by cervical dislocation.

The eyes were enucleated and hemisected in carboxygenated (95% O₂, 5% CO₂) artificial cerebrospinal fluid (ACSF) solution containing (in mM): 125 NaCl, 2.5 KCl, 2 CaCl₂, 1 MgCl₂, 1.25 NaH₂PO₄, 26 NaHCO₃, 20 glucose, and 0.5 l-glutamine (pH 7.4). Next, the retina was flat-mounted onto an Anodisc (#13, 0.2 μ m pore size, GE Healthcare, Germany) with the ganglion cell layer

(GCL) facing up. To uniformly label the GCL cells, bulk electroporation was performed with the fluorescent Ca²⁺ indicator Oregon-Green BAPTA-1 (OGB-1; Invitrogen, Germany), as described earlier (44, 73), using 4-mm plate electrodes (CUY700P4E/L, Xceltis, Germany) and 9 pulses (~ 9.2 V, 100 ms pulse width at 1 Hz). After electroporation, the tissue was immediately moved to the microscope's recording chamber, where it was continuously perfused with carboxygenated ACSF at $\sim 36^\circ\text{C}$ and left to recover for ~ 30 min before recordings started. Additionally, Sulforhodamine-101 (SR101, Invitrogen, Germany) was added to the ACSF (~ 0.1 μ M final concentration) to visualize blood vessels and identify damaged cells.

Two-photon Ca²⁺ recordings and light stimulation. We recorded light stimulus-evoked Ca²⁺ signals in GCL cells of the explanted mouse retina using a MOM-type two-photon (2P) microscope (74, 75) from Sutter Instruments (purchased from Science Products, Germany), as described earlier (43, 44). In brief, the microscope was powered by a mode-locked Ti: Sapphire laser (MaiTai-HP DeepSee, Newport Spectra-Physics, Germany) at 927 nm. Two detection pathways allowed simultaneously recording of OGB-1 and SR101 fluorescence (HQ 510/84 and HQ 630/60, respectively; both Chroma/AHF, Germany) through a 16x water immersion objective (CFI75 LWD \times 16 /0.8W, DIC N2, Nikon, Germany). A custom-written software (ScanM, by M. Müller and T.E.) running under IGOR Pro 6.3 for Windows (Wavemetrics, USA) was used to acquire time-lapsed (64x64 pixels) image scans at a frame rate of 7.8125 Hz. Higher resolution images were acquired using 512x512 pixel scans. Additionally, to register the scan field positions, the outline of the retina and the optic disc were traced.

The retinas were presented with color noise stimulus using a visual stimulator tuned to the spectral sensitivities of mice (76). This stimulus consisted of independent binary dense noise (28x28 pixel frames, each pixel covering $(0.83^\circ)^2$ of visual angle) in the UV and green stimulator channels at 5 or 30 Hz. The stimulus contained 5 different training sequences (96 s each) interspersed with 6 repeats of a 10 s test sequence (Suppl. Fig. S1a).

In total, we used three data sets for modeling: (i) responses of $n=96$ GCL neurons to 5-Hz noise recorded in dorsal retina ($n=2$ eyes); (ii) responses of $n=427$ GCL neurons to 5-Hz noise recorded ventrally ($n=5$ eyes); in this dataset, we also presented two other stimuli: a full-field chirp (700 μ m in diameter) and a moving bar stimulus (300x1,000 μ m bright bar moving at 8 directions at 1 mm/s). The responses to these latter stimuli were used to functionally classify the recorded GCL neurons (43). (iii) $n=64$ GCL neurons to 30-Hz noise recorded ventrally ($n=2$ eyes). Note that all cell numbers are after quality control (see below).

Data preprocessing and analysis. For each cell, we calculated a quality index (QI , with $0 \leq QI \leq 1$) for its responses to each stimulus type as follows:

$$QI = \text{Var}[E[C]_t]/E[\text{Var}[C]_t] \quad (1)$$

where C is a t -by- r response matrix (time samples, t , by repetitions, r). The higher QI , the more reliable the response and the higher the signal-to-noise ratio. For the noise stimulus, QI_{noise} was determined based on the test sequence responses. For the following analysis, we only used cells with $QI_{noise} > 0.25$; in case chirp and moving bar responses were also recorded, neurons had to fulfill $QI_{chirp} > 0.35$ or $QI_{bar} > 0.6$ to be included.

In case of the noise stimulus, we preprocessed each cell's Ca^{2+} signal by Z-normalizing the raw traces and matching sampling frequency of the recording (7.8125 Hz) to the stimulus frequency (5 or 30 Hz) via linear interpolation. Then, the traces were detrended using a high-pass filter (> 0.1 Hz) and their 1st order derivatives were calculated, with negative values set to zero. We used the average of a cell's responses to the 6 test sequence repeats as ground truth. Excluding the test sequences, we had per cell a total of 480 s of data, of which we used 440 s ($\sim 91\%$) for training and the remaining 40 s ($\sim 9\%$) for validation (i.e., to pick the hyperparameters of the SI model, see below).

For chirp and moving bar responses, we first detrended the traces and then normalized them to $[0, 1]$ (44). Using these responses, the cells were classified to different functional groups (43) using RGC type classifier (see below).

To estimate the directional tuning from the moving bar responses, we first performed singular value decomposition (SVD) on the mean response matrix, resulting in a temporal and a directional component. We then summed the directional vectors in 2D planes and used the resulting vector length as direction selectivity index. Next, by shuffling trial labels and computing the tuning curve for 1,000 times (permutation test), we got the null distribution (no directional tuning). The percentile of true vector length was used as p-value of directional tuning (43). Here, we considered cells with $p < 0.05$ as direction-selective (DS) and the remaining ones as non-DS.

RGC type classifier. To predict the functional type of GCL cells, we used a Random Forest Classifier (RFC; (77)), which was trained on a published mouse dataset (43). In that study, features were extracted from the responses to different visual stimuli (e.g., chirp and moving bar) and used to cluster GCL cells into 32 RGC types and 14 additional dAC types. Here, we learned a mapping f from response features (20 features from responses to chirp, ϕ_{chirp} and 8 features from responses to moving bar stimulus, ϕ_{mb}) and two additional parameters $\Theta = \{\theta_{soma}, \theta_{DS}\}$ to functional cell type labels L by training a RFC for the dataset from (43):

$$f : (\phi_{chirp}, \phi_{bar}, \Theta) \mapsto L \quad (2)$$

where θ_{soma} denotes soma size to distinguish between alpha and non-alpha RGC types and θ_{DS} denotes p-value of permutation test for direction selectivity to distinguish between DS and non-DS RGC types.

We fit the RFC on a subset of data from (43) and validated its performance on a held-out test dataset. The classifier had a prediction accuracy of $\sim 76\%$ on a held-out test dataset (Suppl. Fig. S5). To apply the trained classifier to our newly recorded dataset, we projected the RGC responses (normalized to $[-1, 1]$) into the feature space described in (43) by computing the dot product between the response and the feature matrices. We used the RFC implementation provided by the python package `scikit-learn` (78) to train the classifier.

2D models.

Stand-alone SI model (2D). As baseline model to predict the responses of neurons to the noise stimulus, we employed a stand-alone SI model (supervised learning), in which we used factorized spatial and temporal convolutional filters (Fig. 2a; (79, 80)). This SI model consisted of one spatial convolutional layer (16x2x1x9x9, output channels x input channels x depth x image width x image height), one temporal convolutional layer (16x16x8x1x1, with 8 stimulus frames preceding an event), and — after flattening the spatial dimension — one fully connected layer (FC; 96x6,400, output x input channels), followed by an exponential function. No padding was used. The loss function was defined as:

$$L_{SI} = \sum_i (\hat{r}_i - \bar{r}_i \log \hat{r}_i) + \alpha_1 \|\bar{w}_{cs}\|_2 + \alpha_2 \|\bar{w}_{ct}\|_2 + \beta \|\bar{w}_f\|_1 \quad (3)$$

Here, the first term is the Poisson loss between predicted responses (\hat{r}_i) and ground truth (\bar{r}_i) (with i denoting the neuron index), the second term is the L2 penalty on the weights of the spatial convolutional filters (\bar{w}_{cs}) with hyperparameter α_1 , the third term is the L2 penalty on the weights of temporal convolutional filters (\bar{w}_{ct}) with hyperparameter α_2 , and the last term is the L1 penalty on the FC layer (\bar{w}_f) with hyperparameter β .

After performing a grid search for the three hyperparameters, we picked $\alpha_1 = 10, \alpha_2 = 10, \beta = 1/16$ which yielded the best performance on the validation data. After training, we estimated the neurons' spatio-temporal RF filters by computing gradients for each neuron, starting with a blank image sequence as input. These gradients represent the first-order approximation of the input that maximizes the neuron's activation (6). For visualization, we extracted the spatial and temporal RFs via SVD.

As a metric of biological plausibility, we calculated the coefficient of determination (R-squared; $[0, 1]$) of fitting 2D Gaussian distributions to the spatial (component of) the convolutional filters. We set the R-squared value to 0 if the sigma of the fitted Gaussian was larger than the size of the filter (i.e., 9 pixels). We calculated this fit quality for the filter of the chromatic channel with the dominant response. Because the mouse retina is divided into a more green-sensitive dorsal and a more UV-sensitive ventral retina (e.g., (44)), this meant that for dorsal neurons

839 we only determined the R-squared for filters for the green
840 stimulus channel, and for ventral neurons for the UV stim-
841 ulus channel.

842 **SI-PCA model (2D).** The spatial convolutional filters of the
843 SI-PCA model were composed from PCA basis functions
844 (W). The model was trained to learn the weights of these
845 basis functions. The filters were produced by performing
846 PCA transformation on natural images recorded in mouse
847 habitats (38):

$$W = U^T \quad (4)$$

848 where U contains the eigenvectors of the covariance matrix
849 of the centered data in each column.

850 For example, when using 4 PCA bases, the shape of learn-
851 able weight matrix was 16x4 (channel number x basis
852 number), the shape of PCA bases was 4x2x1x9x9 (basis
853 number x chromatic channel x depth x image width x im-
854 age height), and the resulted spatial filter had the shape of
855 16x2x1x9x9. We varied the number of used basis (hyper-
856 parameter) and selected the one which achieved the best
857 performance on validation data (Suppl. Fig. S1b; Suppl.
858 Fig. S3b).

859 **SI-DCT model (2D).** For the SI-DCT model, its spatial con-
860 volutional filters were composed from DCT basis func-
861 tions, which were defined as:

$$F(u, v) = \alpha(u)\alpha(v) \cos\left[\frac{(2i+1)\pi}{2N}u\right] \cos\left[\frac{(2j+1)\pi}{2N}v\right] \quad (5)$$

$$\alpha(u) = \begin{cases} \sqrt{\frac{1}{N}} & u = 0 \\ \sqrt{\frac{2}{N}} & u \neq 0 \end{cases} \quad (6)$$

$$\alpha(v) = \begin{cases} \sqrt{\frac{1}{N}} & v = 0 \\ \sqrt{\frac{2}{N}} & v \neq 0 \end{cases} \quad (7)$$

862 where i and j denote pixel index of the input image (size
863 (N, N)); u and v denote DCT coefficient index of the DCT
864 filter. Here, we employed DCT basis functions for one-
865 channel gray images and thus used different bases for each
866 chromatic channel. For example, when using 4 DCT bases,
867 the shape of learnable weight matrix was 16x4x2 (channel
868 number x basis number x chromatic channel), the shape of
869 basis function was 4x1x9x9 (basis number x depth x image
870 width x image height), and the resulted spatial filter had
871 the shape of 16x2x1x9x9. Like for SI-PCA, we varied the
872 number of used basis and picked the one which achieved
873 the best performance on validation data (Suppl. Fig. S1b).

874 **Stand-alone EC model (2D).** We used a similar EC model
875 architecture (convolutional autoencoder) and loss function
876 as in (38). The model's encoder contained a single con-
877 volutional layer (with weights denoted \vec{w}_c) followed by a
878 rectified linear unit (ReLU) function, one fully-connected
879 (FC) layer, and another ReLU function. The decoder con-
880 tained one FC layer, one ReLU function, a single deconvo-
881 lutional layer (with weights denoted \vec{w}_d), and a hyperbolic

882 tangent (tanh) function to map back to the original data
883 range ($[-1, 1]$).

884 As a measure of reconstruction quality, we used mean
885 squared error (MSE; (37, 38)). Gaussian noise was added
886 to the encoder output for redundancy reduction (37, 81, 82)
887 and an L1 penalty (hyperparameter β) was imposed to its
888 activation (\vec{h}) for sparse readouts (37, 81, 83). We also ap-
889 plied L2 regularization on the convolutional and deconvo-
890 lutional layers to encourage the learning of smooth filters
891 (42, 84, 85). We used 16 9x9 convolutional and deconvo-
892 lutional filters. The activation tensor (16x28x28, out-
893 put channel x image width x image height) following the
894 first convolutional layer was flattened to a one-dimensional
895 vector with 12,544 inputs before feeding into the FC layer.
896 The loss function for the EC model was:

$$L_{EC} = \sum_i (\vec{x}_i - \hat{x}_i)^2 + \alpha(\|\vec{w}_c\|_2 + \|\vec{w}_d\|_2) + \beta\|\vec{h}\|_1 \quad (8)$$

897 where the first term is the MSE error between the predic-
898 tion \hat{x}_i and ground truth x_i with image index i , and the
899 next two terms denote the L2 and L1 penalties.

900 **Hybrid model (2D).** The hybrid (semi-supervised) model
901 consisted of a SI and an EC branch (for details on the two
902 models' architectures, see above). These branches were
903 trained simultaneously, sharing the spatial convolutional
904 filters (\vec{w}_{cs}). The total loss function of the hybrid model
905 was derived from the loss functions of the two branches as
906 follows:

$$L_{Hybrid} = wL_{SI} + (1-w)L_{EC} \quad (9)$$

$$L_{SI} = \left(\sum_i (\vec{r}_i - \hat{r}_i \log \hat{r}_i) + \alpha_1 \|\vec{w}_{cs}\|_2 + \alpha_2 \|\vec{w}_{ct}\|_2 / w \right. \\ \left. + \beta_1 \|\vec{w}_f\|_1 / w \right) / N_1 \quad (10)$$

$$L_{EC} = \left(\sum_j (\vec{x}_j - \hat{x}_j)^2 + \alpha_3 \|\vec{w}_{cs}\|_2 + \alpha_3 \|\vec{w}_d\|_2 / (1-w) \right. \\ \left. + \beta_2 \|\vec{h}\|_1 / (1-w) \right) / N_2 \quad (11)$$

907 Here, i and j denote neuron and image index, respectively;
908 N_1 and N_2 the number of neurons and images, respec-
909 tively. The weight (w , with $0 \leq w \leq 1$) controlled the
910 impact of each branch's loss function on the shared spa-
911 tial filters. Practically, we used $w = 10^{-8}$ for L_{SI} and
912 $w = (1 - 10^{-8})$ for L_{EC} when $w = 0$ and $w = 1$, respec-
913 tively. Note that we added w to the denominator of the last
914 two terms to maintain the same regularization for \vec{w}_{ct} and
915 \vec{w}_f in a stand-alone SI model when varying w . For L_{EC} ,
916 similar to L_{SI} , we added $(1-w)$ to the denominator of
917 the last two terms to keep the same regularization for \vec{w}_d
918 and \vec{h} in a stand-alone EC model when varying w . We
919 used different data to train the EC branch of the hybrid
920 model: natural images, phase-scrambled natural images

921 and noise. All hybrid models were trained for a maximum
922 of 100 epochs (Suppl. Fig. S1c,d); training was stopped
923 early when the validation loss started decreasing.

924 Tuning all hyperparameters jointly in a grid search was
925 computationally prohibitive. Hence, for the SI branch,
926 we varied the hyperparameters around those determined
927 for the stand-alone configuration ($\alpha_1 = 10, \alpha_2 = 10, \beta_1 =$
928 $1/16$; see above), while for the EC branch, we varied the
929 hyperparameters systematically around the values ($\alpha_3 =$
930 $10^3, \beta_2 = 1/16$) used in (38). To tune w , we devised a
931 linear search approach by normalizing the loss functions
932 (using N_1 and N_2).

933 After training the hybrid model, we estimated the spatio-
934 temporal RFs of all neurons using a gradient ascent algo-
935 rithm (6). We visualized the spatial and temporal compo-
936 nent of RFs using SVD (cf. Fig. 3b), and the magnitude of
937 the RF was indicated in the spatial component.

938 We trained 2D models using all training data (440 s) with a
939 learning rate of $\mu = 10^{-4}$. In case less data were used (i.e.,
940 to evaluate data efficiency), we kept all hyperparameters
941 the same as for the full data case but doubled the learning
942 rate. This was done because the stand-alone SI model and
943 the hybrid model could not reach the minimum of valida-
944 tion loss within 100 epochs (when less data were used).

945 3D models.

946 **Stand-alone SI model (3D).** The 3D SI model consisted
947 of one spatio-temporal convolutional layer (16x2x8x9x9,
948 output channels x input channels x depth x image width
949 x image height; depth varied with the frequency of noise
950 stimuli, $n=8$ and $n=30$ for 5-Hz and 30-Hz noise, respec-
951 tively), and — after flattening all dimension — one FC
952 layer (96x6,400, output channels x input channels; output
953 channel varied with cell numbers $n=96, 64$ or 427 for dif-
954 ferent data sets; see above), followed by an exponential
955 function. No padding was used. The loss function was
956 defined as:

$$L_{SI} = \sum_i (\hat{r}_i - \bar{r}_i \log \hat{r}_i) + \alpha \|\vec{w}_c\|_2 + \beta \|\vec{w}_f\|_1 \quad (12)$$

957 This equation differs from Equation () with respect to the
958 L2 penalty, which is here on the weights of the spatio-
959 temporal convolutional filters (\vec{w}_c) with hyperparameter α
960 for the second term. After performing a grid search for the
961 two hyperparameters, we picked $\alpha = 100, \beta = 1/4$ which
962 yielded the best performance on the validation data. After
963 training, we estimated and extracted the cells' spatial and
964 temporal RFs via SVD for visualization.

965 **SI-PCA model (3D).** For the 3D SI-PCA models, we applied
966 Equation () to the movie clips (2x8x9x9, chromatic chan-
967 nel x depth x image width x image height; depth varied
968 with the frequency of noise stimuli, $n=8$ and $n=30$ for 5-
969 Hz and 30-Hz noise, respectively). Like for 2D SI-PCA
970 models, we varied the number of used bases and picked
971 the number for which the model achieved the best perfor-
972 mance on the validation data (Suppl. Fig. S3a).

973 **Stand-alone EC model (3D).** The 3D EC models used a se-
974 quence of frames from a movie clip as input and featured
975 3D spatio-temporal convolutional layers (with weights de-
976 noted \vec{w}_c) in the encoder. The decoder contained deconvol-
977 utional layers with weights \vec{w}_d . In the past-encoding case,
978 we fed an 8-frame clip (frames at $t-7$ to t) to the model
979 and aimed at reconstructing the 7th frame (at $t-1$). In the
980 future-prediction case, the goal was to predict the 8th frame
981 (at t) with the input being the first 7 frames ($t-7$ to $t-1$)
982 of the clip. The loss functions was similar to that given
983 by Equation () except that (i) \vec{w}_c features different a shape
984 (16x2x8x9x9, output channel x chromatic channel x filter
985 depth x filter width x filter height), and (ii) x_i denotes the
986 7th frame for the past encoding and the 8th frame for the
987 future prediction model (Suppl. Fig. S2b,c,d).

988 **Hybrid model (3D).** The 3D hybrid models consisted of a
989 SI branch and an EC branch with shared spatio-temporal
990 convolutional filters (\vec{w}_c ; see above). Like for the 2D hy-
991 brid models, the total loss function was a weighted sum of
992 losses for the two branches as follows:

$$L_{Hybrid} = wL_{SI} + (1-w)L_{EC} \quad (13)$$

$$L_{SI} = \left(\sum_i (\hat{r}_i - \bar{r}_i \log \hat{r}_i) + \alpha_1 \|\vec{w}_c\|_2 + \beta_1 \|\vec{w}_f\|_1 / w \right) / N_1 \quad (14)$$

$$L_{EC} = \left(\sum_j (\hat{x}_j - \bar{x}_j)^2 + \alpha_2 \|\vec{w}_c\|_2 + \alpha_2 \|\vec{w}_d\|_2 / (1-w) + \beta_2 \|\vec{h}\|_1 / (1-w) \right) / N_2 \quad (15)$$

993 Here, i denotes neuron index, j movie clip index, N_1 neu-
994 ron number, and N_2 the number of movie clips. Again,
995 instead of tuning all hyperparameters jointly via a grid
996 search, we varied the hyperparameters around the val-
997 ues determined for the stand-alone SI configuration ($\alpha_1 =$
998 $100, \beta_1 = 1/4$) for the SI branch. For the EC branch,
999 we varied the hyperparameters systematically around the
1000 values ($\alpha_2 = 10^4, \beta_2 = 1/16$) used in the stand-alone EC
1001 models. We then tuned w linearly after normalizing the
1002 loss functions (using N_1 and N_2). We also visualized the
1003 spatial and temporal RF components using SVD (Fig. 5a,
1004 bottom).

1005 Acknowledgments

1006 We thank Dylan Paiton and Katrin Franke for helpful
1007 discussions, and Merle Harrer for excellent technical as-
1008 sistance. This work was supported by the German Re-
1009 search Foundation (DFG): SFB 1233, Robust Vision: In-
1010 ference Principles and Neural Mechanisms, projects 10
1011 and 12, project number: 276693517; and under Ger-
1012 many's Excellence Strategy EXC 2064/1 (project number
1013 390727645); the European Union's Horizon 2020 research
1014 and innovation programme under the Marie Skłodowska-
1015 Curie grant (agreement No 674901); the Max Planck So-
1016 ciety (M.FE.A.KYBE0004); and the German Ministry of

1017 Education and Research (BMBF; FKZ: 01GQ1002), and
1018 the Tübingen AI Center (FKZ: 01IS18039A). The funders
1019 had no role in study design, data collection and analysis,
1020 decision to publish, or preparation of the manuscript.

1021 Author Contributions

1022 Conceptualization: Y.Q.; Methodology: Y.Q., D.K., K.S.,
1023 D.G., L.H., T.S., and T.E.; Data acquisition & curation:
1024 K.S.; Formal analysis: Y.Q. with input from D.K., M.B.,
1025 L.B., and T.E.; Investigation: Y.Q. with input from D.K.,
1026 K.S., L.B., M.B., and T.E.; Writing – original Draft: Y.Q.,
1027 D.K., L.B., and T.E.; Writing – review & editing: all au-
1028 thors; Visualization: Y.Q.; D.G. (confusion matrix); Soft-
1029 ware: Y.Q.; L.H. and D.G. (classifier); Resources: T.S. and
1030 T.E.; Supervision: M.B., L.B., and T.E.; Funding acquisi-
1031 tion: L.B., M.B., and T.E.

1032 Declaration of Interests

1033 The authors declare no competing interests.

1034 Data and Code Availability

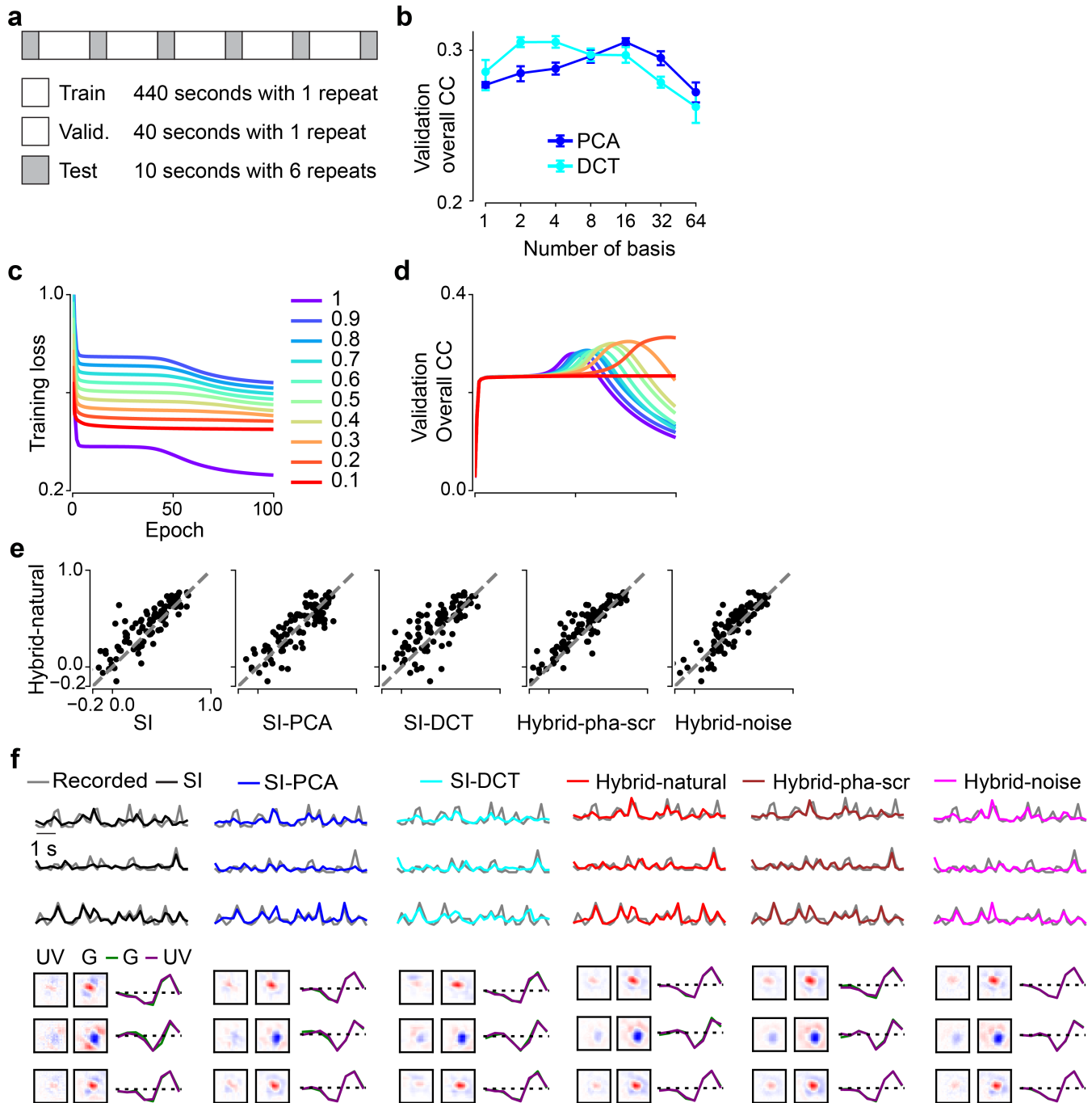
1035 Data and code would be available upon publication.

1036 Bibliography

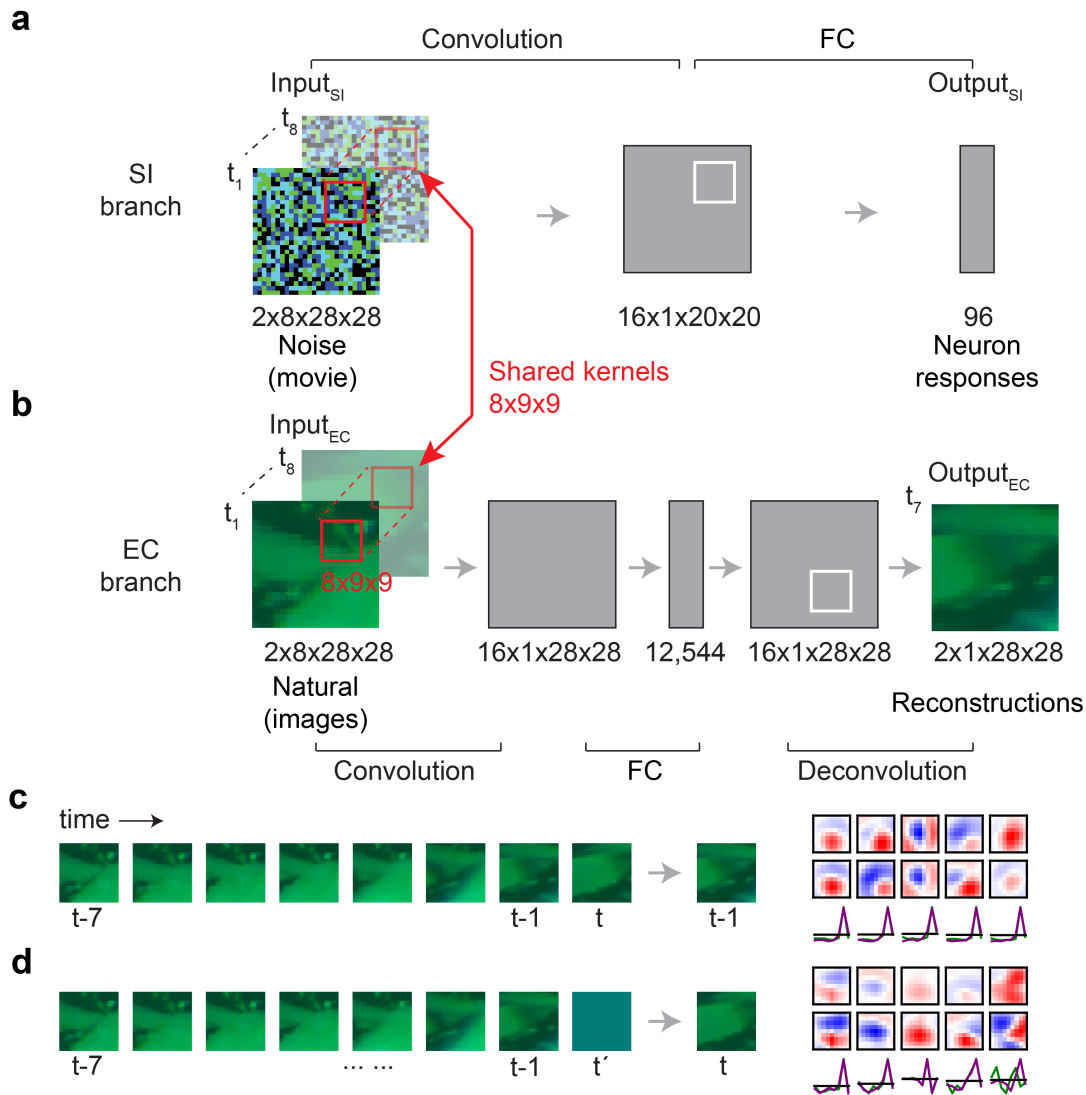
- 1037 1. Ian H Stevenson and Konrad P Kording. How advances in neural recording affect data
1038 analysis. *Nature neuroscience*, 14(2):139–142, 2011.
- 1039 2. EJ Chichilnisky. A simple white noise analysis of neuronal light responses. *Network:
1040 Computation in Neural Systems*, 12(2):199–213, 2001.
- 1041 3. Michael C-K Wu, Stephen V David, and Jack L Gallant. Complete functional character-
1042 ization of sensory neurons by system identification. *Annu. Rev. Neurosci.*, 29:477–505,
1043 2006.
- 1044 4. Jonathan W Pillow, Jonathon Shlens, Liam Paninski, Alexander Sher, Alan M Litke,
1045 EJ Chichilnisky, and Eero P Simoncelli. Spatio-temporal correlations and visual sig-
1046 nalling in a complete neuronal population. *Nature*, 454(7207):995–999, 2008.
- 1047 5. Vasilis Marmarelis. *Analysis of physiological systems: The white-noise approach*.
1048 Springer Science & Business Media, 2012.
- 1049 6. Melinda E Koelling and Duane Q Nykamp. Computing linear approximations to non-
1050 linear neuronal response. *Network: Computation in Neural Systems*, 19(4):286–313,
1051 2008.
- 1052 7. Tim Gollisch and Markus Meister. Eye smarter than scientists believed: neural compu-
1053 tations in circuits of the retina. *Neuron*, 65(2):150–164, 2010.
- 1054 8. Esteban Real, Hiroki Asari, Tim Gollisch, and Markus Meister. Neural circuit inference
1055 from function to structure. *Current Biology*, 27(2):189–198, 2017.
- 1056 9. Ben Wilmore, Ryan J Prenger, Michael C-K Wu, and Jack L Gallant. The berkeley
1057 wavelet transform: a biologically inspired orthogonal wavelet transform. *Neural compu-
1058 tation*, 20(6):1537–1564, 2008.
- 1059 10. Niru Maheswaranathan, David B Kastner, Stephen A Baccus, and Surya Ganguli. In-
1060 ferring hidden structure in multilayered neural circuits. *PLoS computational biology*, 14
1061 (8):e1006291, 2018.
- 1062 11. Daniel LK Yamins, Ha Hong, Charles F Cadieu, Ethan A Solomon, Darren Seibert, and
1063 James J DiCarlo. Performance-optimized hierarchical models predict neural responses
1064 in higher visual cortex. *Proceedings of the National Academy of Sciences*, 111(23):
1065 8619–8624, 2014.
- 1066 12. Umut Güçlü and Marcel AJ van Gerven. Deep neural networks reveal a gradient in
1067 the complexity of neural representations across the ventral stream. *Journal of Neuro-
1068 science*, 35(27):10005–10014, 2015.
- 1069 13. Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553):
1070 436–444, 2015.
- 1071 14. Demis Hassabis, Dharshan Kumaran, Christopher Summerfield, and Matthew
1072 Botvinick. Neuroscience-inspired artificial intelligence. *Neuron*, 95(2):245–258, 2017.
- 1073 15. Maxwell H Turner, Luis Gonzalo Sanchez Giraldo, Odelia Schwartz, and Fred Rieke.
1074 Stimulus- and goal-oriented frameworks for understanding natural vision. *Nature neuro-
1075 science*, 22(1):15–24, 2019.
- 1076 16. Blake A Richards, Timothy P Lillicrap, Philippe Beaudoin, Yoshua Bengio, Rafal Bogacz,
1077 Amelia Christensen, Claudia Clopath, Rui Ponte Costa, Archy de Berker, Surya Ganguli,
1078 et al. A deep learning framework for neuroscience. *Nature neuroscience*, 22(11):1761–
1079 1770, 2019.
- 1080 17. Daniel LK Yamins and James J DiCarlo. Using goal-driven deep learning models to
1081 understand sensory cortex. *Nature neuroscience*, 19(3):356–365, 2016.

- 1082 18. Lane McIntosh, Niru Maheswaranathan, Aran Nayebi, Surya Ganguli, and Stephen Bac-
1083 cus. Deep learning models of the retinal response to natural scenes. *Advances in neural
1084 information processing systems*, 29:1369–1377, 2016.
- 1085 19. David Klindt, Alexander S Ecker, Thomas Euler, and Matthias Bethge. Neural system
1086 identification for large populations separating “what” and “where”. In *Advances in Neural
1087 Information Processing Systems*, pages 3506–3516, 2017.
- 1088 20. Pouya Bashivan, Kohitij Kar, and James J DiCarlo. Neural population control via deep
1089 image synthesis. *Science*, 364(6439), 2019.
- 1090 21. Carlos R Ponce, Will Xiao, Peter F Schade, Till S Hartmann, Gabriel Kreiman, and
1091 Margaret S Livingstone. Evolving images for visual neurons using a deep generative
1092 network reveals coding principles and neuronal preferences. *Cell*, 177(4):999–1009,
1093 2019.
- 1094 22. Edgar Y Walker, Fabian H Sinz, Erick Cobos, Taliah Muhammad, Emmanouil
1095 Froudarakis, Paul G Fahey, Alexander S Ecker, Jacob Reimer, Xaq Pitkow, and An-
1096 dreas S Toulas. Inception loops discover what excites neurons most using deep predic-
1097 tive models. *Nature neuroscience*, 22(12):2060–2065, 2019.
- 1098 23. Tom Baden, Thomas Euler, and Philipp Berens. Understanding the retinal basis of
1099 vision across species. *Nature Reviews Neuroscience*, 21(1):5–20, 2020.
- 1100 24. Horace B Barlow et al. Possible principles underlying the transformation of sensory
1101 messages. *Sensory communication*, 1(01), 1961.
- 1102 25. Eero P Simoncelli and Bruno A Olshausen. Natural image statistics and neural repre-
1103 sentation. *Annual review of neuroscience*, 24(1):1193–1216, 2001.
- 1104 26. Eugene Switkes, Melanie J Mayer, and Jeffrey A Sloan. Spatial frequency analysis of
1105 the visual environment: Anisotropy and the carpentered environment hypothesis. *Vision
1106 research*, 18(10):1393–1399, 1978.
- 1107 27. Xiangmin Xu, Christine E Collins, Ilya Khaytin, Jon H Kaas, and Vivien A Casagrande.
1108 Unequal representation of cardinal vs. oblique orientations in the middle temporal visual
1109 area. *Proceedings of the National Academy of Sciences*, 103(46):17490–17495, 2006.
- 1110 28. Ahna R Girshick, Michael S Landy, and Eero P Simoncelli. Cardinal rules: visual orien-
1111 tation perception reflects knowledge of environmental statistics. *Nature neuroscience*,
1112 14(7):926–932, 2011.
- 1113 29. Simon Laughlin. A simple coding procedure enhances a neuron’s information capacity.
1114 *Zeitschrift für Naturforschung c*, 36(9-10):910–912, 1981.
- 1115 30. J Hans van Hateren and Dan L Ruderman. Independent component analysis of natural
1116 image sequences yields spatio-temporal filters similar to simple cells in primary visual
1117 cortex. *Proceedings of the Royal Society of London. Series B: Biological Sciences*, 265
1118 (1412):2315–2320, 1998.
- 1119 31. Suva Roy, Na Young Jun, Emily L Davis, John Pearson, and Greg D Field. Inter-mosaic
1120 coordination of retinal receptive fields. *Nature*, 592(7854):409–413, 2021.
- 1121 32. Joseph J Atick and A Norman Redlich. Towards a theory of early visual processing.
1122 *Neural computation*, 2(3):308–320, 1990.
- 1123 33. Joseph J Atick. Could information theory provide an ecological theory of sensory pro-
1124 cessing? *Network: Computation in neural systems*, 3(2):213–251, 1992.
- 1125 34. Christina Enroth-Cugell and John G Robson. The contrast sensitivity of retinal ganglion
1126 cells of the cat. *The Journal of physiology*, 187(3):517–552, 1966.
- 1127 35. Dana H Ballard. Modular learning in neural networks. In *AAAI*, pages 279–284, 1987.
- 1128 36. Geoffrey E Hinton and Ruslan R Salakhutdinov. Reducing the dimensionality of data
1129 with neural networks. *science*, 313(5786):504–507, 2006.
- 1130 37. Samuel Ocko, Jack Lindsey, Surya Ganguli, and Stephane Deny. The emergence of
1131 multiple retinal cell types through efficient coding of natural movies. In *Advances in
1132 Neural Information Processing Systems*, pages 9389–9400, 2018.
- 1133 38. Yongrong Qiu, Zhijian Zhao, David Klindt, Magdalena Kautzky, Klaudia P Szatko, Frank
1134 Schaeffel, Katharina Rifai, Katrin Franke, Laura Busse, and Thomas Euler. Natural
1135 environment statistics in the upper and lower visual field are reflected in mouse retinal
1136 specializations. *Current Biology*, 2021.
- 1137 39. Dylan M Paiton, Charles G Frye, Sheng Y Lundquist, Joel D Bowen, Ryan Zarccone, and
1138 Bruno A Olshausen. Selectivity and robustness of sparse coding networks. *Journal of
1139 Vision*, 20(12):10–10, 2020.
- 1140 40. Jan Eichhorn, Fabian Sinz, and Matthias Bethge. Natural image coding in v1: how much
1141 use is orientation selectivity? *PLoS computational biology*, 5(4):e1000336, 2009.
- 1142 41. Wiktor Młynarski, Michal Hledik, Thomas R Sokolowski, and Gašper Tkačik. Statistical
1143 analysis and optimality of neural systems. *Neuron*, 109(7):1227–1241, 2021.
- 1144 42. Benjamin T Vincent and Roland J Baddeley. Synaptic energy efficiency in retinal pro-
1145 cessing. *Vision research*, 43(11):1285–1292, 2003.
- 1146 43. Tom Baden, Philipp Berens, Katrin Franke, Miroslav Román Rosón, Matthias Bethge,
1147 and Thomas Euler. The functional diversity of retinal ganglion cells in the mouse. *Nature*,
1148 529(7586):345–350, 2016.
- 1149 44. Klaudia P Szatko, Maria M Korympidou, Yanli Ran, Philipp Berens, Deniz Dalkara, Timm
1150 Schubert, Thomas Euler, and Katrin Franke. Neural circuits in the mouse retina support
1151 color vision in the upper visual field. *Nature communications*, 11(1):1–14, 2020.
- 1152 45. Cassandra L Schlamp, Angela D Montgomery, Caitlin E Mac Nair, Claudia Schuartz,
1153 Daniel J Willmer, and Robert W Nickells. Evaluation of the percentage of ganglion cells
1154 in the ganglion cell layer of the rodent retina. *Molecular vision*, 19:1387, 2013.
- 1155 46. Gerald H Jacobs, Gary A Williams, and John A Fenwick. Influence of cone pigment
1156 coexpression on spectral sensitivity and color vision in the mouse. *Vision research*, 44
1157 (14):1615–1622, 2004.
- 1158 47. Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional net-
1159 works. In *European conference on computer vision*, pages 818–833. Springer, 2014.
- 1160 48. Katrin Franke, Philipp Berens, Timm Schubert, Matthias Bethge, Thomas Euler, and
1161 Tom Baden. Inhibition decorrelates visual feature representations in the inner retina.
1162 *Nature*, 542(7642):439–444, 2017.
- 1163 49. Robert E Soodak. Two-dimensional modeling of visual receptive fields using gaussian
1164 subunits. *Proceedings of the National Academy of Sciences*, 83(23):9259–9263, 1986.
- 1165 50. Rajesh PN Rao and Dana H Ballard. Predictive coding in the visual cortex: a functional
1166 interpretation of some extra-classical receptive-field effects. *Nature neuroscience*, 2(1):
1167 79–87, 1999.

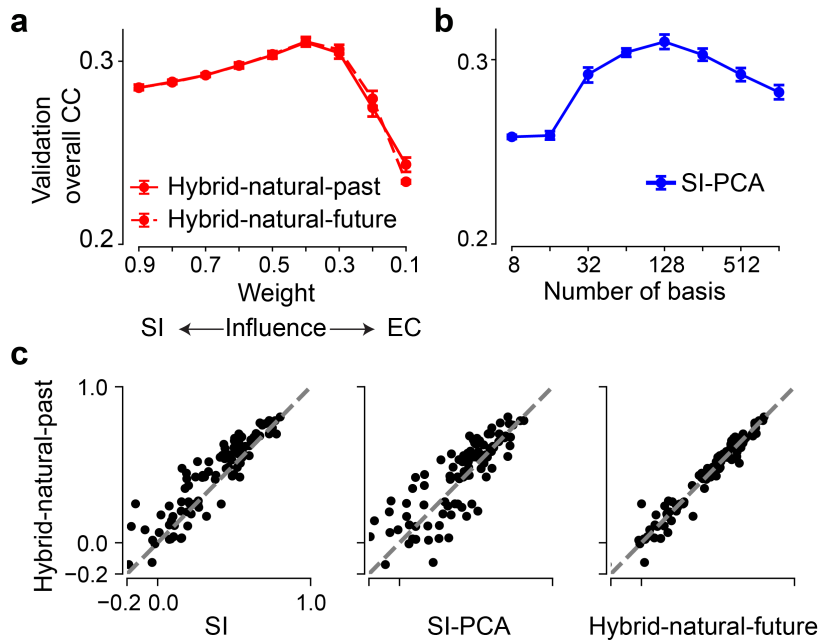
- 1168 51. Toshihiko Hosoya, Stephen A Baccus, and Markus Meister. Dynamic predictive coding
1169 by the retina. *Nature*, 436(7047):71–77, 2005.
- 1170 52. Jamie Johnston, Sofie-Helene Seibel, Léa Simone Adele Darnet, Sabine Renninger,
1171 Michael Orger, and Leon Lagnado. A retinal circuit generating a dynamic predictive
1172 code for oriented features. *Neuron*, 102(6):1211–1222, 2019.
- 1173 53. J Hans van Hateren. Real and optimal neural images in early vision. *Nature*, 360(6399):
1174 68–70, 1992.
- 1175 54. Joseph J Atick and A Norman Redlich. What does the retina know about natural
1176 scenes? *Neural computation*, 4(2):196–210, 1992.
- 1177 55. Matthew Chalk, Olivier Marre, and Gašper Tkačik. Toward a unified theory of efficient,
1178 predictive, and sparse coding. *Proceedings of the National Academy of Sciences*, 115
1179 (1):186–191, 2018.
- 1180 56. Jerome Y Lettvin, Humberto R Maturana, Warren S McCulloch, and Walter H Pitts.
1181 What the frog's eye tells the frog's brain. *Proceedings of the IRE*, 47(11):1940–1951,
1182 1959.
- 1183 57. J Alexander Bae, Shang Mu, Jinseop S Kim, Nicholas L Turner, Ignacio Tartavull, Nico
1184 Kemnitz, Chris S Jordan, Alex D Norton, William M Silversmith, Rachel Prentki, et al.
1185 Digital museum of retinal ganglion cells with dense anatomy and physiology. *Cell*, 173
1186 (5):1293–1306, 2018.
- 1187 58. Nicholas M Tran, Karthik Shekhar, Irene E Whitney, Anne Jacobi, Inbal Benhar, Gu-
1188 osong Hong, Wenjun Yan, Xian Adiconis, McKinzie E Arnold, Jung Min Lee, et al.
1189 Single-cell profiles of retinal ganglion cells differing in resilience to injury reveal neuro-
1190 protective genes. *Neuron*, 104(6):1039–1055, 2019.
- 1191 59. Jillian Goetz, Zachary F Jessen, Anne Jacobi, Adam Mani, Sam Cooler, Devon Greer,
1192 Sabah Kadri, Jeremy Segal, Karthik Shekhar, Joshua Sanes, et al. Unified classifica-
1193 tion of mouse retinal ganglion cells using function, morphology, and gene expression.
1194 *Morphology, and Gene Expression*, 2021.
- 1195 60. Horace B Barlow and Richard M Hill. Selective sensitivity to direction of movement in
1196 ganglion cells of the rabbit retina. *Science*, 139(3553):412–412, 1963.
- 1197 61. Fabian H Sinz, Alexander S Ecker, Paul G Fahey, Edgar Y Walker, Erick Cobos, Em-
1198 manouil Froudarakis, Dimitri Yatsenko, Xaq Pitkow, Jacob Reimer, and Andreas S Tol-
1199 lias. Stimulus domain transfer in recurrent models for large scale cortical population
1200 prediction on video. *BioRxiv*, page 452672, 2018.
- 1201 62. Jon Touryan, Gidon Felsen, and Yang Dan. Spatial structure of complex cell receptive
1202 fields measured with natural images. *Neuron*, 45(5):781–791, 2005.
- 1203 63. Andrew Saxe, Stephanie Nelli, and Christopher Summerfield. If deep learning is the
1204 answer, what is the question? *Nature Reviews Neuroscience*, 22(1):55–67, 2021.
- 1205 64. Alexander Heitman, Nora Brackbill, Martin Greschner, Alexander Sher, Alan M Litke,
1206 and EJ Chichilnisky. Testing pseudo-linear models of responses to natural scenes in
1207 primate retina. *bioRxiv*, page 045336, 2016.
- 1208 65. Nicole C Rust and J Anthony Movshon. In praise of artifice. *Nature neuroscience*, 8
1209 (12):1647–1650, 2005.
- 1210 66. Yifeng Zhang, In-Jung Kim, Joshua R Sanes, and Markus Meister. The most numerous
1211 ganglion cell type of the mouse retina is a selective feature detector. *Proceedings of the
1212 National Academy of Sciences*, 109(36):E2391–E2398, 2012.
- 1213 67. Horace B Barlow. Summation and inhibition in the frog's retina. *The Journal of physiol-
1214 ogy*, 119(1):69–88, 1953.
- 1215 68. Sophie Deneve and Matthew Chalk. Efficiency turns the table on neural encoding,
1216 decoding and noise. *Current Opinion in Neurobiology*, 37:141–148, 2016.
- 1217 69. Michael Teti, Emily Meyer, and Garrett Kenyon. Can lateral inhibition for sparse coding
1218 help explain v1 neuronal responses to natural stimuli? In *2020 IEEE Southwest Sym-
1219 posium on Image Analysis and Interpretation (SSIAI)*, pages 120–124. IEEE, 2020.
- 1220 70. Horace Barlow. Redundancy reduction revisited. *Network: computation in neural sys-
1221 tems*, 12(3):241, 2001.
- 1222 71. Brett Vintch, J Anthony Movshon, and Eero P Simoncelli. A convolutional subunit model
1223 for neuronal responses in macaque v1. *Journal of Neuroscience*, 35(44):14829–14841,
1224 2015.
- 1225 72. Santiago A Cadena, George H Denfield, Edgar Y Walker, Leon A Gatys, Andreas S
1226 Tolias, Matthias Bethge, and Alexander S Ecker. Deep convolutional models improve
1227 predictions of macaque v1 responses to natural images. *PLoS computational biology*,
1228 15(4):e1006897, 2019.
- 1229 73. Kevin L Briggman and Thomas Euler. Bulk electroporation and population calcium
1230 imaging in the adult mammalian retina. *Journal of neurophysiology*, 105(5):2601–2609,
1231 2011.
- 1232 74. Thomas Euler, Susanne E Hausselt, David J Margolis, Tobias Breuning, Xavier
1233 Castell, Peter B Detwiler, and Winfried Denk. Eyecup scope—optical recordings of light
1234 stimulus-evoked fluorescence signals in the retina. *Pflügers Archiv-European Journal
1235 of Physiology*, 457(6):1393–1414, 2009.
- 1236 75. Thomas Euler, Katrin Franke, and Tom Baden. Studying a light sensor with light: mul-
1237 tiphoton imaging in the retina. In *Multiphoton Microscopy*, pages 225–250. Springer,
1238 2019.
- 1239 76. Katrin Franke, André Maia Chagas, Zhijian Zhao, Maxime JY Zimmermann, Philipp
1240 Bartel, Yongrong Qiu, Klaudia P Szatko, Tom Baden, and Thomas Euler. An arbitrary-
1241 spectrum spatial visual stimulator for vision research. *elife*, 8:e48779, 2019.
- 1242 77. Leo Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001.
- 1243 78. F Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blon-
1244 del, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau,
1245 M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python.
1246 *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- 1247 79. Lin Sun, Kui Jia, Dit-Yan Yeung, and Bertram E Shi. Human action recognition using
1248 factorized spatio-temporal convolutional networks. In *Proceedings of the IEEE interna-
1249 tional conference on computer vision*, pages 4597–4605, 2015.
- 1250 80. Du Tran, Heng Wang, Lorenzo Torresani, Jamie Ray, Yann LeCun, and Manohar Paluri.
1251 A closer look at spatiotemporal convolutions for action recognition. In *Proceedings of
1252 the IEEE conference on Computer Vision and Pattern Recognition*, pages 6450–6459,
1253 2018.
- 1254 81. Eizaburo Doi and Michael S Lewicki. A theory of retinal population coding. *Advances in
1255 neural information processing systems*, 19:353, 2007.
- 1256 82. MCW Van Rossum, Brendan J O'Brien, and Robert G Smith. Effects of noise on the
1257 spike timing precision of retinal ganglion cells. *Journal of neurophysiology*, 89(5):2406–
1258 2419, 2003.
- 1259 83. David J Field. What is the goal of sensory coding? *Neural computation*, 6(4):559–601,
1260 1994.
- 1261 84. David H Hubel and Torsten N Wiesel. Receptive fields of single neurones in the cat's
1262 striate cortex. *The Journal of physiology*, 148(3):574–591, 1959.
- 1263 85. David Marr and Ellen Hildreth. Theory of edge detection. *Proceedings of the Royal
1264 Society of London. Series B. Biological Sciences*, 207(1167):187–217, 1980.



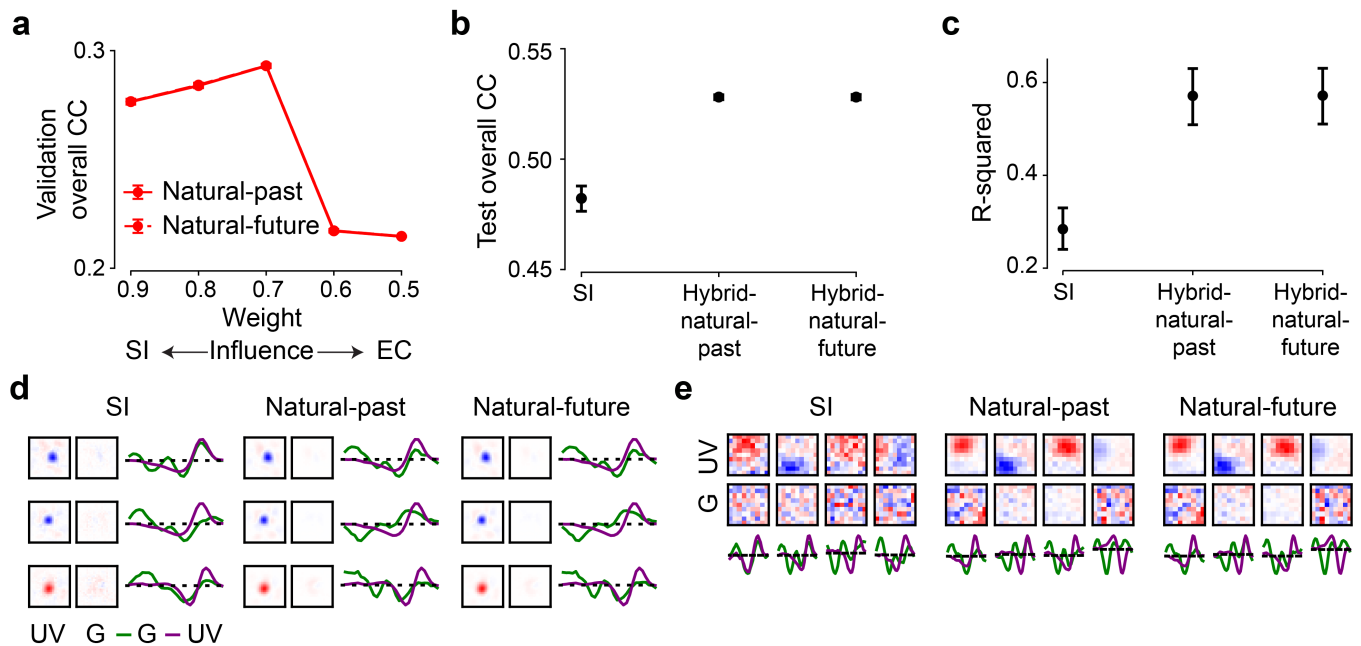
Supplemental Fig. S1. Training of 2D models. **a.** The noise stimulus (9 minutes in total) containing training and validation data (1 repeat) and test data (6 repeats). **b.** Model performance (mean) based on validation data for SI-PCA and SI-DCT with different numbers of basis. SI-PCA and SI-DCT yielded best performance when using 16 and 4 bases, respectively (each model for $n=10$ random seeds; error bars represent 2.5 and 97.5 percentiles with bootstrapping). **c.** Training loss as a function of training epochs for the hybrid model (Input_{EC}, natural scenes) with different weights (w), indicated by color (right). As weight decreased from 1 to 0.2, more training epochs were needed to reach the best performance. The hybrid model performed best for $w = 0.2$. Note that the hybrid model showed a slower change in correlation coefficient (CC) around the peak at $w = 0.2$ (compared to $w = 1$), demonstrating the regularization effects of the EC branch on the hybrid model. **d.** Model performance based on validation data (with linear correlation coefficient as metric) during the hybrid-natural model training with different weights (colors as in (c)). As weight decreased from 1 to 0.2, more training epochs were needed to reach the best performance. The hybrid model performed best for $w = 0.2$. Note that the hybrid model showed a slower change in correlation coefficient (CC) around the peak at $w = 0.2$ (compared to $w = 1$), demonstrating the regularization effects of the EC branch on the hybrid model. **e.** Scatter plots for model predictions based on test data at a particular seed (each dot representing one neuron). Hybrid with natural scenes as input_{EC} ($w = 0.2$) vs. SI, SI with PCA basis (16 bases), SI with DCT basis (4 bases), hybrid-pha-scr ($w = 0.3$) and hybrid-noise ($w = 0.4$). **f.** Upper: Three representative GCL cell responses (gray traces) to noise stimulus together with predictions of the best performing models on test data (black, SI; blue, SI with PCA basis; cyan, SI with DCT basis; red, hybrid w/ natural scenes as input in EC path; brown, hybrid w/ phase-scrambled scenes as input in EC path; magenta, hybrid w/ noise as input in EC path). Lower: Learned spatio-temporal RFs of the example cells, visualized by SVD. Same random seed as in (e).



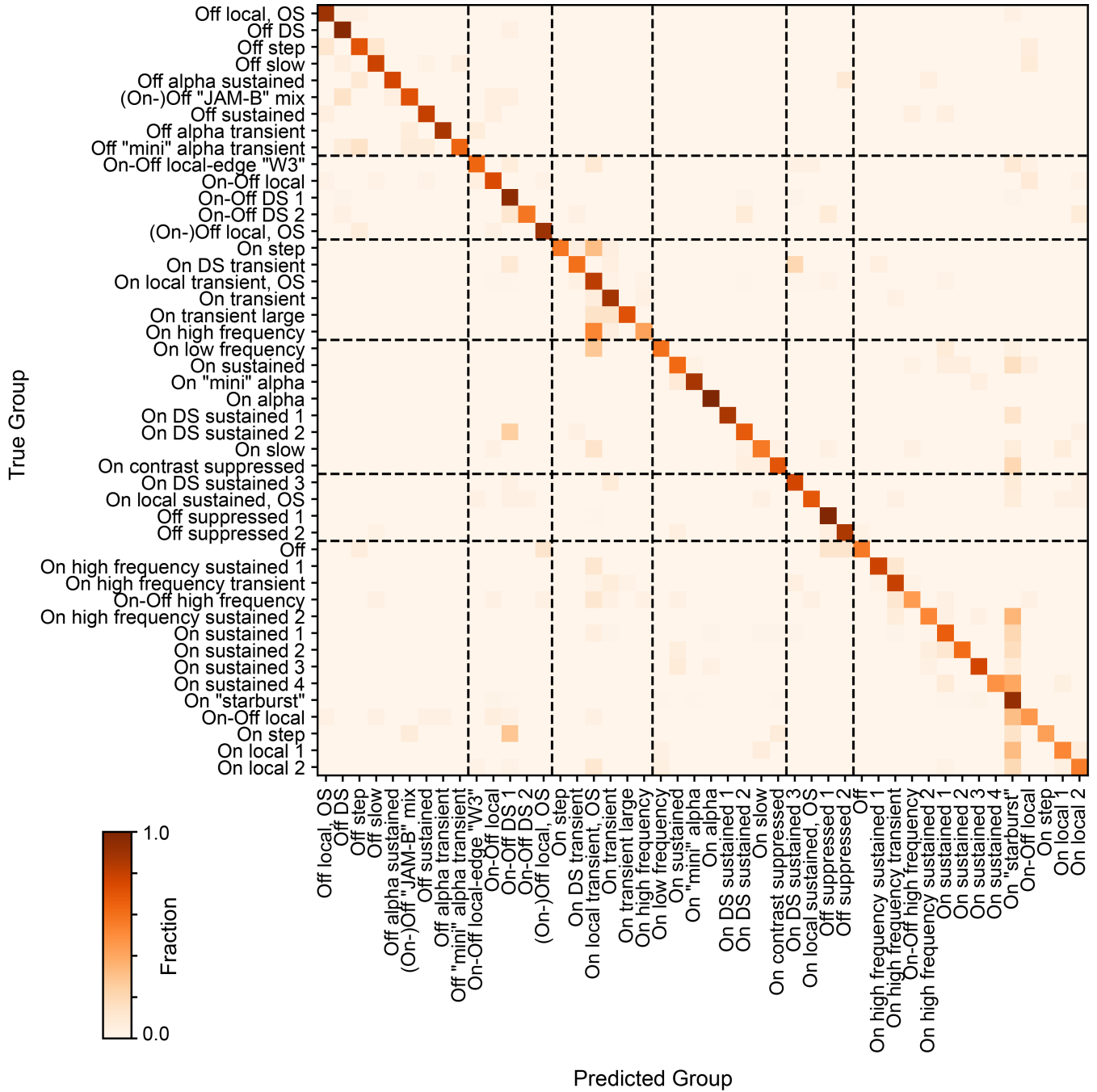
Supplemental Fig. S 2. Three-dimensional hybrid networks embedding natural movies. **a,b.** Illustration of SI network (a) with 3D spatio-temporal convolutional filter, and EC network (b), reconstructing the 7th frame (at $t - 1$) based on 8 continuous frames ($t - 7$ to t ; encoding the past, c). Combined as a hybrid network, the two branches were trained in parallel with shared 3D filters (Input_{EC}, 8-frame UV-green movie clip; Output_{EC}, reconstruction of the 7th frame of Input_{EC}). **c.** Example for input/output of the EC model for encoding the past (left; also see b) and exemplary spatio-temporal convolutional filters when using natural movies as input to train the EC model alone (right). **d.** Example for input/output of the EC model for predicting the future, i.e., predicting the 8th frame from the first 7 frames ($t - 7$ to $t - 1$) of the clip, and exemplary spatio-temporal filters when using natural movies as input to train the EC model alone. During preprocessing, the 8th frame of input was set to the mean of the first 7 frames, for UV and green channel, respectively. Note that for stand-alone EC models, all temporal components of filters for past encoding were very similar while those for future prediction were much more diverse.



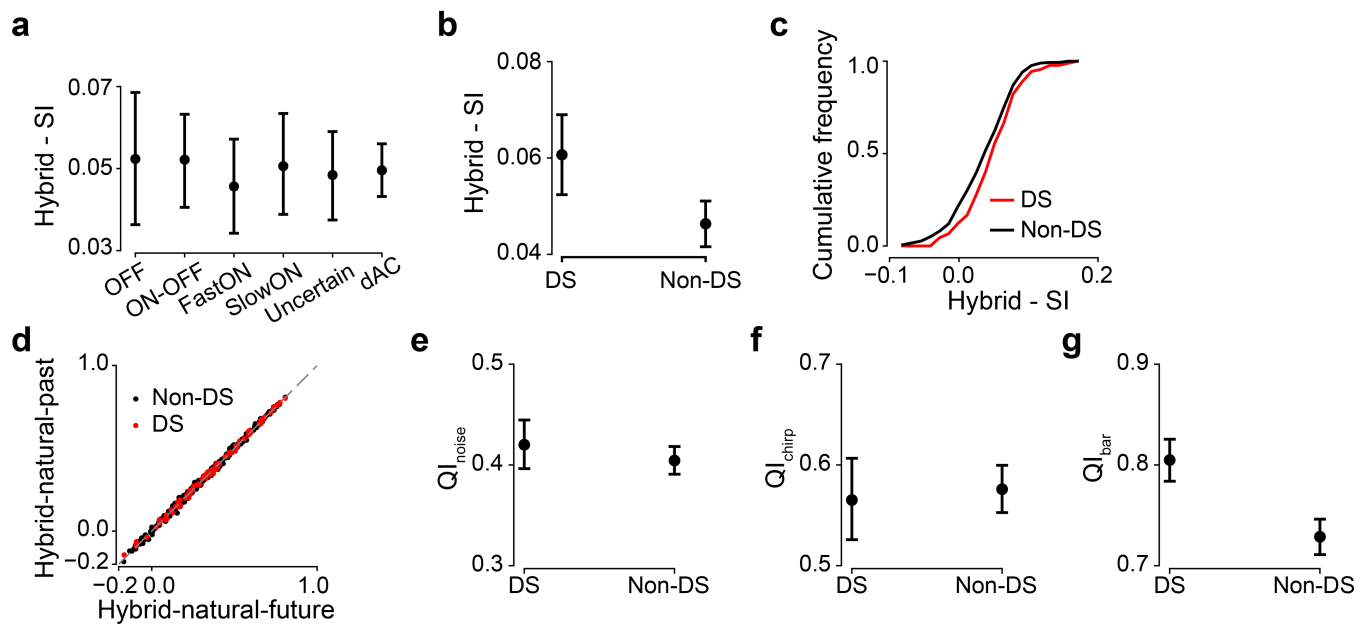
Supplemental Fig. S3. Training of 3D hybrid models. a,b. Model performance (mean) based on validation data for hybrid models w/ natural movies as input_{EC} (a), applying past encoding (hybrid-natural-past) or future prediction (hybrid-natural-future) and for different weights, and for the SI-PCA model (b) with different numbers of basis (each model for n=10 random seeds). c. Scatter plots for model predictions based on test data at a particular seed (each dot representing one neuron). hybrid-natural-past ($w = 0.4$) vs. SI, SI-PCA (128 PCA bases) and hybrid-natural-future ($w = 0.4$). Error bars in (a)–(b) represent 2.5 and 97.5 percentiles with bootstrapping. Both 3D hybrid models performed similarly, with a peak in predictive performance on the validation data at around $w = 0.4$ (a). This value of w was higher than for the 2D hybrid models ($w = 0.2$; cf. Fig. 3c). We also examined the low-pass filtering effects on the 3D SI model by using PCA filters (3D SI-PCA) and varying the number of basis (b). Like for the 2D case when varying the number of basis, we found a maximum in performance on the validation data at 128 bases, which was larger than the 16 bases in the 2D case (cf. Suppl. Fig. S1b).



Supplemental Fig. S4. Hybrid model for encoding neuronal responses to 30-Hz dense noise. To test hybrid models for different stimuli, we recorded neuronal responses to the 30-Hz dense noise in the ventral retina. We yielded n=64 neurons after quality control (Methods), which were used to train the SI and hybrid networks. a. Model performance (mean) based on validation data for hybrid models (w/ natural movies as input_{EC}), applying encoding-past (hybrid-natural-past) or predicting-future (hybrid-natural-future) and for different weights. Each model for n=10 random seeds. Both models with similar performance for all weights, peaking at $w = 0.7$. b. Model performance (mean) based on test data for SI, hybrid-natural-past ($w = 0.7$) and hybrid-natural-future ($w = 0.7$). Each model for n=10 random seeds. The two hybrid models had better performance with smaller standard deviation compared the SI model ($p < 0.0001$ for SI and hybrid-natural-past, $p = 0.9992$ for hybrid-natural-past and hybrid-natural-future; two-sided permutation test, n=10,000 repeats). c. R-squared (mean) of fitting a 2D Gaussian to all the spatial filters in UV stimulus channel (each model for n=10 random seeds; $p < 0.0001$ for SI and hybrid-natural-past, $p = 0.9888$ for hybrid-natural-past and hybrid-natural-future; two-sided permutation test, n=10,000 repeats). d. Learned spatio-temporal filters of the three representative cells, visualized by SVD. Note that because all neurons in this data set were recorded in the ventral retina, their responses were dominated by the UV channel. Different temporal filters in the UV channel were observed for these neurons (cf. the very similar temporal filters in the green channel for neurons' responses to 5-Hz noise in Fig. 3b, Fig. 5a lower). e. Exemplary shared spatial and temporal filters of 3D models, visualized by SVD and for one random seed. Temporal: UV and green channels indicated by purple and green lines, respectively. Error bars in (a)–(c) represent 2.5 and 97.5 percentiles with bootstrapping.



Supplemental Fig. S 5. Confusion matrix for a trained random forest classifier. Normalized confusion matrix (true cell types against predicted cell types) for a trained random forest classifier evaluated on a test dataset (for details, see Methods). Dotted line indicates separation of 6 broad functional cell groups (43).



Supplemental Fig. S 6. Hybrid model for different cell types. **a.** Performance difference (mean) between hybrid-natural-future and SI based on test data for different cell types (each model for $n=10$ random seeds). **b.** Performance difference (mean) between hybrid-natural-future and SI based on test data for DS and non-DS cells (each model for $n=10$ random seeds). **c.** Cumulative histogram of model prediction difference between hybrid-natural-future ($w = 0.7$) and SI on test data, for DS (red) and non-DS cells, at one particular seed. **d.** Scatter plots for model predictions based on test data at a particular seed (each dot representing one neuron) for DS and non-DS cells and hybrid-natural-past ($w = 0.7$) vs. hybrid-natural-future ($w = 0.7$). Note that the predictions of two hybrid models were similar for most of neurons. **e.** Quality index (mean) for DS and non-DS cells based on responses to the repeated test sequences in the noise stimuli ($p=0.2881$, two-sided permutation test, $n=10,000$ repeats; for details, see Methods). **f.** Like (e) but for chirp responses ($p=0.6714$, two-sided permutation test, $n=10,000$ repeats). **g.** Like (e) but for bar stimulus responses ($p<0.0001$, two-sided permutation test, $n=10,000$ repeats). Error bars in (a),(b),(e)-(g) represent 2.5 and 97.5 percentiles with bootstrapping.