
DECONTEXTUALIZED LEARNING FOR INTERPRETABLE HIERARCHICAL REPRESENTATIONS OF VISUAL PATTERNS

R. Ian Etheredge^{1, 2, 3, *}, Manfred Schartl^{4, 5, 6, 7}, and Alex Jordan^{1, 2, 3}

¹Department of Collective Behaviour, Max Planck Institute of Animal Behavior, Konstanz, Germany

²Centre for the Advanced Study of Collective Behaviour, University of Konstanz, Konstanz, Germany

³Department of Biology, University of Konstanz, Konstanz, Germany

⁴Centro de Investigaciones Científicas de las Huastecas Aguazarca, A.C., Calnali, Hidalgo, Mexico

⁵Developmental Biochemistry, Biocenter, University of Würzburg, Würzburg, Bavaria, Germany

⁶Hagler Institute for Advanced Study, Texas A&M University, College Station, TX, USA

⁷Xiphophorus Genetic Stock Center, Texas State University San Marcos, San Marcos, TX, USA

August 25, 2020

SUMMARY

1 Apart from discriminative models for classification and object detection tasks, the application of deep
2 convolutional neural networks to basic research utilizing natural imaging data has been somewhat
3 limited; particularly in cases where a set of interpretable features for downstream analysis is needed,
4 a key requirement for many scientific investigations. We present an algorithm and training paradigm
5 designed specifically to address this: decontextualized hierarchical representation learning (DHRL).
6 By combining a generative model chaining procedure with a ladder network architecture and latent
7 space regularization for inference, DHRL address the limitations of small datasets and encourages
8 a disentangled set of hierarchically organized features. In addition to providing a tractable path for
9 analyzing complex hierarchal patterns using variation inference, this approach is generative and can
10 be directly combined with empirical and theoretical approaches. To highlight the extensibility and
11 usefulness of DHRL, we demonstrate this method in application to a question from evolutionary
12 biology.

13 **Keywords** Generative Modeling · Interpretable AI · Disentangled Representation Learning ·
14 Hierarchical Features · Image Analysis · Small Data

*Corresponding author: ietheredge@ab.mpg.de

15 **1 Introduction**

16 The application of deep convolutional neural networks (CNNs¹) to supervised tasks is quickly becoming ubiquitous,
17 even outside of standardized visual classification tasks.² In the life sciences, researchers are leveraging these powerful
18 models for a broad range of domain-specific discriminative tasks such as automated tracking of animal movement,³⁻⁶
19 the detection and classification of cell lines,⁷⁻⁹ and mining genomics data.¹⁰

20 A key motivation for the expanded use of deep feed-forward networks lies in their capacity to capture increasingly
21 abstract and robust representations. However, outside of the objective function they have been optimized on, building
22 interpretability into these representations is often difficult as networks naturally absorb all correlations found in the
23 sample data and the features which are useful for defining class boundaries can become highly complex (Figure S1). For
24 many investigations the main objective falls outside of a clearly defined detection or classification task, e.g. identifying
25 a set of descriptive features for downstream analysis, and interpretability and generalizability is much more important.
26 Because of this, in contrast to many traditional computer vision algorithms,¹¹⁻¹⁴ the application of more expressive
27 approaches built on CNNs and other deep networks to research has been limited¹⁵ (Figure 2).

28 Unsupervised learning, a family of algorithms designed to uncover unknown patterns in data without the use of labeled
29 samples, offers an alternative for compression, clustering, and feature extraction using deep networks. Generative
30 modeling techniques have been especially effective in capturing the complexity of natural images, i.e. generative
31 adversarial networks (GANs¹⁶) and variational autoencoders (VAEs,^{17,18}). VAEs in particular offer an intuitive way
32 for analyzing data. As an extension of variational inference, VAEs combine an inference model, which performs
33 amortized inference (typically a CNN) to approximate the true posterior distribution and encode samples into a set of
34 latent variables ($q_\phi(z|x)$), and a generative model which generates new samples from those latent variables ($p_\theta(x|z)$).
35 Instead of optimizing on a discriminative task, the objective function in VAEs is less strictly defined but typically
36 seeks to minimize the reconstruction error between inputs x and outputs $p_\theta(q_\phi(x))$ (reconstruction loss) as well as the
37 divergence between the distribution of latent variables $q_\phi(z|x)$ and the prior distribution $p(z)$ (latent regularization).

38 **1.1 Overcoming Hurdles to Application**

39 In VAEs, two problems often arise which are of primary concern to researchers using natural imaging data. 1) The
40 mutual information between x and z can become vanishingly small, resulting in an uninformative latent code and overfit
41 to sample data, the information preference problem;^{22,24} this is particularly true when using powerful convolutional
42 decoders which are needed to create realistic model output.^{20,23,24} 2) In contrast to the hierarchical representations
43 produced by deep feed-forward networks used for discriminative tasks, in generative models local feature contexts
44 become emphasized at the cost of large-scale spatial relationships. This is a product of the restrictive mean-field
45 assumption of pixel-wise comparisons and produces generative models capable of reproducing complex image features
46 while using only local feature contexts without capturing higher-order spatial relationships within the latent encoding.²²

Decontextualized learning for interpretable hierarchical representations of visual patterns

Table 1: Desired characteristics of an integrative tool for investigations of natural image data and general representation learning meta-prior enforcement strategies.

Desired Characteristic	Representation Learning Meta-Prior ¹⁹	Example Approach
Disentangling factors of variation	Limited number of shared factors of variation	Latent regularization ^{20,21}
Capturing spatial relationships	Hierarchical organization of representation	Hierarchical model architecture ²²
Incorporating existing knowledge	Local variation on manifolds	Structured latent codes ²³
Connect analyses and experiments	Local variation on manifolds	Generative models ^{16,17}
Inference	Probability mass and local variation on manifolds	Variational inference ¹⁷

47 The basis of a more expressive and robust approach for investigating natural image data has some key requirements: 1) provide a useful representation which disentangles factors of variation along a set of interpretable axes; 2) capture feature contexts and hierarchical relationships; 3) incorporate existing knowledge of feature importance and relationships between samples; 4) allow for statistical inference of complex traits; and 5) provide direct connections between analytical, virtual and experimental approaches. Here we integrate meta-prior enforcement strategies taken from representation learning¹⁹ to specifically address the requirements of researchers using natural image data (Table 1).

53 Here we propose to address the limitations of existing approaches and incorporate the specific requirements of researchers using a combination of meta-prior enforcement strategies. VAEs with a ladder network architecture has been shown to better capture a hierarchy of feature by mitigating the explain away problem of lower level feature, allowing for bottom-up and top-down feedback.²² Additionally, combining pixel-wise error with a perceptual loss function²⁵ adapted from neural style transfer,^{26,27} may also reduce the restrictive assumptions of amortized inference and pixel-wise reconstruction error by balancing them against abstract measures of visual similarity.

59 In terms of the latent regularization, a disentangled representation of causal factors requires an information-preserving latent code. Choosing a regularization techniques which mitigate the trade off between inference and data fit²¹ can encourage the disentanglement of generative factors along a set of variables in an interpretable way. We also propose a novel training paradigm inspired by GAN chaining that further relaxes the natural covariances in the data: *decontextualized learning* and actually uses the restrictive assumptions of GAN generator networks to our advantage to overcome the limitations of small datasets, typical for many studies in the natural sciences and further increase the disentanglement of generative factors (Figure 1, Methods 4.2).

66 While several metrics have been proposed for assessing interpretability and disentanglement,²⁸⁻³⁰ these metrics rely heavily on the associated labels, well defined features or stipulations from classification of detection competitions, e.g.³¹ In addition to being highly domain specific, for most practical investigations in the natural sciences, these types of

Decontextualized learning for interpretable hierarchical representations of visual patterns

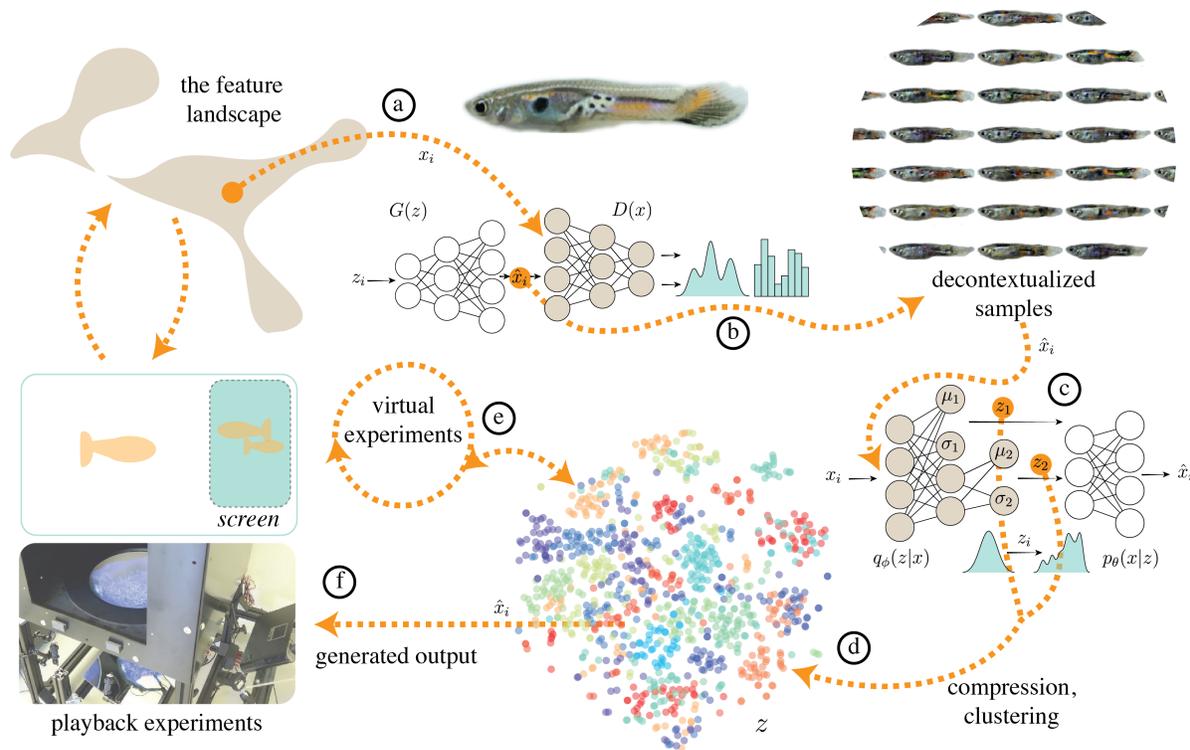


Figure 1: *Overview.* a) Many patterns (e.g. male guppy ornaments) consist of combinations of several elements which have hierarchical relationships, spatial dependence, and feature contexts which may hold distinct biological importance. In our proposed framework, small sample sizes are supplemented using a generative (GAN) model which learns images statistics sufficient to produce novel out of sample examples (b). This model can be used to produce an unlimited number of novel samples and also reduces the covariances within sample data, which is advantageous for disentangling generative features. We use these "decontextualized" generated outputs as input (c) to a variational auto encoder (VAE). Via a specific combination of meta-prior enforcement strategies and network architectures, we capture the hierarchical structure, disentangling factors of variation across multiple scales of increasing abstraction (z_1 through z_n). Using the learned distribution over these variables, the latent representation, parameterized by a mean and variance term, we (d) define a color-pattern space. Using this low dimensional representation we can (e) interface with downstream models such as evolutionary algorithms and (f) produce photo-realistic outputs to be used in playback experiments and immersive VR. Interpolations through the color-pattern space with animated models and VR allows researchers to manipulate generated output for experimental tests. Techniques represented with dashed lines: (d) capturing a hierarchy of visual features, (e) combining a low-dimensional latent representation with virtual experiments and (f) playback experiments represent the current gaps in our analytical and experimental framework. Our approach directly address these shortcomings to span these gaps, creating a robust, integrated framework for investigating natural visual stimuli.

69 labels do not exist and we must often rely on fundamentally qualitative assessments. In many cases, labeled data is not
 70 available and interpreting traversals of the latent code (Figure S2) may introduce our own perceptual biases. Here, we
 71 adapt an approach from explainable AI: integrated gradients³² in application to latent variable exploration too provide a
 72 direct assessment of latent variables, quantifying latent feature attributions without the necessity of labeled data and
 73 allows for exploring latent variables without adding additional human biases (Methods 4.3).

74 We demonstrate the proposed framework using two example datasets: male guppy ornamentation and butterfly wing
 75 patterns from the discipline of sensory ecology and evolution (see Appendix A for motivation and background on
 76 existing approaches).

77 2 Results

78 While biological datasets are typically small, they are usually highly structured and standardized compared to large
79 classification datasets (e.g. ImageNet³³). This provides an advantage for controlling noise and uninformative covariates
80 in the data. Using a modified infoGAN²³ architecture, we incorporate prior knowledge about the structure of our sample
81 data to generate realistic samples from the complex image distribution conditioned on a set of latent variables. Here, we
82 incorporate prior knowledge about our samples of male guppy ornamentation images by providing a 32-class categorical
83 latent code (Figure 3b, top right). These 32-classes represent the 32 individual tanks, unique subsets of the overall
84 sample, with shared traits related to guppy ornamentation patterns. The categories learned by the trained model possess
85 unique features which also covary in the sample data, e.g. a distinct black bar and orange stripe which characterizes
86 one guppy species, *P. wengei* (Figure S2, a). While generated samples share characteristics and even resemble known
87 varieties, generated samples possess decontextualized combinations of features across examples (Figure S2, a). We use
88 these, decontextualized samples as input to our variational (VAE) model for our "decontextualized" training paradigm.
89 GAN training and VAE training are performed in separate steps so that models are not jointly optimized. The generated
90 samples from the trained GAN model are used as training data to a variational model (Figure 1) with a hierarchical model
91 architecture²² which consists of 10 latent variables across four codes (z_1, \dots, z_4) with increasing expressivity, (Methods
92 4.2.2). We observe distinct clusters in the latent space of the trained model which correspond to sample categories
93 and differs qualitatively from two existing methods (raw pixel and perceptual loss embeddings using tSNE,^{34,35} Figure
94 2). The unique latent space of the four latent encodings capture unique factors of variation in the sample data in a
95 scale-dependent way (Figure S2, Figure S3). In this model, z_1 , the latent code with the lowest capacity captures local
96 traits such as the color and intensity of discrete patches, e.g. z_{11} encodes variation in the intensity of an orange spot (S2
97 4b, left). At higher levels (z_2, \dots, z_4), latent variables encode complex traits which combine multiple elements, (S2 4b,
98 right). We use this same latent representation to describe the relationship between samples and calculate likelihood
99 estimates. Samples with rare traits, e.g. such as the "Tr5" strain in our sample data which are distinctly melanated,
100 cluster together in the embedded space, and have a low sample likelihood (3).

101 Embedding the 4, 10-dimensional, latent codes reveals scale-dependent relationships between elements. In z_1 (Figure
102 S3, left) color values and local features dominate the relationship between points (Figure S3, left). Nearest neighbor
103 samples (Minkowski distance³⁶ in the 10-dimensional space, Figure S3, b) show color similarity whereas higher order
104 features, e.g. patterning and morphology, determine the relationships between samples in the more expressive latent
105 spaces (z_2, \dots, z_4). Though we find strong covariance between features across scales, in some cases the nearest
106 neighbors samples differ greatly depending on the scale and feature context (Figure S3, b).

107 We assess the level of disentanglement of our trained variational models using the metric established in³⁰ using known
108 class labels as attribute classes (butterfly species, learned class from infoGAN pre-training, and guppy strain varieties).
109 Across models, we find the most expressive latent codes (z_4) provide the highest degree of disentanglement between
110 known classes with the highest disentanglement score overall using our decontextualized, DHRL method (see Table 2).

Decontextualized learning for interpretable hierarchical representations of visual patterns

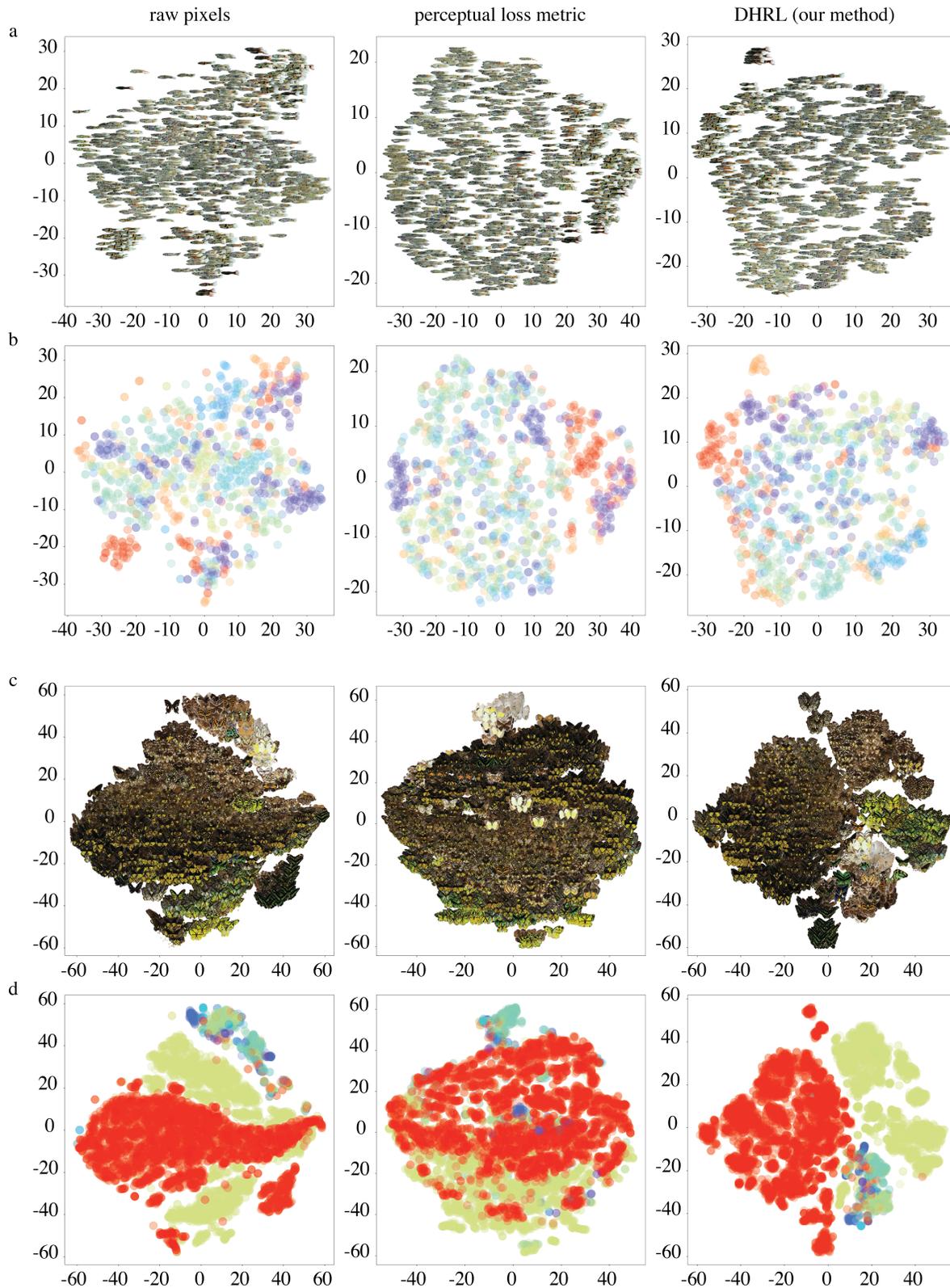


Figure 2: *Comparing with existing techniques.* 2-dimensional embedding of (left) raw pixel distributions, (middle) using a perceptual similarity score^(15,25), and (right) our framework. a) guppies, b) butterflies. Colors indicate unique sub groups for each sample (guppy variety and butterfly species).

Decontextualized learning for interpretable hierarchical representations of visual patterns

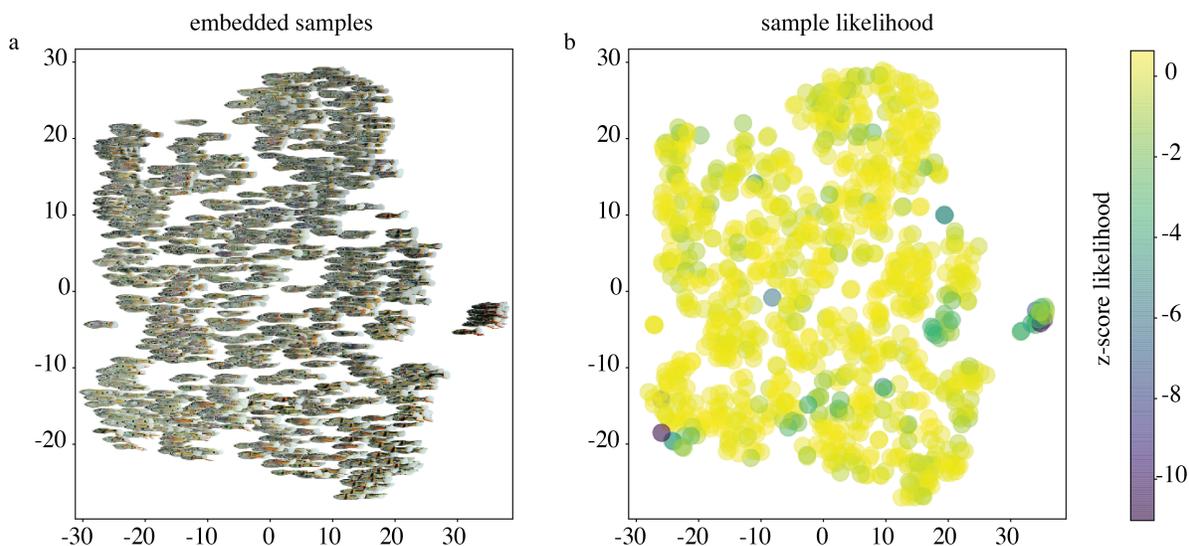


Figure 3: *Sample likelihood estimates*. a) Embedded samples b) Normalized (standard score) likelihood estimates for each sample.

Table 2: Disentanglement and completeness metrics for of VLAE inference network across datasets and when using our decontextualized learning approach (DHRL).

VLAE Training Data	$D(z_1), C(z_1)$	$D(z_2), C(z_2)$	$D(z_3), C(z_3)$	$D(z_4), C(z_4)$
Butterflies (n=9531)	0.64, 0.60	0.67, 0.55	0.63, 0.51	0.88, 0.60
Guppies (n=987)	0.29, 0.32	0.12, 0.13	0.13, 0.16	0.56, 0.66
Guppies (gen., n=19k)	0.12, 0.16	0.13, 0.18	0.32, 0.42	0.62, 0.75
Guppies DHRL	0.14, 0.16	0.18, 0.23	0.31, 0.39	0.90, 0.95

111 We also provide a qualitative approach for attributing latent variables to image features using network gradients
 112 (Methods 4.3); when labels are unknown. In Figure 4, a-d we visualize one variable of z_1 , the least expressive
 113 latent variable space (z_{13}) of the DHRL-trained guppy latent variable model. We find that the same latent variable
 114 controls the relative intensity of green color patches across individuals. Looking at a single variable of more expressive
 115 latent codes z_{27} of the trained butterfly model (4, e-h) we find that this latent variable controls the size of yellow
 116 patches on the lower wings relative to the size of yellow patches on the upper wings (when patches are not present
 117 this variable has no effect (4, f). Further investigation of latent variables can be performed using the provided tool
 118 (<https://github.com/ietheredge/VisionEngine/notebooks/IntegratedGradients.ipynb>).

119 Using the latent representation, z , of our DHRL trained variational model of guppy ornaments as input, we apply an
 120 evolutionary algorithm (Figure 5), defined by a fitness function from the guppy literature: oranger, higher contrast males
 121 are preferred by females.³⁷ Starting from a parent population initialized by our sample embedding (900 samples), we
 122 simulate 500 generations under these selective forces. We observe exaggerated and more numerous orange and black

Decontextualized learning for interpretable hierarchical representations of visual patterns

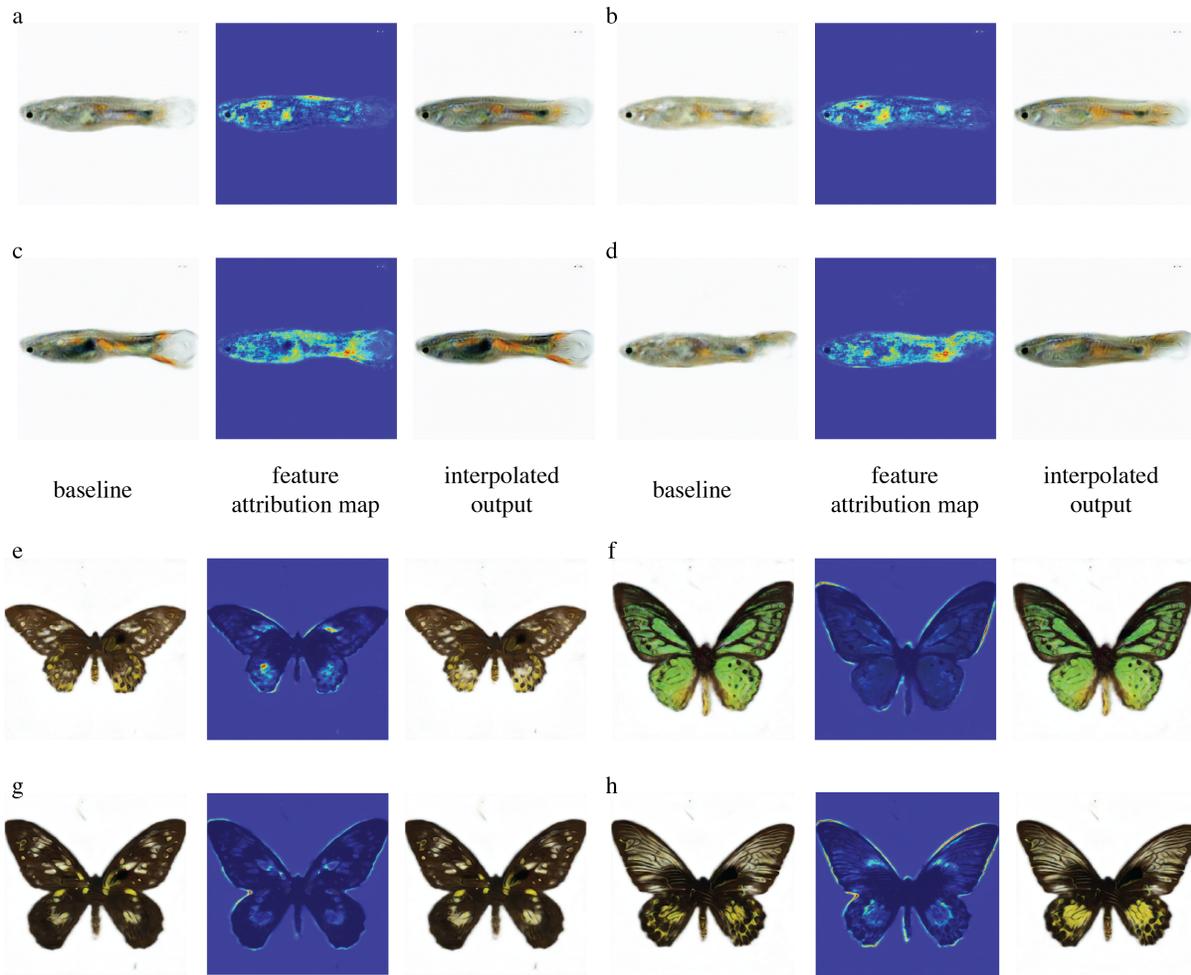


Figure 4: *Latent variable feature attribution*. a-d) Four samples of generated guppy images performing integrated gradients feature attribution (see Methods 4.3). a-d) guppy images visualizing the feature attributions of the latent variable z_{13} of the DHRL trained variational model. e-h) butterfly images using latent feature z_{37} . Heatmap values have been normalized using a standard score. Images to the left are generated with the latent feature set to its lowest value in the sample, to the right with the highest value in the sample.

123 patches in novel configurations compared to the initial population (Figure 5, b). Projecting the latent representation of
124 generations 1, 250, and 500, we find that instead of a single peak, after several generations, many novel solutions are
125 optimized (Figure 5 a). Investigating the values of the latent variables over generations reveals two distinct latent factors
126 driven to fixation in the population under these selective forces (S4). We also observe to population optimization of
127 latent factors over time in Movie S5. Using a single Titan Xp GPU with 12GB memory we could simulate a population
128 size of 1000 individuals in an average of 19.5 seconds per generation.

Decontextualized learning for interpretable hierarchical representations of visual patterns

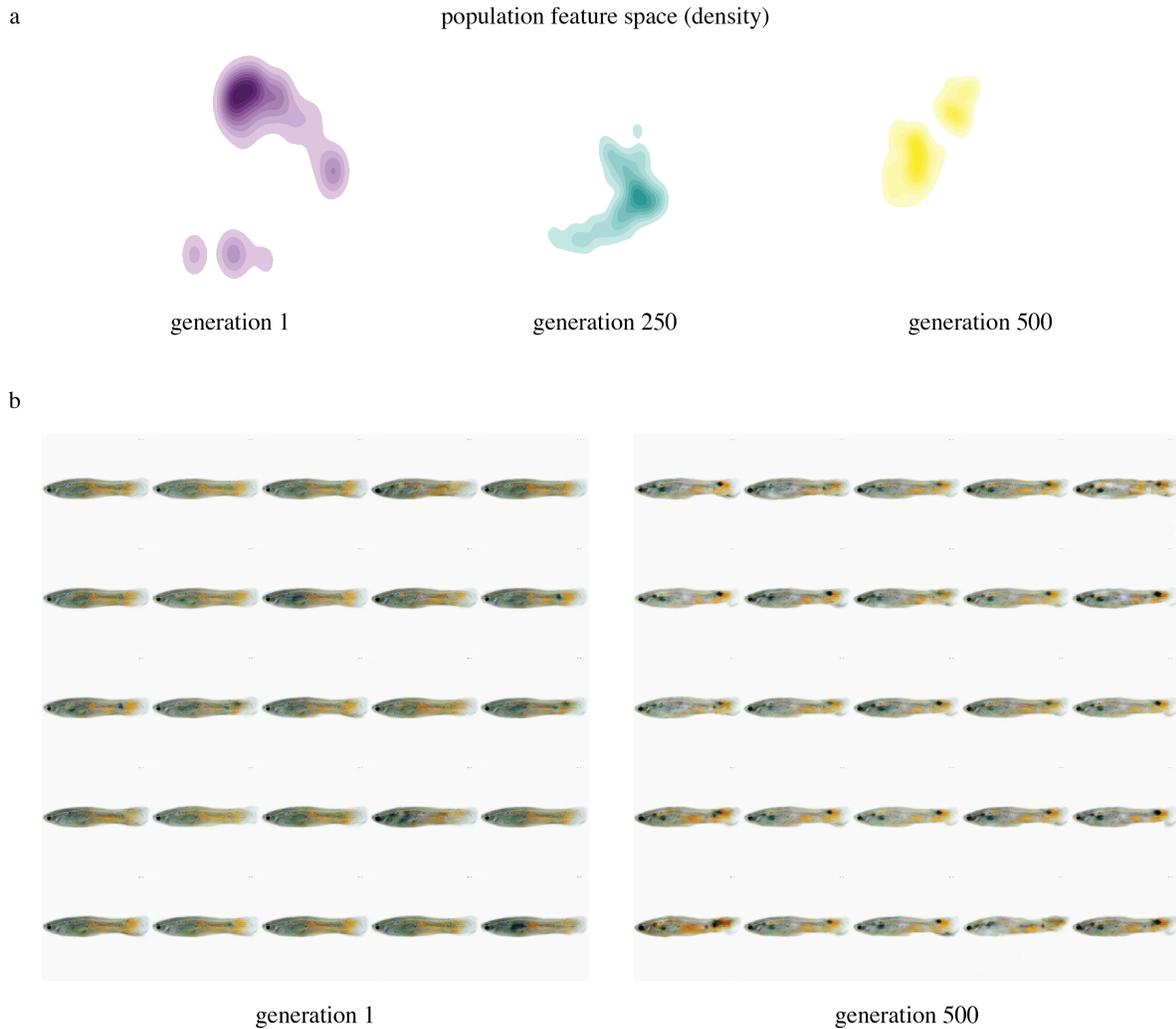


Figure 5: *Virtual Experiments*. a) Kernel density plot of samples over generations 1, 250, 500 selecting orange ornaments and contrast. After 500 generations the population has shifted from the initial sample distribution, finding two peaks which maximize the fitness function. b) Samples of initial parent population, left, with the highest fitness, compared to those with the highest fitness after 500 generations. Samples in later generations show higher numbers of brighter orange and dark melanated patches and increased within body contrast.

129 **3 Discussion**

130 Supervised discriminative learning algorithms are already becoming an integral tool for researchers across disciplines
131 whereas unsupervised generative modeling approaches remain a relatively young and active area of machine learning
132 research. Already, the highly expressive generative models like the ones presented here are transforming the way we
133 interact with image data. By solving problems in a more general way, generative modeling approaches provide more
134 direct connections to hypothesis testing and connecting observations. Here, we demonstrate how these approaches may
135 serve as an engine for more integrative studies of animal coloration patterns, and natural image data more generally,
136 directly connecting approaches.

137 Analytically, our approach captures important hierarchical features across spatial scales that existing approaches do not
138 account for (Figure 2, Figure S3, Appendix A.1), it removes the inherent biases of predefined filters by learning features
139 directly from the sample data, and it disentangles complex factors of variation into a useful, meaningful representation
140 (Figure S2, Figure 4). More than compressing data into a low dimensional space, this approach is generative and can
141 create novel out-of-sample examples with high fidelity. This is a potentially transformative extension for researchers in
142 the natural sciences which is not offered by existing approaches, allowing researchers to test analytical results with
143 virtual experiments, and empirically, by using virtual reality playback experiments or observational studies (see Movie
144 S6).

145 These techniques can be adapted to many domain specific questions (see A.1 for a specific discussion regarding the
146 potential impact of this approach on the study of color pattern evolution). As the latency between input and output
147 decreases in video playback experiments, integrating instantaneous behavioral feedback and in-the-loop methods for
148 hypothesis testing may be used to design complex real-time assays. More sophisticated virtual experiments may also
149 incorporate agent based models and evolutionary algorithms working directly on the latent representation to create
150 complex simulations (e.g as in,³⁸ Figure 5). In our demonstration, we are able to simulate 1000 individuals in under 20
151 seconds per generation with very little optimization and asynchronous approaches may already be possible. Analytically,
152 as research in machine learning aimed at understanding how information is organized and used by algorithms advances,
153 a growing theoretical framework with a basis in statistical mechanics³⁹ and information theory⁴⁰ may provide additional
154 avenues for investigating the statistical properties of color pattern spaces and their evolution.

155 **4 Experimental Procedures**

156 **4.1 Materials Availability**

157 Guppy images were collected from a maintained stock at the University of Wuerzburg under authorization 568/300-
158 1870/13 of the Veterinary Office of the District Government of Lower Franconia, Germany, in accordance with the
159 German Animal Protection Law (TierSchG). Individuals were imaged on a white background with fixed lighting

160 conditions⁴¹ using a Cannon D600 digital camera. Images were down sampled and center cropped to final size of 256 x
161 256 pixels. The dataset consists of 977 standardized RGB images across three species and 13 individual strains.

162 Butterfly images were downloaded from the Natural History Museum, London under a creative commons license (DOIs:
163 <https://doi.org/10.5519/qd.gvq3p7xq>, <https://doi.org/10.5519/qd.pw8srv43>). This dataset consists of 9531 RGB images.

164 For each dataset, we segmented samples from the background using a customized object segmentation network adapted
165 from.⁴² For each dataset we annotated 8 samples to train the segmentation network. All samples were cropped and
166 resized to 256 x 256 and placed on a transparent background (RGBA). For calculating the perceptual loss during
167 training, images were translated to 3-channel images with a white background using alpha blending.

168 Updated links to original data repositories can be accessed here: [https://github.com/ietheredge/
169 VisionEngine/README.md](https://github.com/ietheredge/VisionEngine/README.md).

170 4.1.1 Data and Code Availability

171 All models were implemented using Tensorflow 2.2 and can be accessed here: [https://github.com/ietheredge/
172 VisionEngine](https://github.com/ietheredge/VisionEngine), including installation and evaluation scripts to reproduce our results. Instructions for creating new
173 data loaders for training new datasets using this method can be found at [https://github.com/ietheredge/
174 VisionEngine/data_loaders/datasets/README.md](https://github.com/ietheredge/VisionEngine/data_loaders/datasets/README.md).

175 4.2 Key Methods

176 DHRL relies on a three-step process of sequential training where first a generative adversarial network is trained to
177 transform a noise sample into realistic out of sample examples. Next, a variational autoencoder is pre-trained on the
178 generated samples. Then finally, the pretrained variational model is fine-tuned on the original samples.

179 4.2.1 InfoGAN

180 We use an unsupervised approach to disentangle discrete and continuous latent factors adapted from²³ (InfoGAN) which
181 modifies the minimax game typically used for training GANs such that:

$$\min_{G,Q} \max_D V_I(D, G, Q) = V(D, G) - \lambda L_I(G, Q) \quad (1)$$

182 where $V(D, G)$ is the original GAN objective introduced in¹⁶ and $L_I(G, Q)$ approximates the lower bound of the
183 mutual information $I(c; G(z, c))$ using Monte Carlo sampling such that $L_I(G, Q) \leq I(c; G(z, c))$.²³ Like the generator
184 G and discriminator D , Q is parameterized as a neural network and shares all convolutional layers with D .

185 Both discrete $Q(c_d|x)$ and continuous latent codes $Q(c_c|x)$ are provided with continuous latent codes treated as a
186 factored Gaussian distributions. Importantly, InfoGAN does not require supervision and no labels are provided, e.g.²⁹

Decontextualized learning for interpretable hierarchical representations of visual patterns

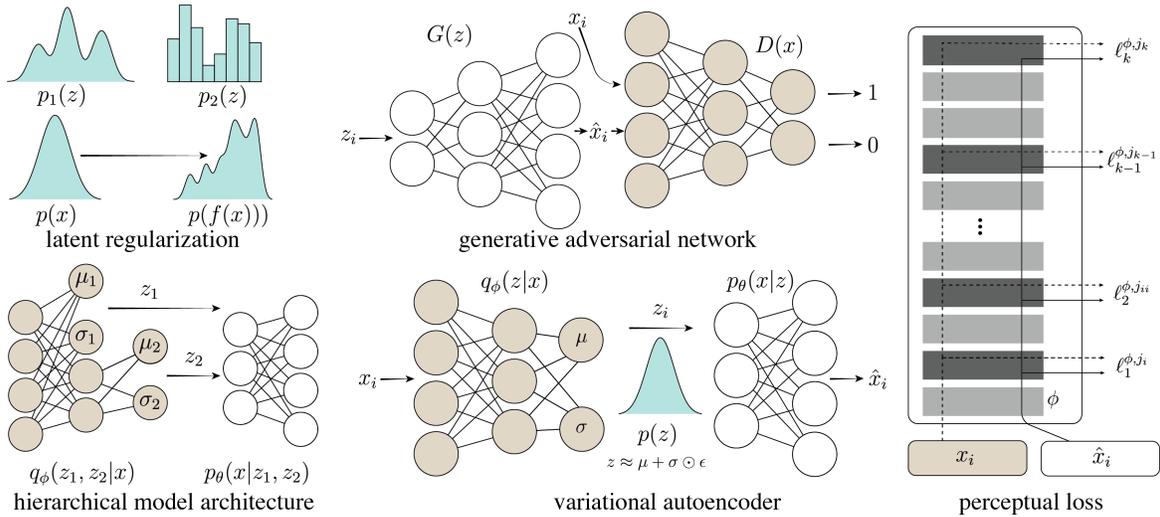


Figure 6: *Key Methods* *Top left*: the distributions of our latent representation may be parameterized by a number of continuous or discrete variables. In infoGAN, a categorical latent code is combined with continuous latent codes which allows for disentangling substructures in the sample data without labeled samples. *Top right*, example structure of a generative adversarial network. Here, a noise vector, z_i is input to the generator network $G(z)$ which produces a reconstructed output \hat{x}_i . A real sample, x_i , and generated sample \hat{x}_i are subsequently passed through a separate discriminator network $D(x)$ which determines if the sample is real (1) or generated (0). In infoGAN, the latent encoding of generated samples is optimized an additional network Q which shares all convolutional layers with D . *Bottom left*, the generic architecture of a variational ladder autoencoder (VLAE). Multiple latent spaces (z_1, z_2, \dots, z_k) are learned with each successive input layer having increasing expressivity and abstraction (ability to combine features across spatial scales). *Bottom middle*, structure of a variational autoencoder (VAE). x_i and \hat{x}_i are an example input and its reconstructed output, the probabilistic encoder or inference model, $q_\phi(z|x)$, performs posterior inference learning shared model parameters, ϕ , across samples, approximating the true posterior distribution. The probabilistic decode, $p_\theta(Z|X)$, $p_\theta(X|Z)$, learns a joint distribution of the encoded space, Z , and the data space x . The low dimensional bottleneck, Z is a distribution of latent variables capable of reconstructing sample inputs, parameterized by a vector of means μ and standard deviations σ . The noise term ϵ allows for the parameters of this multivariate distribution to be optimized using back propagation, known as the reparametrization trick. *Right*, Perceptual loss models use a pretrained network, ϕ , e.g. VGG-16.⁴³ Two samples, the original input and reconstructed output, are input to the model and the maxpooling layer activations for each are used as outputs. The distance between these functions emphasizes higher-level similarity than standard pixel-wise differences. Outputs from the shallow layer ℓ_1 represents low-level local features where as output from the deeper layer ℓ_k contains information from across spatial scales and more abstract representations. The euclidean distance between these activation outputs gives a metric for the similarity of the two inputs as "perceived" by a network pre-trained on a much broader dataset. Perceptual loss functions can be used as a stand-alone transfer-learning approach to finding perceptual differences between samples or as part of any network as an additional or alternative reconstruction loss (see 2).

187 We substitute the original generator and discriminator models from²³ with the architecture described in⁴⁴ and increase
188 the flexibility of the latent code, providing additional continuous and discrete latent codes. For guppy experiments, we
189 provide two continuous and 19 discrete codes (samples were drawn from 19 paternal lines). For the basis noise vector
190 input to the generator, we used 100-unit random noise vector.

191 4.2.2 Variational Ladder Autoencoder

192 In contrast to hierarchical architectures, e.g.,^{45,46} we learn a hierarchy of features by using multiple latent codes with
193 increasing levels of abstraction from,²² i.e. $q_\phi(z_1, \dots, z_L|x)$. The expressivity of z_i is determined by its depth. The
194 encoder $q_\phi(z_1, \dots, z_L|x)$ consists of four blocks such that:

$$H_\ell = G_\ell(H_{\ell-1}) \quad (2)$$

$$z_\ell \sim \mathcal{N}(\mu_\ell(H_\ell), \mathbf{I}) \quad (3)$$

195 where H_ℓ , G_ℓ , and μ_ℓ are neural networks. For our encoder model, G_ℓ is a stack of convolutional, batch normalization,
196 and leaky rectified linear unit activation (Conv-BN-LeakyReLU), we stack four Conv-BN-LeakyReLU blocks for
197 each G_ℓ with increasing number of channels for each subsequent convolutional layer, i.e. N-channels/2, N-channels,
198 N-channels, N-channels*2 where N-channels is 16, 64, 256, 1024 for G_1, G_2, G_3, G_4 respectively. We apply spectral
199 normalization to all convolutional layers (see below). Because we want to preserve feature localization, we use average
200 pooling followed by a squeeze-excite block to apply a context-aware weighting to each channel (see below).

201 Similarly, the decoder, $p_\theta(x|z_1, \dots, z_L)$, is composed of blocks such that:

$$\tilde{z}_\ell = U_\ell([\tilde{z}_{\ell+1}; V_\ell(z_\ell)]) \quad (4)$$

202 where $[\cdot; \cdot]$ denotes channel-wise concatenation. Parallel to G_ℓ , blocks in the decoder: U_ℓ are composed of Conv-BN-
203 ReLU blocks (note the use of ReLU and not LeakyReLU in the decoder) with decreasing number of channels in each
204 convolutional layer, i.e. N-channels*2, N-channels, N-channels, N-channels/2 where N-channels is 1024, 256, 64, 16.
205 No spectral normalization wrappers or squeeze-excite layers are applied in the decoder.

206 4.2.3 Squeeze-Excite Layers

207 Squeeze-and-Excitation Networks⁴⁷ were proposed to improve feature interdependence by adaptively weighting each
208 channel within a feature map based on the filter relevance by applying a channel-wise recalibration. Here we
209 apply squeeze-excite (SE) layers prior to the variational layer such that each embedding z_i captures features with
210 cross-channel dependencies. Each SE layer consists of a global average pooling layer which averages channel-wise
211 features followed by two fully connected layers with relu activations, the first with size channels/16 and the second with

212 the same size as the number of input channels. Finally a sigmoid, "excite," layer assigns channel wise probabilities
 213 which are then multiplied channel wise with the original inputs.

214 4.2.4 Reconstruction Loss

215 We minimize the negative log likelihood of the sample data by minimizing the mean squared error between input and
 216 output, jointly optimizing the reconstruction loss for each sample x :

$$\begin{aligned}\mathcal{L}_{\text{pixel-wise}} &= \mathbb{E}_{p_{\text{data}}}(x) \mathbb{E}_{q_{\phi}}(z | x) [\log p_{\theta}(x | z)] \\ &= \frac{1}{n} \sum_{i=1}^n (x_i - p_{\theta}(q_{\phi}(x_i)))^2\end{aligned}\quad (5)$$

217 To relax the restrictive mean-field assumption which is implicit in minimizing the pixel-wise error, we jointly optimize
 218 the similarity between inputs and outputs using intermediate layers of a pretrained network, VGG16,⁴³ as feature
 219 maps.²⁵⁻²⁷ Here we calculate the Gram matrices of feature maps, which match the feature distributions of real and
 220 generated outputs for each layer as:

$$G_{ab}^{\ell} = \frac{\sum_{cd} F_{cda}^{\ell}(x) F_{cdb}^{\ell}(x)}{CD} \quad (6)$$

$$\mathcal{L}_{\text{perceptual}} = \sum_{\ell=1}^L \frac{\frac{1}{n} \sum_{i=1}^n (G_{ab}^{\ell}(x_i) - G_{cd}^{\ell}(p_{\theta}(q_{\phi}(x_i))))^2}{L} \quad (7)$$

221 for feature maps F_a and F_b in layer ℓ across locations c, d . This measures the correlation between image filters and is
 222 equivalent to minimizing the distance between the distribution of features across feature maps, independently of feature
 223 position.⁴⁸

224 The combined reconstruction loss is a weighted sum of the perceptual loss and pixel-wise error:

$$\mathcal{L}_{\text{reconstruction}} = \alpha \mathcal{L}_{\text{perceptual}} + \beta \mathcal{L}_{\text{pixel-wise}} \quad (8)$$

225 where α and β are Lagrange multipliers controlling the influence of each loss term. Here we set $\alpha = 1\text{e-}6$ and $\beta = 1\text{e}5$
 226 to balance the contribution of reconstruction terms with variational loss (see below).

227 4.2.5 Maximum Mean Discrepancy

228 We use the maximum mean discrepancy approach (MMD)²¹ to maximize the similarity between the statistical moments
229 of $p(z)$ and $q_\phi(x)$ using the kernel embedding trick:

$$\text{MMD}(p(z)||q_\phi(z)) = \mathbb{E}_{p(z),p(z')} [k(z, z')] + \mathbb{E}_{q_\phi(z),q_\phi(z')} [k(z, z')] - 2\mathbb{E}_{p(z),q_\phi(z')} [k(z, z')] \quad (9)$$

230 using a Gaussian kernel, $k(z, z')$, such that

$$k(z, z') = e^{-\frac{\|z-z'\|^2}{2\sigma^2}} \quad (10)$$

231 to measure the similarity between $p_\theta(z)$ and $q_\phi(z)$ in Euclidean space. We measured similarity using multiple kernels
232 with varying degrees of smoothness, controlled by the value of σ^2 , i.e. multi-kernel MMD,⁴⁹ with varying bandwidths:
233 $\sigma^2 = 1e-6, 1e-5, 1e-4, 1e-3, 1e-2, 1e-1, 1, 5, 10, 15, 20, 25, 30, 35, 100, 1e3, 1e4, 1e5, 1e6$.

234 Weighing the influence of MMD kernel differences on the combined objective function is controlled by the Lagrange
235 multiplier λ applied across each latent code. Giving the combined objective:

$$\mathcal{L}_{\text{total}} = \left(\sum_i^L \lambda \text{MK-MMD}(q_\phi(z_i)||p(z_i)) \right) + \mathcal{L}_{\text{reconstruction}} \quad (11)$$

236 where L is the number of hierarchical latent codes and z_i is the n-dimensional latent code and the prior, $p(z_i) = \mathcal{N}(0, \mathbf{I})$
237 and $\mathcal{L}_{\text{reconstruction}}$ define above. Here, we set $\lambda = 1$.

238 4.2.6 Denoise Training

239 In addition to further relaxing the contribution of pixel-wise error, adding a denoising criterion has been shown to yield
240 better sample likelihood by learning to map both training data and corrupted inputs to the true posterior, providing more
241 robust training for out of sample data.⁵⁰ We implement this with the addition of noise layer which samples a corrupted
242 input \tilde{x} from input x before passing \tilde{x} to the encoder $q_\phi(z|\tilde{x})$. We use apply random binomial noise (salt and pepper) to
243 ten percent of pixels.

244 4.2.7 Spectral Normalization

245 Spectral normalization has been proposed as a method to prevent exploding gradients when using rectified linear units
246 to stabilize GAN training via a global regularization on the weight matrix of each layer as opposed to gradient clipping
247 to provide bounded first derivatives (the Lipschitz constraint).⁵¹

248 4.3 Latent Feature Attribution and Disentanglement

249 Understanding the importance of features for model predictions is an active area of research. Integrated gradients,
 250 introduced by,³² assigns feature importance, determining causal relationships between predictions and image features
 251 by summing the gradients along paths between x' and x .

$$\text{IG}_i(x) ::= (x_i - x'_i) \times \int_{\alpha=0}^1 \frac{\partial \mathcal{P}(x' + \alpha \times (x - x'))}{\partial x_i} d\alpha \quad (12)$$

252 We adapt this procedure to investigate the contribution of each latent variable parameter z_i where we use a baseline z ,
 253 an encoding of a single sample x and iterate z_j while holding all other z_l constant and summing the gradients of the
 254 decoder $p_\theta(x|z)$ such that:

$$\text{IG}_i^{\text{approx}}(p_\theta(x|z^j)) ::= (p_\theta(x|z^j)_i - p_\theta(x|z^{j'})_i) \times \sum_{k=1}^m \frac{\partial \mathcal{P}(p_\theta(x|z^{j'}) : z_j^{j'} = z_j^{j'} + \frac{k}{m} \times (z_j^j - z_j^{j'}))}{\partial p_\theta(x|z^j)_i} \times \frac{1}{m} \quad (13)$$

255 where j is the axis of latent code being interpolated, i is the individual feature (pixel), $p_\theta(x|z)$ is the reconstructed
 256 output, $p_\theta(x|z')$ is the baseline reconstructed output, k is the perturbation constant, and m is the number of steps in
 257 the approximation of the integral. We use the Riemann sum approximation of the integral over the interpolated path
 258 \mathcal{P} which involves computing the gradient in a loop over the inputs for $k = 1, \dots, m$. Here, we use $m = 300$ and
 259 $k = 2 \max(|z|)$ for each z^j starting from a baseline $p_\theta(x|z^{j'}) : z_j = -\max(|z|)$.

260 We use the technique developed in³⁰ for assessing disentanglement, measuring the relative entropy of latent factors for
 261 predicting class labels. We measure disentanglement of D_i of each latent code is measured by $D_i = (1 - H_K(P_i))$
 262 where H_K is the entropy and P_i is the relative importance of the generative factor. We also include a metric of
 263 completeness C_i , approximating the degree to which the generative factor is captured by a single latent variable, where
 264 $C_j = (1 - H_D(P_j))$ where P_j is the unweighted contribution of generative factors.³⁰ Here, in the absence of labeled
 265 features, we use species (butterflies), breeding line variants (guppies), and predicted class of the generative model
 266 (generated guppies, 4.2.1, above) for each model as approximate class labels (one class). This approximation naturally
 267 overestimates D_i and underestimates C_j as there is some overlap between classes in terms of visual features (see Figure
 268 2, Figure S3). While³⁰ proposes a third term to evaluate representations I to measure the relative informativeness,
 269 we found that this value was highly coupled to the choice of the Lagrange multiplier λ used for latent regularization
 270 (above).

271 4.4 Simulating Evolution on the Latent Space

272 For demonstrating an example virtual experiment, we use a genetic algorithm, with a parent population of 1000 random
 273 samples, evolved over 500 generations. Parent samples are random initialized across the the latent variables of each

274 latent code. Fitness was calculate as an equally weighted sum of the total percentage of pixels within two ranges (orange
275 $\text{rgb}(0.9, 0.55, 0.) > \text{rgb}(1., 0.75, 0.1)$ and black $\text{rgb}(0., 0., 0.) < \text{rgb}(0.2, 0.2, 0.2)$) measured on the generated output, a
276 simplification of empirical results from the literature.^{37,52} During each generation predicted fitness for each sample in
277 the population was measured by the fitness of the nearest neighboring value in the reference table (for processing speed).
278 To simulate weak selective pressure on the fitness function, we drew 500 random parent subsamples weighted by their
279 proportional fitness. An additional 200 samples were drawn, without the proportional fitness weighting. Together,
280 from the 700 subsamples in each generation we drew 300 random pairs, the "alleles" from each sample (the specific
281 latent variable values) were chosen randomly with equal probability to create a combined offspring between the two
282 samples. Each combined offspring then had two alleles randomly mutated, one by drawing from a random normal
283 distribution and the other by replacing an existing value with zero (similar to destabilizing and stabilizing mutations).
284 The next generation thus consisted of 100 samples, 700 parent samples + 300 offspring. This process was repeated for
285 500 generations.

286 **5 Acknowledgments**

287 We would like to thank members of the Dept. of Collective Behavior, Max Planck Institute of Animal Behavior and
288 Centre for the Advanced Study of Collective Behaviour, University of Konstanz for comments on earlier versions of the
289 manuscript as well as the Max Planck Computing and Data Facility for use of computational resources.

290 **5.1 Author contributions**

291 RIE conceived the approach and designed the methodology; MS and RIE collected sample data. RIE wrote the
292 manuscript. AJ secured funding. All authors contributed to editing and approving the manuscript.

293 **5.2 Declaration of Interests**

294 The Authors have no financial or non-financial competing interest.

295 **References**

- 296 ¹ Yann LeCun, Koray Kavukcuoglu, and Clément Faret. Convolutional networks and applications in vision. In
297 *Proceedings of 2010 IEEE international symposium on circuits and systems*, pages 253–256. IEEE, 2010.
- 298 ² Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual
299 object classes (voc) challenge. *International journal of computer vision*, 88(2):303–338, 2010.
- 300 ³ Katarzyna Bozek, Laetitia Hebert, Alexander S Mikheyev, and Greg J Stephens. Towards dense object tracking in
301 a 2d honeybee hive. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages
302 4185–4193, 2018.

Decontextualized learning for interpretable hierarchical representations of visual patterns

- 303 ⁴ Alexander Mathis, Pranav Mamidanna, Kevin M. Cury, Taiga Abe, Venkatesh N. Murthy, Mackenzie W. Mathis,
304 and Matthias Bethge. Deeplabcut: markerless pose estimation of user-defined body parts with deep learning. *Nature*
305 *Neuroscience*, 2018.
- 306 ⁵ Talmo D. Pereira, Diego E. Aldarondo, Lindsay Willmore, Mikhail Kislin, Samuel S. H Wang, Mala Murthy, and
307 Joshua W. Shaevitz. Fast animal pose estimation using deep neural networks. *Nature Methods*, 16(1):117–125, 2019.
- 308 ⁶ Jacob M Graving, Daniel Chae, Hemal Naik, Liang Li, Benjamin Koger, Blair R Costelloe, and Iain D Couzin.
309 Deepposekit, a software toolkit for fast and robust animal pose estimation using deep learning. *Elife*, 8:e47994,
310 2019.
- 311 ⁷ Thorsten Falk, Dominic Mai, Robert Besch, Özgün Çiçek, Ahmed Abdulkadir, Yassine Marrakchi, Anton Böhm,
312 Jan Deubner, Zoe Jäckel, Katharina Seiwald, et al. U-net: deep learning for cell counting, detection, and morphometry.
313 *Nature methods*, 16(1):67–70, 2019.
- 314 ⁸ Julian Riba, Jonas Schoendube, Stefan Zimmermann, Peter Koltay, and Roland Zengerle. Single-cell dispensing
315 and ‘real-time’ cell classification using convolutional neural networks for higher efficiency in single-cell cloning.
316 *Scientific reports*, 10(1):1–9, 2020.
- 317 ⁹ Claire McQuin, Allen Goodman, Vasilii Chernyshev, Lee Kamensky, Beth A Cimini, Kyle W Karhohs, Minh Doan,
318 Liya Ding, Susanne M Rafelski, Derek Thirstrup, et al. Cellprofiler 3.0: Next-generation image processing for
319 biology. *PLoS biology*, 16(7):e2005970, 2018.
- 320 ¹⁰ Ryan Poplin, Pi-Chuan Chang, David Alexander, Scott Schwartz, Thomas Colthurst, Alexander Ku, Dan Newburger,
321 Jojo Dijamco, Nam Nguyen, Pegah T Afshar, et al. A universal snp and small-indel variant caller using deep neural
322 networks. *Nature biotechnology*, 36(10):983–987, 2018.
- 323 ¹¹ Paul VC Hough. Method and means for recognizing complex patterns, December 18 1962. US Patent 3,069,654.
- 324 ¹² Christopher G Harris, Mike Stephens, et al. A combined corner and edge detector. In *Alvey vision conference*,
325 volume 15, pages 10–5244. Citeseer, 1988.
- 326 ¹³ D. G. Lowe. Object recognition from local scale-invariant features. In *Proceedings of the Seventh IEEE International*
327 *Conference on Computer Vision*, volume 2, pages 1150–1157 vol.2, Sep. 1999.
- 328 ¹⁴ David G. Lowe. Distinctive image features from scale-invariant keypoints. *Int. J. Comput. Vision*, 60(2):91–110,
329 November 2004.
- 330 ¹⁵ Briana D Ezray, Drew C Wham, Carrie E Hill, and Heather M Hines. Unsupervised machine learning reveals
331 mimicry complexes in bumblebees occur along a perceptual continuum. *Proceedings of the Royal Society B*,
332 286(1910):20191501, 2019.
- 333 ¹⁶ Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron C.
334 Courville, and Yoshua Bengio. Generative adversarial networks. *CoRR*, abs/1406.2661, 2014.

Decontextualized learning for interpretable hierarchical representations of visual patterns

- 335 ¹⁷ Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. In Yoshua Bengio and Yann LeCun, editors,
336 *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014,*
337 *Conference Track Proceedings*, 2014.
- 338 ¹⁸ Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate
339 inference in deep generative models. In Eric P. Xing and Tony Jebara, editors, *Proceedings of the 31st International*
340 *Conference on Machine Learning*, volume 32 of *Proceedings of Machine Learning Research*, pages 1278–1286,
341 Beijing, China, 22–24 Jun 2014. PMLR.
- 342 ¹⁹ Y. Bengio, A. Courville, and P. Vincent. Representation learning: A review and new perspectives. *IEEE Transactions*
343 *on Pattern Analysis and Machine Intelligence*, 35(8):1798–1828, Aug 2013.
- 344 ²⁰ Shengjia Zhao, Jiaming Song, and Stefano Ermon. Infovae: Balancing learning and inference in variational
345 autoencoders. 33:5885–5892, 2019.
- 346 ²¹ Arthur Gretton, Karsten Borgwardt, Malte Rasch, Bernhard Schölkopf, and Alex J Smola. A kernel method for the
347 two-sample-problem. In *Advances in neural information processing systems*, pages 513–520, 2007.
- 348 ²² Shengjia Zhao, Jiaming Song, and Stefano Ermon. Learning hierarchical features from deep generative models. In
349 *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 4091–4099. JMLR. org,
350 2017.
- 351 ²³ Xi Chen, Yan Duan, Rein Houthoofd, John Schulman, Ilya Sutskever, and Pieter Abbeel. Infogan: Interpretable
352 representation learning by information maximizing generative adversarial nets. In *Advances in neural information*
353 *processing systems*, pages 2172–2180, 2016.
- 354 ²⁴ Tian Qi Chen, Xuechen Li, Roger B Grosse, and David K Duvenaud. Isolating sources of disentanglement in
355 variational autoencoders. In *Advances in Neural Information Processing Systems*, pages 2610–2620, 2018.
- 356 ²⁵ Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution.
357 In *European Conference on Computer Vision*, 2016.
- 358 ²⁶ Leon A. Gatys, Alexander S. Ecker, and Matthias Bethge. Texture synthesis using convolutional neural networks. In
359 *Proceedings of the 28th International Conference on Neural Information Processing Systems - Volume 1, NIPS'15*,
360 pages 262–270, Cambridge, MA, USA, 2015. MIT Press.
- 361 ²⁷ Leon Gatys, Alexander Ecker, and Matthias Bethge. A neural algorithm of artistic style. *Journal of Vision*,
362 16(12):326–326, Sep 2016.
- 363 ²⁸ Irina Higgins, Loïc Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed,
364 and Alexander Lerchner. beta-vae: Learning basic visual concepts with a constrained variational framework. In *5th*
365 *International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference*
366 *Track Proceedings*, 2017.

Decontextualized learning for interpretable hierarchical representations of visual patterns

- 367 ²⁹Hyunjik Kim and Andriy Mnih. Disentangling by factorising. In Jennifer Dy and Andreas Krause, editors,
368 *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine*
369 *Learning Research*, pages 2649–2658, Stockholmsmässan, Stockholm Sweden, 10–15 Jul 2018. PMLR.
- 370 ³⁰Cian Eastwood and Christopher K. I. Williams. A framework for the quantitative evaluation of disentangled
371 representations. In *International Conference on Learning Representations*, 2018.
- 372 ³¹Grégoire Mesnil Yann Dauphin, Xavier Glorot, Salah Rifai, Yoshua Bengio, Ian Goodfellow, Erick Lavoie, Xavier
373 Muller, Guillaume Desjardins, David Warde-Farley, Pascal Vincent, et al. Unsupervised and transfer learning
374 challenge: a deep learning approach. In *Proceedings of ICML Workshop on Unsupervised and Transfer Learning*,
375 pages 97–110, 2012.
- 376 ³²Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In Doina Precup and
377 Yee Whye Teh, editors, *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney,*
378 *NSW, Australia, 6-11 August 2017*, volume 70 of *Proceedings of Machine Learning Research*, pages 3319–3328.
379 PMLR, 2017.
- 380 ³³Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej
381 Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual
382 Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015.
- 383 ³⁴Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*,
384 9(Nov):2579–2605, 2008.
- 385 ³⁵Dmitry Kobak and Philipp Berens. The art of using t-sne for single-cell transcriptomics. *Nature communications*,
386 10(1):1–14, 2019.
- 387 ³⁶Richard Dubes and Anil K. Jain. Clustering methodologies in exploratory data analysis. In *Advances in computers*,
388 volume 19, pages 113–228. Elsevier, 1980.
- 389 ³⁷Anne E Houde. Mate choice based upon naturally occurring color-pattern variation in a guppy population. *Evolution*,
390 41(1):1–10, 1987.
- 391 ³⁸David Ha and Jürgen Schmidhuber. Recurrent world models facilitate policy evolution. In *Advances in Neural*
392 *Information Processing Systems 31*, pages 2451–2463. Curran Associates, Inc., 2018. [https://worldmodels.](https://worldmodels.github.io)
393 [github.io](https://worldmodels.github.io).
- 394 ³⁹Yasaman Bahri, Jonathan Kadmon, Jeffrey Pennington, Sam S Schoenholz, Jascha Sohl-Dickstein, and Surya
395 Ganguli. Statistical mechanics of deep learning. *Annual Review of Condensed Matter Physics*, 2020.
- 396 ⁴⁰Naftali Tishby, Fernando C. Pereira, and William Bialek. The information bottleneck method. pages 368–377, 1999.
- 397 ⁴¹Darrell J Kemp. Female mating biases for bright ultraviolet iridescence in the butterfly eureka hecabe (pieridae).
398 *Behavioral Ecology*, 19(1):1–8, 2008.

Decontextualized learning for interpretable hierarchical representations of visual patterns

- 399 ⁴² Sergi Caelles, Kevis-Kokitsi Maninis, Jordi Pont-Tuset, Laura Leal-Taixé, Daniel Cremers, and Luc Van Gool.
400 One-shot video object segmentation. In *Proceedings of the IEEE conference on computer vision and pattern*
401 *recognition*, pages 221–230, 2017.
- 402 ⁴³ Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. 2015.
- 403 ⁴⁴ Joseph Redmon, Santosh Kumar Divvala, Ross B. Girshick, and Ali Farhadi. You only look once: Unified, real-time
404 object detection. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas,*
405 *NV, USA, June 27-30, 2016*, pages 779–788. IEEE Computer Society, 2016.
- 406 ⁴⁵ Philip Bachman. An architecture for deep, hierarchical generative models. In *Advances in Neural Information*
407 *Processing Systems*, pages 4826–4834, 2016.
- 408 ⁴⁶ Casper Kaae Sønderby, Tapani Raiko, Lars Maaløe, Søren Kaae Sønderby, and Ole Winther. Ladder variational
409 autoencoders. In *Advances in neural information processing systems*, pages 3738–3746, 2016.
- 410 ⁴⁷ Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *Proceedings of the IEEE conference on*
411 *computer vision and pattern recognition*, pages 7132–7141, 2018.
- 412 ⁴⁸ Yanghao Li, Naiyan Wang, Jiaying Liu, and Xiaodi Hou. Demystifying neural style transfer. In Carles Sierra, editor,
413 *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI 2017, Melbourne,*
414 *Australia, August 19-25, 2017*, pages 2230–2236. ijcai.org, 2017.
- 415 ⁴⁹ Arthur Gretton, Dino Sejdinovic, Heiko Strathmann, Sivaraman Balakrishnan, Massimiliano Pontil, Kenji Fukumizu,
416 and Bharath K Sriperumbudur. Optimal kernel choice for large-scale two-sample tests. In *Advances in neural*
417 *information processing systems*, pages 1205–1213, 2012.
- 418 ⁵⁰ Daniel Im Jiwoong Im, Sungjin Ahn, Roland Memisevic, and Yoshua Bengio. Denoising criterion for variational
419 auto-encoding framework. In *Thirty-First AAAI Conference on Artificial Intelligence*, 2017.
- 420 ⁵¹ Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization for generative
421 adversarial networks. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC,*
422 *Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018.
- 423 ⁵² John A Endler and Anne E Houde. Geographic variation in female preferences for male traits in *poecilia reticulata*.
424 *Evolution*, 49(3):456–468, 1995.
- 425 ⁵³ Henry Walter Bates. *The naturalist on the river Amazons*. London; Toronto: JM Dent; New York: EP Dutton,[1910,
426 reprinted 1921], 1863.
- 427 ⁵⁴ Charles Darwin. *The descent of man: and selection in relation to sex*. John Murray, Albemarle Street., 1871.
- 428 ⁵⁵ Alfred Russel Wallace. The colors of animals and plants. *The American Naturalist*, 11(11):641–662, 1877.
- 429 ⁵⁶ Fritz Müller. *Über die vorteile der mimicry bei schmetterlingen*. 1878.
- 430 ⁵⁷ Edward Bagnall Poulton. *The colours of animals: their meaning and use, especially considered in the case of insects*.
431 D. Appleton, 1890.

Decontextualized learning for interpretable hierarchical representations of visual patterns

- 432 ⁵⁸ Gerald Handerson Thayer. *Concealing-coloration in the animal kingdom: an exposition of the laws of disguise*
433 *through color and pattern: being a summary of Abbott H. Thayer's discoveries*. Macmillan Company, 1918.
- 434 ⁵⁹ Hugh Bamford Cott. Adaptive coloration in animals. 1940.
- 435 ⁶⁰ John Maynard Smith. Natural selection and the concept of a protein space. *Nature*, 225(5232):563–564, 1970.
- 436 ⁶¹ Jakob Von Uexküll. A stroll through the worlds of animals and men: A picture book of invisible worlds. *Semiotica*,
437 89(4):319–391, 1992.
- 438 ⁶² Cedric P van den Berg, Jolyon Troscianko, John A Endler, N Justin Marshall, and Karen L Cheney. Quantitative
439 colour pattern analysis (qcpa): A comprehensive framework for the analysis of colour patterns in nature. *Methods in*
440 *Ecology and Evolution*, 11(2):316–332, 2020.
- 441 ⁶³ Rafael Maia, Hugo Gruson, John A. Endler, and Thomas E. White. pavo 2: New tools for the spectral and spatial
442 analysis of colour in r. *Methods in Ecology and Evolution*, 0(0), 2019.
- 443 ⁶⁴ Mary Caswell Stoddard, Rebecca M. Kilner, and Christopher Town. Pattern recognition algorithm reveals how birds
444 evolve individual egg pattern signatures. *Nature Communications*, 5:4117 EP –, Jun 2014. Article.
- 445 ⁶⁵ Felipe M. Gawryszewski. Color vision models: Some simulations, a general n-dimensional model, and the
446 colourvision r package. *Ecology and Evolution*, 8(16):8159–8170, 2018.
- 447 ⁶⁶ Cynthia Tedore and Sönke Johnsen. Using RGB displays to portray color realistic imagery to animal eyes. *Current*
448 *Zoology*, 63(1):27–34, 06 2016.
- 449 ⁶⁷ John A. Endler. Variation in the appearance of guppy color patterns to guppies and their predators under different
450 visual conditions. *Vision Research*, 31(3):587 – 608, 1991.
- 451 ⁶⁸ John A. Endler and Paul W. Mielke JR. Comparing entire colour patterns as birds see them. *Biological Journal of*
452 *the Linnean Society*, 86(4):405–431, 2005.
- 453 ⁶⁹ John A. Endler. A framework for analysing colour pattern geometry: adjacent colours. *Biological Journal of the*
454 *Linnean Society*, 107(2):233–253, 2012.
- 455 ⁷⁰ Jolyon Troscianko and Martin Stevens. Image calibration and analysis toolbox – a free software suite for objectively
456 measuring reflectance, colour and pattern. *Methods in Ecology and Evolution*, 6(11):1320–1331, 2015.
- 457 ⁷¹ Eleanor M. Caves and Sönke Johnsen. Acuityview: An r package for portraying the effects of visual acuity on
458 scenes observed by an animal. *Methods in Ecology and Evolution*, 9(3):793–797, 2018.
- 459 ⁷² John A. Endler, Gemma L. Cole, and Alexandra M. Kranz. Boundary strength analysis: Combining colour pattern
460 geometry and coloured patch visual properties for use in predicting behaviour and fitness. *Methods in Ecology and*
461 *Evolution*, 9(12):2334–2348, 2018.
- 462 ⁷³ Mary Caswell Stoddard and Daniel Osorio. Animal coloration patterns: Linking spatial vision to quantitative
463 analysis. *The American Naturalist*, 193(2):164–186, 2019.
- 464 ⁷⁴ H.F. Nijhout. Elements of butterfly wing patterns. *Journal of Experimental Zoology*, 291(3):213–225, 2001.

Decontextualized learning for interpretable hierarchical representations of visual patterns

- 465 ⁷⁵ G Th Fechner. Ueber die subjectiven nachbilder und nebenbilder. *Annalen der Physik*, 126(7):427–470, 1840.
- 466 ⁷⁶ Ray Fuller and Jorge A Santos. *Human factors for highway engineers*. Pergamon Amsterdam, The Netherlands,
467 2002.
- 468 ⁷⁷ Gemma L Cole and John A Endler. Male courtship decisions are influenced by light environment and female
469 receptivity. *Proceedings of the Royal Society B: Biological Sciences*, 283(1839):20160861, 2016.
- 470 ⁷⁸ Andreas Nieder, David J Freedman, and Earl K Miller. Representation of the quantity of visual items in the primate
471 prefrontal cortex. *Science*, 297(5587):1708–1711, 2002.
- 472 ⁷⁹ S Fujita, T Kitayama, N Mizoguchi, Y Oi, N Koshikawa, and M Kobayashi. Spatiotemporal profiles of transcallosal
473 connections in rat insular cortex revealed by in vivo optical imaging. *Neuroscience*, 206:201–211, 2012.
- 474 ⁸⁰ Scott Cheng-Hsin Yang, Mate Lengyel, and Daniel M Wolpert. Active sensing in the categorization of visual
475 patterns. *Elife*, 5:e12215, 2016.
- 476 ⁸¹ Laura A Kelley and Jennifer L Kelley. Animal visual illusion and confusion: the importance of a perceptual
477 perspective. *Behavioral Ecology*, 25(3):450–463, 2014.
- 478 ⁸² Sami Merilaita, Nicholas E. Scott-Samuel, and Innes C. Cuthill. How camouflage works. *Philosophical Transactions
479 of the Royal Society B: Biological Sciences*, 372(1724):20160341, 2017.
- 480 ⁸³ Clelia Gasparini, Giovanna Serena, and Andrea Pilastro. Do unattractive friends make you look better? context-
481 dependent male mating preferences in the guppy. *Proceedings of the Royal Society B: Biological Sciences*,
482 280(1756):20123072, 2013.
- 483 ⁸⁴ David Marr. Vision: A computational investigation into the human representation and processing of visual
484 information, henry holt and co. Inc., New York, NY, 2(4.2), 1982.
- 485 ⁸⁵ David H Hubel and Torsten N Wiesel. Receptive fields, binocular interaction and functional architecture in the cat's
486 visual cortex. *The Journal of physiology*, 160(1):106–154, 1962.
- 487 ⁸⁶ Diego E Pafundo, Mark A Nicholas, Ruilin Zhang, and Sandra J Kuhlman. Top-down-mediated facilitation in the
488 visual cortex is gated by subcortical neuromodulation. *Journal of Neuroscience*, 36(10):2904–2914, 2016.
- 489 ⁸⁷ Spencer J Ingley, Mohammad Rahmani Asl, Chengde Wu, Rongfeng Cui, Mahmoud Gadelhak, Wen Li, Ji Zhang,
490 Jon Simpson, Chelsea Hash, Trisha Butkowski, et al. anyfish 2.0: an open-source software platform to generate and
491 share animated fish models to study behavior. *SoftwareX*, 3:13–21, 2015.
- 492 ⁸⁸ John R Stowers, Maximilian Hofbauer, Renaud Bastien, Johannes Griessner, Peter Higgins, Sarfarazhussain Farooqui,
493 Ruth M Fischer, Karin Nowikovsky, Wulf Haubensak, Iain D Couzin, et al. Virtual reality for freely moving animals.
494 *Nature methods*, 14(10):995, 2017.
- 495 ⁸⁹ H. Naik, R. Bastien, N. Navab, and I. D. Couzin. Animals in virtual environments. *IEEE Transactions on
496 Visualization and Computer Graphics*, 26(5):2073–2083, 2020.

Decontextualized learning for interpretable hierarchical representations of visual patterns

- 497 ⁹⁰ Kuo-Hua Huang, Peter Rupperecht, Thomas Frank, Koichi Kawakami, Tewis Bouwmeester, and Rainer W Friedrich.
498 A virtual reality system to analyze neural activity and behavior in adult zebrafish. *Nature Methods*, pages 1–9, 2020.
- 499 ⁹¹ John R. Stowers, Anton L. Fuhrmann, Maximilian Hofbauer, Martin Streinzer, Axel Schmid, Michael H. Dickinson,
500 and Andrew D. Straw. Reverse engineering animal vision with virtual reality and genetics. *Computer*, 47:38–45,
501 2014.
- 502 ⁹² Charles Darwin. *On the origin of species*. John Murray, 1859.
- 503 ⁹³ Ronald A Fisher. The evolution of sexual preference. *The Eugenics Review*, 7(3):184, 1915.
- 504 ⁹⁴ Russell Lande. Models of speciation by sexual selection on polygenic traits. *Proceedings of the National Academy
505 of Sciences*, 78(6):3721–3725, 1981.
- 506 ⁹⁵ Mark Kirkpatrick. Sexual selection and the evolution of female choice. *Evolution*, 36(1):1–12, 1982.
- 507 ⁹⁶ Yoh Iwasa and Andrew Pomiankowski. The evolution of mate preferences for multiple sexual ornaments. *Evolution*,
508 48(3):853–867, 1994.
- 509 ⁹⁷ Richard O Prum. The lande–kirkpatrick mechanism is the null model of evolution by intersexual selection:
510 implications for meaning, honesty, and design in intersexual signals. *Evolution: International Journal of Organic
511 Evolution*, 64(11):3085–3100, 2010.
- 512 ⁹⁸ Eleanor M. Caves, Nicholas C. Brandley, and Sönke Johnsen. Visual acuity and the evolution of signals. *Trends in
513 Ecology Evolution*, 33(5):358 – 372, 2018.
- 514 ⁹⁹ Mathieu Joron and James LB Mallet. Diversity in mimicry: paradox or paradigm? *Trends in Ecology & Evolution*,
515 13(11):461–466, 1998.
- 516 ¹⁰⁰ William A Searcy and Stephen Nowicki. *The evolution of animal communication: reliability and deception in
517 signaling systems*. Princeton University Press, 2005.
- 518 ¹⁰¹ Derek A Roff. The evolution of mate choice: a dialogue between theory and experiment. *Annals of the New York
519 Academy of Sciences*, 1360(1):1–15, 2015.
- 520 ¹⁰² John A Endler. Predation, light intensity and courtship behaviour in poecilia reticulata (pisces: Poeciliidae). *Animal
521 Behaviour*, 35(5):1376–1385, 1987.

522 **Appendix A Example Application to the Evolution of Color Patterns: Background**

523 The incredible variety of color patterns seen in nature evolved under the selective forces imposed by the environment, and
524 the visual experience of their receivers.^{53–59} Quantifying this diversity, and reliably testing the functional significance of
525 these traits is fundamental to understanding fitness landscapes⁶⁰ and underlies many subdisciplines of sensory ecology,
526 cognitive neuroscience, collective behavior, and evolution.

527 Creating quantitative descriptions of color patterns which take into account the unique sensory and semiotic worlds
528 of their receivers⁶¹ has been a central challenge in visual ecology. Many tools have been developed: Quantitative

Decontextualized learning for interpretable hierarchical representations of visual patterns

529 Colour Pattern Analysis,⁶² PAVO,⁶³ Natural Pattern Match,⁶⁴ among others.^{65–73} Each of these tools uses one or an
530 ensemble of complimentary metrics from image analysis and computer vision, e.g. image statistics, edge detection, and
531 landmark-based filters.¹⁴

532 Still, fundamental gaps remain. One of these gaps is the difficulty in building quantitative descriptions of complex
533 features with multiple subelements. Most existing approaches fail to capture the full complexity of many of color
534 patterns; the algorithms themselves are insufficiently expressive. This is particularly true when spatial or scale
535 dependent relationships between features exist, e.g. the irregular patterns of male guppy ornamentation or butterfly
536 wing patterns where similar sets of elements are arranged in species-specific configurations.⁷⁴ Recently, researchers
537 have begun employing machine learning algorithms such as non-linear dimensionality reduction, e.g. t-distribute
538 stochastic neighbor embedding (t-SNE,^{34,35} Figure 2), and deep neural networks (Figure 2,¹⁵ Figure S1). Still, while
539 these techniques can better represent more complex relationships between pixel values within an image, current
540 implementations do not disentangle features across scales or provide extensions to downstream experiments.

541 While complex trait may be difficult to quantify, they are nonetheless biologically relevant in terms of feature context^{75–80}
542 and the perceptions of shape, motion, and attention.^{57–59,81–83} And in the brain, we know that perception is hierarchically
543 organized,⁸⁴ and representations made at higher levels of the visual cortex and its homologs heavily influence the
544 perception of low-level features.^{85,86} While measuring local features across an image provides important insight on
545 regularity and the nature of wide-field variation, a collection of local feature descriptions across space is fundamentally
546 different to a feature description built across scales.

547 Another gap is in building direct connections between approaches. Establishing spectral sensitivity, acuity, and feature
548 importance is typically done using stimulus playback experiments or behavioral assays. However, beyond using
549 statistical descriptions of features to guide researchers in the creation of stimuli there are few explicit connections
550 between analysis and experiment. The current state of the art: immersive virtual reality (VR) and low-latency playback
551 experiment—with fully animated, photo-realistic, 3D models, provide a rich experimental basis for investigating
552 the relationship between visual inputs, neural activity, and behavior.^{87–90} VR systems are also beginning to better
553 account for species-specific sensory biases including photoreceptor sensitivity, flicker fusion rate, acuity, and depth
554 perception.^{89,91} Still, currently these approaches rely on human-in-the-loop interventions for creating stimulus with
555 even moderate complexity.

556 Additionally, because color pattern traits have evolved under selective pressure from multiple receivers, establishing
557 these types of evolutionary trade-offs is important to our understanding. However, experimental approaches often require
558 large, highly disruptive manipulations such as translocation experiments or large scale crossbreeding experiments.
559 Simulations and virtual experiments may better allows researchers to be explicit about the stimulus that is being tested
560 and greatly reduce the number of subjects needed (Methods 4.4).

561 **A.1 The potential impacts of this approach on the study of evolution**

562 This platform may be used to address many outstanding questions regarding the functional significance of color pattern
563 traits; here, we discuss some of these questions. 1) What are the constraints on the evolvability of a given trait? By
564 identifying the topographical relationship between different traits within the color pattern space we can test predictions
565 about the selective forces acting on them related to their geometric relationships, e.g. the axes of variation in traits meant
566 to communicate viability should show increased orthogonality compared to co-occurring traits which have evolved
567 under a Fisherian process.⁹²⁻⁹⁷ 2) Categorical perception is an important perceptual mechanism for understanding the
568 evolution of color signals.⁹⁸ But in systems where color patterns are used for mimicry^{53,55,99} or novelty, investigating
569 the boundaries between complex traits is fundamental. By performing traversals across the distribution of the latent
570 variables, interpolating between samples can allow for tests of continuous¹⁰⁰ versus categorical perception¹⁰¹ of complex
571 traits. 3) Many color pattern traits have evolved under selective pressure from multiple receivers, e.g. both females and
572 predators shape the diversity of male guppy ornaments.¹⁰² Establishing these types of evolutionary trade-offs is difficult
573 and often requires large, highly disruptive manipulations such as translocation experiments. Using evolutionary models
574 similar to the ones presented here researchers can simulate multiple fitness landscapes and evolutionary trajectories
575 simultaneously to perform a broad range of virtual experiments. Importantly, while each of these examples place either
576 analytical, experimental, or virtual results at the center, by using the platform presented here, they maintain direct
577 connections across approaches. Furthermore, they can incorporate existing techniques⁶⁷⁻⁷³ as image preprocessing
578 routines, during playback, or constraints on virtual experiments.

579 **Appendix B Supplemental Figures**

Decontextualized learning for interpretable hierarchical representations of visual patterns

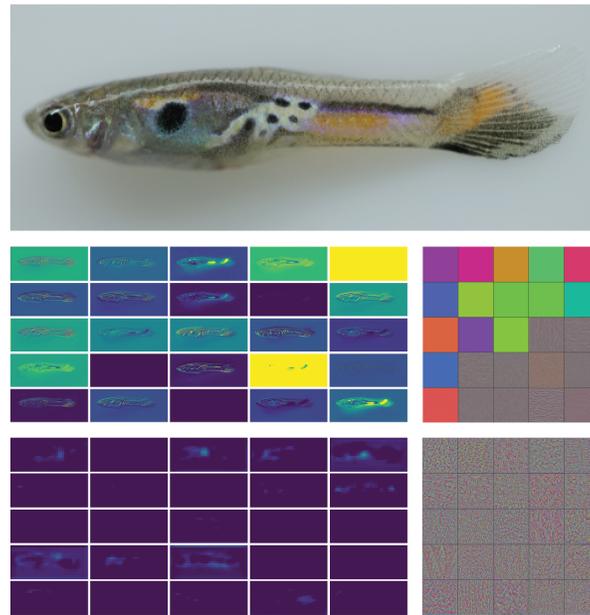


Figure S1: *Convolutional layers*. In typical *supervised* discriminative models, the objective being optimized is well defined, e.g. accurate classification or localization. As such, the representations provided by the downstream convolutional layers of deep networks take on characteristics optimized for task performance. At higher and higher network layers, the boundaries between classes can become complex and specialized to this objective because of the usefulness of such representations to identifying complex boundaries. Middle: Features learned at lower layers relate to color patches or gestures whereas at higher levels (bottom) features become complex and interpretability can be difficult. Left, image pixels which are activated by pretrained image filters (yellow represents higher activations). Right, the maximally activating image feature for each filter.

Decontextualized learning for interpretable hierarchical representations of visual patterns

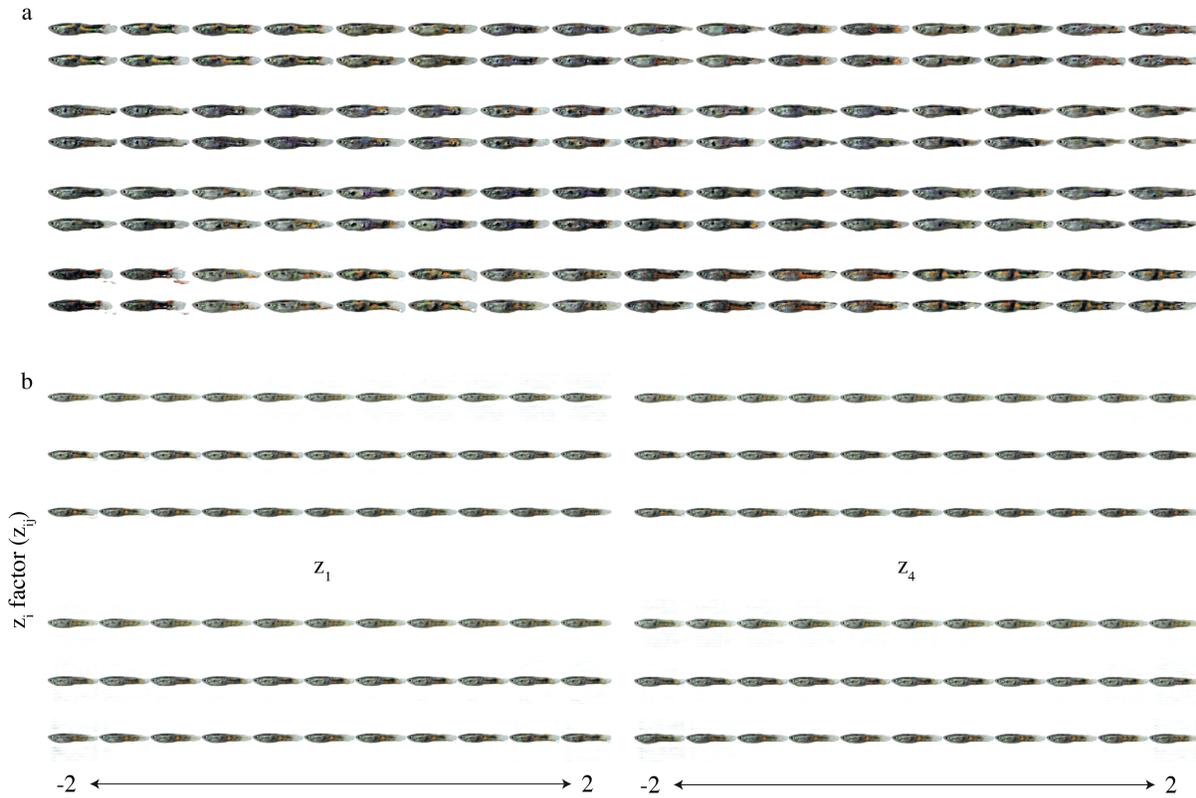


Figure S2: *Exploring latent variables.* a) We incorporated knowledge of our sample data by providing a 32-class categorical latent code and were able to generate examples from distinct classes learned by the model which capture meaningful combinations of features in our sample data. b) Latent traversal of 3 latent variables from z_1, \dots, z_4 . Top and middle are two embedded samples and bottom a latent code initialized at zero. For each latent variable (rows) we traverse values between -2 and 2 for the generated output. We see that each latent code has consistent effects. All samples in both a and b are generated from the generative models (b and c in Figure 1)

Decontextualized learning for interpretable hierarchical representations of visual patterns

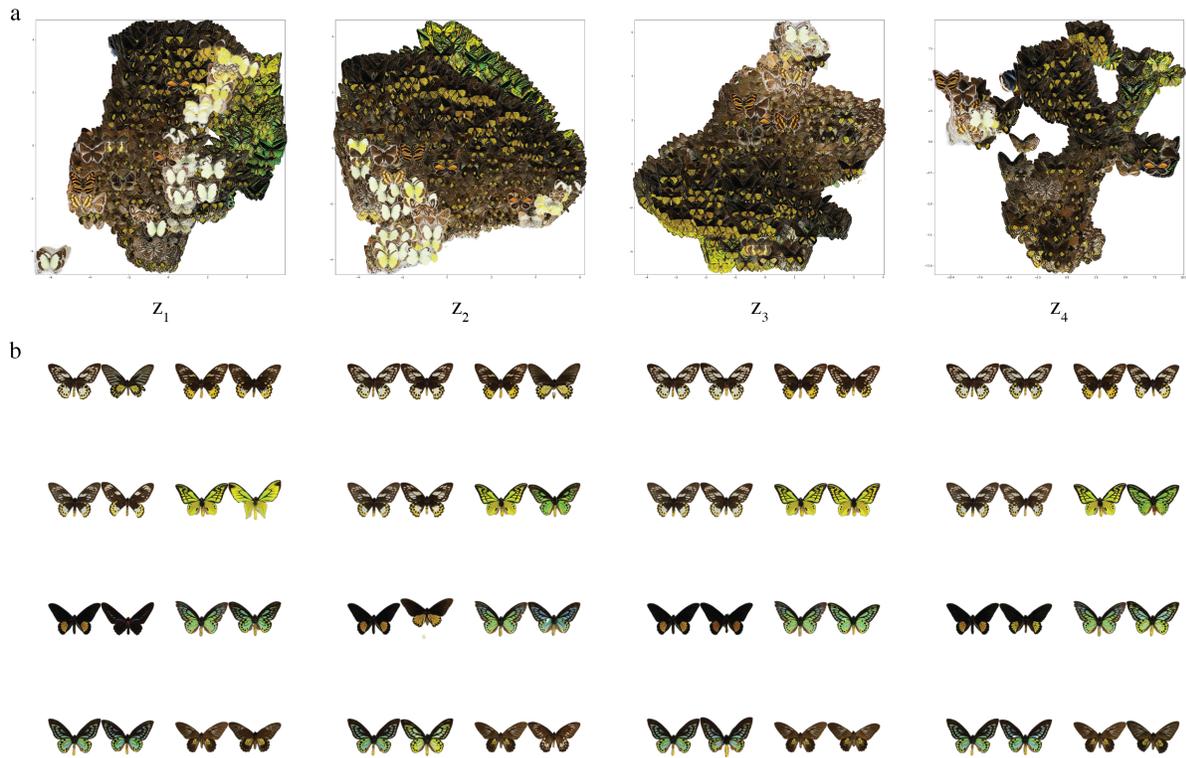


Figure S3: *Exploring latent representations.* a) 2D embedding (using tSNE^{34,35}) of butterfly images using the 4 hierarchical latent encodings. The relationship between images at lower levels are dominated by color value similarities whereas at higher layers pattern elements at increasing spatial scales define the relationship between samples b) Nearest neighbors of 8 random samples based in the Minkowski distance³⁶ between the 10-dimensional space of each latent code z_1, \dots, z_4

Decontextualized learning for interpretable hierarchical representations of visual patterns

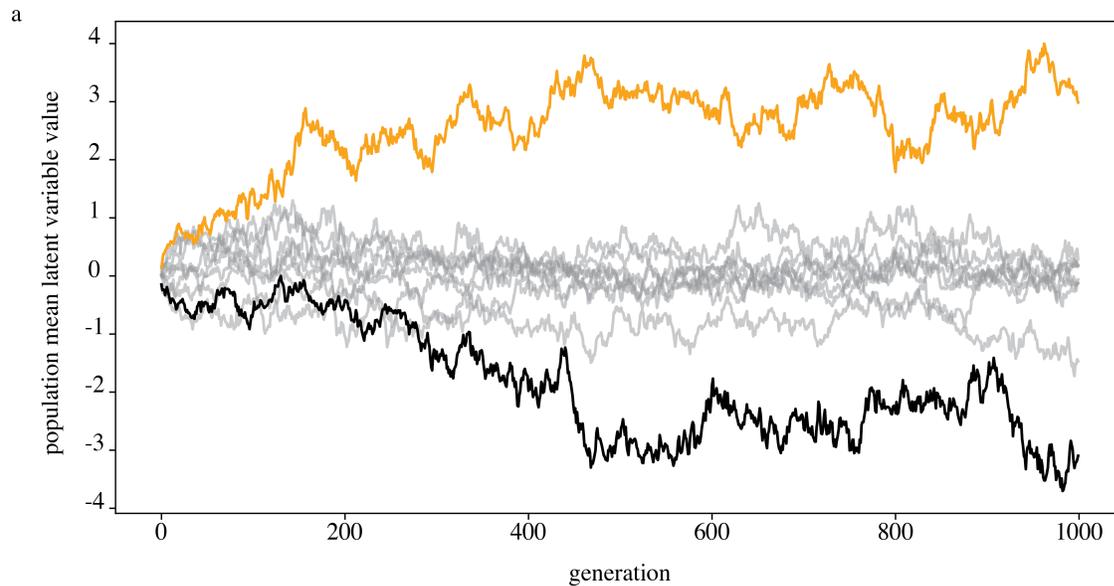


Figure S4: *Latent variable "alleles" over generations* a) We find two alleles are driven to fixation in the population after several generations selecting for oranger and higher contrast color patterns in guppies.

Figure S5: *Movie 1: The combined pattern space over 500 generations, visualized in 2D using tSNE*

Figure S6: *Movie 2: VR animation of learned coloration pattern models to an animated guppy for virtual playback.*