

# 1 **Computer vision, machine learning, and the promise** 2 **of phenomics in ecology and evolutionary biology**

3  
4 Moritz Lürig<sup>1</sup>; Seth Donoughe<sup>2</sup>; Erik I. Svensson<sup>1</sup>; Arthur Porto<sup>3,4</sup>; Masahito Tsuboi<sup>1</sup>

5 <sup>1</sup> Department of Biology, Lund University, 22362 Lund, Sweden

6 <sup>2</sup> Department of Molecular Genetics and Cell Biology, University of Chicago, Chicago, US

7 <sup>3</sup> Department of Biological Sciences, Louisiana State University, Baton Rouge, US

8 <sup>4</sup> Center for Computation & Technology, Louisiana State University, Baton Rouge, US

## 9 **Abstract**

10 For centuries, ecologists and evolutionary biologists have used images such as drawings,  
11 paintings, and photographs to record and quantify the shapes and patterns of life. With the  
12 advent of digital imaging, biologists continue to collect image data at an ever-increasing rate.  
13 This immense body of data provides insight into a wide range of biological phenomena,  
14 including phenotypic trait diversity, population dynamics, mechanisms of divergence and  
15 adaptation and evolutionary change. However, the rate of image acquisition frequently  
16 outpaces our capacity to manually extract meaningful information from the images. Moreover,  
17 manual image analysis is low-throughput, difficult to reproduce, and typically measures only a  
18 few traits at a time. This has proven to be an impediment to the growing field of phenomics -  
19 the study of many phenotypic dimensions together. Computer vision (CV), the automated  
20 extraction and processing of information from digital images, is a way to alleviate this  
21 longstanding analytical bottleneck. In this review, we illustrate the capabilities of CV for fast,  
22 comprehensive, and reproducible image analysis in ecology and evolution. First, we briefly  
23 review phenomics, arguing that ecologists and evolutionary biologists can most effectively  
24 capture phenomic-level data by using CV. Next, we describe the primary types of image-based  
25 data, and review CV approaches for extracting them (including techniques that entail machine  
26 learning and others that do not). We identify common hurdles and pitfalls, and then highlight  
27 recent successful implementations of CV in the study of ecology and evolution. Finally, we  
28 outline promising future applications for CV in biology. We anticipate that CV will become a

29 basic component of the biologist's toolkit, further enhancing data quality and quantity, and  
30 sparking changes in how empirical ecological and evolutionary research will be conducted.

### 31 **From phenotypes to phenomics: measuring traits at scale**

32 Faced with the overwhelming complexity of the living world, most life scientists confine their  
33 efforts to a small set of observable traits. Although a drastic simplification of organismal  
34 complexity, the focus on single phenotypic attributes often provides a tractable, operational  
35 approach to understand biological phenomena, e.g. phenotypic trait diversity, population  
36 dynamics, mechanisms of divergence and adaptation and evolutionary change. However,  
37 there are also obvious limitations in how much we can learn from studying small numbers of  
38 *phenotypes* in isolation. Evolutionary and conservation biologist Michael Soulé was one of the  
39 first to demonstrate the value of collecting and analyzing many phenotypes at once in his early  
40 study of the side-blotched lizard (*Uta stansburiana*; [Soulé 1967]; reviewed in Houle et al. ).  
41 While doing so, he defined the term "*phenome*" as "the phenotype as a whole" (Soulé 1967).  
42 *Phenomics*, by extension, is the comprehensive study of phenomes. In practice, this entails  
43 collecting and analyzing multidimensional phenotypes with a wide range of quantitative and  
44 high-throughput methods (Houle et al. 2010, Bilder et al. 2009). Given that biologists are now  
45 attempting to understand increasingly complex and high dimensional relationships (Walsh  
46 2007), it is surprising that phenomics still remains underutilized (Fig. 1), both as  
47 methodological approach and as an overarching conceptual and analytical framework (Houle  
48 et al. 2010).

49 Phenomic datasets are essential if we are to understand some of the most compelling  
50 but challenging questions in the study of ecology and evolution. For instance, natural selection  
51 typically does not operate on single traits, but on multiple traits simultaneously (Lande and  
52 Arnold 1983, Phillips and Arnold 1999b). Such correlational selection can bias evolutionary  
53 change and shape genetic covariance patterns by building up and maintaining linkage  
54 disequilibrium (Schluter 1996, Phillips and Arnold 1999a, Sinervo and Svensson 2002,

55 Svensson et al. 2021), and it is a phenomenon that can most effectively be uncovered in  
56 multidimensional phenotypic datasets. Another example is pleiotropy, which generates  
57 patterns of covariation among traits that are impossible to predict if only a few simple traits are  
58 measured (Visscher and Yang 2016, Saltz et al. 2017). Phenotypic plasticity, which is  
59 increasingly recognized in mediating evolutionary trajectories (Pfennig et al. 2010), is also an  
60 inherently multivariate phenomenon involving many traits and interactions between traits, so  
61 it should be quantified as such (Morel-Journel et al. 2020). Finally, community stability  
62 depends on species interactions and ecological niches of organisms; niches are typically  
63 influenced by many traits at once (Blonder 2018, Laughlin et al. 2020). Put simply: if we are to  
64 draw a complete picture of biological processes and aim to understand their causal  
65 relationships at various levels of biological organization, we need to measure more traits.  
66 Phenomic datasets will make our conclusions and inferences more robust if underpinned by  
67 more complete information without systematic biases.

68 High dimensional phenotypic data are also needed for uncovering the causal links  
69 between genotypes, environmental factors, and phenotypes, i.e. to understand the genotype-  
70 phenotype map (Houle et al. 2010, Orgogozo et al. 2015). The advent of genomics - high  
71 throughput molecular methods to analyze the structure, function or evolution of an organism's  
72 genome in parts or as a whole (Church and Gilbert 1984, Feder and Mitchell-Olds 2003) - has  
73 already improved our understanding of many biological phenomena. This includes the  
74 emergence and maintenance of biological diversity (Seehausen et al. 2014), the inheritance  
75 and evolution of complex traits (Pitchers et al. 2019), and the evolutionary origin of key  
76 metabolic traits (Ishikawa et al. 2019). Thus, accessible molecular tools have lowered the  
77 hurdles for discovery-based genomic research and shifted the focus away from the study of  
78 observable organismal traits and phenotypes towards their molecular basis. However, a  
79 similar “moonshot-program” for the phenotype, i.e. an ensemble of phenomics methods that  
80 matches genomics in their comprehensiveness, is still lacking (Freimer and Sabatti 2003). The  
81 growing mismatch in how efficiently molecular and phenotypic data are collected may hamper

82 further scientific progress in ecological and evolutionary research (Houle et al. 2010, Orgogozo  
83 et al. 2015, Lamichhane et al. 2019).

84       Following previous calls for phenomic research programs (Bilder et al. 2009, Houle et  
85 al. 2010, Furbank and Tester 2011), some recent studies have collected phenotypic data with  
86 high dimensionality, for example, in plants (Ubbens and Stavness 2017), animals (Cheng et  
87 al. 2011, Kühl and Burghardt 2013, Pitchers et al. 2019) and microbes (Zackrisson et al. 2016,  
88 French et al. 2018). The methods in such studies included 2D- and 3D-scanners, camera  
89 traps, robotic imaging platforms, computational image analysis, morphometrics,  
90 transcriptomics, metabolomics, and automated data loggers to record physiological and  
91 behavioral data from organisms (Houle et al. 2010). Many of these techniques produce image-  
92 based data. In general, ecologists and evolutionary biologists use digital imaging to quantify  
93 the external phenotype of an organism (i.e. its visually observable characteristics), to count  
94 organisms (e.g. cells on microscope slides), or to detect the presence of an organism (e.g. in  
95 images collected by camera traps). Existing work has supplied us with an immense body of  
96 image data that has provided insight into a wide range of biological phenomena, yet, when  
97 biologists manually extract phenotypes from images for phenomic-scale research, they  
98 confront several main bottlenecks.

99       A major constraint when working with large amounts of images (~1000 or more) is  
100 processing time and cost. Manual extraction of phenotypic data from images is slow and it  
101 requires trained domain experts whose work is extremely expensive. Moreover, the collection  
102 of such metrics in a manual fashion entails subjective decisions by the researcher, which may  
103 make it prone to error, and certainly makes reproducibility difficult. Last, manually measured  
104 traits tend to be low-dimensional measurements of higher dimensional traits. For example,  
105 external color traits, such as eye color phenotypes, are often scored as discrete categories  
106 (e.g. red vs blue phenotypes), whereas *pixel* level information (number of red vs. blue pixels)  
107 can provide a continuous phenotypic metric (Liu et al. 2010). Such quantitative, high-  
108 dimensional data can provide insight into previously hidden axes of variation in natural  
109 phenotypes. In this review we extol *computer vision* (CV), the automatic extraction of

110 meaningful information from images, as a promising toolbox to collect phenotypic information  
111 on a massive scale. The field has blossomed in recent years, producing a diverse array of  
112 computational tools to increase analytic efficiency, data dimensionality, and reproducibility.  
113 We argue that CV is poised to become a basic component of the data analysis toolkit in  
114 ecology and evolution, enabling researchers to collect and explore phenomic-scale data.

## 115 **Digital images as data**

### 116 **The structure of digital images**

117 A two-dimensional image is an intuitive way to record, store, and analyze organismal  
118 phenotypes. In the pre-photography era, ecologists and evolutionary biologists used drawings  
119 to capture the shapes and patterns of life, later to be replaced by analog photography, which  
120 allowed for qualitative assessment and simple, often only qualitative analysis of phenotypic  
121 variation. With the advent of digital photography, biologists could collect phenotypic data at  
122 unprecedented rates using camera stands, camera traps, microscopes, scanners, video  
123 cameras, or any other instrument with semiconductor image sensors (hereafter “image  
124 sensors”). Image sensors produce two-dimensional raster images (also known as bitmap  
125 images), which store incoming visible light or other electromagnetic signals into discrete,  
126 locatable picture elements - in short: pixels (Fig. 2). Each pixel contains quantitative  
127 phenotypic information that is organized as an array of rows and columns, whose dimensions  
128 are also referred to as “pixel resolution” or just “resolution”. An image with 1000 rows and 1500  
129 columns has a resolution of 1000 x 1500 (= 1 500 000 pixels, or 1.5 megapixels). The same  
130 applies for digital videos, which are simply a series of digital images displayed in succession,  
131 where the frame rate (measured as frames per second = fps) describes the speed of that  
132 succession.

133         At the pixel level, images or video frames can store variable amounts of information,  
134 depending on the *bit depth*, which refers to the number of distinct values that a pixel can  
135 represent (Fig. 2). In binary images, pixels contain information as a single bit, which can take

136 exactly two values - typically black or white ( $2^1$  values = 2 intensity values). Grayscale images  
137 from typical consumer cameras have a bit depth of 8, thus each pixel can take a value between  
138 0-255 ( $2^8$  values = 256 intensity values), which typically represents a level of light intensity,  
139 also referred to as pixel intensity. Color images are typically composed of at least three sets  
140 of pixel arrays, also referred to as channels, each of which contain values for either red, green  
141 or blue (RGB; Fig. 2). Each channel, when extracted from an RGB image, is a grayscale  
142 representation of the intensities for a single-color channel. Through the combination of pixel  
143 values at each location into triplets, colors are numerically represented. Today the industrial  
144 standard for color images is 24-bit depth, in which each color channel has a bit depth of 8 and  
145 thus can represent 256 colors (Fig. 1). Thus, 24-bit RGB images can represent over 16 million  
146 color variations in each pixel ( $2^{24} = 256 \times 256 \times 256 = 16\,777\,216$  intensity values), which  
147 already greatly surpasses the estimated 2.28 million of color variations that humans can  
148 perceive (Pointer and Attridge 1998).

149 Today, high resolution image sensors are an affordable way to store externally visible  
150 phenotypic information, like color and shape. However, advanced image sensors can also  
151 combine information from different spectra other than the visible light, like infrared radiation,  
152 which can be used to quantify individual body temperatures. With thermal image sensors,  
153 biologists can estimate body surface temperatures, which are correlated with internal (core)  
154 body temperatures (Tattersall and Cadena 2010), particularly in small animals like insects  
155 (Tsubaki et al. 2010, Svensson et al. 2020). Thermal imaging, or thermography, offers new  
156 opportunities for ecophysiological evolutionary research of how animals cope with heat or cold  
157 stress in their natural environments (Fig. 4B; Tattersall et al. 2009, Tattersall and Cadena  
158 2010, Svensson and Waller 2013). Fluorescence spectroscopy is another way to quantify  
159 phenotypes in high throughput and with high detail. For example, plate readers typically used  
160 in microbial and plankton research, can combine light in the visible spectrum with images  
161 containing information of cell fluorescence or absorbance to an “image stack” (Roeder et al.  
162 2012). Image stacks and the inclusion of multiple spectral channels provide a promising

163 avenue of research towards capturing a more complete representation of the phenotype (Fig.  
164 4A; Hense et al. 2008, Di et al. 2014).

165

## 166 **A brief introduction to computer vision**

167 CV is an interdisciplinary field at the intersection of *signal processing* and *machine learning*  
168 (Fig. 6; Mitchell 1997), which is concerned with the automatic and semiautomatic extraction of  
169 information from digital images. After image acquisition, a typical CV workflow involves three  
170 main steps: preprocessing, *segmentation*, and measurement (Fig. 3).

171 **Preprocessing** - Independent of how much care has been taken during image acquisition,  
172 preprocessing is an important step to prepare images for the CV routines to follow. There is a  
173 wealth of image processing techniques that can be applied at this stage, such as  
174 transformations to reduce noise (e.g. gaussian blur) or enhance contrast (e.g. histogram  
175 adjustment). Images can also be masked or labeled to filter the image so that subsequent  
176 steps are applied to the intended portions of each image. Defining the appropriate coordinate  
177 space (i.e. pixel-to-mm ratios) is also part of preprocessing. This step is highly specific to the  
178 respective study system or image dataset and may initially require some fine-tuning by the  
179 scientist to ensure data quality, which, however, can typically be automated afterwards.

180 **Segmentation** The central step in all phenotyping related CV-pipelines is the segmentation  
181 of images into pixels that contain the desired signal or information (*foreground*) and all other  
182 pixels (background). In its most basic form, segmentation of grayscale images can be done  
183 by simple signal processing algorithms, such as a *threshold* or *watershed*. Similarly, *feature*  
184 *detection* algorithms examine pixels and their adjacent region for specific characteristics or  
185 key points, e.g. whether groups of pixel form edges, corners, ridges or blobs. Videos or  
186 multiple images of the same scene provide an additional opportunity for segmentation:  
187 foreground detection can detect changes in image sequences to determine the pixels of  
188 interest (e.g. a specimen placed in an arena, or animals moving against a static background),  
189 while subsequent background subtraction isolates the foreground for further processing.  
190 Finally, *object detection* describes the high-level task of finding instances of *semantic objects*

191 (organisms, organs, structures, etc.) in an image, which is typically addressed through  
192 classical machine learning or *deep learning* (see section Computer vision methods). In  
193 classical machine learning, *features* must be first engineered or extracted from a training  
194 dataset using *feature detectors*, then used to train a classifier, and finally applied to the actual  
195 dataset (Mitchell 1997). Deep learning algorithms are a family of machine learning methods  
196 based on artificial *neural networks* that “learn” what constitutes the object of interest during  
197 the training phase. With sufficient training from *labelled* images, deep learning-powered object  
198 detection algorithms can be highly accurate and often greatly outperform pre-existing object  
199 recognition methods - in some cases even human experts (Buetti-Dinh et al. 2019).

200 **Measurement** What sorts of data can we extract from images? CV can retrieve a multitude of  
201 phenotypic traits from digital images in a systematic and repeatable fashion. In the simplest  
202 case such traits have been measured traditionally and are established in each study system,  
203 such as body size (e.g. length or diameter) or color (e.g. “red phenotype” vs “blue phenotype”).  
204 In such cases, switching from a manual approach to a semi- or fully automatic CV approach  
205 is straightforward, because the target traits are well embedded in existing statistical and  
206 conceptual frameworks. The main benefits from CV are that costly manual labor is reduced  
207 and that the obtained data becomes more reproducible, because the applied CV analysis  
208 pipeline can be stored and re-executed. It is also possible to increase the number of  
209 dimensions without much extra effort and without discarding the traditionally measured traits  
210 (Table 1). For example, in addition to body size, one could extract body shape traits, i.e. the  
211 outline of the body itself (i.e. contour coordinates of the foreground), and texture (i.e. all pixel  
212 intensities within the foreground). Such high dimensional traits can be directly analyzed using  
213 multivariate statistics, or transformed into continuous low dimensional traits, such as  
214 continuous shape features (circularity or area), texture features (color intensity or variation,  
215 pixel distribution), or moments of the raw data (Table 1).



## 216 **A history of computer vision methods**

217 The field of CV is now close to celebrating its 6th decade. It first emerged in the late 1950s  
218 and early 1960s, in the context of artificial intelligence research (Rosenblatt 1958). At the time,  
219 it was widely considered a stepping-stone in our search for understanding human intelligence  
220 (Minsky 1961). Given its long history, a wide variety of CV techniques have emerged since its  
221 inception, but they all contain variations of the same basic mechanism. CV is, from the  
222 methodological standpoint, the process of extracting meaningful features from image data and  
223 then the use of such features to perform tasks, which, as described above, may include  
224 classification, segmentation, recognition, detection, among others. In this section, we will not  
225 aim at presenting an all-encompassing review of all CV methods, but rather to identify the  
226 major trends in the field and highlight the techniques that have proved useful in the context of  
227 biological research. It is worth noting that even classical CV approaches are still routinely used  
228 in the modern literature, either in isolation or, most commonly, in combination with others. In  
229 a large part, methodological choices in CV are highly domain-specific (see section *Practical*  
230 *considerations for computer vision*, Fig. 4, and Fig. 6).

231

### 232 **First wave - Hand-crafted features**

233 The first wave of CV algorithms is also the closest one to the essence of CV, namely, the  
234 process of extracting features from images. Starting with the work of Larry Roberts, which  
235 aimed at deriving 3D information from 2D images (Roberts 1963), researchers in the 1970s  
236 and 1980s developed different ways to perform feature extraction from raw pixel data. Such  
237 features tended to be low-level features, such as lines, edges, texture or lighting, but provided  
238 us with the initial basic geometric understanding of the data contained in images. A notable  
239 example of such algorithms is the watershed algorithm. First developed in 1979 (Beucher  
240 1979), the watershed algorithm treats images as a topographic map, in which pixel intensity  
241 represents its height, and attempts to segment the image into multiple separate 'drainage  
242 basins'. This algorithm is still routinely used in signal processing techniques (Fig. 6) and can

243 be effectively used to process biological images such as those obtained through animal or  
244 plant cell microscopy (McQuin et al. 2018). Other initial low-level hand-crafted approaches  
245 that achieved popularity include the Canny and Sobel filters (edge detectors; Canny 1986,  
246 Kanopoulos et al. 1988) and Hough transforms (ridge detection; Duda and Hart 1972).

247 Another approach that gained popularity in the CV literature in the early 1990s was  
248 principal component analysis (PCA). In a PCA, independent, aggregate statistical features are  
249 extracted from multidimensional datasets. These can be used, for example, in classification.  
250 One of the most notable uses of PCA in the context of CV was the eigenfaces approach (Turk  
251 and Pentland 1991). Essentially, Turk and Pentland (1991) noted that one could decompose  
252 a database of face images into eigenvectors (or characteristic images) through PCA. These  
253 eigenvectors could then be linearly combined to reconstruct any image in the original dataset.  
254 A new face could be decomposed into statistical features and further compared to other known  
255 images in a multidimensional space. PCA has notably found many other uses in biology (e.g.,  
256 (Ringnér 2008).

257 In the late 1990s and early 2000s, Scale Invariant Feature Transform (SIFT; Lowe  
258 1999, 2004) and Histogram of Oriented Gradients (HOG; Dalal and Triggs 2005) were  
259 developed. Both SIFT and HOG represent intermediate-level local features that can be used  
260 to identify keypoints that are shared across images. In both approaches, the first step is the  
261 extraction of these intermediate-level features from image data, followed by a feature matching  
262 step that tries to identify those features in multiple images. Finding keypoints across images  
263 is an essential step in many CV applications, such as object detection, landmarking, and  
264 image registration. These intermediate-level features have several advantages over the lower-  
265 level features mentioned above, most notably the ability to be detected in a wide-variety of  
266 scales, noise and illumination. Another key aspect of SIFT and HOG features is that they are  
267 generally invariant to certain geometric transformations, such as uniform scaling and simple  
268 affine distortions.

269

## 270 **Second wave - Initial machine-learning approaches**

271 While the use of hand-crafted features spurred much of the initial work in CV, soon it became  
272 apparent that without image standardization, those low- and intermediate-level features will  
273 often fall short of producing sufficiently robust CV algorithms. For example, images belonging  
274 to the same class can often look very different and the identification of a common set of shared  
275 low-level features can prove to be quite challenging. Consider, for instance, the task of  
276 classifying animal images. Two cat breeds can look quite different, despite belonging to the  
277 same general class (cat). As such, while the initial feature-engineering approaches were  
278 essential for the development of the field, it was only with the advent of machine-learning that  
279 CV acquired more generalizable applications.

280 Machine learning algorithms for CV can be divided in two main categories (but see Box  
281 2): supervised and unsupervised (Geoffrey Hinton, Terrence J. Sejnowski 1999).  
282 Unsupervised algorithms attempt to identify previously unidentified patterns on unlabeled  
283 data. In other words, no supervision is applied to the algorithm during learning. While it can  
284 be argued that PCA was one of the first successful unsupervised learning algorithms applied  
285 directly to CV, here we group PCA with “first wave” tools due to its use as a feature extractor.  
286 Other unsupervised learning algorithms commonly used in CV include clustering techniques,  
287 such as k-means (Lloyd 1982) and gaussian mixture models (GMM; Reynolds and Rose  
288 1995). Clustering algorithms represented some of the first machine learning approaches for  
289 CV. Their aim is to find an optimal set of objects (or components) that are more similar to each  
290 other than to those in other sets. This type of approach allowed researchers to find hidden  
291 patterns embedded in multidimensional data, proving useful for classification and  
292 segmentation tasks.

293 However, it is in the supervised domain that machine learning for CV has been most  
294 successful (Heileman and Myler 1989). In supervised learning approaches, the user supplies  
295 labeled training data in the form of input-output pairs (Box 2). The ML algorithm iteratively  
296 “learns” a function that maps input into output for the labeled training data. Among the initial  
297 supervised learning approaches for CV, Support Vector Machines (SVM) were by far the most

298 common approach (Cortes and Vapnik 1995). Given a certain image dataset and their  
299 corresponding labels (e.g., classes in a classification task), SVMs find the hyperplane (in  
300 feature space) that maximizes the separation between the classes of interest. An essential  
301 aspect of SVMs is that such learned decision boundaries separating the classes can be  
302 nonlinear in the original feature space, allowing the model to separate classes that would not  
303 be separable by a purely linear technique (Cortes and Vapnik 1995).

304

### 305 **Third wave - Ensemble methods**

306 While SVMs were extremely successful in CV and spurred much of the supervised work that  
307 happened afterwards, it became clear by the early 2000s that single estimators often  
308 underperformed approaches combining the predictions of several independent estimators, an  
309 approach known as ensemble methods (Dietterich 2000). Ensemble methods represent a  
310 slightly different philosophical approach to machine learning, in which multiple models are  
311 trained to solve the same task and their individual results are combined to obtain an even  
312 better model performance. Several ensemble methods have been developed in the literature,  
313 but they are generally divided in two main families: bagging and boosting.

314         Bagging approaches combine several models that were trained in parallel through an  
315 averaging process (Bauer and Kohavi 1999). Each underlying model is trained independently  
316 of the others based on a bootstrap resample of the original dataset. As a consequence, each  
317 model is trained with slightly different and (almost) independent data, greatly reducing the  
318 variance in the combined model predictions. A classic example of bagging approach is the  
319 random forest algorithm (Breiman 2001), in which multiple learning trees are fitted to bootstrap  
320 resamples of the data and posteriorly combined through mean averaging (or majority vote).

321         Boosting, on the other hand, combines learners sequentially rather than in parallel  
322 (Bauer and Kohavi 1999). Among boosting algorithms, gradient boosting (Friedman 2000) is  
323 one of the most widely used in CV. In gradient boosting, models are combined in a cascade  
324 fashion, such that a downstream model is fitted to the residuals of upstream models. As a  
325 consequence, while each individual model in the cascade is only weakly related to the overall

326 task, the combined algorithm (i.e., the entire cascade) represents a strong learner that is  
327 directly related to the task of interest. Since this approach, if unchecked, will lead the final  
328 model to overfit the training data, regularization procedures are usually applied when using  
329 gradient boosting.

330

### 331 **Fourth wave – Deep learning**

332 Deep learning approaches are, at the time of this writing, the state-of-the-art in CV and have  
333 recently become more accessible through the community-wide adoption of code-sharing  
334 practices (e.g. via <https://github.com/> or <https://stackoverflow.com/>). *Deep learning* refers to a  
335 family of machine learning methods based on artificial neural networks with multiple steps that  
336 perform *convolutions* or other mathematical operations on the input data, each of which is  
337 referred to as a *hidden layer*. Networks with dozens or hundreds of hidden layers (i.e. deep  
338 networks) allow for the extraction of high-level features from raw image data (LeCun et al.  
339 2015). While they have only recently become widespread, the history of artificial neural  
340 networks is at least as old as the field of CV itself. One of first successful attempts in the study  
341 of artificial neural networks was the perceptron (Rosenblatt 1958), a computer whose  
342 hardware design was inspired by neurons, and which was used to classify a set of inputs into  
343 two categories. This early work, while successful, was largely restricted to linear functions and  
344 therefore could not deal with non-linearity, such as XOR functions (Minsky et al. 1969). As a  
345 consequence, artificial neural network research remained rather understudied until the early  
346 80s when training procedures for multi-layer perceptrons were introduced (i.e.,  
347 backpropagation; Rumelhart and McClelland 1987). Even then, multi-layer approaches were  
348 computationally taxing, and the hardware requirements represented an important bottleneck  
349 to research in neural network-based CV, which remained disfavored compared to much lighter  
350 approaches, such as SVMs.

351 It eventually became clear that a neural network approach to CV represented a  
352 fundamental leap for CV. When compared to the hand-crafted features that dominated the  
353 field for most of its history, neural networks learn features from the data itself, therefore

354 eliminating the need for feature engineering (LeCun et al. 2015). In a large part, deep learning  
355 approaches for CV have only emerged in force due to two major developments at the  
356 beginning of the 21st century. On one side, hardware capability greatly increased due to high  
357 consumer demand for personal computing and gaming. On the other, there was a widespread  
358 adoption of the internet, leading to an exponential increase in data availability through shared  
359 image databases and labelled data. Today, deep learning is a general term that encompasses  
360 a wide variety of approaches that share an architectural commonality of relying on training  
361 neural networks with multiple hidden layers. However, this superficial similarity hides a  
362 considerable array of differences between different algorithms and one could say that the field  
363 of deep learning is as diverse as the domains in which CV is applied. We present some of the  
364 most relevant classes of deep learning approaches in Box 2.

## 365 **Practical considerations for computer vision**

### 366 **Measurement theory: define your traits thoughtfully**

367 Defining meaningful phenotypes is deceptively challenging. Traditionally, biologists relied on  
368 intuition and natural history conventions to define phenotypes, but this approach can obscure  
369 the fact that phenomes are exceedingly high-dimensional, and many dimensions are infinitely  
370 divisible. When deciding what to measure, we suggest that researchers consider  
371 *measurement theory*, a qualitative formalization of the relationship between actual  
372 measurements and the entity that the measurements are intended to represent (Houle et al.  
373 2011). In phenomics using CV, we recommend that researchers adhere to the following three  
374 principles: i) Ensure that the measurements are meaningful in the theoretical context of  
375 research questions. ii) Remember that all measurements are estimates. Measurements  
376 without uncertainties should always be avoided. iii) Be careful with units and scale types,  
377 particularly when composite values, such as the proportion of one measurement over another,  
378 are used as a measurement. Wolman (2006) and Houle et al. (2011) give details of  
379 measurement theory and practical guidelines for its use in ecology and evolutionary biology.

380

381 **Image quality and pertinent metadata: collect images that are maximally useful**

382 As a general rule of thumb, images taken for any CV analysis should have a *signal-to-noise*  
383 ratio (SNR) sufficiently high so that the signal (i.e. the phenotypic information) is detectable  
384 from the image background. High SNR can be achieved by using high resolution imaging  
385 devices (e.g. DSLR cameras or flatbed scanners), ensuring that the object is in focus (e.g.  
386 automatically or by fixing the distance between camera and object), and by creating a high  
387 contrast between object and background (e.g. by using backgrounds that are of contrasting  
388 color or brightness to the organism or object). We recommend to iteratively assess suitability  
389 of imaging data early on in a project and adjust if necessary. This means taking pilot datasets,  
390 processing them, measuring traits, estimating measurement errors, and then updating the  
391 image collection process. Moreover, it is good practice to include a color or size reference  
392 whenever possible. It helps researchers to assess if the image has sufficient SNR, increases  
393 reproducibility, and helps to evaluate measurement bias as we discuss in the next section.

394

395 **Every measurement is an estimate: dealing with measurement error**

396 Another important aspect to consider when measuring phenotypes from images is the  
397 (in)accuracy of measurements. Formally, measurement inaccuracy is composed of  
398 imprecision and bias, corresponding to random and systematic differences between measured  
399 and true values, and can be expressed as the following relationship

400

$$401 \text{ inaccuracy} = \text{imprecision} + \text{bias}^2$$

402

403 (Grabowski and Porto 2017, Tsuboi et al. 2020). These two sources of errors characterize  
404 distinct aspects of a measurement: precise measurements may still be inaccurate if biased,  
405 and unbiased measurements may still be inaccurate if imprecise (Fig. 5). Measurement  
406 imprecision can be evaluated by the coefficient of variation (standard deviation divided by the  
407 mean) of repeated measurements. Bias requires a knowledge of true values.

408 We ultimately need to understand if a measurement is sufficiently accurate to address  
409 the research question at hand. Repeatability is a widely used estimator of measurement  
410 accuracy in ecology and evolutionary biology (Wolak et al. 2012), which in our notation could  
411 be expressed as

412

$$413 \text{ repeatability} = 1 - \frac{\text{inaccuracy}}{\text{total variance}}$$

414

415 This expression clarifies that the repeatability depends both on measurement inaccuracy and  
416 total variance in the data. For example, volume estimates of deer antler from 3D  
417 photogrammetry have an average inaccuracy of 8.5%, which results in repeatabilities of 67.8-  
418 99.7% depending on the variance in antler volume that a dataset contains (Tsuboi et al. 2020).  
419 In other words, a dataset with little variation requires more accurate measurement to achieve  
420 the same repeatability as a dataset with more variation. Therefore, the impact of measurement  
421 error has to be evaluated in the specific context of data analysis.

422 CV-based phenomics is extremely useful in this regard because it allows researchers  
423 to identically repeat a measurement process, and thereby evaluate inaccuracies in order to  
424 improve measurement precision. In the aforementioned example of volume estimated from  
425 3D photogrammetry (Tsuboi et al. 2020), it was found that 70% of the total inaccuracy arose  
426 from the error in scaling arbitrary voxel units into real volumetric units. Therefore, by using the  
427 mean of two estimates obtained from two copies of an image that are scaled twice  
428 independently as a representative measurement, the inaccuracy dropped to 5.5%. However,  
429 the opportunity to improve accuracy by repeated measurements is limited if a majority of error  
430 arises from the stored images themselves. For this reason, we recommend always taking  
431 repeated images of the same subject at least for a subset of data. This will allow evaluating  
432 the magnitude of error due to images relative to the error due to acquisition of measurements  
433 from images. If the error caused by images is large compared to the error caused by data



434 acquisition, it may be necessary to modify imaging and/or preprocessing protocol to increase  
435 SNR.

436         Assessing measurement bias requires separate treatments. When linear (length) or  
437 chromatic (color) measurements are obtained from images, it is a good general practice to  
438 include size and color scales as part of images to estimate bias as the difference between  
439 known values of imaged scales and measurements obtained through CV (i.e. the reference  
440 card in Fig. 3). Knowing the true value may be difficult in some cases, such as domain area  
441 or circularity (Hoffmann et al. 2018), since they are hard to characterize without a CV. When  
442 multiple independent methods to measure the same character exist, we recommend using  
443 them on sample data to determine the bias of one method relative to the other.

444

#### 445 **Selecting a CV pipeline: as simple as possible, as complex as necessary**

446 When using CV tools there are usually many different ways to collect a specific type of  
447 phenotypic information from images (Fig. 6). Therefore, one of the first hurdles to overcome  
448 when considering the use of CV is selecting the appropriate technique from among a large  
449 and growing set of choices. The continued emergence of novel algorithms to collect, process  
450 and analyze image-derived data may sometimes make us believe that any “older” technology  
451 is immediately outdated. Deep learning, specifically CNNs, is a prominent example of an  
452 innovation in CV that was frequently communicated as so “revolutionary” and “transformative”  
453 that many scientists believed it would replace all existing methods. However, despite the  
454 success of CNNs, there are many cases where they are inappropriate or unfeasible, e.g. due  
455 to small sample sizes, hardware or time constraints, or because of the complexity that deep  
456 learning implementations entail, despite many efforts to make this technology more tractable.  
457 We discourage readers from defaulting to using the newest technology stacks; rather, we  
458 suggest that researchers be pragmatic as to which is the fastest and simplest way to get the  
459 phenotypic information of desire from any given set of images.

460         Begin by considering the size of a given image dataset, whether it is complete, e.g.  
461 after an experiment, or whether there will be continued future additions, e.g. as part of a long-

462 term experiment or field survey. As a rough rule of thumb, if a dataset encompasses only a  
463 thousand images or fewer, consider it “small”; if a dataset has thousands to tens of thousands  
464 images, consider it “large” (see Fig. 6 for methodological suggestions for each case). The next  
465 assessment should be about the SNR in your images: images taken in the laboratory typically  
466 have a high degree of standardization, e.g. controlled light environment or background, and  
467 thus a very high SNR. Field images can also have a high SNR, for example, if they are taken  
468 against the sky or if the trait of question is very distinct from the background through bright  
469 coloration. If the dataset is “small” and/or has high SNR, it may not be necessary to use the  
470 more sophisticated CV tools; instead, signal processing, e.g. threshold or watershed  
471 algorithms, may already be sufficient for segmentation although typically some pre- and post-  
472 processing is typically still required (e.g. blurring to remove noise, “morphology”-operations to  
473 close gaps, or masking false positives).

474 For large datasets, images with low SNR, or if the information of interest is variable  
475 across images (e.g. traits are photographed from different angles or partially covered up),  
476 machine learning approaches are probably more suitable. In contrast to signal processing,  
477 where segmentation results are immediately available, all machine learning image analysis  
478 pipelines include iterative training and validation phases, followed by a final testing phase.  
479 Such a workflow can be complex to initiate but pays off in the long run by providing  
480 segmentation results that become increasingly robust if more training data is supplied over  
481 time. Classic machine learning algorithms often require an intermediate amount of training  
482 data (500-1000 or more images) before they can produce satisfactory results. In this category,  
483 SVM or HOG algorithms are a good choice when areas of interest do not contrast sufficiently  
484 from the surrounding area, for example, when automatically detecting landmarks (Porto and  
485 Lysne Voje 2020). Deep learning algorithms require much larger training datasets (several  
486 1000s to 10000s) but are less sensitive to noise and idiosyncrasies of the foreground. Thus,  
487 for large and continuously growing data sets, or for recurring image analysis tasks, deep  
488 learning has become the standard approach for segmentation. Deeper networks may  
489 increase model accuracy, and thus improve the segmentation results, but have an increasing

490 risk of overfitting the contained information - i.e. the model is less generalizable to input data.  
491 While the implementation of deep learning pipelines requires some expertise, they can be  
492 retrained and are typically less domain specific than classic machine learning pipelines  
493 (O'Mahony et al. 2020).

## 494 **Recent examples of computer vision to collect phenomic data**

495 "Phenomics" as a term has not yet gained widespread attention in the ecological and  
496 evolutionary biology research communities (Fig. 1), but many biologists are engaged in  
497 research programs that are collecting phenomic data, even though it is not called as such.  
498 Some of them are already using automatic or semi-automatic CV to collect phenotypic data.  
499 Here we present small a selection of promising applications of CV to answer ecological or  
500 evolutionary research questions (points matching panels in Fig. 4):

501 **A. Resource competition traits** - Species diversity within ecological communities is  
502 often thought to be governed by competition for limiting resources (Chesson 2000).  
503 However, the exact traits that make species or individuals the best competitors under  
504 resource limitation conditions are difficult to identify among all other traits. In this  
505 example, the phenotypic space underlying resource competition was explored by  
506 implementing different limitation scenarios for experimental phytoplankton  
507 communities. Images were taken with a plate reader that used a combination of visible  
508 light and fluorometry measurements (Hense et al. 2008). The images were analyzed  
509 using signal processing, which allowed the rapid segmentation of several 1000 images  
510 by combining information from multiple fluorescence emission excitation spectra to an  
511 image stack. As a result, over 100 traits related to morphology (shape, size, and  
512 texture) and internal physiology (pigment content, distribution of pigments within each  
513 cell) were obtained at the individual cell level. (Gallego et al., unpublished data)

514 **B. Thermal adaptation and thermal reaction norms** - Variation in body temperature  
515 can be an important source of fitness variation (Kingsolver and Huey 2008, Svensson

516 et al. 2020). Quantifying body temperature and thermal reaction norms in response to  
517 natural and sexual selection allows us to test predictions from evolutionary theory  
518 about phenotypic plasticity and canalization (Lande 2009, Chevin et al. 2010).  
519 However, body temperature is an internal physiological trait that is difficult to quantify  
520 in a non-invasive way on many individuals simultaneously and under natural  
521 conditions. Thermal imaging is an efficient and non-invasive method to quantify such  
522 physiological phenotypes on a large scale and can be combined with thermal loggers  
523 to measure local thermal environmental conditions in the field (Svensson and Waller  
524 2013, Svensson et al. 2020).

525 **C. Stochastically patterned morphological traits** - In contrast to homologous,  
526 landmark-based morphological traits, tissues also form emergent patterns that are  
527 unique to every individual. The arrangement of veins on the wings of damselflies is  
528 one such example. By measuring the spacing, angles, and connectivities within the  
529 adult wing tissue, researchers have proposed hypotheses about the mechanisms of  
530 wing development and physical constraints on wing evolution (Hoffmann et al. 2018,  
531 Salcedo et al. 2019).

532 **D. Morphometrics and shape of complex structures** - Landmark-based  
533 morphometrics has become a popular tool used to characterize morphological  
534 variation in complex biological structures. Despite its popularity, landmark data is still  
535 collected mainly through manual annotation, a process which represents a significant  
536 bottleneck for phenomic studies. However, machine-learning-based CV can be used  
537 to accurately automate landmark data collection in morphometric studies not only in  
538 2D (McPeck et al. 2008, Porto and Voje 2020), but also in 3D (Porto et al. 2020).

539 **E. Volumes of morphologically complex traits.** Many topics in evolutionary ecology  
540 concerns investment of resources into a particular trait. However, measuring energetic  
541 investment, either as mass or volume of the target traits, has been challenging  
542 because many traits are morphologically complex, making it difficult to estimate  
543 investment from a combination of linear measurements. Photogrammetry is a low-cost

544 and fast technique to create 3D surface images from a set of images. Using a simple  
545 protocol and a free proprietary software, Tsuboi et al. (2020) demonstrated that  
546 photogrammetry can accurately measure the volume of antler in deer family Cervidae.  
547 The protocol is still relatively low-throughput due primarily to the necessity of high  
548 number of images (> 50) per sample, but it allows extensive sampling (*sensu* Houle et  
549 al. 2010) of linear, area and volumetric measurements of antler structures.

## 550 **Outlook**

551 In this review we provided a broad overview of various CV techniques and gave some recent  
552 examples of their application in ecological and evolutionary research. We presented CV as a  
553 promising toolkit to overcome the image analysis bottleneck in phenomics. However, to be  
554 clear, we do not suggest that biologists discontinue the collection of univariate traits like body  
555 size or discrete colors. Such measures are undoubtedly useful if they contain explanatory  
556 value and predictive power. Instead, we propose that CV can help to i) collect them with higher  
557 throughput, ii) in a more reproducible fashion, and to iii) collect additional traits so we can  
558 interpret them in the context of trait combinations. We argue that CV is not bound to  
559 immediately replace existing methods, but it simply opens the opportunity to place empirical  
560 research of phenotypes on a broader base. We also note that CV based phenomics can be  
561 pursued in a deductive or inductive fashion. In the former case, scientists would simply  
562 conduct hypothesis driven research including a wider array of traits into causal models (Houle  
563 et al. 2011); in the latter, they would engage in discovery-based data-mining approaches that  
564 allow scientists to form hypotheses a posteriori based on the collected data (Kell and Oliver  
565 2004).

566 Although CV based phenomics provides new opportunities for many areas of study,  
567 we identify several specific fields that will profit most immediately from CV. First, evolutionary  
568 quantitative genetics will benefit tremendously from increased sample sizes that CV-  
569 phenomics entails, because the bottleneck of the field has been the difficulty in accurately

570 estimating key statistics such as genetic variance covariance matrices and selection gradients.  
571 The recent discovery of tight matches between mutational, genetic, and macroevolutionary  
572 variances in *Drosophila* wing shape (Houle et al. 2017) is exemplary of a successful phenomic  
573 project. Second, large-scale empirical studies of the genotype-phenotype map will finally  
574 become possible, because of the availability of high-throughput phenotypic data and analytical  
575 framework to deal with big data (Pitchers et al. 2019, Zheng et al. 2019, Maeda et al. 2020).  
576 Third, studies of fossil time-series will gain opportunities to document and analyze the  
577 dynamics of long-term phenotypic evolution with unprecedented temporal resolution (Liow et  
578 al. 2017, Brombacher et al. 2017). The ever-growing technology of CV indicates that these  
579 are likely a small subset of unforeseen future applications of CV phenomics in our field. Just  
580 like the technological advancements in DNA-sequencing has revolutionized our view of  
581 genomes, development and molecular evolution in the past decades, we anticipate that the  
582 way we look at phenotypic data will be changing in the coming years.

583         Just as CV is changing what it means to measure a trait, there is a complementary  
584 change in what can be considered scientific image data in the first place. Large, publicly  
585 available image datasets are fertile ground for ecology and evolutionary research. Such  
586 databases include both popular and non-scientific social media (e.g. <https://www.flickr.com/>  
587 or <https://www.instagram.com/>), but also quality-controlled and vetted natural history and  
588 species identification resources with global scope and ambitions (e.g.  
589 <https://www.inaturalist.org/>). Successful examples of how such public image databases can  
590 be useful are in studies aiming to quantify the frequencies variation of discrete traits, such as  
591 color polymorphism frequencies in different geographic regions (Leighton et al. 2016). These  
592 manual efforts in mining available public image resources can potentially be replaced in the  
593 future using more automated machine learning or CV approaches. Similarly, the corpus of  
594 published scientific literature is full of image data that can be combined and re-analyzed in  
595 order to address larger-scale questions (Hoffmann et al. 2018, Church et al. 2019a, 2019b).

596         Previous calls for phenomics argued that, to make phenomics a successful endeavor,  
597 it has to be extensive, aiming at measuring many different aspects of the phenotypes, as well

598 as intensive, aiming at characterizing each measurement accurately with large sample size  
599 and with high temporal resolution (Bilder et al. 2009, Houle et al. 2010, Furbank and Tester  
600 2011). We agree with this view, but we also emphasize that phenomics is nothing conceptually  
601 new in this respect. As we discussed, many researchers in our field have already adopted  
602 phenomic pipelines, studying high-dimensional phenotypic data acquired by high-throughput  
603 measuring technologies without using the term phenomics. If so, what is the conceptual value  
604 of phenomics? In our opinion, phenomics is more than just a rigorous version of conventional  
605 research of organismal phenotypes, but also a dedication towards phenotypic data.  
606 Phenomics shifts us from viewing phenotypes as given entities towards viewing them as part  
607 of the phenome at the whole organismal level.

608

## 609 **Acknowledgements**

610 The publication of this study was funded through the Swedish Research Council International  
611 Postdoc Grant (2016-06635) to MT. ML was supported by a Swiss National Science  
612 Foundation Early Postdoc.Mobility grant (SNSF: P2EZP3\_191804). EIS was funded by a grant  
613 from the Swedish Research Council (VR: grant no. 2016-03356). SD was supported by the  
614 Jane Coffin Childs Memorial Fund.

## 615 **Authors contributions**

616 ML conceived the idea for this review and initiated its writing. In the process, all authors  
617 contributed equally to the development and discussion of ideas, and to the writing of the  
618 manuscript.

## 619 **Conflict of interest**

620 The authors declare that the research was conducted in the absence of any commercial or  
621 financial relationships that could be construed as a potential conflict of interest.

622 **References**

- 623 Bauer, E., and R. Kohavi. 1999. An Empirical Comparison of Voting Classification  
624 Algorithms: Bagging, Boosting, and Variants. *Machine learning* 36:105–139.
- 625 Beucher, S. 1979. Use of watersheds in contour detection. *Proceedings of the International*  
626 *Workshop on Image*.
- 627 Bilder, R. M., F. W. Sabb, T. D. Cannon, E. D. London, J. D. Jentsch, D. S. Parker, R. A.  
628 Poldrack, C. Evans, and N. B. Freimer. 2009. Phenomics: the systematic study of  
629 phenotypes on a genome-wide scale. *Neuroscience* 164:30–42.
- 630 Blonder, B. 2018. Hypervolume concepts in niche- and trait-based ecology. *Ecography*  
631 41:1441–1455.
- 632 Breiman, L. 2001. Random Forests. *Machine learning* 45:5–32.
- 633 Brombacher, A., P. A. Wilson, I. Bailey, and T. H. G. Ezard. 2017. The Breakdown of Static  
634 and Evolutionary Allometries during Climatic Upheaval. *The American naturalist*  
635 190:350–362.
- 636 Bruijning, M., M. D. Visser, C. A. Hallmann, and E. Jongejans. 2018. trackdem : Automated  
637 particle tracking to obtain population counts and size distributions from videos in r.  
638 *Methods in ecology and evolution / British Ecological Society* 9:965–973.
- 639 Buetti-Dinh, A., V. Galli, S. Bellenberg, O. Ilie, M. Herold, S. Christel, M. Boretska, I. V.  
640 Pivkin, P. Wilmes, W. Sand, M. Vera, and M. Dopson. 2019. Deep neural networks  
641 outperform human expert’s capacity in characterizing bioleaching bacterial biofilm  
642 composition. *Biotechnology reports (Amsterdam, Netherlands)* 22:e00321.
- 643 Canny, J. 1986. A Computational Approach to Edge Detection. *IEEE transactions on pattern*  
644 *analysis and machine intelligence* PAMI-8:679–698.
- 645 Cheng, K. C., X. Xin, D. P. Clark, and P. La Riviere. 2011. Whole-animal imaging, gene  
646 function, and the Zebrafish Phenome Project. *Current opinion in genetics &*  
647 *development* 21:620–629.
- 648 Chesson, P. 2000. Mechanisms of Maintenance of Species Diversity. *Annual review of*



649 ecology and systematics 31:343–366.

650 Chevin, L.-M., R. Lande, and G. M. Mace. 2010. Adaptation, plasticity, and extinction in a  
651 changing environment: towards a predictive theory. *PLoS biology* 8:e1000357.

652 Church, G. M., and W. Gilbert. 1984. Genomic sequencing. *Proceedings of the National  
653 Academy of Sciences of the United States of America* 81:1991–1995.

654 Church, S. H., S. Donoughe, B. A. S. de Medeiros, and C. G. Extavour. 2019a, July. Insect  
655 egg size and shape evolve with ecology but not developmental rate.

656 Church, S. H., S. Donoughe, B. A. S. de Medeiros, and C. G. Extavour. 2019b. A dataset of  
657 egg size and shape from more than 6,700 insect species. *Scientific data* 6:104.

658 Cortes, C., and V. Vapnik. 1995. Support-Vector Networks. *Machine learning* 20:273–297.

659 Dalal, N., and B. Triggs. 2005. Histograms of oriented gradients for human detection. Pages  
660 886–893 vol. 1 2005 IEEE Computer Society Conference on Computer Vision and  
661 Pattern Recognition (CVPR'05).

662 Dietterich, T. G. 2000. Ensemble Methods in Machine Learning. Pages 1–15 *Proceedings of  
663 the First International Workshop on Multiple Classifier Systems*. Springer-Verlag, Berlin,  
664 Heidelberg.

665 Di, Z., M. J. D. Klop, V.-M. Rogkoti, S. E. Le Dévédec, B. van de Water, F. J. Verbeek, L. S.  
666 Price, and J. H. N. Meerman. 2014. Ultra high content image analysis and phenotype  
667 profiling of 3D cultured micro-tissues. *PloS one* 9:e109688.

668 Duda, R. O., and P. E. Hart. 1972. Use of the Hough transformation to detect lines and  
669 curves in pictures. *Communications of the ACM* 15:11–15.

670 Feder, M. E., and T. Mitchell-Olds. 2003. Evolutionary and ecological functional genomics.  
671 *Nature reviews. Genetics* 4:651–657.

672 Freimer, N., and C. Sabatti. 2003. The human phenome project. *Nature genetics* 34:15–21.

673 French, S., B. E. Coutts, and E. D. Brown. 2018. Open-Source High-Throughput Phenomics  
674 of Bacterial Promoter-Reporter Strains. *Cell systems* 7:339–346.e3.

675 Friedman, J. H. 2000. Greedy Function Approximation: A Gradient Boosting Machine.  
676 *Annals of Statistics*.

677 Furbank, R. T., and M. Tester. 2011. Phenomics--technologies to relieve the phenotyping  
678 bottleneck. *Trends in plant science* 16:635–644.

679 Gehan, M. A., N. Fahlgren, A. Abbasi, J. C. Berry, S. T. Callen, L. Chavez, A. N. Doust, M. J.  
680 Feldman, K. B. Gilbert, J. G. Hodge, J. S. Hoyer, A. Lin, S. Liu, C. Lizárraga, A.  
681 Lorence, M. Miller, E. Platon, M. Tessman, and T. Sax. 2017. PlantCV v2: Image  
682 analysis software for high-throughput plant phenotyping. *PeerJ* 5:e4088.

683 Geoffrey Hinton, Terrence J. Sejnowski. 1999. *Unsupervised Learning: Foundations of*  
684 *Neural Computation*. MIT Press.

685 Gerum, R. C., S. Richter, B. Fabry, and D. P. Zitterbart. 2017. ClickPoints: an expandable  
686 toolbox for scientific image annotation and analysis. *Methods in ecology and evolution /*  
687 *British Ecological Society* 8:750–756.

688 Grabowski, M., and A. Porto. 2017. How many more? Sample size determination in studies  
689 of morphological integration and evolvability. *Methods in ecology and evolution / British*  
690 *Ecological Society* 8:592–603.

691 Hakim, A., Y. Mor, I. A. Toker, A. Levine, M. Neuhof, Y. Markovitz, and O. Rechavi. 2018.  
692 WorMachine: machine learning-based phenotypic analysis tool for worms. *BMC biology*  
693 16:8.

694 Heileman, G. L., and H. R. Myler. 1989. *Theoretical and experimental aspects of supervised*  
695 *learning in artificial neural networks*. phd, University of Central Florida, USA.

696 Hense, B. A., P. Gais, U. Jutting, H. Scherb, and K. Rodenacker. 2008. Use of fluorescence  
697 information for automated phytoplankton investigation by image analysis.

698 Hoffmann, J., S. Donoughe, K. Li, M. K. Salcedo, and C. H. Rycroft. 2018. A simple  
699 developmental model recapitulates complex insect wing venation patterns. *Proceedings*  
700 *of the National Academy of Sciences of the United States of America* 115:9905–9910.

701 Houle, D., G. H. Bolstad, K. van der Linde, and T. F. Hansen. 2017. Mutation predicts 40  
702 million years of fly wing evolution. *Nature* 548:447–450.

703 Houle, D., D. R. Govindaraju, and S. Omholt. 2010. Phenomics: the next challenge. *Nature*  
704 *Reviews. Genetics* 11:855–866.

705 Houle, D., J. Mezey, P. Galpern, and A. Carter. 2003. Automated measurement of  
706 *Drosophila* wings. *BMC evolutionary biology* 3:25.

707 Houle, D., C. Pélabon, G. P. Wagner, and T. F. Hansen. 2011. Measurement and Meaning  
708 in Biology.

709 Hsiang, A. Y., K. Nelson, L. E. Elder, E. C. Sibert, S. S. Kahanamoku, J. E. Burke, A. Kelly,  
710 Y. Liu, and P. M. Hull. 2018. AutoMorph : Accelerating morphometrics with automated  
711 2D and 3D image processing and shape extraction. *Methods in ecology and evolution* /  
712 British Ecological Society 9:605–612.

713 Ishikawa, A., N. Kabeya, K. Ikeya, R. Kakioka, J. N. Cech, N. Osada, M. C. Leal, J. Inoue,  
714 M. Kume, A. Toyoda, A. Tezuka, A. J. Nagano, Y. Y. Yamasaki, Y. Suzuki, T. Kokita, H.  
715 Takahashi, K. Lucek, D. Marques, Y. Takehana, K. Naruse, S. Mori, O. Monroig, N.  
716 Ladd, C. J. Schubert, B. Matthews, C. L. Peichel, O. Seehausen, G. Yoshizaki, and J.  
717 Kitano. 2019. A key metabolic gene for recurrent freshwater colonization and radiation  
718 in fishes. *Science* 364:886–889.

719 Kanopoulos, N., N. Vasanthavada, and R. L. Baker. 1988. Design of an image edge  
720 detection filter using the Sobel operator. *IEEE journal of solid-state circuits* 23:358–367.

721 Kell, D. B., and S. G. Oliver. 2004. Here is the evidence, now what is the hypothesis? The  
722 complementary roles of inductive and hypothesis-driven science in the post-genomic  
723 era. *BioEssays: news and reviews in molecular, cellular and developmental biology*  
724 26:99–105.

725 Kingsolver, J. G., and R. B. Huey. 2008. Size, temperature, and fitness: three rules.  
726 *Evolutionary ecology research* 10:251–268.

727 Kühl, H. S., and T. Burghardt. 2013. Animal biometrics: quantifying and detecting phenotypic  
728 appearance. *Trends in ecology & evolution* 28:432–441.

729 Lamichhaney, S., D. C. Card, P. Grayson, J. F. R. Tonini, G. A. Bravo, K. Näpflin, F.  
730 Termignoni-Garcia, C. Torres, F. Burbrink, J. A. Clarke, T. B. Sackton, and S. V.  
731 Edwards. 2019. Integrating natural history collections and comparative genomics to  
732 study the genetic architecture of convergent evolution. *Philosophical transactions of the*

733 Royal Society of London. Series B, Biological sciences 374:20180248.

734 Lande, R. 2009. Adaptation to an extraordinary environment by evolution of phenotypic  
735 plasticity and genetic assimilation. *Journal of evolutionary biology* 22:1435–1446.

736 Lande, R., and S. J. Arnold. 1983. The Measurement of Selection on Correlated Characters.  
737 *Evolution* 37:1210–1226.

738 Laughlin, D. C., J. R. Gremer, P. B. Adler, R. M. Mitchell, and M. M. Moore. 2020. The Net  
739 Effect of Functional Traits on Fitness. *Trends in ecology & evolution* 35:1037–1047.

740 LeCun, Y., Y. Bengio, and G. Hinton. 2015. Deep learning. *Nature* 521:436–444.

741 Leighton, G., P. S. Hugo, A. Roulin, and A. Amar. 2016. Just Google it: assessing the use of  
742 Google Images to describe geographical variation in visible traits of organisms 7.

743 Le, V.-L., M. Beurton-Aimar, A. Zemmari, A. Marie, and N. Parisey. 2020. Automated  
744 landmarking for insects morphometric analysis using deep neural networks. *Ecological  
745 informatics* 60:101175.

746 Liow, L. H., E. Di Martino, M. Krzeminska, M. Ramsfjell, S. Rust, P. D. Taylor, and K. L.  
747 Voje. 2017. Relative size predicts competitive outcome through 2 million years. *Ecology  
748 letters* 20:981–988.

749 Liu, F., A. Wollstein, P. G. Hysi, G. A. Ankra-Badu, T. D. Spector, D. Park, G. Zhu, M.  
750 Larsson, D. L. Duffy, G. W. Montgomery, D. A. Mackey, S. Walsh, O. Lao, A. Hofman,  
751 F. Rivadeneira, J. R. Vingerling, A. G. Uitterlinden, N. G. Martin, C. J. Hammond, and  
752 M. Kayser. 2010. Digital Quantification of Human Eye Color Highlights Genetic  
753 Association of Three New Loci.

754 Lloyd, S. 1982. Least squares quantization in PCM. *IEEE transactions on information theory  
755 / Professional Technical Group on Information Theory* 28:129–137.

756 Lowe, D. G. 1999. Object recognition from local scale-invariant features. Pages 1150–1157  
757 vol.2 *Proceedings of the Seventh IEEE International Conference on Computer Vision*.

758 Lowe, D. G. 2004. Distinctive image features from scale-invariant keypoints. *International  
759 journal of computer vision*.

760 Lürig, M. D. 2018. phenopype - a phenotyping pipeline for python.

761 Maeda, T., J. Iwasawa, H. Kotani, N. Sakata, M. Kawada, T. Horinouchi, A. Sakai, K.  
762 Tanabe, and C. Furusawa. 2020. High-throughput laboratory evolution reveals  
763 evolutionary constraints in *Escherichia coli*. *Nature communications* 11:5970.

764 McPeck, M. A., L. Shen, J. Z. Torrey, and H. Farid. 2008. The Tempo and Mode of Three-  
765 Dimensional Morphological Evolution in Male Reproductive Structures.

766 McQuin, C., A. Goodman, V. Chernyshev, L. Kametsky, B. A. Cimini, K. W. Karhohs, M.  
767 Doan, L. Ding, S. M. Rafelski, D. Thirstrup, W. Wiegnaebe, S. Singh, T. Becker, J. C.  
768 Caicedo, and A. E. Carpenter. 2018. CellProfiler 3.0: Next-generation image processing  
769 for biology. *PLoS biology* 16:e2005970.

770 Minsky, M. 1961. Steps toward Artificial Intelligence. *Proceedings of the IRE* 49:8–30.

771 Minsky, M, Papert, and S. 1969. *Perceptrons*.

772 Mitchell, T. M. 1997. *Machine learning*. 1997. Burr Ridge, IL: McGraw Hill 45:870–877.

773 Morel-Journel, T., V. Thuillier, F. Pennekamp, E. Laurent, D. Legrand, A. S. Chaine, and N.  
774 Schtickzelle. 2020. A multidimensional approach to the expression of phenotypic  
775 plasticity. *Functional ecology*.

776 O'Mahony, N., S. Campbell, A. Carvalho, S. Harapanahalli, G. V. Hernandez, L. Krpalkova,  
777 D. Riordan, and J. Walsh. 2020. Deep Learning vs. Traditional Computer Vision.

778 Orgogozo, V., B. Morizot, and A. Martin. 2015. The differential view of genotype-phenotype  
779 relationships. *Frontiers in genetics* 6:179.

780 Pfennig, D. W., M. A. Wund, E. C. Snell-Rood, T. Cruickshank, C. D. Schlichting, and A. P.  
781 Moczek. 2010. Phenotypic plasticity's impacts on diversification and speciation. *Trends*  
782 *in ecology & evolution* 25:459–467.

783 Phillips, C. P., and S. J. Arnold. 1999a. Hierarchical comparison of genetic variance  
784 covariance matrices. I. Using the Flury hierarchy. *Evolution; international journal of*  
785 *organic evolution* 53:1506–1515.

786 Phillips, P. C., and S. J. Arnold. 1999b. Hierarchical Comparison of Genetic Variance-  
787 Covariance Matrices. I. Using the Flury Hierarchy. *Evolution; international journal of*  
788 *organic evolution* 53:1506–1515.

789 Pitchers, W., J. Nye, E. J. Márquez, A. Kowalski, I. Dworkin, and D. Houle. 2019. A  
790 Multivariate Genome-Wide Association Study of Wing Shape in *Drosophila*  
791 *melanogaster*. *Genetics* 211:1429–1447.

792 Pointer, M. R., and G. G. Attridge. 1998. The number of discernible colours. *Color Research*  
793 & Application: Endorsed by Inter-Society Color Council, The Colour Group (Great  
794 Britain), Canadian Society for Color, Color Science Association of Japan, Dutch Society  
795 for the Study of Color, The Swedish Colour Centre Foundation, Colour Society of  
796 Australia, *Centre Français de la Couleur* 23:52–54.

797 Porto, A., and K. Lysne Voje. 2020. ML-morph: A fast, accurate and general approach for  
798 automated detection and landmarking of biological structures in images. *Methods in*  
799 *ecology and evolution* / British Ecological Society.

800 Porto, A., S. M. Rolfe, and A. Murat Maga. 2020, September 19. ALPACA: a fast and  
801 accurate approach for automated landmarking of three-dimensional biological  
802 structures.

803 Porto, A., and K. L. Voje. 2020. ML-morph: A fast, accurate and general approach for  
804 automated detection and landmarking of biological structures in images. *Methods in*  
805 *ecology and evolution* / British Ecological Society 11:500–512.

806 Reynolds, D. A., and R. C. Rose. 1995. Robust text-independent speaker identification using  
807 Gaussian mixture speaker models. *IEEE transactions on audio, speech, and language*  
808 *processing* 3:72–83.

809 Ringnér, M. 2008. What is principal component analysis? *Nature biotechnology* 26:303–304.

810 Roberts, L. G. 1963. *Machine perception of three-dimensional solids*. Massachusetts  
811 Institute of Technology.

812 Roeder, A. H. K., A. Cunha, M. C. Burl, and E. M. Meyerowitz. 2012. A computational image  
813 analysis glossary for biologists. *Development* 139:3071–3080.

814 Rosenblatt, F. 1958. The perceptron: a probabilistic model for information storage and  
815 organization in the brain. *Psychological review* 65:386–408.

816 Rumelhart, D. E., and J. L. McClelland. 1987. *Learning Internal Representations by Error*

817 Propagation. Pages 318–362 *Parallel Distributed Processing: Explorations in the*  
818 *Microstructure of Cognition: Foundations*. MIT Press.

819 Salcedo, M. K., J. Hoffmann, S. Donoughe, and L. Mahadevan. 2019. Computational  
820 analysis of size, shape and structure of insect wings. *Biology open* 8.

821 Saltz, J. B., F. C. Hessel, and M. W. Kelly. 2017. Trait Correlations in the Genomics Era.  
822 *Trends in ecology & evolution*.

823 Schindelin, J., I. Arganda-Carreras, E. Frise, V. Kaynig, M. Longair, T. Pietzsch, S.  
824 Preibisch, C. Rueden, S. Saalfeld, B. Schmid, J.-Y. Tinevez, D. J. White, V.  
825 Hartenstein, K. Eliceiri, P. Tomancak, and A. Cardona. 2012. Fiji: an open-source  
826 platform for biological-image analysis. *Nature methods* 9:676–682.

827 Schluter, D. 1996. Adaptive radiation along genetic lines of least resistance. *Evolution;*  
828 *international journal of organic evolution* 50:1766–1774.

829 Seehausen, O., R. K. Butlin, I. Keller, C. E. Wagner, J. W. Boughman, P. A. Hohenlohe, C.  
830 L. Peichel, G. P. Saetre, C. Bank, A. Brannstrom, A. Brelsford, C. S. Clarkson, F.  
831 Eroukhmanoff, J. L. Feder, M. C. Fischer, A. D. Foote, P. Franchini, C. D. Jiggins, F. C.  
832 Jones, A. K. Lindholm, K. Lucek, M. E. Maan, D. A. Marques, S. H. Martin, B.  
833 Matthews, J. I. Meier, M. Most, M. W. Nachman, E. Nonaka, D. J. Rennison, J.  
834 Schwarzer, E. T. Watson, A. M. Westram, and A. Widmer. 2014. Genomics and the  
835 origin of species. *Nature reviews. Genetics* 15:176–192.

836 Sinervo, B., and E. Svensson. 2002. Correlational selection and the evolution of genomic  
837 architecture. *Heredity* 89:329–338.

838 Soulé, M. 1967. PHENETICS OF NATURAL POPULATIONS I. PHENETIC  
839 RELATIONSHIPS OF INSULAR POPULATIONS OF THE SIDE-BLOTCHED LIZARD.  
840 *Evolution; international journal of organic evolution* 21:584–591.

841 Svensson, E. I., S. J. Arnold, R. Bürger, K. Csilléry, J. Draghi, J. M. Henshaw, A. G. Jones,  
842 S. De Lisle, D. A. Marques, K. McGuigan, M. N. Simon, and A. Runemark. 2021.  
843 Correlational selection in the age of genomics. *Nature Ecology and Evolution*.

844 Svensson, E. I., M. Gomez-Llano, and J. Waller. 2020, June 12. Natural and sexual

845 selection on phenotypic plasticity favour thermal canalization.

846 Svensson, E. I., and J. T. Waller. 2013. Ecology and sexual selection: evolution of wing  
847 pigmentation in calopterygid damselflies in relation to latitude, sexual dimorphism, and  
848 speciation. *The American naturalist* 182:E174–95.

849 Tattersall, G. J., D. V. Andrade, and A. S. Abe. 2009. Heat exchange from the toucan bill  
850 reveals a controllable vascular thermal radiator. *Science* 325:468–470.

851 Tattersall, G. J., and V. Cadena. 2010. Insights into animal temperature adaptations  
852 revealed through thermal imaging. *The Imaging Science Journal* 58:261–268.

853 Tsubaki, Y., Y. Samejima, and M. T. Siva-Jothy. 2010. Damselfly females prefer hot males:  
854 higher courtship success in males in sunspots. *Behavioral ecology and sociobiology*  
855 64:1547–1554.

856 Tsuboi, M., B. T. Kopperud, C. Syrowatka, M. Grabowski, K. L. Voje, C. Pélabon, and T. F.  
857 Hansen. 2020. Measuring Complex Morphological Traits with 3D Photogrammetry: A  
858 Case Study with Deer Antlers. *Evolutionary biology* 47:175–186.

859 Turk, M., and A. Pentland. 1991. Eigenfaces for recognition. *Journal of cognitive*  
860 *neuroscience* 3:71–86.

861 Ubbens, J. R., and I. Stavness. 2017. Deep Plant Phenomics: A Deep Learning Platform for  
862 Complex Plant Phenotyping Tasks. *Frontiers in plant science* 8:1190.

863 Visscher, P. M., and J. Yang. 2016. A plethora of pleiotropy across complex traits.

864 Walsh, B. 2007, January. Escape from flatland.

865 Weinstein, B. G. 2015. MotionMeerkat: integrating motion video detection and ecological  
866 monitoring. *Methods in ecology and evolution / British Ecological Society* 6:357–362.

867 Wolak, M. E., D. J. Fairbairn, and Y. R. Paulsen. 2012. Guidelines for estimating  
868 repeatability. *Methods in ecology and evolution / British Ecological Society* 3:129–137.

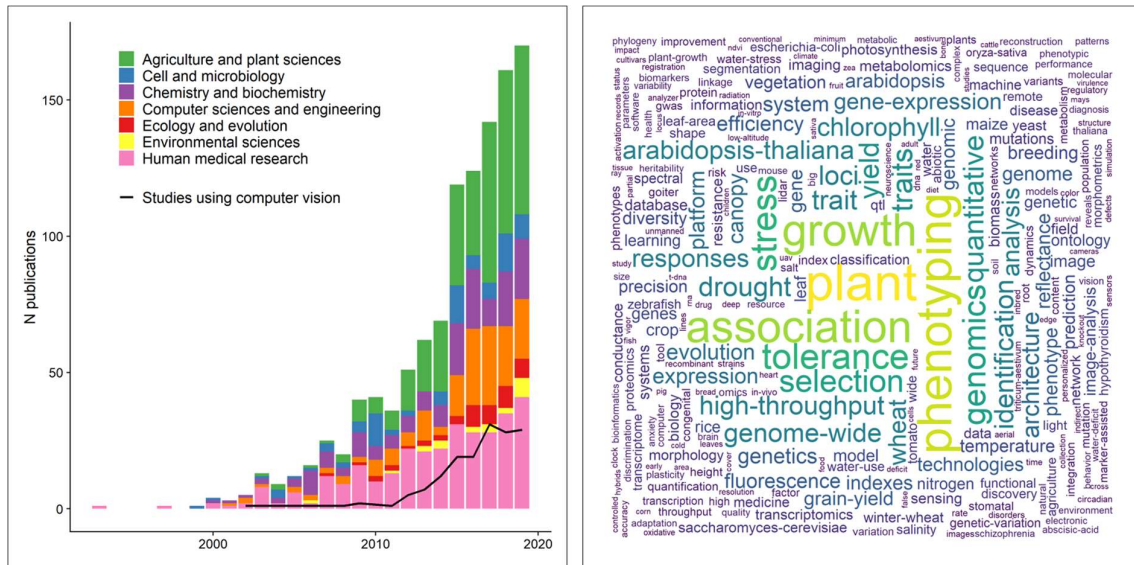
869 Wolman, A. G. 2006. Measurement and meaningfulness in conservation science.  
870 *Conservation biology: the journal of the Society for Conservation Biology* 20:1626–  
871 1634.

872 Zackrisson, M., J. Hallin, L.-G. Ottosson, P. Dahl, E. Fernandez-Parada, E. Ländström, L.



873            Fernandez-Ricaud, P. Kaferle, A. Skyman, S. Stenberg, S. Omholt, U. Petrovič, J.  
874            Warringer, and A. Blomberg. 2016. Scan-o-matic: High-Resolution Microbial Phenomics  
875            at a Massive Scale. *G3* 6:3003–3014.  
876            Zheng, J., J. L. Payne, and A. Wagner. 2019. Cryptic genetic variation accelerates evolution  
877            by opening access to diverse adaptive peaks. *Science* 365:347–353.  
  
878

879 **Figures**



880

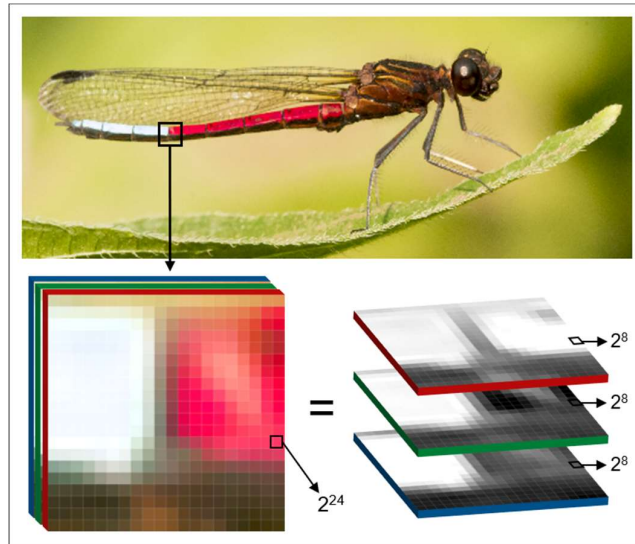
881 Figure 1 - The current state of phenomics research. Left panel: a literature survey using the  
 882 search key “phenomic\*” on in a Web of Science topic search (title, keywords, abstract) resulted  
 883 in 1323 papers (on 23/10/2020). Here we show only papers published between 1990 and  
 884 before 2020 (1125 papers) for better visual inference. Traditionally phenomics approaches are  
 885 used in agricultural sciences and crop research to maximize yield, as well as in human  
 886 medicine to study drug responses and disorder phenotypes. The black line denotes the studies  
 887 that used computer vision or some sort of image analysis (acquired with a topic search using  
 888 the strings "computer vision", "image analysis", "image based", "image processing"), indicating  
 889 that only a small subset of the studies uses image analysis. Right panel: a word cloud that  
 890 was constructed using the 500 most used keywords from the papers presented in the left  
 891 panel.

892

893

894

895



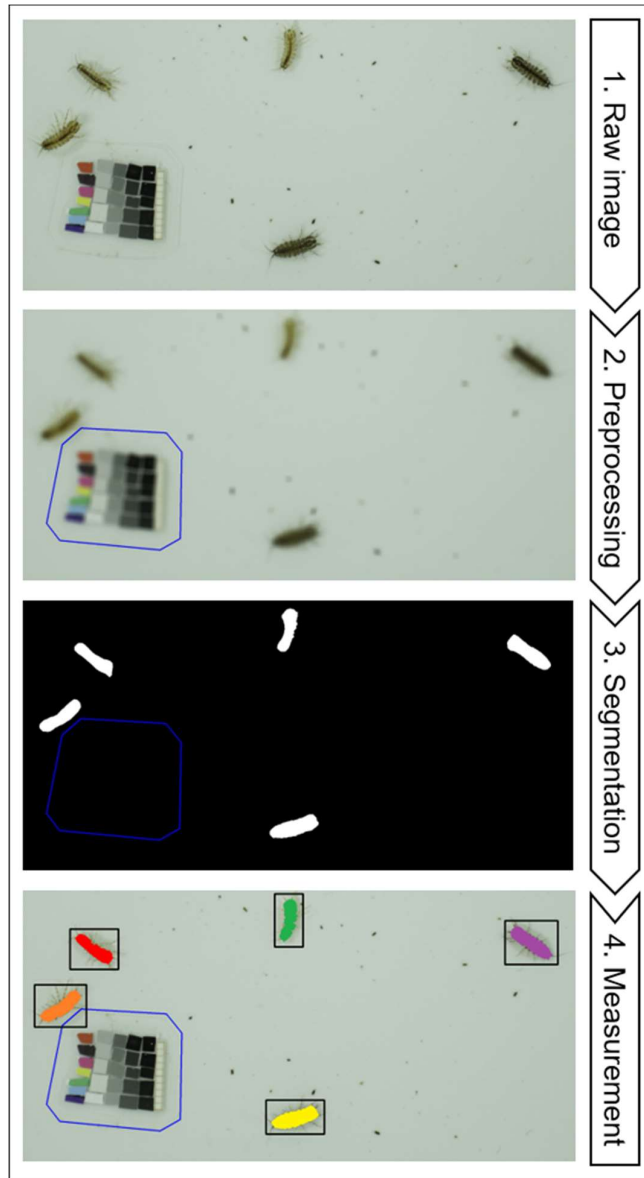
896

897 Figure 2 - The structure of digital images. Two-dimensional raster images, as produced by  
898 most commercially available cameras, are composed of three color channels red, green, blue  
899 = RGB), each of which by itself is a grayscale image. The industrial standard for color  
900 representation on the pixel level is 24 bit ( $2^{24} = 16\,777\,216$  possible color variations per pixel),  
901 which is achieved through additive mixing of each of the 8 bit channels ( $2^8 + 2^8 + 2^8$ ). This  
902 enormous range of color intensities among several million pixels is a potentially very high-  
903 resolution representation of organismal traits, or the organism as a whole. Therefore, digital  
904 images are a useful medium for phenomics research, as they offer an inexpensive, memory  
905 efficient and standardizable way to capture, store and analyze complex phenotypes. The  
906 photograph shows a blue-tip jewel damselfly (*Chlorocypha curta*) in Cameroon (Africa) - image  
907 by Erik Svensson.

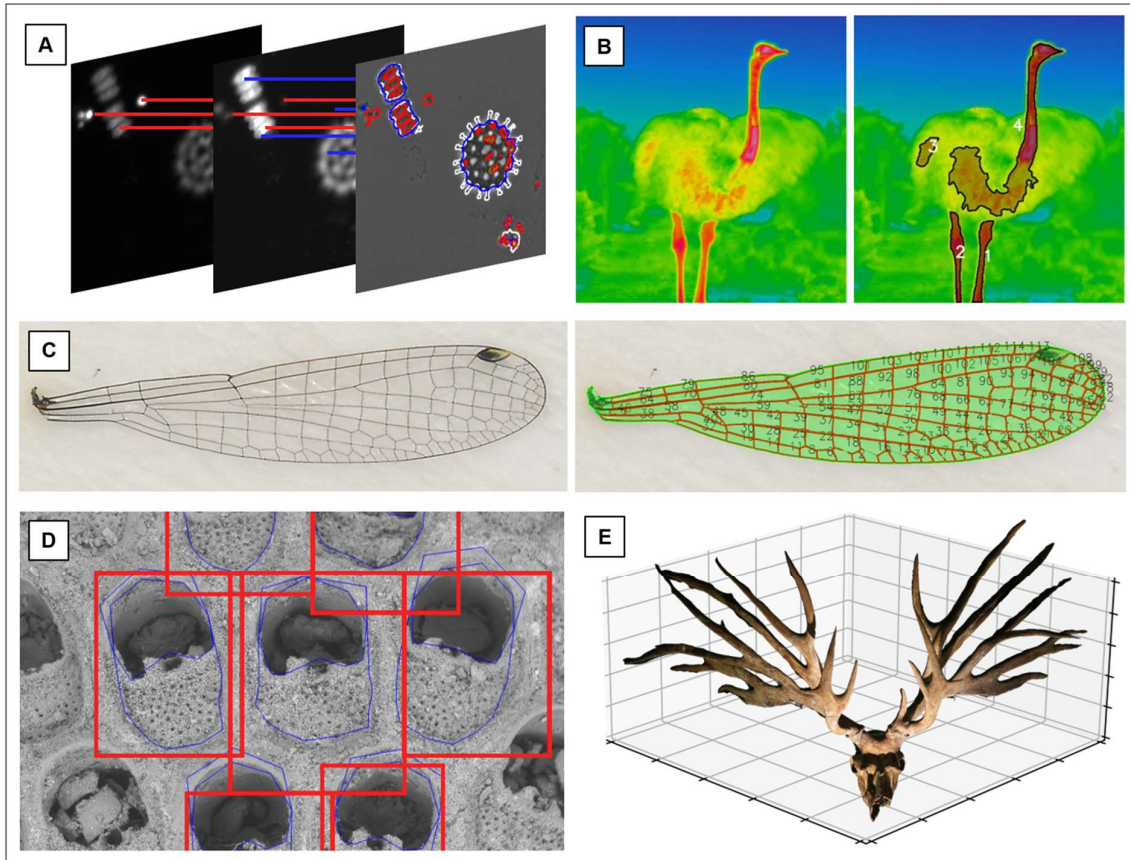
908

909

910

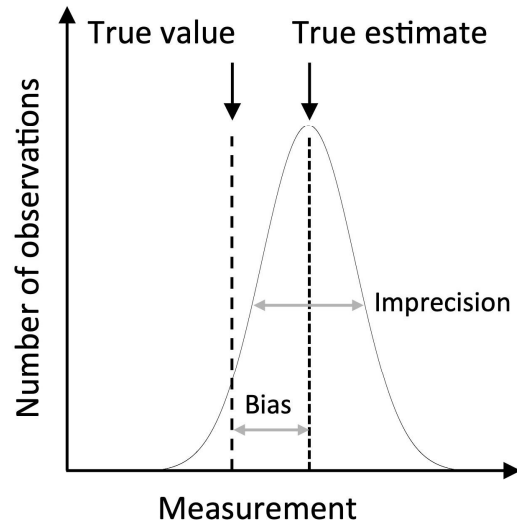


913 Figure 3 - A typical computer vision workflow using signal processing. 1) Raw image - The  
914 goal is to detect, count and measure freshwater isopods (*Asellus aquaticus*, image by Moritz  
915 Lürig) from the raw image that was taken under controlled laboratory conditions. 2)  
916 Preprocessing - The operating principle of most signal processing workflows is that the objects  
917 of interest are made to contrast strongly from all other pixels, meaning that images should  
918 have a high *signal-to-noise ratio* (SNR. In this specific case a high SNR is already present,  
919 because the isopods are much darker than the tray they are sitting on and much larger than  
920 the fecal pellets and other detritus around them. To further increase the SNR, gaussian blur  
921 blends pixels in a given neighborhood (=kernel size), which effectively removes the smaller  
922 dark objects. The reference card gets excluded manually and can be used to convert pixels to  
923 millimeters and to correct the color space. 3) Segmentation - Using a thresholding algorithm  
924 all connected pixels that are above a specific grayscale value and larger than a specified area  
925 are designated foreground (white) and all pixels become background (black). The output from  
926 this step is referred to as a "binary mask". 4) Measurement - Now the white pixels from the  
927 binary mask can be used to locate the areas of interest in the raw image and to extract  
928 information from them. Discrimination between multiple instances of the same class is referred  
929 to as instance segmentation.



930

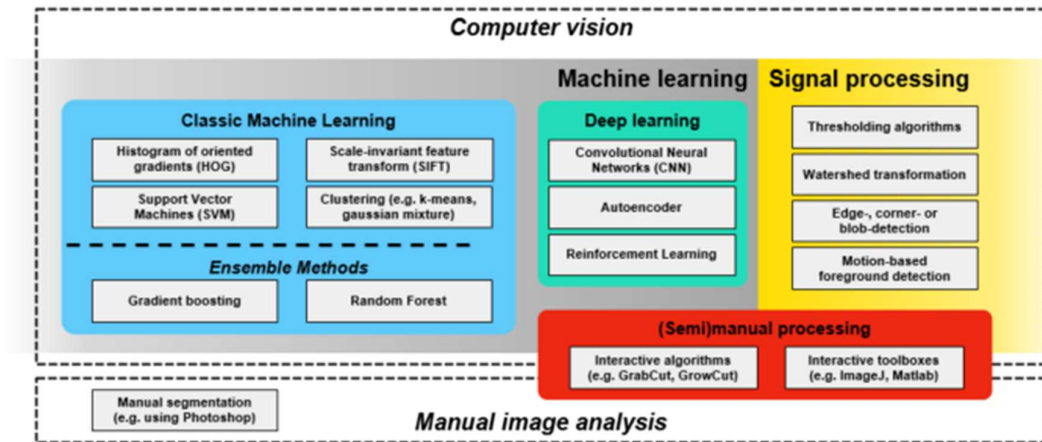
931 Figure 4 - Different types of high dimensional phenotypic data that are collected using a fully  
 932 or semi-automatic computer vision approach. A) Morphology and fluorescence traits of  
 933 phytoplankton communities are represented through a combination of shape features (e.g.  
 934 circularity, perimeter length, area) and texture features (e.g. blob intensity and distribution  
 935 within the cell) from images showing fluorescence intensity (pictograms on the left; images by  
 936 Irene Gallego). B) In ostriches (*Struthio camelus*), surface temperatures of bare body parts  
 937 without feathers (necks and legs) are detected using signal processing (image by Erik  
 938 Svensson). C) Signal processing approach that captures individual domains of a damselfly  
 939 wing via thresholding (image by Masahito Tsuboi). D) Ensemble-based approach to shape  
 940 prediction of individual zooids within a bryozoan colony (image by Arthur Porto) E) 3D image  
 941 of the skull of extinct deer *Eucladoceros dicranios* from which we can measure linear, area,  
 942 and volumetric measurements of antler features (image by Masahito Tsuboi).



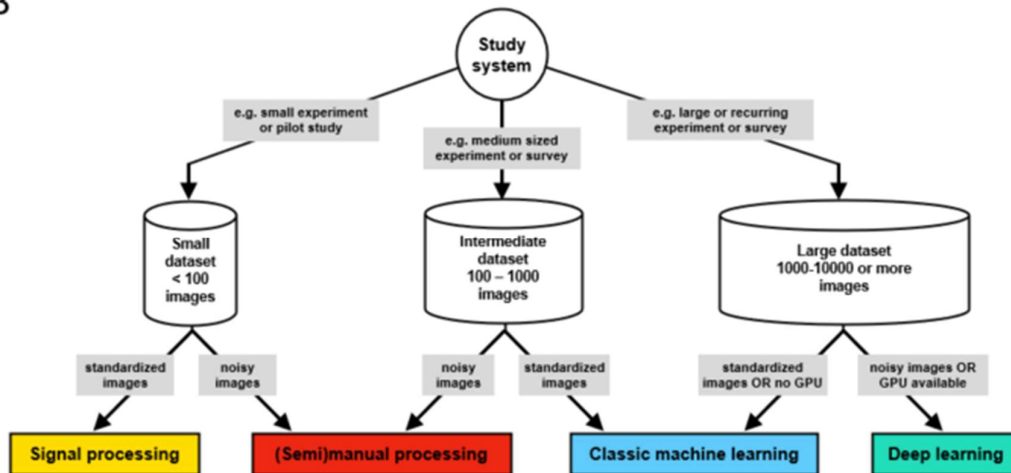
943

944 Figure 5 - Schematic illustration of bias and imprecision. X-axis represents phenotypic values  
 945 and Y-axis represents number of observations. The gaussian curve shows the distribution of  
 946 repeated measurements of the same specimen. Dashed line is the true estimate, and the  
 947 variance of measurements around the true estimate is the imprecision. The true value may  
 948 deviate systematically from the true estimate (long-dashed line). The difference between true  
 949 estimate and true value is the bias.

A



B



950



951 Figure 6 - Computer vision (CV) methods overview - which is the right one for my data? A) CV  
952 is a field at the intersection of machine learning and signal processing which is concerned with  
953 the automatic and semiautomatic extraction of information from digital images. B) Decision  
954 tree for CV methods: begin by considering the size of a given image dataset, whether it is  
955 complete, e.g. after an experiment, or whether there will be continued future additions, e.g. as  
956 part of a long-term experiment or field survey. The next assessment should be about the  
957 signal-to-noise ratio (SNR) in your images: images taken in the laboratory typically have a high  
958 degree of standardization and thus a very high SNR, which makes them suitable for a signal  
959 processing approach. In contrast to signal processing, where segmentation results are  
960 immediately available, all machine learning image analysis pipelines include iterative training  
961 and validation phases, followed by a final testing phase. Such a workflow can be complex to  
962 initiate but pays off in the long run by providing segmentation results that become increasingly  
963 robust if more training data is supplied over time. Deep learning algorithms require large  
964 training datasets (several 1000s to 10000s) but are less sensitive to noise and idiosyncrasies  
965 of the foreground. Thus, for large and continuously growing data sets, or for recurring image  
966 analysis tasks, deep learning has become the standard approach for segmentation.

967

968

969

970 **Tables**

971 Table 1 - Classes of phenotypic data. Depending on the research question, scientists define  
972 their phenotypes of interest using specific or abstract, low or high dimensional traits (see  
973 section *Measurement theory*). The human eye excels at rapidly recognizing externally visible  
974 phenotypes (e.g. benthic vs. limnetic morphotypes of fish), but has difficulties discerning what  
975 constitutes such phenotypes. Computer vision offers an objective way to collect any data type  
976 with high efficiency and reproducibility. For instance, by breaking down low dimensional traits  
977 (e.g. red vs. blue phenotype) into continuous low or high dimensional metrics (e.g. degree of  
978 red- or blueness), the decision of what constitutes a phenotype becomes more reproducible.

Trait type	Low dimensional	High dimensional
Specific / directly measurable	Size, discrete color (“red phenotype” vs. “blue phenotype”) and morphotype scoring (e.g. benthic vs limnetic)	Shape coordinates, texture maps, landmarks
Abstract / derived	Shape (e.g. circularity, area) and texture features (e.g. mean, SD, uniformity), moments, principal components, hypervolumes	Matrices, activation maps

979

980

981

982 Table 2 - Select examples of recent open-source computer vision libraries with a biology-  
 983 context. Although typically first developed for a particular study system or organism (e.g.  
 984 PlantCV or WorMachine), most CV applications apply techniques that are generally applicable  
 985 to any type of phenotypic data contained in digital images.

Name	Year	Reference	Repository	Purpose	Application type	Description	Techniques
AutoMorph	2018	(Hsiang et al. 2018)	<a href="https://github.com/HullLab/AutoMorph">https://github.com/HullLab/AutoMorph</a>	object detection and feature extraction	Python package	High throughput segmentation	Signal processing
ClickPoints	2017	(Gerum et al. 2017)	<a href="https://github.com/fabrylab/clickpoints">https://github.com/fabrylab/clickpoints</a>	labelling, label evaluation	Python package	Interactive labelling tool	signal processing
DeepMerkat	2018	(Weinstein 2015)	<a href="https://github.com/bw4sz/DeepMerkat">https://github.com/bw4sz/DeepMerkat</a>	object detection, classification	Python	Background subtraction and image classification for stationary cameras in ecological videos	Signal processing, deep learning
EB-Net	2020	(Le et al. 2020)	<a href="https://github.com/linhlevandlu/CNN_Beetles_Landmarks">https://github.com/linhlevandlu/CNN_Beetles_Landmarks</a>	keypoint and feature detection	Python	Insect morphometrics	deep learning
ImageJ	2012	(Schindelin et al. 2012)	<a href="https://fiji.sc/">https://fiji.sc/</a> <a href="https://imagej.nih.gov/ij/download.html">https://imagej.nih.gov/ij/download.html</a>	multi-purpose	standalone	Comprehensive, multi-purpose image processing library	manual processing, signal processing, classic machine learning, feature extraction
ML-morph	2020	(Porto and Lysne Voje 2020)	<a href="https://github.com/agport/ml-morph">https://github.com/agport/ml-morph</a>	landmark detection; geometric morphometrics	Python package	High throughput morphometrics	Classic machine learning, ensemble Methods
MotionMeerkat	2015	(Weinstein 2015)	<a href="https://github.com/bw4sz/DeepMerkat">https://github.com/bw4sz/DeepMerkat</a>	motion tracking	Python package/standalone	Deep learning driven motion detection	Signal processing, deep learning
Phenotype	2020	(Lürig 2018)	<a href="https://github.com/mlueurig/phenotype">https://github.com/mlueurig/phenotype</a>	object detection, feature extraction, motion tracking	Python package	Computer vision library with high throughput workflows	signal processing

PlantCV	2017	(Gehan et al. 2017)	<a href="https://github.com/danforthcenter/plantcv">https://github.com/danforthcenter/plantcv</a>	object detection and feature extraction; spectral analysis	Python package	Plant phenotyping library	signal processing, classic machine learning
Scan-o-matic	2016	(Zackrisson et al. 2016)	<a href="https://github.com/Scan-o-Matic/scano-matic">https://github.com/Scan-o-Matic/scano-matic</a>	object detection and feature extraction	Python package	Microbial phenotyping platform	Signal processing
Trackdem	2017	(Bruijning et al. 2018)	<a href="https://github.com/marioleinbruijning/trackdem">https://github.com/marioleinbruijning/trackdem</a>	motion tracking and blob counting	R package	Behavioral analysis pipeline	Signal processing
WingMachine	2003	(Houle et al. 2003)	<a href="https://www.bio.fsu.edu/~dhoule/Software/">https://www.bio.fsu.edu/~dhoule/Software/</a>	keypoint and feature detection	standalone	Drosophila wing morphometrics	Signal processing, feature extraction
WorMachine	2018	(Hakim et al. 2018)	<a href="https://github.com/adamhak/WorMachineClient">https://github.com/adamhak/WorMachineClient</a>	object detection and feature extraction	Matlab	Integrated image processing and feature extraction	Signal processing, classic machine learning; deep learning

986

987

988

989

990 **Boxes**

991 Box 1 - Glossary of terms relevant for computer vision and machine learning in ecology and  
 992 evolution used in this review. Terms in this list are printed in *italic* when first mentioned in the  
 993 main text.  
 994

bit depth	number of values a pixel can take (e.g. 8 bit = $2^8 = 256$ values)
computer vision	technical domain at the intersection of signal processing, machine learning, robotics and other scientific areas that is concerned with the automated extraction of information from digital images and videos.
convolution	mathematical operation by which information contained in images are abstracted. Each convolutional layer produces a feature map, which is passed on to the next layer.
deep learning	machine learning methods based on neural networks. supervised learning = algorithm learns input features from input-output pairs (e.g. labelled images). unsupervised = algorithm looks for undetected patterns (e.g. images without labelling)
feature	a measurable property or pattern. can be specific (e.g. edges, corners, points) or abstract (e.g. convolution via kernels), and combined to vectors and matrices (feature maps)
feature detection	methods for making pixel-level or pixel-neighborhood decisions on whether parts of an image are a feature or not
foreground	all pixels of interest in a given image, whereas the background constitutes all other pixels. the central step in computer vision is the segmentation of all pixels into foreground and background
hidden layer	a connected processing step in neural networks during which information is received, processed (e.g. convolved), and passed on to the next layer
kernel	a small mask or matrix to perform operations on images, for example, blurring, sharpening or edge detection. the kernel operation is performed pixel wise, sliding across the entire image.
labelling	typically manual markup of areas of interest in an image by drawing bounding boxes or polygons around the contour. can be multiple objects and multiple classes of objects per image. can also refer to assigning whole images to a class (e.g. relevant for species identification)
machine learning	subset of artificial intelligence: the study and implementation of computer algorithms that improve automatically through experience. (Mitchell 1997)
measurement theory	concerns the relationship between measurements and nature so that inferences from measurements reflect the underlying reality intended to be represent (Houle et al. 2011).
neural network	deep learning algorithms that use multi layered ("deep") abstractions of information to extract higher level features from input via convolution

object detection	methods for determining whether a pixel region constitutes an object that belongs to the foreground or not, based on its features
phenome	the phenotype as a whole (sensu Soulé 1967)
phenomics	the acquisition of high-dimensional phenotypic data on an organism-wide scale
phenotype	a single trait or a specific set of traits that are part of the phenome
pixel	short for picture element; the smallest accessible unit of a digital raster image. Pixels have finite values (=intensities), e.g. 256 in an 8-bit grayscale image.
segmentation	the classification of all pixels in an image into foreground and background, either manually by labelling the area of interest, or automatically, by means of signal processing or machine learning algorithms. semantic segmentation = all pixels of a class, instance segmentation = all instances of a class
signal processing	technically correct: digital image processing (not to be confused with image analysis or image editing). subfield of engineering that is concerned with the filtering or modification of digital images by means of algorithms and filter matrices (kernels),
signal-to-noise ratio (SNR)	describes the level of the pixels containing the desired signal (i.e. the phenotypic information) to all other pixels. Lab images typically have a high SNR, field images a low SNR.
threshold algorithm	pixel-intensity based segmentation of images, e.g. based on individual pixel intensity (binary thresholding) or their intensity with respect to their neighborhood (adaptive thresholding). creates a binary mask which contains only black or white pixels
training data	representative image dataset to train a machine learning algorithm. can be created manually by labelling images, or semi-automatic by using signal processing for segmentation. can contain single or multiple classes
watershed algorithm	the segmentation of images by treating the pixels as a topographic map of basins, where bright pixels have high elevation and dark pixels have low elevation.

995

996

997

998

999

1000

1001

1002

1003 **Box 2 - An overview of the main deep learning architectures and approaches.**

1004

1005 **Families of network topologies**

1006

1007 A. **Deep convolutional network** - A large and common family of neural networks  
1008 composed an input layer, an output layer and multiple hidden layers. These networks  
1009 feature convolution kernels that process input data and pooling layers that simplify the  
1010 information processed through the convolutional kernels. For certain tasks, the input  
1011 can be a window of the image, rather than the entire image.

1012 B. **Deconvolutional Network** - A smaller family of neural networks that perform the  
1013 reverse process when compared to convolutional networks. It starts with the processed  
1014 data (i.e., the output of the convolutional network) and it aims to separate what has  
1015 been convoluted. Essentially, it constructs upwards from processed data (e.g.,  
1016 reconstructs an image from a label).

1017 C. **Generative Adversarial Network** - A large family of networks composed of two  
1018 separate networks, a generator and a discriminator. The generator is trained to  
1019 generate realistic data, while the discriminator is trained to differentiate between  
1020 generated data from actual samples. Essentially, in this approach, the objective is for  
1021 the generator to generate such realistic data that the discriminator cannot tell it apart  
1022 from samples.

1023 D. **Autoencoders** - A family of networks is trained in an unsupervised manner. The  
1024 autoencoder aims to learn how to robustly represent the original dataset, oftentimes in  
1025 smaller dimensions, even in the presence of noise. Autoencoders are composed of  
1026 multiple layers, and it can be divided into two main parts: the encoder and the decoder.  
1027 The encoder maps the input into the representation and the decoder uses the  
1028 representation to reconstruct the original input.

1029 E. **Deep Belief Network** - A family of generative networks that are composed of multiple  
1030 layers of hidden units, in which there can be connections between layers but not within  
1031 units within layers. Deep belief networks can be conceived as being composed of  
1032 multiple simpler networks, where each subnetwork's hidden layer acts as a visible layer  
1033 to another subnetwork.

1034

1035 **Learning Classes**

1036

1037 A. **Supervised Learning** - Training data is provided when fitting the model. The training  
1038 dataset is composed of inputs and expected outputs. Models are tested by making  
1039 predictions based on inputs and comparing them with expected outputs.

1040 B. **Unsupervised Learning** - No training data is provided to the model. Unsupervised  
1041 learning relies exclusively on inputs. Models trained using unsupervised learning are  
1042 used to describe or extract relationships in image data, such as clustering or  
1043 dimensionality reduction.

1044 C. **Reinforcement Learning** - The learning process occurs in a supervised manner, but  
1045 not through the use of static training datasets. Rather, in reinforcement learning, the  
1046 model is directed towards a goal, with a limited set of actions it may perform, and model  
1047 improvement is obtained through feedback. The learning itself occurs exclusively  
1048 through feedback obtained based on past action. This feedback can be quite noisy and  
1049 delayed.

1050 D. **Hybrid Learning Problems**

1051 **Semi-Supervised Learning** - Semi supervised learning relies on training datasets  
1052 where only a small percentage of the training dataset is labeled, with the  
1053 remaining images having no label. It is a hybrid in between supervised and  
1054 unsupervised learning, since the model has to make effective use of unlabeled data  
1055 while relying only partially on labeled ones.

1056 **Self-Supervised Learning** - Self supervised learning uses a combination of  
1057 unsupervised and supervised learning. In this approach, supervised learning is used

1058 to solve a pretext task for which training data is available (or can be artificially  
1059 provided), and whose representation can be used to solve an unsupervised learning  
1060 problem. Generative adversarial networks rely on this technique to learn how to  
1061 artificially generate image data.  
1062

### 1063 **Other learning Techniques**

- 1064
- 1065 A. **Active Learning** - During active learning, the model can query the user during the  
1066 learning process to require labels for new data points. It requires human interaction,  
1067 and it aims to being more efficient about what training data is used by the model
  - 1068 B. **Online Learning** - Online learning techniques are often used in situations where  
1069 observations are streamed through time and in which the probability distribution of the  
1070 data might drift over time. In this technique, the model is updated as more data  
1071 becomes available, allowing the model itself to change through time.
  - 1072 C. **Transfer Learning** - Transfer learning is a useful technique when training a model for  
1073 a task that is related to another task for which a robust model is already available.  
1074 Essentially, it treats the already robust model as a starting point from which to train a  
1075 new model. It greatly diminishes the training data needs of supervised models and it  
1076 is, therefore, used when the available training data is limited.
  - 1077 D. **Ensemble Learning** - As mentioned in the main text, ensemble learning refers to a  
1078 learning technique in which multiple models are trained either in parallel or sequentially  
1079 and the final prediction is the result of the combination of the predictions generated by  
1080 each component.  
1081

1082

1083