

# 1 DeepPoseKit: a software toolkit for 2 fast and robust animal pose 3 estimation using deep learning

4 **Jacob M. Graving**<sup>1,2,3,\*</sup>, **Daniel Chae**<sup>4</sup>, **Hemal Naik**<sup>1,2,3,5</sup>, **Liang Li**<sup>1,2,3</sup>, **Benjamin**  
5 **Koger**<sup>1,2,3</sup>, **Blair R. Costelloe**<sup>1,2,3</sup>, **Iain D. Couzin**<sup>1,2,3,\*</sup>

\*For correspondence:

[jgraving@gmail.com](mailto:jgraving@gmail.com); [icouzin@ab.mpg.de](mailto:icouzin@ab.mpg.de)

6 <sup>1</sup>Department of Collective Behaviour, Max Planck Institute of Animal Behavior, 78464  
7 Konstanz, Germany; <sup>2</sup>Department of Biology, University of Konstanz, 78464 Konstanz,  
8 Germany; <sup>3</sup>Centre for the Advanced Study of Collective Behaviour, University of Konstanz,  
9 78464 Konstanz, Germany; <sup>4</sup>Department of Computer Science, Princeton University,  
10 08544 Princeton, NJ, USA; <sup>5</sup>Chair for Computer Aided Medical Procedures, Technische  
11 Universität München, 80333 Munich, Germany

---

13 **Abstract** Quantitative behavioral measurements are important for answering questions across scientific  
14 disciplines—from neuroscience to ecology. State-of-the-art deep-learning methods offer major advances in  
15 data quality and detail by allowing researchers to automatically estimate locations of an animal's body parts  
16 directly from images or videos. However, currently-available animal pose estimation methods have limitations  
17 in speed and robustness. Here we introduce a new easy-to-use software toolkit, *DeepPoseKit*, that addresses  
18 these problems using an efficient multi-scale deep-learning model, called *Stacked DenseNet*, and a fast  
19 GPU-based peak-detection algorithm for estimating keypoint locations with subpixel precision. These advances  
20 improve processing speed >2x with no loss in accuracy compared to currently-available methods. We  
21 demonstrate the versatility of our methods with multiple challenging animal pose estimation tasks in laboratory  
22 and field settings—including groups of interacting individuals. Our work reduces barriers to using advanced  
23 tools for measuring behavior and has broad applicability across the behavioral sciences.

---

## 25 Introduction

26 Understanding the relationships between individual behavior, brain activity (reviewed by *Krakauer et al. 2017*),  
27 and collective and social behaviors (*Rosenthal et al., 2015; Strandburg-Peshkin et al., 2013; Jolles et al., 2017;*  
28 *Klibaite et al., 2017; Klibaite and Shaevitz, 2019*) is a central goal of the behavioral sciences—a field that spans  
29 disciplines from neuroscience to psychology, ecology, and genetics. Measuring and modelling behavior is key to  
30 understanding these multiple scales of complexity, and, with this goal in mind, researchers in the behavioral  
31 sciences have begun to integrate theory and methods from physics, computer science, and mathematics  
32 (*Anderson and Perona, 2014; Berman, 2018; Brown and De Bivort, 2018*). A cornerstone of this interdisciplinary  
33 revolution is the use of state-of-the-art computational tools, such as computer vision algorithms, to automatically  
34 measure locomotion and body posture (*Dell et al., 2014*). Such a rich description of animal movement then  
35 allows for modeling, from first principles, the full behavioral repertoire of animals (*Stephens et al., 2011;*  
36 *Berman et al., 2014a, 2016; Wiltchko et al., 2015; Johnson et al., 2016b; Todd et al., 2017; Klibaite et al., 2017;*  
37 *Markowitz et al., 2018; Klibaite and Shaevitz, 2019; Costa et al., 2019*). Tools for automatically measuring  
38 animal movement represent a vital first step toward developing unified theories of behavior across scales  
39 (*Berman, 2018; Brown and De Bivort, 2018*). Therefore, technical factors like scalability, robustness, and usability  
40 are issues of critical importance, especially as researchers across disciplines begin to increasingly rely on these  
41 methods.

42 Two of the latest contributions to the growing toolbox for quantitative behavioral analysis are from *Mathis*  
43 *et al. (2018)* and *Pereira et al. (2019)*, who make use of a popular type of machine learning model called  
44 *convolutional neural networks*, or *CNNs* (*LeCun et al. 2015*; Appendix 1), to automatically measure detailed  
45 representations of animal posture—structural *keypoints*, or *joints*, on the animal's body—directly from images

46 and without markers. While these methods offer a major advance over conventional methods with regard  
47 to data quality and detail, they have disadvantages in terms of speed and robustness, which may limit their  
48 practical applications. To address these problems, we introduce a new software toolkit, called *DeepPoseKit*, with  
49 methods that are fast, robust, and easy-to-use. We run experiments using multiple datasets to compare our new  
50 methods with those from *Mathis et al. (2018)* and *Pereira et al. (2019)*, and we find that our approach offers  
51 considerable improvements. These results also demonstrate the flexibility of our toolkit for both laboratory and  
52 field situations and exemplify the wide applicability of our methods across a range of species and experimental  
53 conditions.

## 54 **Measuring animal movement with computer vision**

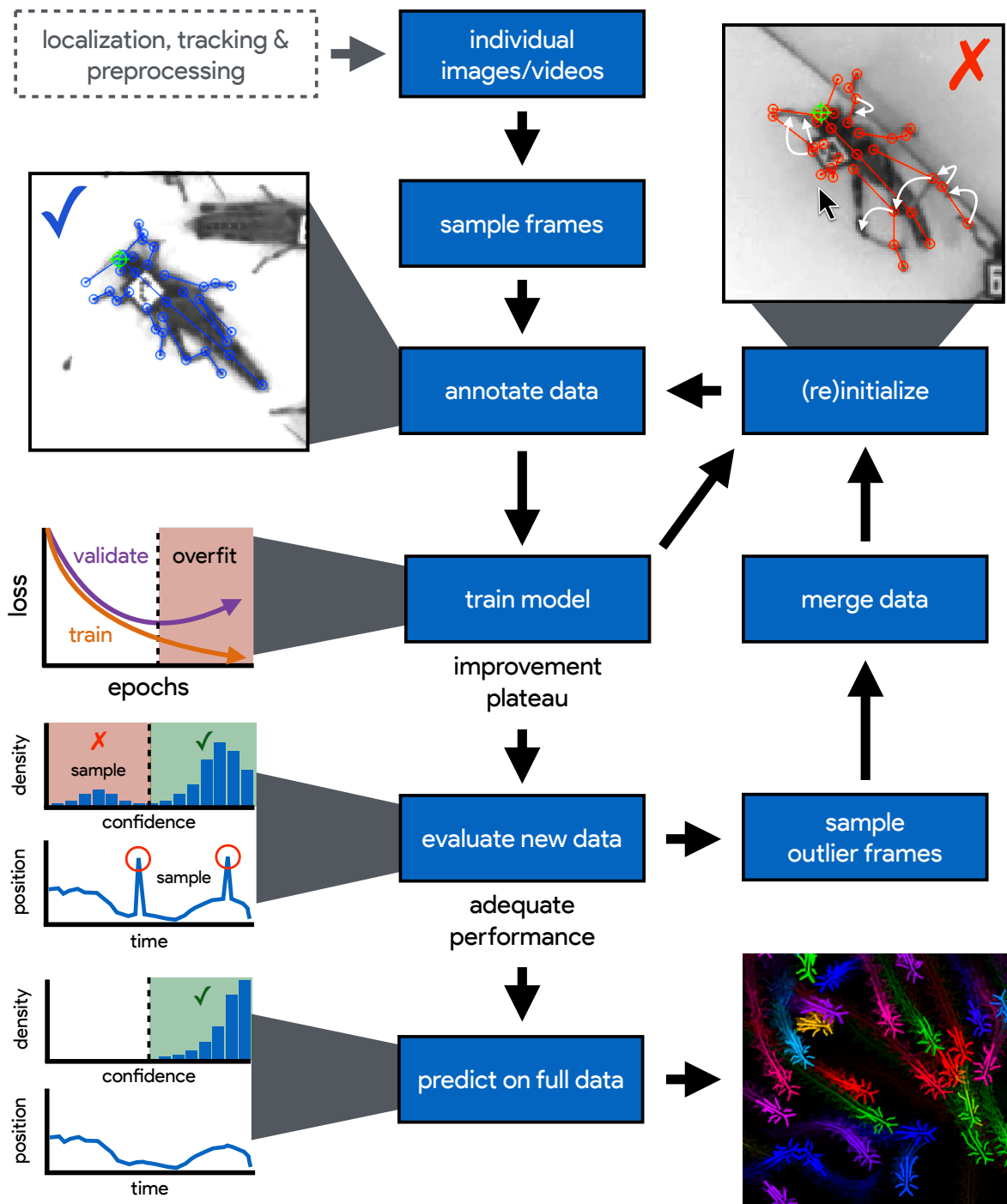
55 Collecting high-quality behavioral data is a challenging task, and while direct observations are important for  
56 gathering qualitative data about a study system, a variety of automated methods for quantifying movement have  
57 become popular in recent years (*Dell et al., 2014; Anderson and Perona, 2014; Kays et al., 2015*). Methods like  
58 video monitoring and recording help to accelerate data collection and reduce the effects of human intervention,  
59 but the task of manually scoring videos is time consuming and suffers from the same limitations as direct  
60 observation, namely observer bias and mental fatigue. Additionally, due to limitations of human observers'  
61 ability to process information, many studies that rely on manual scoring use relatively small datasets to estimate  
62 experimental effects, which can lead to increased rates of statistical errors. Studies that lack the statistical  
63 resolution to robustly test hypotheses (commonly called "power" in frequentist statistics) also raise concerns  
64 about the use of animals for research, as statistical errors caused by sparse data can impact researchers' ability  
65 to accurately answer scientific questions. These limitations have led to the development of automated methods  
66 for quantifying behavior using advanced imaging technologies (*Dell et al., 2014*) as well as sophisticated tags and  
67 collars with GPS, accelerometry, and acoustic-recording capabilities (*Kays et al., 2015*). Tools for automatically  
68 measuring the behavior of individuals now play a central role in our ability to study the neurobiology and  
69 ecology of animals, and reliance on these technologies for studying animal behavior will only increase in the  
70 future.

71 The rapid development of computer vision hardware and software in recent years has allowed for the use of  
72 automated image-based methods for measuring behavior across many experimental contexts (*Dell et al., 2014*).  
73 Early methods for quantifying movement with these techniques required highly-controlled laboratory conditions.  
74 However, because animals exhibit different behaviors depending on their surroundings (*Strandburg-Peshkin*  
75 *et al., 2017; Francisco et al., 2019; Akhund-Zade et al., 2019*), laboratory environments are often less than ideal  
76 for studying many natural behaviors. Most conventional computer vision methods are also limited in their ability  
77 to accurately track groups of individuals over time, but nearly all animals are social at some point in their life and  
78 exhibit specialized behaviors when in the presence of conspecifics (*Strandburg-Peshkin et al., 2013; Rosenthal*  
79 *et al., 2015; Jolles et al., 2017; Francisco et al., 2019; Versace et al., 2019*). These methods also commonly track  
80 only the animal's center of mass, which reduces the behavioral output of an individual to a two-dimensional or  
81 three-dimensional particle-like trajectory. While trajectory data are useful for many experimental designs, the  
82 behavioral repertoire of an animal cannot be fully described by its aggregate locomotory output. For example,  
83 stationary behaviors, like grooming and antennae movements, or subtle differences in walking gaits cannot be  
84 reliably detected by simply tracking an animal's center of mass (*Berman et al., 2014a; Wiltschko et al., 2015*).

85 Together these factors have driven the development of software that can accurately track the positions of  
86 marked (*Crall et al., 2015; Graving, 2017; Wild et al., 2018; Boenisch et al., 2018*) or unmarked (*Pérez-Escudero*  
87 *et al., 2014; Romero-Ferrero et al., 2019*) individuals as well as methods that can quantify detailed descriptions  
88 of an animal's posture over time (*Stephens et al., 2011; Berman et al., 2014a; Wiltschko et al., 2015; Mathis*  
89 *et al., 2018; Pereira et al., 2019*). Recently these advancements have been further improved through the  
90 use of deep learning, a class of machine learning algorithms that learn complex statistical relationships from  
91 data (*LeCun et al., 2015*). Deep learning has opened the door to accurately tracking large groups of marked  
92 (*Wild et al., 2018; Boenisch et al., 2018*) or unmarked (*Romero-Ferrero et al., 2019*) individuals and has made  
93 it possible to measure the body posture of animals in nearly any context—including in the wild (*Nath et al.,*  
94 *2019*)—by tracking the positions of user-defined body parts (*Mathis et al., 2018; Pereira et al., 2019*). These  
95 advances have drastically increased the quality and quantity, as well as the diversity, of behavioral data that are  
96 potentially available to researchers for answering scientific questions.

## 97 **Animal pose estimation using deep learning**

98 In the past, conventional methods for measuring posture with computer vision relied on species-specific algo-  
99 rithms (*Uhlmann et al., 2017*), highly-specialized or restrictive experimental setups (*Mendes et al., 2013; Kain*  
100 *et al., 2013*), attaching intrusive physical markers to the study animal (*Kain et al., 2013*), or some combination  
101 thereof. These methods also typically required expert computer-vision knowledge to use, were limited in the



**Figure 1.** An illustration of the workflow for DeepPoseKit. Multi-individual images are localized, tracked, and preprocessed into individual images, which is not required for single-individual image datasets. An initial image set is sampled, annotated, and then iteratively updated using the active learning approach described by *Pereira et al. (2019)* (see Appendix 2). As annotations are made, the model is trained (Figure 2) with the current training set and keypoint locations are initialized for unannotated data to reduce the difficulty of further annotations. This is repeated until there is a noticeable improvement plateau for the initialized data—where the annotator is providing only minor corrections—and for the validation error when training the model (Appendix Figure 8). New data from the full dataset are evaluated with the model, and the training set is merged with new examples that are sampled based on the model's predictive performance, which can be assessed with techniques described by *Mathis et al. (2018)*; *Nath et al. (2019)* for identifying outlier frames and minimizing extreme prediction errors—shown here as the distribution of confidence scores predicted by the model and predicted body part positions with large temporal derivatives—indicating extreme errors. This process is repeated as necessary until performance is adequate when evaluating new data. The pose estimation model can then be used to make predictions for the full data set, and the data can be used for further analysis.

**Figure 1–video 1.** A visualization of the posture data output for a group of locusts (5x speed).

102 number or type of body parts that could be tracked (*Mendes et al., 2013*), involved capturing and handling  
103 the study animals to attach markers (*Kain et al., 2013*)—which is not possible for many species—and despite  
104 best efforts to minimize human involvement, often required manual intervention to correct errors (*Uhlmann  
105 et al., 2017*). All of these methods were built to work for a small range of conditions and typically required  
106 considerable effort to adapt to novel contexts.

107 In contrast to conventional computer-vision methods, modern deep-learning-based methods can be used  
108 to achieve near human-level accuracy in almost any scenario by manually annotating data (Figure 1)—known as  
109 a *training set*—and training a general-purpose image-processing algorithm—a convolutional neural network or  
110 CNN—to automatically estimate the locations of an animal's body parts directly from images (Figure 2). State-of-  
111 the-art machine learning methods, like CNNs, use these training data to parameterize a model describing the  
112 statistical relationships between a set of input data—i.e., images—and the desired output distribution—i.e.,  
113 posture keypoints. After adequate training, a model can be used to make predictions on previously-unseen  
114 data from the same dataset—inputs that were not part of the training set, which is known as *inference*. In other  
115 words, these models are able to generalize human-level expertise at scale after having been trained on only a  
116 relatively small number of examples. We provide more detailed background information on using CNNs for  
117 pose estimation in Appendices 1–5.

118 Similar to conventional pose estimation methods, the task of implementing deep-learning models in software  
119 and training them on new data is complex and requires expert knowledge. However, in most cases, once the  
120 underlying model and training routine are implemented, a high-accuracy pose estimation model for a novel  
121 context can be built with minimal modification—often just by changing the training data. With a simplified  
122 toolkit and high-level software interface designed by an expert, even scientists with limited computer-vision  
123 knowledge can begin to apply these methods to their research. Once the barriers for implementing and training  
124 a model are sufficiently reduced, the main bottleneck for using these methods becomes collecting an adequate  
125 training set—a labor-intensive task made less time-consuming by techniques described in Appendix 2.

126 *Mathis et al. (2018)* and *Pereira et al. (2019)* were the first to popularize the use of CNNs for animal pose  
127 estimation. These researchers built on work from the human pose estimation literature (e.g., *Andriluka et al.  
128 2014; Insafutdinov et al. 2016; Newell et al. 2016*) using a type of *fully-convolutional neural network* or *F-CNN*  
129 (*Long et al. 2015*; Appendix 3) often referred to as an *encoder-decoder* model (Appendix 3 Box 1). These models  
130 are used to measure animal posture by training the network to transform images into probabilistic estimates  
131 of keypoint locations, known as *confidence maps* (shown in Figure 2), that describe the body posture for one  
132 or more individuals. These confidence maps are processed to produce the 2-D spatial coordinates of each  
133 keypoint, which can then be used for further analysis.

134 While these models typically need large amounts of training data, both *Mathis et al. (2018)* and *Pereira  
135 et al. (2019)* have demonstrated that near human-level accuracy can be achieved with few training examples  
136 (Appendix 2). In order to ensure generalization to large datasets, both groups of researchers introduced ideas  
137 related to iteratively refining the training set used for model fitting (*Mathis et al., 2018; Pereira et al., 2019*). In  
138 particular, *Pereira et al. (2019)* describe a technique known as *active learning* where a trained model is used  
139 to initialize new training data and reduce annotation time (Appendix 2). *Mathis et al. (2018)* describe multiple  
140 techniques that can be used to further refine training data and minimize errors when making predictions on the  
141 full dataset. Simple methods to accomplish this include filtering data or selecting new training examples based  
142 on confidence scores or the entropy of the confidence maps from the model output. *Nath et al. (2019)* also  
143 introduced the use temporal derivatives (i.e., speed and acceleration) and autoregressive models to identify  
144 outlier frames, which can then be labeled to refine the training set or excluded from further analysis on the final  
145 dataset (Figure 1).

## 146 **Pose estimation models and the speed-accuracy trade-off**

147 *Mathis et al. (2018)* developed their pose estimation model, which they call *DeepLabCut*, by modifying a  
148 previously-published model called *DeeperCut* (*Insafutdinov et al., 2016*). The *DeepLabCut* model (*Mathis  
149 et al., 2018*), like the *DeeperCut* model, is built on the popular *ResNet* architecture (*He et al., 2016*)—a state-of-  
150 the-art deep-learning model used for image classification. This choice is advantageous because the use of a  
151 popular architecture allows for incorporating a pre-trained encoder to improve performance and reduce the  
152 number of required training examples (*Mathis et al., 2018*), known as *transfer learning* (*Pratt 1993*; Appendix  
153 2)—although, as will be seen, our results suggest that transfer learning offers only a small improvement over a  
154 randomly-initialized model. However, this choice of of a pre-trained architecture is also disadvantageous as  
155 the model is *overparameterized* with >25 million parameters. Overparameterization allows the model to make  
156 accurate predictions, but this may come with the cost of slow inference. To alleviate these effects, work from  
157 *Mathis and Warren (2018)* showed that inference speed for the *DeepLabCut* model (*Mathis et al., 2018*) can be  
158 improved by decreasing the resolution of input images, but this is achieved at the expense of accuracy.



159 With regard to model design, *Pereira et al. (2019)* implement a modified version of a model called *SegNet*  
160 (*Badrinarayanan et al., 2015*), which they call *LEAP* (LEAP Estimates Animal Pose), that attempts to limit model  
161 complexity and overparameterization with the goal of maximizing inference speed (see Appendix 5)—however,  
162 the comparisons we make in this paper suggest this strategy achieved only limited success compared to the  
163 DeepLabCut model (*Mathis et al., 2018*). The LEAP model is advantageous because it is explicitly designed for  
164 fast inference but has disadvantages such as a lack of robustness to data variance, like rotations or shifts in  
165 lighting, and an inability to generalize to new experimental setups. Additionally, to achieve maximum perform-  
166 ance, the training routine for the LEAP model introduced by *Pereira et al. (2019)* requires computationally  
167 expensive preprocessing that is not practical for many datasets, which makes it unsuitable for a wide range of  
168 experiments (see Appendix 5 for more details).

169 Together the methods from *Mathis et al. (2018)* and *Pereira et al. (2019)* represent the two extremes of a  
170 phenomenon known as the *speed-accuracy trade-off* (*Huang et al., 2017b*)—an active area of research in the  
171 machine learning literature. *Mathis et al. (2018)* prioritize accuracy over speed by using a large overparam-  
172 eterized model (*Insafutdinov et al., 2016*), and *Pereira et al. (2019)* prioritize speed over accuracy by using a  
173 smaller less-robust model. While this speed-accuracy trade-off can limit the capabilities of CNNs, there has been  
174 extensive work to make these models more efficient without impacting accuracy (e.g., *Chollet 2017; Huang et al.*  
175 *2017a; Sandler et al. 2018*). To address the limitations of this trade-off, we apply recent developments from the  
176 machine learning literature and provide an effective solution to the problem.

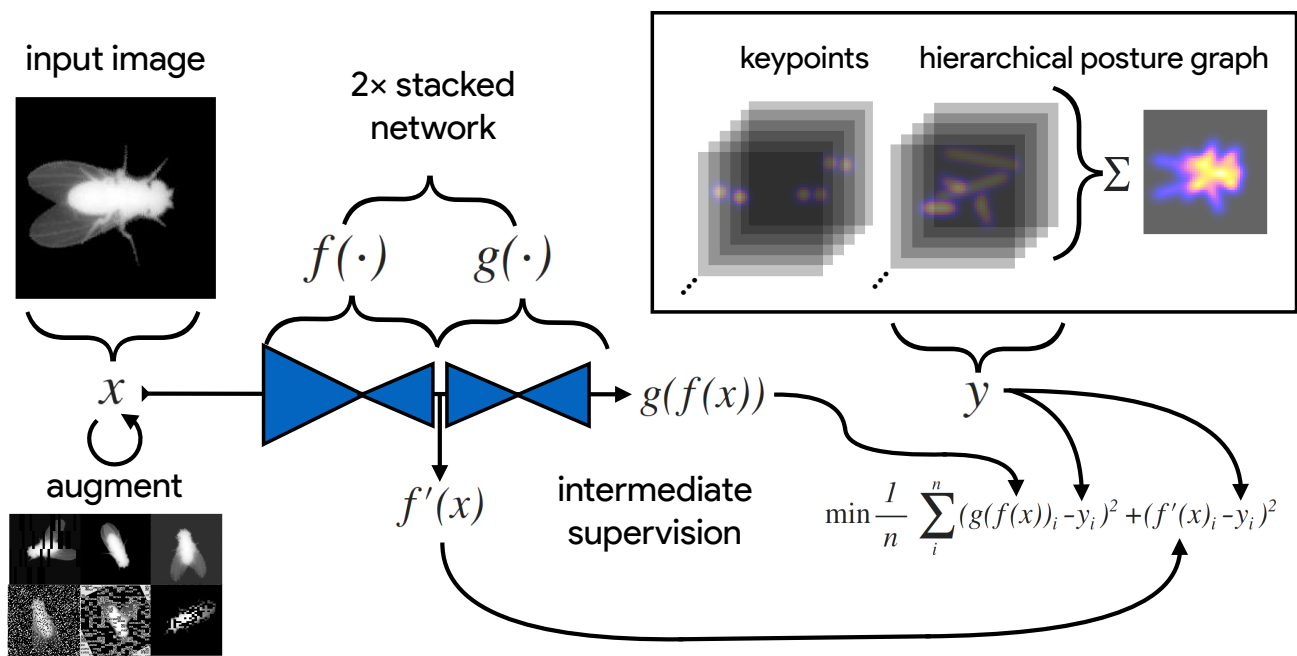
177 In the case of F-CNN models used for pose estimation, improvements in efficiency and robustness have  
178 been made through the use of *multi-scale inference* (Appendix 3 Box 1) by increasing connectivity between the  
179 model's many layers across multiple spatial scales (Appendix 3 Figure 1). Multi-scale inference implicitly allows  
180 the model to simultaneously integrate large-scale global information, such as the lighting, image background,  
181 or the orientation of the focal individual's body trunk; information from intermediate scales like anatomical  
182 geometry related to cephalization and bilateral symmetry; and fine-scale local information that could include  
183 differences in color, texture, or skin patterning for specific body parts. This multi-scale design gives the model  
184 capacity to learn the hierarchical relationships between different spatial scales and efficiently aggregate them  
185 into a joint representation when solving the posture estimation task (see Box 1 and Appendix 3 Figure 1 for  
186 further discussion)

## 187 **Individual vs. multiple pose estimation**

188 Most work on human pose estimation now focuses on estimating the pose of multiple individuals in an  
189 image (e.g., *Cao et al. 2017*). For animal pose estimation, the methods from *Pereira et al. (2019)* are limited  
190 to estimating posture for single individuals—known as *individual pose estimation*—while the methods from  
191 *Mathis et al. (2018)* can also be used to estimate posture for multiple individuals simultaneously—known as  
192 *multiple pose estimation*. However, the majority of work on multiple pose estimation, including *Mathis et al.*  
193 *(2018)*, has not adequately solved the tracking problem of linking individual posture data across frames in a  
194 video, especially after visual occlusions—although recent work has attempted to address this problem (*Iqbal*  
195 *et al., 2017; Andriluka et al., 2018*). Additionally, as the name suggests, the task of multiple pose estimation  
196 requires exhaustively annotating images of multiple individuals—where every individual in the image must be  
197 annotated to prevent the model from learning conflicting information. This type of annotation task is even  
198 more laborious and time consuming than annotations for individual pose estimation and the amount of labor  
199 increases proportionally with the number of individuals in each frame, which makes this approach intractable  
200 for many experimental systems.

201 Reliably tracking the position of individuals over time is important for most behavioral studies, and there  
202 are a number of diverse methods already available for solving this problem (*Pérez-Escudero et al., 2014; Crall*  
203 *et al., 2015; Graving, 2017; Romero-Ferrero et al., 2019; Wild et al., 2018; Boenisch et al., 2018*). Therefore, to  
204 avoid solving an already-solved problem of tracking individuals, and to circumvent the cognitively complex task  
205 of annotating data for multiple pose estimation, the work we describe in this paper is purposefully limited to  
206 individual pose estimation where each image contains only a single focal individual, which may be cropped from  
207 a larger multi-individual image after localization and tracking. We introduce a top-down posture estimation  
208 framework that can be readily adapted to existing behavioral analysis workflows, which could include any  
209 method for localizing and tracking individuals.

210 The additional step of localizing and tracking individuals naturally increases the processing time for producing  
211 posture data from raw image data, which varies depending on the algorithms being used and the number of  
212 individuals in each frame. While tracking and localization may not be practical for all experimental systems,  
213 which could make our methods difficult to apply "out-of-the-box", the increased processing time from automated  
214 tracking algorithms is a reasonable trade-off for most systems given the costly alternative of increased manual  
215 labor when annotating data. This trade-off seems especially practical when considering that the posture data



**Figure 2.** An illustration of the model training process for our Stacked DenseNet model in DeepPoseKit (see Appendix 1 for details about training models). Input images  $x$  (top-left) are augmented (bottom-left) with various spatial transformations (rotation, translation, scale, etc.) followed by noise transformations (dropout, additive noise, blurring, contrast, etc.) to improve the robustness and generalization of the model. The ground truth annotations are then transformed with matching spatial augmentations (not shown for the sake of clarity) and used to draw the confidence maps  $y$  for the keypoints and hierarchical posture graph (top-right). The images  $x$  are then passed through the network to produce a multidimensional array  $g(f(x))$ —a stack of images corresponding to the keypoint and posture graph confidence maps for the ground truth  $y$ . Mean squared error between the outputs for both networks  $g(f(x))$  and  $f'(x)$  and the ground truth data  $y$  is then minimized (bottom-right), where  $f'(x)$  indicates a subset of the output from  $f(x)$ —only those feature maps being optimized to reproduce the confidence maps for the purpose of intermediate supervision (Appendix 4). The loss function is minimized until the validation loss stops improving—indicating that the model has converged or is starting to overfit to the training data.

216 produced by most multiple pose estimation algorithms still need to be linked across video frames to maintain  
 217 the identity of each individual, which is effectively a bottom-up method for achieving the same result. Limiting  
 218 our methods to individual pose estimation also simplifies the pose detection problem as processing confidence  
 219 maps produced by the model does not require computationally-expensive local peak detection and complex  
 220 methods for grouping keypoints into individual posture graphs (e.g., *Insafutdinov et al. 2016*; *Cao et al. 2017*;  
 221 Appendix 3). Additionally, because individual pose estimation is such a well-studied problem in computer vision,  
 222 we can readily build on state-of-the-art methods for this task (see Appendices 3 and 4 for details).

## 223 Methods and Results

224 Here we introduce fast, flexible, and robust pose estimation methods, with a software interface—a high-level  
 225 programming interface (API) and graphical user-interface (GUI) for annotations—that emphasizes usability.  
 226 Our methods build on the state-of-the-art for individual pose estimation (*Newell et al. 2016*; Appendix 4),  
 227 convolutional regression models (*Jégou et al. 2017*; Appendix 3 Box 1), and conventional computer vision  
 228 algorithms (*Guizar-Sicairos et al., 2008*) to improve model efficiency and achieve faster, more accurate results  
 229 on multiple challenging pose estimation tasks. We developed two model implementations—including a new  
 230 model architecture that we call *Stacked DenseNet*—and a new method for processing confidence maps called  
 231 *subpixel maxima* that provides fast and accurate peak detection for estimating keypoint locations with subpixel  
 232 precision—even at low spatial resolutions. We also discuss a modification to incorporate a hierarchical posture  
 233 graph for learning the multi-scale geometry between keypoints on the animal's body, which increases accuracy  
 234 when training pose estimation models. We ran experiments to optimize our approach and compared our new  
 235 models to the models from *Mathis et al. (2018)* (DeepLabCut) and *Pereira et al. (2019)* (LEAP) in terms of speed,  
 236 accuracy, training time, and generalization ability. We benchmark these models using three image datasets  
 237 recorded in the laboratory and the field—including multiple interacting individuals that were first localized and  
 238 cropped from larger, multi-individual images (see "Datasets" for details).

## 239 **An end-to-end pose estimation framework**

240 We provide a full-featured, extensible, and easy-to-use software package that is written entirely in the Python  
241 programming language (Python Software Foundation) and is built on the popular Keras deep-learning package  
242 (Chollet et al., 2015)—using TensorFlow as a backend (Abadi et al., 2015). Our software is a complete, end-to-  
243 end pipeline (Figure 1) with a custom GUI for creating annotated training data with active learning similar to  
244 Pereira et al. (2019; Appendix 2), as well as a flexible pipeline for data augmentation (Jung 2018; Appendix 2;  
245 shown in Figure 2), model training and evaluation (Figure 2; Appendix 1), and running inference on new data.  
246 We designed our high-level programming interface using the same guidelines from Keras (Chollet et al., 2015)  
247 to allow the user to go from idea to result as quickly as possible, and we organized our software into a Python  
248 module called *DeepPoseKit*. The code, documentation, and examples for our entire software package are freely  
249 available at <https://github.com/jgraving/deepposekit> under a permissive open-source license.

## 250 **Our pose estimation models**

251 To achieve the goal of “fast animal pose estimation” introduced by Pereira et al. (2019), while maintaining  
252 the robust predictive power of models like DeepLabCut (Mathis et al., 2018), we implemented two fast pose  
253 estimation models that extend the state-of-the-art model for individual pose estimation introduced by Newell  
254 et al. (2016) and the current state-of-the-art for convolutional regression from Jégou et al. (2017). Our model  
255 implementations use fewer parameters than both the DeepLabCut model (Mathis et al., 2018) and LEAP model  
256 (Pereira et al., 2019) while simultaneously removing many of the limitations of these architectures.

257 In order to limit overparameterization while minimizing performance loss, we designed our models to allow  
258 for multi-scale inference (Appendix 3 Box 1) while optimizing our model hyperparameters for efficiency. Our  
259 first model is a novel implementation of *FC-DenseNet* from Jégou et al. (2017; Appendix 3 Box 1) arranged in a  
260 stacked configuration similar to Newell et al. (2016; Appendix 4). We call this new model Stacked DenseNet, and  
261 to the best of our knowledge, this is the first implementation of this model architecture in the literature—for  
262 pose estimation or otherwise. Further details for this model are available in Appendix 7. Our second model  
263 is a modified version of the *Stacked Hourglass* model from Newell et al. (2016; Appendix 4) with hyperparam-  
264 eters that allow for changing the number of filters in each convolutional block to constrain the number of  
265 parameters—rather than using 256 filters for all layers as described in Newell et al. (2016).

## 266 **Subpixel keypoint prediction on the GPU**

267 In addition to implementing our efficient pose estimation models, we developed a new method to process model  
268 outputs to allow for faster, more accurate predictions. When using a fully-convolutional posture estimation  
269 model, the confidence maps produced by the model must be converted into coordinate values for the predictions  
270 to be useful, and there are typically two choices for making this conversion. The first is to move the confidence  
271 maps out of GPU memory and post-process them on the CPU. This solution allows for easy, flexible, and accurate  
272 calculation of the coordinates with subpixel precision (Insafutdinov et al., 2016; Mathis et al., 2018). However,  
273 CPU processing is not ideal because moving large arrays of data between the GPU and CPU can be costly, and  
274 computation on the CPU is generally slower. The other option is to directly process the confidence maps on the  
275 GPU and then move the coordinate values from the GPU to the CPU. This approach usually means converting  
276 confidence maps to integer coordinates based on the row and column index of the global maximum for each  
277 confidence map (Pereira et al., 2019). However, this means that, to achieve a precise estimation, the confidence  
278 maps should be predicted at the full resolution of the input image, or larger, which slows down inference speed.

279 As an alternative to these two strategies, we introduce a new GPU-based convolutional layer that we call  
280 *subpixel maxima*. This layer uses the fast, efficient, image registration algorithm introduced by Guizar-Sicairos  
281 et al. (2008) to translationally align a centered two-dimensional Gaussian filter to each confidence map via  
282 Fourier-based convolution. The translational shift between the filter and each confidence map allows us to  
283 calculate the coordinates of the global maxima with high speed and subpixel precision. This technique allows  
284 for accurate predictions of keypoint locations even if the model’s confidence maps are dramatically smaller than  
285 the resolution of the input image.

## 286 **Learning multi-scale relationships between keypoints**

287 Minimizing extreme prediction errors is important to prevent downstream effects on any further behavioral  
288 analysis (Seethapathi et al., 2019)—especially in the case of analyses based on time-frequency transforms like  
289 those from Beraman et al. (2014a, 2016); Klibaite et al. (2017); Todd et al. (2017); Klibaite and Shaevitz (2019)  
290 and Pereira et al. (2019) where high magnitude errors can cause inaccurate behavioral classifications. One way  
291 to minimize extreme errors when estimating posture is to incorporate multiple spatial scales when making  
292 predictions. Our pose estimation models are implicitly capable of using information from multiple scales (see

293 Appendix 3 Box 1), but there is no explicit signal that optimizes the model to take advantage of this information  
294 when making predictions.

295 To remedy this, we modified the model's output to predict, in addition to the keypoint locations, a hierarchical  
296 graph of edges describing the multi-scale geometry between keypoints—similar to the part affinity fields  
297 described by *Cao et al. (2017)*. This was achieved by adding an extra set of confidence maps to the output where  
298 edges in the postural graph are represented by Gaussian-blurred lines the same width as the Gaussian peaks in  
299 the keypoint confidence maps. Our posture graph output then consists of four levels: (1) a set of confidence  
300 maps for the smallest limb segments in the graph (e.g., foot to ankle, knee to hip, etc.; Figure 2), (2) a set of  
301 confidence maps for individual limbs (e.g., left leg, right arm, etc.; Figure 3), (3) a map with the entire postural  
302 graph, and (4) a fully-integrated map that incorporates the entire posture graph and confidence peaks for all of  
303 the joint locations (Figure 2). Each level of the hierarchical graph is built from lower levels in the output, which  
304 forces the model to learn correlated features across multiple scales when making predictions.

## 305 Experiments and model comparisons

306 We ran three main experiments to test and optimize our approach. First, we compared our new subpixel  
307 maxima layer to an integer-based global maxima with downsampled outputs ranging from  $1\times$  to  $\frac{1}{16}\times$  the input  
308 resolution using our Stacked DenseNet model. Next, we tested if training our Stacked DenseNet model to  
309 predict the multi-scale geometry of the posture graph improves accuracy. Finally, we compared our model  
310 implementations of Stacked Hourglass and Stacked DenseNet to the models from *Pereira et al. (2019)* (LEAP)  
311 and *Mathis et al. (2018)* (DeepLabCut), which we also implemented in our framework (see Appendix 7 for  
312 details). We assessed both the inference speed and prediction accuracy of each model as well as training time  
313 and generalization ability. When comparing these models we incorporated the relevant improvements from  
314 our experiments—including subpixel maxima and predicting multi-scale geometry between keypoints—unless  
315 otherwise noted (see Appendix 7).

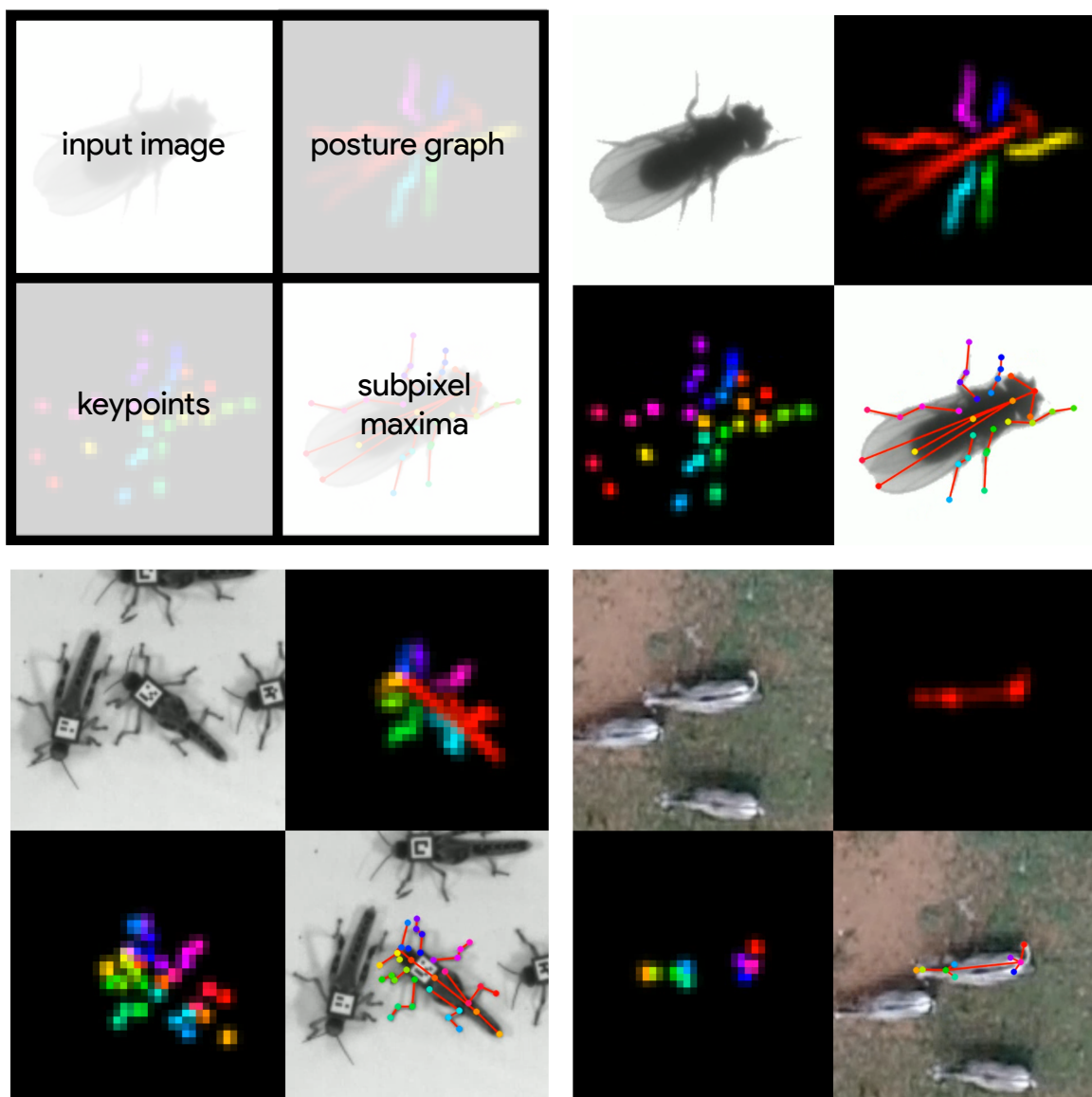
316 While we do make comparisons to the DeepLabCut model (*Mathis et al., 2018*) we do not use the same  
317 training routine as *Mathis et al. (2018)* and *Nath et al. (2019)*, who use binary cross-entropy loss for optimizing  
318 the confidence maps in addition to the location refinement maps described by *Insafutdinov et al. (2016)*. We  
319 made this modification in order to hold the training routine for each model constant, while only varying the  
320 model itself. However, we find that these differences between training routines effectively have no impact on  
321 performance when the models are trained using the same dataset and data augmentations (Appendix 7 Figure  
322 1). We also provide qualitative comparisons to demonstrate that, when trained with our DeepPoseKit framework,  
323 our implementation of the DeepLabCut model (*Mathis et al., 2018*) appears to produce fewer prediction errors  
324 than the original implementation from *Mathis et al. (2018)*; *Nath et al. (2019)* when applied to a novel video  
325 (Appendix 7 Figure 1-Figure supplement 1; Appendix 7 Figure 1-video 1).

## 326 Datasets

327 We performed experiments using the vinegar or "fruit" fly (*Drosophila melanogaster*) dataset (Figure 3-video 1)  
328 provided by *Pereira et al. (2019)*, and to demonstrate the versatility of our methods we also compared model  
329 performance across two previously unpublished posture data sets from groups of desert locusts (*Schistocerca*  
330 *gregaria*) recorded in a laboratory setting (Figure 3-video 2), and herds of Grévy's zebras (*Equus grevyi*) recorded  
331 in the wild (Figure 3-video 3). The locust and zebra datasets are particularly challenging for pose estimation as  
332 they feature multiple interacting individuals—with focal individuals centered in the frame—and the latter with  
333 highly-variable environments and light conditions. These datasets are freely-available from [https://github.com/  
334 jgraving/deepposekit-data](https://github.com/jgraving/deepposekit-data) (*Graving et al., 2019*).

335 Our locust dataset consisted of a group of 100 locusts in a circular plastic arena 1-m in diameter. The  
336 locust group was recorded from above using a high-resolution camera (Basler ace acA2040-90umNIR) and  
337 video recording system (Motif, loopbio GmbH). Our zebra dataset consisted of variably sized groups in the  
338 wild recorded from above using a commercially-available quadcopter drone (DJI Phantom 4 Pro). Locusts were  
339 localized and tracked using 2-D barcode markers (*Graving, 2017*) attached to the thorax with cyanoacrylate glue,  
340 and any missing localizations (<0.02% of the total dataset) between successful barcode reads were interpolated  
341 with linear interpolation. Individual zebra were localized using custom deep-learning software based on Faster  
342 R-CNN (*Ren et al., 2015*) for predicting bounding boxes. The positions of each zebra were then tracked across  
343 frames using a linear assignment algorithm (*Munkres, 1957*) and data were manually verified for accuracy.

344 After positional tracking, the videos were then cropped using the egocentric coordinates of each individual  
345 and saved as separate videos—one for each individual. The images used for each training set were randomly  
346 selected using the k-means sampling procedure (with  $k=10$ ) described by *Pereira et al. (2019)* (Appendix 2).  
347 After annotating the images with keypoints, we rotationally and translationally aligned the images and keypoints  
348 using the central body axis of the animal in each labeled image. This step allowed us to more easily perform



**Figure 3.** A visualization of the datasets we used to evaluate our methods (Table 1). For each dataset, confidence maps for the keypoints (bottom-left) and posture graph (top-right) are illustrated using different colors for each map. These outputs are from our Stacked DenseNet model at  $\frac{1}{4}\times$  resolution.

**Figure 3–video 1.** A video of a behaving fly from *Pereira et al. (2019)* with pose estimation outputs visualized.

**Figure 3–video 2.** A video of a behaving locust with pose estimation outputs visualized.

**Figure 3–video 3.** A video of a behaving Grévy's zebra with pose estimation outputs visualized.



**Table 1.** Datasets used for model comparisons.

Name	Species	Resolution	# Images	# Keypoints	Individuals	Source
Vinegar fly	<i>Drosophila melanogaster</i>	192×192	1500	32	Single	<i>Pereira et al. (2019)</i>
Desert locust	<i>Schistocerca gregaria</i>	160×160	800	35	Multiple	This paper
Grévy's zebra	<i>Equus grevyi</i>	160×160	900	9	Multiple	This paper

349 data augmentations (see "Model training") that allow the model to make accurate predictions regardless of the  
350 animal's body size and orientation (see Appendix 5). However, this preprocessing step is not a strict requirement  
351 for training, and there is no requirement for this preprocessing step when making predictions on new unlabeled  
352 data, such as with the methods described by *Pereira et al. (2019)* (Appendix 5). Before training each model we  
353 split each annotated dataset into randomly selected training and validation sets with 90% training examples  
354 and 10% validation examples, unless otherwise noted. The details for each dataset are described in Table 1.

### 355 Model training

356 For each experiment, we set our model hyperparameters to the same configuration for our Stacked DenseNet  
357 and Stacked Hourglass models. Both models were trained with  $\frac{1}{4}\times$  resolution outputs and a stack of two  
358 networks with two outputs where loss was applied (see Figure 2). Although our model hyperparameters could  
359 be infinitely adjusted to trade off between speed and accuracy, we compared only one configuration for each of  
360 our model implementations. These results are not meant to be an exhaustive search of model configurations as  
361 the best configuration will depend on the application. The details of the hyperparameters we used for each  
362 model are described in Appendix 7.

363 To make our posture estimation tasks closer to realistic conditions and properly demonstrate the robustness  
364 of our methods to rotation, translation, scale, and noise, we applied various augmentations to each data set  
365 during training (Figure 2). All models were trained using data augmentations that included random flipping,  
366 or mirroring, along both the horizontal and vertical image axes with each axis being independently flipped by  
367 drawing from a Bernoulli distribution (with  $p = 0.5$ ), random rotations around the center of the image drawn  
368 from a uniform distribution in the range  $[-180^\circ, +180^\circ]$ , random scaling drawn from a uniform distribution in  
369 the range [90%, 110%] for flies and locusts and [75%, 125%] for zebras (to account for greater size variation  
370 in the data set), and random translations along the horizontal and vertical axis independently drawn from a  
371 uniform distribution with the range  $[-5\%, +5\%]$ —where percentages are relative to the original image size. After  
372 performing these spatial augmentations we also applied a variety of noise augmentations that included additive  
373 noise—i.e., adding or subtracting randomly-selected values to pixels; dropout—i.e., setting individual pixels  
374 or groups of pixels to a randomly-selected value; blurring or sharpening—i.e., changing the composition of  
375 spatial frequencies; and contrast ratio augmentations—i.e., changing the ratio between the highest value and  
376 lowest value in the image. These augmentations help to further ensure robustness to shifts in lighting, noise,  
377 and occlusions. See Appendix 2 for further discussion on data augmentation.

378 We trained our models (Figure 2) using mean squared error loss optimized using the ADAM optimizer  
379 (*Kingma and Ba, 2014*) with a learning rate of  $1 \times 10^{-3}$  and a batch size of 16. We lowered the learning rate by a  
380 factor of 5 each time the validation loss did not improve by more than  $1 \times 10^{-3}$  for 10 epochs. We considered  
381 models to be converged when the validation loss stopped improving for 50 epochs, and we calculated validation  
382 error as the Euclidean distance between predicted and ground-truth image coordinates for only the best  
383 performing version of the model, which we evaluated at the end of each epoch during optimization. We  
384 performed this procedure five times for each experiment and randomly selected a new validation set for each  
385 replicate.

### 386 Model evaluation

387 Machine learning models are typically evaluated for their ability to generalize to new data, known as *predictive*  
388 *performance*, using a held-out *test set*—a subsample of annotated data that is not used for training or validation.  
389 However, when fitting and evaluating a model on a small dataset, using an adequately-sized validation and test  
390 set can lead to erroneous conclusions about the predictive performance of the model if the training set is too  
391 small (*Kuhn and Johnson, 2013*). Therefore, to maximize the size of the training set, we elected to use only a  
392 validation set for model evaluation.

393 Generally a test set is used to avoid biased performance measures caused by overfitting the model hyper-  
394 parameters to the validation set. However, we did not adjust our model architecture to achieve better  
395 performance on our validation set—only to achieve fast inference speeds. While we did use validation error to  
396 decide when to lower the learning rate during training and when to stop training, lowering the learning rate

397 in this way should have no effect on the generalization ability of the model, and because we heavily augment  
398 our training set during optimization—forcing the model to learn a much larger image distribution than what is  
399 included in the training and validation sets—overfitting to the validation set is unlikely. We also demonstrate  
400 the generality of our results for each experiment by randomly selecting a new validation set with each replicate.  
401 All of these factors make the Euclidean error for the unaugmented validation set a reasonable measure of the  
402 predictive performance for each model.

403 The inference speed for each model was assessed by running predictions on 100,000 randomly generated  
404 images with a batch size of 1 for real-time speeds and a batch size of 100 for offline speeds, unless otherwise  
405 noted. Our hardware consisted of a Dell Precision Tower 7910 workstation (Dell, Inc.) running Ubuntu Linux  
406 v18.04 with 2x Intel Xeon E5-2623 v3 CPUs (8 cores, 16 threads at 3.00GHz), 64GB of RAM, a Quadro P6000 GPU  
407 and a Titan Xp GPU (NVIDIA Corporation). We used both GPUs for training models and evaluating predictive  
408 performance, and we only used the faster Titan Xp GPU for benchmarking inference speeds and training time.  
409 While the hardware we used for development and testing is on the high-end of the performance spectrum, there  
410 is no requirement for this level of performance, and our software can easily be run on lower-end hardware. We  
411 evaluated inference speeds on multiple consumer-grade desktop computers and found similar performance  
412 ( $\pm 10\%$ ) when using the same GPU.

### 413 **Assessing prediction accuracy with Bayesian inference**

414 To more rigorously assess performance differences between models, we parameterized the Euclidean error  
415 distribution for each experiment by fitting a Bayesian linear model with a Gamma-distributed likelihood function.  
416 This model takes the form:

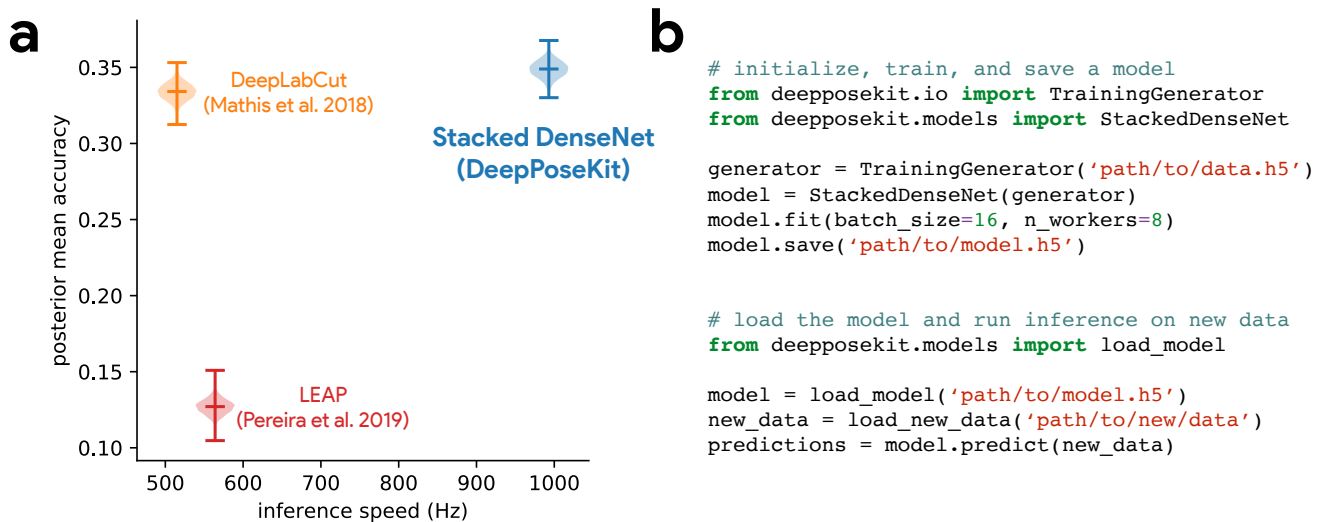
$$\begin{aligned} p(y|X, \theta_\mu, \theta_\phi) &\sim \text{Gamma}(\alpha, \beta) \\ \alpha &= \mu^2 \phi^{-1} \\ \beta &= \mu \phi^{-1} \\ \mu &= h(X\theta_\mu) \\ \phi &= h(X\theta_\phi) \end{aligned}$$

417 where  $X$  is the design matrix composed of binary indicator variables for each pose estimation model,  $\theta_\mu$   
418 and  $\theta_\phi$  are vectors of intercepts,  $h(\cdot)$  is the softplus function (*Dugas et al., 2001*)—or  $h(x) = \log(1 + e^x)$ —used to  
419 enforce positivity of  $\mu$  and  $\phi$ , and  $y$  is the Euclidean error of the pose estimation model. Parameterizing our  
420 error distributions in this way allows us to calculate the posterior distributions for the mean  $E[y] = \alpha\beta^{-1} \equiv \mu$  and  
421 variance  $\text{Var}[y] = \alpha\beta^{-2} \equiv \phi$ . This parameterization then provides us with a statistically rigorous way to assess  
422 differences in model accuracy in terms of both central tendency and spread—accounting for both epistemic  
423 uncertainty (unknown unknowns; e.g., parameter uncertainty) and aleatoric uncertainty (known unknowns; e.g.,  
424 data variance). Details of how we fitted these models can be found in Appendix 6.

### 425 **Subpixel prediction allows for fast and accurate inference**

426 We compared the accuracy of our subpixel maxima layer to an integer-based maxima layer using the fly dataset.  
427 We found significant accuracy improvements across every downsampling configuration (Appendix Figure 5).  
428 Even with confidence maps at  $\frac{1}{8}$  the resolution of the original image, error did not drastically increase compared  
429 to full-resolution predictions. Making predictions for confidence maps at such a downsampled resolution allows  
430 us to achieve very fast inference >1000 Hz while maintaining high accuracy. We also provide speed comparisons  
431 with the other models we tested and find that our Stacked DenseNet model is faster than the DeepLabCut  
432 model (*Mathis et al., 2018*) for both offline (batch size = 100) and real-time speeds (batch size = 1). While we  
433 find that our Stacked DenseNet model is faster than the LEAP model (*Pereira et al., 2019*) for offline processing  
434 (batch size = 100), the LEAP model (*Pereira et al., 2019*) is significantly faster for real-time processing (batch  
435 size = 1). Our Stacked Hourglass model (*Newell et al., 2016*) is about the same or slightly faster than Stacked  
436 DenseNet for offline speeds (batch size = 100), but is much slower for real-time processing (batch size = 1).

437 Achieving fast pose estimation using CNNs typically relies on massively parallel processing on the GPU with  
438 large batches of data or requires downsampling the images to increase speed, which increases error (*Mathis and*  
439 *Warren, 2018*). These factors make fast and accurate real-time inference challenging to accomplish. Our Stacked  
440 DenseNet model, with a batch size of one, can run inference at  $\sim 30$ – $110$  Hz—depending on resolution (Appendix  
441 Figure 5b). These speeds are faster than the DeepLabCut model (*Mathis et al., 2018*) and could be further  
442 improved by downsampling the input image resolution or reconfiguring the model with fewer parameters. This



**Figure 4.** Our Stacked DenseNet model estimates posture at approximately 2x—or greater—the speed of the LEAP model (Pereira et al., 2019) and the DeepLabCut model (Mathis et al., 2018) while also achieving similar accuracy to the DeepLabCut model (Mathis et al., 2018)—shown here as mean accuracy  $(1 + \text{Euclidean error})^{-1}$  for our most challenging dataset of multiple interacting Grévy's zebras (*E. grevyi*) recorded in the wild (a). See Figure 4 supplement 1 for further details. Our software interface is designed to be straightforward but flexible. We include many options for expert users to customize model training with sensible default settings to make pose estimation as easy as possible for beginners. For example, training a model and running inference on new data requires writing only a few lines of code and specifying some basic settings (b).

**Figure 4—Figure supplement 1.** Euclidean error distributions for each model across our three datasets. Letter-value plots (left) show the raw error distributions for each model. Violinplots of the posterior distributions for the mean and variance (right) show statistical differences between the error distributions. Overall the LEAP model (Pereira et al., 2019) was the worst performer on every dataset in terms of both mean and variance. Our Stacked DenseNet model was the best performer for the fly dataset, while our Stacked DenseNet model and the DeepLabCut model (Mathis et al., 2018) both performed equally well on the locust and zebra datasets. The posteriors for the DeepLabCut model (Mathis et al., 2018) and our Stacked DenseNet model are highly overlapping for these datasets, which suggests they are not statistically discernible from one another. Our Stacked Hourglass model (Newell et al., 2016) performed equally to the DeepLabCut model (Mathis et al., 2018) and our Stacked DenseNet model for the locust dataset but performed slightly worse for the fly and zebra datasets.

443 allows our methods to be flexibly used for real-time closed-loop behavioral experiments with prediction errors  
 444 similar to current state-of-the-art methods.

445 **Predicting multi-scale geometry improves accuracy and reduces extreme errors**  
 446 We find that training our Stacked DenseNet model to predict a hierarchical posture graph reduces keypoint  
 447 prediction error (Appendix Figure 6), and because the feature maps for the posture graph can be removed  
 448 from the final output during inference, this effectively improves prediction accuracy for free. Both the mean and  
 449 variance of the error distributions were lower when predicting the posture graph, which suggests that learning  
 450 multi-scale geometry both decreases error on average and helps to reduce extreme prediction errors. The  
 451 overall effect size for this decrease in error is fairly small (<1 pixel average reduction in error), but based on the  
 452 results from the zebra dataset, this modification more dramatically improves performance for datasets with  
 453 higher-variance images and sparse posture graphs. Predicting the posture graph may be especially useful for  
 454 animals with long slender appendages such as insect legs and antennae where prediction errors are likely to  
 455 occur due to occlusions and natural variation in the movement of these body parts. These results also suggest  
 456 that annotating multiple keypoints to incorporate an explicit signal for multi-scale information may help improve  
 457 prediction accuracy for a specific body part of interest.

458 **Stacked DenseNet is fast and robust**  
 459 Finally, we benchmarked our new model implementations against the models (Pereira et al., 2019) and Mathis  
 460 et al. (2018). We find that our Stacked DenseNet model outperforms both the LEAP model (Pereira et al.,  
 461 2019) and the DeepLabCut model (Mathis et al., 2018) in terms of speed while also achieving much higher  
 462 accuracy than the LEAP model (Pereira et al., 2019) with similar accuracy to the DeepLabCut model (Mathis  
 463 et al. 2018; Figure 4a). We found that both the Stacked Hourglass and Stacked DenseNet models outperformed  
 464 the LEAP model (Pereira et al., 2019). Notably our Stacked DenseNet model achieved approximately 2x faster  
 465 inference speeds with 3x higher mean accuracy. Not only were our models' average prediction error significantly  
 466 improved, but also, importantly, the variance was lower—indicating that our models produced fewer extreme  
 467 prediction errors. At  $\frac{1}{4}$  resolution, our Stacked DenseNet model consistently achieved prediction accuracy

468 nearly identical to the DeepLabCut model (*Mathis et al., 2018*) while running inference at nearly 2× the speed  
469 and using only ~5% of the parameters— ~26 million vs. ~1.5 million. Detailed results of our model comparisons  
470 are shown in Figure 4 supplement 1. While our Stacked DenseNet model is already fast, inference speed could  
471 be further improved by using a  $\frac{1}{8}$ × output without much increase in error (Appendix Figure 5) or by further  
472 adjusting the hyperparameters to constrain the size of the model. Our Stacked Hourglass implementation  
473 followed closely behind the performance of our Stacked DenseNet model and the DeepLabCut model (*Mathis*  
474 *et al., 2018*) but consistently performed more poorly than our Stacked DenseNet model in terms of prediction  
475 accuracy, so we excluded this model from further analysis. We were also able to reproduce the results reported  
476 by *Pereira et al. (2019)* that the LEAP model and the Stacked Hourglass model (*Newell et al., 2016*) have similar  
477 average prediction error for the fly dataset. However, we also find that the LEAP model (*Pereira et al., 2019*) has  
478 much higher variance, which suggests it is more prone to extreme prediction errors—a problem for further data  
479 analysis.

## 480 **Stacked DenseNet trains quickly and requires few training examples**

481 To further compare models, we used our zebra dataset to assess the training time needed for our Stacked  
482 DenseNet model, the DeepLabCut model (*Mathis et al., 2018*), and the LEAP model (*Pereira et al., 2019*) to  
483 reach convergence as well as the amount of training data needed for each model to generalize to new data  
484 from outside the training set. We find that our Stacked DenseNet model, the DeepLabCut model (*Mathis et al.,*  
485 *2018*), and the LEAP model (*Pereira et al., 2019*) all fully converge in just a few hours and reach reasonably high  
486 accuracy after only an hour of training (Appendix Figure 7). However, it appears that our Stacked DenseNet  
487 model tends to converge to a good minimum faster than both the DeepLabCut model (*Mathis et al., 2018*)  
488 and the LEAP model (*Pereira et al., 2019*). We also show that our Stacked DenseNet model achieves good  
489 generalization with few training examples and without the use of transfer learning (Appendix Figure 8). These  
490 results demonstrate that, when combined with data augmentation, as few as five training examples can be  
491 used as an initial training set for labelling keypoints with active learning (Figure 1). Additionally, because our  
492 analysis shows that generalization to new data plateaus after approximately 100 labeled training examples,  
493 it appears that 100 training examples is a reasonable minimum size for a training set—although the exact  
494 number will likely change depending on the variance of the image data being annotated. To further examine the  
495 effect of transfer learning on model generalization, we compared performance between the DeepLabCut model  
496 (*Mathis et al., 2018*) initialized with weights pretrained on the ImageNet database (*Deng et al., 2009*) vs. the  
497 same model with randomly-initialized weights (Appendix Figure 8). As postulated by *Mathis et al. (2018)*, we find  
498 that transfer learning does provide some benefit to the DeepLabCut model's ability to generalize. However, the  
499 effect is relatively small with a mean reduction in Euclidean error of <0.5 pixel. Together these results indicate  
500 that transfer learning is helpful, but not required, for deep learning models to achieve good generalization with  
501 limited training data.

## 502 **Discussion**

503 Here we have presented a new software toolkit, called DeepPoseKit, for estimating animal posture using deep  
504 learning models. We built on the state-of-the-art for individual pose estimation using convolutional neural  
505 networks to achieve fast inference without reducing accuracy or generalization ability. Our new pose estimation  
506 model, called Stacked DenseNet, offers considerable improvements (Figure 4a; Figure supplement 1) over the  
507 models from *Mathis et al. (2018)* (DeepLabCut) and *Pereira et al. (2019)* (LEAP), and our software framework  
508 also provides a simplified interface (Figure 4b) for using these advanced tools to measure animal behavior  
509 and locomotion. We tested our methods across a range of datasets from controlled laboratory environments  
510 with single individuals to challenging field situations with multiple interacting individuals and variable lighting  
511 conditions. We found that our methods perform well for all of these situations and require few training examples  
512 to achieve good predictive performance on new data—without the use of transfer learning. We ran experiments  
513 to optimize our approach and discovered that some straightforward modifications can greatly improve speed  
514 and accuracy. Additionally, we demonstrated that these modifications improve not just the average error  
515 but also help to reduce extreme prediction errors—a key determinant for the reliability of subsequent statistical  
516 analysis.

517 While our results offer a good-faith comparison of the available methods for animal pose estimation,  
518 there is inherent uncertainty that we have attempted to account for but may still bias our conclusions. For  
519 example, deep learning models are trained using stochastic optimization algorithms that give different results  
520 with each replicate, and the Bayesian statistical methods we use for comparison are explicitly probabilistic in  
521 nature. There is also great variability across hardware and software configurations when using these models  
522 in practice (*Mathis and Warren, 2018*), so performance may change across experimental setups and datasets.  
523 Additionally, we demonstrated that some models may perform better than others for specific applications

524 (Figure 4 supplement 1), and to account for this, our toolkit offers researchers the ability to choose the model  
525 that best suits their requirements—including the LEAP model (*Pereira et al., 2019*) and the DeepLabCut model  
526 (*Mathis et al., 2018*).

527 We highlighted important considerations when using CNNs for pose estimation and reviewed the progress of  
528 fully-convolutional regression models from the literature. The latest advancements for these models have been  
529 driven mostly by a strategy of adding more connections between layers to increase performance and efficiency  
530 (e.g., *Jégou et al. 2017*). Future progress for this class of models may require better loss functions (*Goodfellow*  
531 *et al., 2014; Johnson et al., 2016a; Chen et al., 2017*) that more explicitly model the spatial dependencies within  
532 a scene, models that incorporate the temporal structure of the data (*Seethapathi et al., 2019*), and more  
533 mathematically-principled approaches (e.g., *Weigert et al. 2018; Roy et al. 2018*) such as the application of  
534 formal probabilistic concepts (*Kendall and Gal, 2017*) and Bayesian inference at scale (*Tran et al., 2018*).

535 Measuring behavior is a critical factor for many studies in neuroscience (*Krakauer et al., 2017*). Understand-  
536 ing the connections between brain activity and behavioral output requires detailed and objective descriptions  
537 of body posture that match the richness and resolution neural measurement technologies have provided for  
538 years (*Anderson and Perona, 2014; Berman, 2018; Brown and De Bivort, 2018*), which our methods and other  
539 deep-learning-based tools provide (*Mathis et al., 2018; Pereira et al., 2019*). We have also demonstrated the  
540 possibility that our toolkit could be used for real-time inference, which allows for closed-loop experiments  
541 where sensory stimuli or optogenetic stimulation are controlled in response to behavioral measurements  
542 (e.g., *Bath et al. 2014; Stowers et al. 2017*). Using real-time measurements in conjunction with optogenetics or  
543 thermogenetics may be key to disentangling the causal structure of motor output from the brain—especially  
544 given that recent work has shown an animal's response to optogenetic stimulation can differ depending on the  
545 behavior it is currently performing (*Cande et al., 2018*). Real-time behavioral quantification is also particularly  
546 important as closed-loop virtual reality is quickly becoming an indispensable tool for studying sensorimotor  
547 relationships in individuals and collectives (*Stowers et al., 2017*).

548 Quantifying individual movement is essential for revealing the genetic (*Kain et al., 2012; Brown et al., 2013;*  
549 *Ayroles et al., 2015*) and environmental (*Bierbach et al., 2017; Akhund-Zade et al., 2019; Versace et al., 2019*)  
550 underpinnings of phenotypic variation in behavior—as well as the phylogeny of behavior (e.g., *Berman et al.*  
551 *2014b*). Measuring individual behavioral phenotypes requires tools that are robust, scaleable, and easy-to-use,  
552 and our approach offers the ability to quickly and accurately quantify the behavior of many individuals in  
553 great detail. When combined with tools for genetic manipulations (*Ran et al., 2013; Doudna and Charpentier,*  
554 *2014*), high-throughput behavioral experiments (*Alisch et al., 2018; Javer et al., 2018; Werkhoven et al., 2019*),  
555 and behavioral analysis (e.g., *Berman et al. 2014a; Wiltschko et al. 2015*), our methods could help to provide  
556 the data resolution and statistical power needed for dissecting the complex relationships between genes,  
557 environment, and behavioral variation.

558 When used together with other tools for localization and tracking (e.g., *Pérez-Escudero et al. 2014; Crall*  
559 *et al. 2015; Graving 2017; Romero-Ferrero et al. 2019; Wild et al. 2018; Boenisch et al. 2018*), our methods  
560 are capable of reliably measuring posture for multiple interacting individuals. The importance of measuring  
561 detailed representations of individual behavior when studying animal collectives has been well established  
562 (*Strandburg-Peshkin et al., 2013; Rosenthal et al., 2015; Strandburg-Peshkin et al., 2015, 2017*). Estimating  
563 body posture is an essential first step for unraveling the sensory networks that drive group coordination, such  
564 as vision-based networks measured via raycasting (*Strandburg-Peshkin et al., 2013; Rosenthal et al., 2015*).  
565 Additionally, using body pose estimation in combination with computational models of behavior (e.g., *Costa*  
566 *et al. 2019, Wiltschko et al. 2015*) and unsupervised behavioral classification methods (e.g., *Berman et al. 2014a,*  
567 *Pereira et al. 2019*) may allow for further dissection of how information flows through groups by revealing the  
568 networks of behavioral contagion across multiple timescales and sensory modalities. While we have provided a  
569 straightforward solution for applying existing pose estimation methods to measure collective behavior, there  
570 still remain many challenging scenarios where these methods would fail. For example, tracking posture in a  
571 densely-packed bee hive or school of fish would require novel solutions to deal with the 3-D nature of individual  
572 movement, which includes maintaining individual identities and dealing with the resulting occlusions that go  
573 along with imaging these types of biological systems.

574 When combined with unmanned aerial vehicles (UAVs; *Schiffman 2014*) or other field-based imaging (*Fran-*  
575 *cisco et al., 2019*), applying these methods to the study of individuals and groups in the wild can provide  
576 high-resolution behavioral data that goes beyond the capabilities of current GPS and accelerometry-based  
577 technologies (*Nagy et al., 2010, 2013; Kays et al., 2015; Strandburg-Peshkin et al., 2015, 2017; Flack et al.,*  
578 *2018*)—especially for species that impractical to study with tags or collars. Additionally, by applying these meth-  
579 ods in conjunction with 3-D habitat reconstruction—using techniques such as photogrammetry—field-based  
580 studies can begin to integrate fine-scale behavioral measurements with the full 3-D environment in which the  
581 behavior evolved (e.g., *Strandburg-Peshkin et al. 2017; Francisco et al. 2019*). Future advances will likely allow



582 for the calibration and synchronizaton of imaging devices across multiple UAVs. This would make it possible to  
583 measure the full 3-D posture of wild animals in scenarios where fixed camera systems (e.g. *Nath et al. 2019*)  
584 would not be tractable, such as during migratory or predation events. When combined, these technologies  
585 could allow researchers to address questions about the behavioral ecology of animals that were previously  
586 impossible to answer.

587 Computer vision algorithms for measuring behavior at the scale of posture have rapidly advanced in a  
588 very short time; nevertheless, the task of pose estimation is far from solved. There are hard limitations to  
589 this current generation of pose estimation methods that are primarily related to the requirement for human  
590 annotations and user-defined keypoints—both in terms of the number of keypoints, the specific body parts  
591 being tracked, and the inherent difficulty of incorporating temporal information into the annotation and training  
592 procedure. Often the body parts chosen for annotation are an obvious fit for the experimental design and have  
593 reliably-visible reference points on the animal's body that make them easy to annotate. However, in many cases  
594 the required number and type of body parts needed for data analysis may not so obvious—such as in the case  
595 of unsupervised behavior classification methods (*Berman et al., 2014a; Pereira et al., 2019*). Additionally, the  
596 reference points for labeling images with keypoints can be hard to define and consistently annotate across  
597 images, which is often the case for soft or flexible-bodied animals like worms and fish. Moreover, due to the  
598 laborious nature of annotating keypoints, the current generation of methods also rarely takes into account  
599 the natural temporal structure of the data, instead treating each video frame as a statistically independent  
600 event, which can lead to extreme prediction errors (reviewed by *Seethapathi et al. 2019*). Extending these  
601 methods to track the full three-dimensional posture of animals also typically requires the use of multiple  
602 synchronized cameras (*Nath et al., 2019; Günel et al., 2019*), which increases the cost and complexity of  
603 creating an experimental setup, as well as the manual labor required for annotating a training set, which must  
604 include labeled data from every camera view.

605 These limitations make it clear that fundamentally-different methods may be required to move the field  
606 forward. Future pose estimation methods will likely replace human annotations with fully-articulated volumetric  
607 3-D models of the animal's body (*Zuffi et al., 2017*), and the 3-D scene will be learned in an unsupervised way  
608 (e.g., *Jaques et al. 2019*), where the shape, position, and posture of the animal's body, the camera position and  
609 lens parameters, and the background environment and lighting conditions will all be jointly learned directly from  
610 2-D images by a deep-learning model (*Valentin et al., 2019*). These *inverse graphics models* (*Kulkarni et al., 2015;*  
611 *Sabour et al., 2017; Valentin et al., 2019*) will likely take advantage of recently-developed differentiable graphics  
612 engines that allow 3-D rendering parameters to be straightforwardly controlled using computationally-efficient  
613 gradient-based optimization methods (*Valentin et al., 2019*). After optimization, the volumetric 3-D timeseries  
614 data predicted by the deep learning model could be used directly for behavioral analysis or specific keypoints or  
615 body parts could be selected for analysis post-hoc. In order to more explicitly incorporate the natural statistical  
616 properties of the data, these models will also likely rely on the use of perceptual (*Johnson et al., 2016a*) and  
617 adversarial (*Goodfellow et al., 2014*) loss functions that incorporate spatial dependencies within the scene  
618 rather than modelling each video frame as a set of statistically independent pixel distributions—as is the case  
619 with current methods when using pixel-wise mean squared error (e.g. *Pereira et al. 2019*) or cross-entropy loss  
620 (e.g. *Mathis et al. 2018*). Because there would be limited or no requirement for human-provided labels, these  
621 models could also be easily modified to incorporate the temporal structure of the data using autoregressive  
622 representations (*Van den Oord et al., 2016; Oord et al., 2016; Kumar et al., 2019*), rather than modeling the  
623 scene in each video frame as a statistically independent event. Together these advances could lead to larger,  
624 higher-resolution, more reliable behavioral datasets that could revolutionize our understanding of relationships  
625 between behavior, the brain, and the environment.

626 In conclusion, we have presented a new toolkit, called DeepPoseKit, for automatically measuring animal  
627 posture from images. We combined recent advances from the literature to create methods that are fast, robust,  
628 and widely applicable to a range of species and experimental conditions. When designing our framework we em-  
629 phasized usability across the entire software interface, which we expect will help to make these advanced tools  
630 accessible to a wider range of researchers. The fast inference and real-time capabilities of our methods should  
631 also help further reduce barriers to previously intractable questions across many scientific disciplines—including  
632 neuroscience, ethology, and behavioral ecology—both in the laboratory and the field.

### 633 **Author contributions**

634 J.M.G and I.D.C conceived the idea for the project. J.M.G. and D.C. developed the software with input from H.N.  
635 J.M.G implemented the pose estimation models and developed the subpixel maxima algorithm. J.M.G. and D.C.  
636 developed the annotation GUI, data augmentation pipeline, and wrote the documentation. J.M.G., D.C. and H.N.  
637 designed the experiments. J.M.G. and D.C. ran the experiments. B.R.C., B.K., J.M.G., and I.D.C. conceived the idea  
638 to apply posture tracking to zebras. B.R.C. and B.K. provided the annotated zebra posture data. B.K., and L.L.

639 helped with initial testing and improvement of the software interface. L.L. also made significant contributions to  
640 an earlier version of the manuscript. J.M.G. fit the linear models and made the figures. J.M.G. wrote the initial  
641 draft of the manuscript with input from H.N. and D.C., and all authors helped revise the manuscript.

## 642 Acknowledgements

643 We are indebted to Talmo Pereira et al. and A. Mathis et al. for making their software open-source and freely-  
644 available—this project would not have been possible without them. We also thank M. Mathis and A. Mathis for  
645 their comments on the manuscript. We thank François Chollet, the Keras and TensorFlow teams, and Alexander  
646 Jung for their open source contributions, which provided the core programming interface for our work. We  
647 thank A. Strandburg-Peshkin, Vivek H. Sridhar, Michael L. Smith, and Joseph B. Bak-Coleman for their helpful  
648 discussions on the project and comments on the manuscript. We also thank M.L.S. for the use of his GPU. We  
649 thank Felicitas Oehler for annotating the zebra posture data and Chiara Hirschhorn for assistance with filming  
650 the locusts and annotating the locust posture data. We thank Alex Bruttel, Christine Bauer, Jayme Weglarski,  
651 Dominique Leo, Markus Miller and loobio GmbH for providing technical support. We acknowledge the NVIDIA  
652 Corporation for their generous donations to our research. This project received funding from the European  
653 Union's Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie grant agreement  
654 No. 748549. B.R.C. acknowledges support from the University of Konstanz Zukunftskolleg's Investment Grant  
655 program. I.D.C. acknowledges support from NSF Grant IOS-1355061, Office of Naval Research Grants N00014-  
656 09-1-1074 and N00014-14-1-0635, Army Research Office Grants W911NG-11-1-0385 and W911NF14-1-0431,  
657 the Struktur-und Innovationsfonds für die Forschung of the State of Baden-Württemberg, the DFG Centre of  
658 Excellence 2117 "Centre for the Advanced Study of Collective Behaviour" (ID: 422037984), and the Max Planck  
659 Society.

## 660 Animal Ethics Statement

661 All procedures for collecting the zebra (*E. grevyi*) dataset were reviewed and approved by Ethikrat, the inde-  
662 pendent Ethics Council of the Max Planck Society. The zebra dataset was collected with the permission of  
663 Kenya's National Commission for Science, Technology and Innovation (NACOSTI/P/17/59088/15489 and NA-  
664 COSTI/P/18/59088/21567) using drones operated by B.R.C. with the permission of the Kenya Civil Aviation  
665 Authority (authorization numbers: KCAA/OPS/2117/4 Vol. 2 (80), KCAA/OPS/2117/4 Vol. 2 (81), KCAA/OPS/2117/5  
666 (86) and KCAA/OPS/2117/5 (87); RPAS Operator Certificate numbers: RPA/TP/0005 AND RPA/TP/000-0009).

## 667 References

- 668 **Abadi M**, Agarwal A, Barham P, Brevdo E, Chen Z, Citro C, Corrado GS, Davis A, Dean J, Devin M, Ghemawat S, Goodfellow I, Harp  
669 A, Irving G, Isard M, Jia Y, Jozefowicz R, Kaiser L, Kudlur M, Levenberg J, et al., TensorFlow: Large-Scale Machine Learning on  
670 Heterogeneous Systems; 2015. <https://www.tensorflow.org/>, software available from tensorflow.org.
- 671 **Akhund-Zade J**, Ho S, O'Leary C, de Bivort BL. The effect of environmental enrichment on behavioral variability depends on  
672 genotype, behavior, and type of enrichment. *bioRxiv*. 2019; p. 557181.
- 673 **Alisch T**, Crall JD, Kao AB, Zucker D, de Bivort BL. MAPLE (modular automated platform for large-scale experiments), a robot for  
674 integrated organism-handling and phenotyping. *eLife*. 2018; 7:e37166.
- 675 **Anderson DJ**, Perona P. Toward a science of computational ethology. *Neuron*. 2014; 84(1):18–31.
- 676 **Andriluka M**, Iqbal U, Insafutdinov E, Pishchulin L, Milan A, Gall J, Schiele B. PoseTrack: A benchmark for human pose estimation  
677 and tracking. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*; 2018. p. 5167–5176.
- 678 **Andriluka M**, Pishchulin L, Gehler P, Schiele B. 2D Human Pose Estimation: New Benchmark and State of the Art Analysis. In:  
679 *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*; 2014. .
- 680 **Ayinde BO**, Zurada JM. Building efficient convnets using redundant feature pruning. *arXiv preprint arXiv:180207653*. 2018; .
- 681 **Ayroles JF**, Buchanan SM, O'Leary C, Skutt-Kakaria K, Grenier JK, Clark AG, Hartl DL, de Bivort BL. Behavioral idiosyncrasy reveals  
682 genetic control of phenotypic variability. *Proceedings of the National Academy of Sciences*. 2015; 112(21):6706–6711.
- 683 **Badrinarayanan V**, Kendall A, Cipolla R. SegNet: A Deep Convolutional Encoder-Decoder Architecture for Image Segmentation.  
684 *CoRR*. 2015; abs/1511.00561. <http://arxiv.org/abs/1511.00561>.
- 685 **Bath DE**, Stowers JR, Hörmann D, Poehlmann A, Dickson BJ, Straw AD. FlyMAD: rapid thermogenetic control of neuronal activity  
686 in freely walking *Drosophila*. *Nature methods*. 2014; 11(7):756.
- 687 **Berman GJ**. Measuring behavior across scales. *BMC biology*. 2018; 16(1):23.
- 688 **Berman GJ**, Bialek W, Shaevitz JW. Predictability and hierarchy in *Drosophila* behavior. *Proceedings of the National Academy of*  
689 *Sciences*. 2016; 113(42):11943–11948.

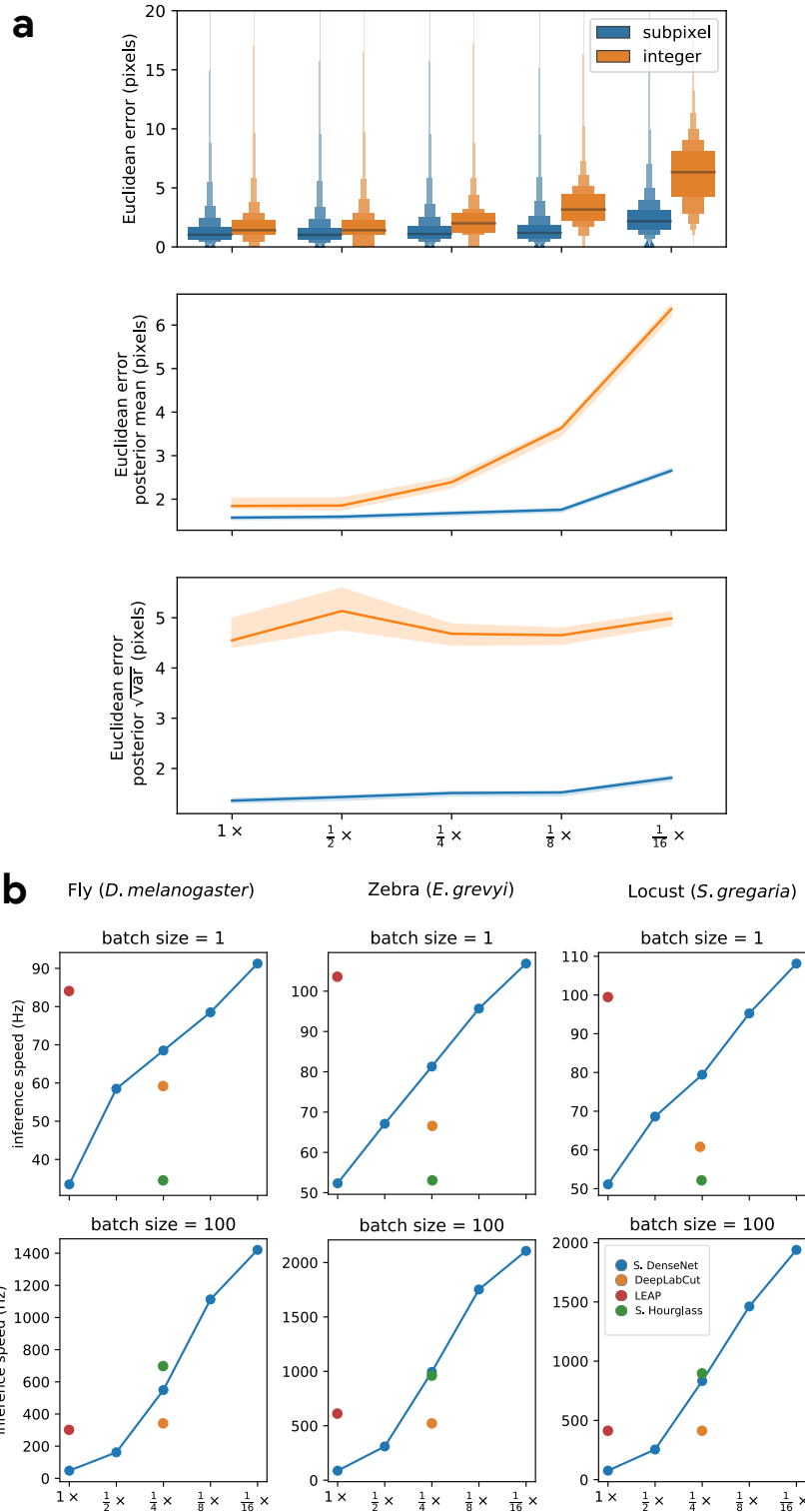
- 690 **Berman GJ**, Choi DM, Bialek W, Shaevitz JW. Mapping the stereotyped behaviour of freely moving fruit flies. *Journal of The Royal Society Interface*. 2014; 11(99):20140672.
- 691
- 692 **Berman GJ**, Choi DM, Bialek W, Shaevitz JW. Mapping the structure of drosophilid behavior. *bioRxiv*. 2014; p. 002873.
- 693 **Bierbach D**, Laskowski KL, Wolf M. Behavioural individuality in clonal fish arises despite near-identical rearing conditions. *Nature communications*. 2017; 8:15361.
- 694
- 695 **Boenisch F**, Rosemann B, Wild B, Dormagen D, Wario F, Landgraf T. Tracking All Members of a Honey Bee Colony Over Their Lifetime Using Learned Models of Correspondence. *Frontiers in Robotics and AI*. 2018; 5:35. <https://www.frontiersin.org/article/10.3389/frobt.2018.00035>, doi: 10.3389/frobt.2018.00035.
- 696
- 697
- 698 **Brown AE**, De Bivort B. Ethology as a physical science. *Nature Physics*. 2018; p. 1.
- 699 **Brown AE**, Yemini EI, Grundy LJ, Jucikas T, Schafer WR. A dictionary of behavioral motifs reveals clusters of genes affecting *Caenorhabditis elegans* locomotion. *Proceedings of the National Academy of Sciences*. 2013; 110(2):791–796.
- 700
- 701 **Cande J**, Namiki S, Qiu J, Korff W, Card GM, Shaevitz JW, Stern DL, Berman GJ. Optogenetic dissection of descending behavioral control in *Drosophila*. *Elife*. 2018; 7:e34275.
- 702
- 703 **Cao Z**, Simon T, Wei SE, Sheikh Y. Realtime multi-person 2d pose estimation using part affinity fields. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*; 2017. p. 7291–7299.
- 704
- 705 **Carpenter B**, Lee D, Brubaker MA, Riddell A, Gelman A, Goodrich B, Guo J, Hoffman M, Betancourt M, Li P. Stan: A Probabilistic Programming Language. *J Stat Softw*. 2017; .
- 706
- 707 **Cauchy A**. Méthode générale pour la résolution des systemes d'équations simultanées. *Comp Rend Sci Paris*. 1847; 25(1847):536–538.
- 708
- 709 **Chen Y**, Shen C, Wei XS, Liu L, Yang J. Adversarial poseNet: A structure-aware convolutional network for human pose estimation. In: *Proceedings of the IEEE International Conference on Computer Vision*; 2017. p. 1212–1221.
- 710
- 711 **Chollet F**, et al., Keras. GitHub; 2015. <https://github.com/fchollet/keras>.
- 712 **Chollet F**. Xception: Deep learning with depthwise separable convolutions. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*; 2017. p. 1251–1258.
- 713
- 714 **Costa AC**, Ahamed T, Stephens GJ. Adaptive, locally linear models of complex dynamics. *Proceedings of the National Academy of Sciences*. 2019; 116(5):1501–1510. <https://www.pnas.org/content/116/5/1501>, doi: 10.1073/pnas.1813476116.
- 715
- 716 **Crall JD**, Gravish N, Mountcastle AM, Combes SA. BEETag: a low-cost, image-based tracking system for the study of animal behavior and locomotion. *PloS one*. 2015; 10(9):e0136487.
- 717
- 718 **Dell AI**, Bender JA, Branson K, Couzin ID, de Polavieja GG, Noldus LP, Pérez-Escudero A, Perona P, Straw AD, Wikelski M, et al. Automated image-based tracking and its application in ecology. *Trends in ecology & evolution*. 2014; 29(7):417–428.
- 719
- 720 **Deng J**, Dong W, Socher R, Li LJ, Li K, Fei-Fei L. Imagenet: A large-scale hierarchical image database. . 2009; .
- 721 **Doudna JA**, Charpentier E. The new frontier of genome engineering with CRISPR-Cas9. *Science*. 2014; 346(6213):1258096.
- 722 **Duane S**, Kennedy AD, Pendleton BJ, Roweth D. Hybrid Monte Carlo. *Phys Lett B*. 1987; 195(2):216–222.
- 723 **Dugas C**, Bengio Y, Bélisle F, Nadeau C, Garcia R. Incorporating second-order functional knowledge for better option pricing. In: *Advances in neural information processing systems*; 2001. p. 472–478.
- 724
- 725 **Flack A**, Nagy M, Fiedler W, Couzin ID, Wikelski M. From local collective behavior to global migratory patterns in white storks. *Science*. 2018; 360(6391):911–914.
- 726
- 727 **Francisco FA**, Nührenberg P, Jordan AL. A low-cost, open-source framework for tracking and behavioural analysis of animals in aquatic ecosystems. *bioRxiv*. 2019; <https://www.biorxiv.org/content/early/2019/03/09/571232>, doi: 10.1101/571232.
- 728
- 729 **Goodfellow I**, Bengio Y, Courville A. Deep learning. MIT press; 2016.
- 730 **Goodfellow I**, Pouget-Abadie J, Mirza M, Xu B, Warde-Farley D, Ozair S, Courville A, Bengio Y. Generative adversarial nets. In: *Advances in neural information processing systems*; 2014. p. 2672–2680.
- 731
- 732 **Graving JM**, pinpoint: behavioral tracking using 2D barcode tags v0.0.1-alpha; 2017. <https://doi.org/10.5281/zenodo.1008970>, doi: 10.5281/zenodo.1008970.
- 733
- 734 **Graving JM**, Chae D, Naik H, Li L, Koger B, Costelloe BR, Couzin ID, Example Datasets for DeepPoseKit; 2019. <https://doi.org/10.5281/zenodo.3366908>, doi: 10.5281/zenodo.3366908.
- 735
- 736 **Guizar-Sicairos M**, Thurman ST, Fienup JR. Efficient subpixel image registration algorithms. *Optics letters*. 2008; 33(2):156–158.

- 737 Günel S, Rhodin H, Morales D, Campagnolo J, Ramdya P, Fua P. DeepFly3D: A deep learning-based approach for 3D limb and  
738 appendage tracking in tethered, adult *Drosophila*. *bioRxiv*. 2019; p. 640375.
- 739 He K, Zhang X, Ren S, Sun J. Deep residual learning for image recognition. In: *Proceedings of the IEEE conference on computer  
740 vision and pattern recognition*; 2016. p. 770–778.
- 741 Hoffman MD, Gelman A. The No-U-Turn sampler: adaptively setting path lengths in Hamiltonian Monte Carlo. *Journal of  
742 Machine Learning Research*. 2014; 15(1):1593–1623.
- 743 Huang G, Liu Z, Van Der Maaten L, Weinberger KQ. Densely connected convolutional networks. In: *Proceedings of the IEEE  
744 conference on computer vision and pattern recognition*; 2017. p. 4700–4708.
- 745 Huang J, Rathod V, Sun C, Zhu M, Korattikara A, Fathi A, Fischer I, Wojna Z, Song Y, Guadarrama S, et al. Speed/accuracy  
746 trade-offs for modern convolutional object detectors. In: *Proceedings of the IEEE conference on computer vision and pattern  
747 recognition*; 2017. p. 7310–7311.
- 748 Insafutdinov E, Pishchulin L, Andres B, Andriluka M, Schiele B. Deepcrut: A deeper, stronger, and faster multi-person pose  
749 estimation model. In: *European Conference on Computer Vision* Springer; 2016. p. 34–50.
- 750 Iqbal U, Milan A, Gall J. PoseTrack: Joint multi-person pose estimation and tracking. In: *Proceedings of the IEEE Conference on  
751 Computer Vision and Pattern Recognition*; 2017. p. 2011–2020.
- 752 Jaques M, Burke M, Hospedales T. Physics-as-Inverse-Graphics: Joint Unsupervised Learning of Objects and Physics from Video.  
753 arXiv preprint arXiv:190511169. 2019; .
- 754 Javer A, Currie M, Lee CW, Hokanson J, Li K, Martineau CN, Yemini E, Grundy LJ, Li C, Ch'ng Q, et al. An open-source platform for  
755 analyzing and sharing worm-behavior data. *Nature methods*. 2018; 15(9):645.
- 756 Jégou S, Drozdal M, Vázquez D, Romero A, Bengio Y. The One Hundred Layers Tiramisu: Fully Convolutional DenseNets for  
757 Semantic Segmentation. *CoRR*. 2017; abs/1611.09326. <http://arxiv.org/abs/1611.09326>.
- 758 Johnson J, Alahi A, Fei-Fei L. Perceptual losses for real-time style transfer and super-resolution. In: *European conference on  
759 computer vision* Springer; 2016. p. 694–711.
- 760 Johnson M, Duvenaud DK, Wiltchko A, Adams RP, Datta SR. Composing graphical models with neural networks for structured  
761 representations and fast inference. In: *Advances in neural information processing systems*; 2016. p. 2946–2954.
- 762 Jolles JW, Boogert NJ, Sridhar VH, Couzin ID, Manica A. Consistent individual differences drive collective behavior and group  
763 functioning of schooling fish. *Current Biology*. 2017; 27(18):2862–2868.
- 764 Jung A, imgaug. GitHub; 2018. <https://github.com/aleju/imgaug>.
- 765 Kain J, Stokes C, Gaudry Q, Song X, Foley J, Wilson R, De Bivort B. Leg-tracking and automated behavioural classification in  
766 *Drosophila*. *Nature communications*. 2013; 4:1910.
- 767 Kain JS, Stokes C, de Bivort BL. Phototactic personality in fruit flies and its suppression by serotonin and white. *Proceedings of  
768 the National Academy of Sciences*. 2012; 109(48):19834–19839.
- 769 Kays R, Crofoot MC, Jetz W, Wikelski M. Terrestrial animal tracking as an eye on life and planet. *Science*. 2015; 348(6240):aaa2478.
- 770 Ke L, Chang MC, Qi H, Lyu S. Multi-Scale Structure-Aware Network for Human Pose Estimation. In: *The European Conference on  
771 Computer Vision (ECCV)*; 2018. .
- 772 Kendall A, Gal Y. What uncertainties do we need in bayesian deep learning for computer vision? In: *Advances in neural  
773 information processing systems*; 2017. p. 5574–5584.
- 774 Kiefer J, Wolfowitz J, et al. Stochastic estimation of the maximum of a regression function. *The Annals of Mathematical Statistics*.  
775 1952; 23(3):462–466.
- 776 Kingma DP, Ba J. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980. 2014; .
- 777 Klambauer G, Unterthiner T, Mayr A, Hochreiter S. Self-normalizing neural networks. In: *Advances in neural information  
778 processing systems*; 2017. p. 971–980.
- 779 Klibaite U, Berman GJ, Cande J, Stern DL, Shaevitz JW. An unsupervised method for quantifying the behavior of paired animals.  
780 *Physical biology*. 2017; 14(1):015006.
- 781 Klibaite U, Shaevitz JW. Interacting fruit flies synchronize behavior. *bioRxiv*. 2019; p. 545483.
- 782 Krakauer JW, Ghazanfar AA, Gomez-Marín A, MacIver MA, Poeppel D. Neuroscience needs behavior: correcting a reductionist  
783 bias. *Neuron*. 2017; 93(3):480–490.
- 784 Kuhn M, Johnson K. *Applied predictive modeling*, vol. 26. Springer; 2013.

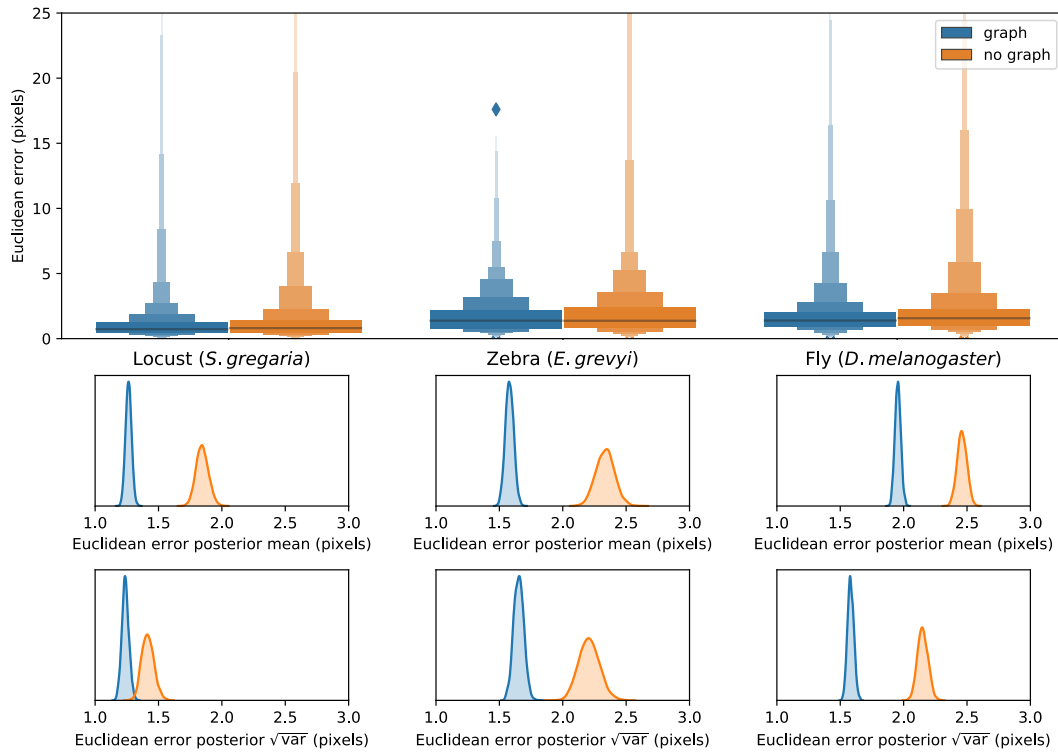
- 785 **Kulkarni TD**, Whitney WF, Kohli P, Tenenbaum J. Deep convolutional inverse graphics network. In: *Advances in neural information*  
786 *processing systems*; 2015. p. 2539–2547.
- 787 **Kumar M**, Babaeizadeh M, Erhan D, Finn C, Levine S, Dinh L, Kingma D. VideoFlow: A flow-based generative model for video.  
788 arXiv preprint arXiv:190301434. 2019; .
- 789 **LeCun Y**, Bengio Y, Hinton G. Deep learning. *nature*. 2015; 521(7553):436.
- 790 **Li H**, Xu Z, Taylor G, Studer C, Goldstein T. Visualizing the loss landscape of neural nets. In: *Advances in Neural Information*  
791 *Processing Systems*; 2018. p. 6391–6401.
- 792 **Long J**, Shelhamer E, Darrell T. Fully convolutional networks for semantic segmentation. In: *Proceedings of the IEEE conference on*  
793 *computer vision and pattern recognition*; 2015. p. 3431–3440.
- 794 **Markowitz JE**, Gillis WF, Beron CC, Neufeld SQ, Robertson K, Bhagat ND, Peterson RE, Peterson E, Hyun M, Linderman SW, et al.  
795 The striatum organizes 3D behavior via moment-to-moment action selection. *Cell*. 2018; 174(1):44–58.
- 796 **Mathis A**, Mamidanna P, Cury KM, Abe T, Murthy VN, Mathis MW, Bethge M. DeepLabCut: markerless pose estimation of user-  
797 defined body parts with deep learning. *Nature Neuroscience*. 2018; <https://www.nature.com/articles/s41593-018-0209-y>.
- 798 **Mathis A**, Warren RA. On the inference speed and video-compression robustness of DeepLabCut. bioRxiv. 2018; <https://www.biorxiv.org/content/early/2018/10/30/457242>, doi: 10.1101/457242.
- 800 **Mendes CS**, Bartos I, Akay T, Márka S, Mann RS. Quantification of gait parameters in freely walking wild type and sensory  
801 deprived *Drosophila melanogaster*. *elife*. 2013; 2:e00231.
- 802 **Munkres J**. Algorithms for the assignment and transportation problems. *Journal of the society for industrial and applied*  
803 *mathematics*. 1957; 5(1):32–38.
- 804 **Nagy M**, Akos Z, Biro D, Vicsek T. Hierarchical group dynamics in pigeon flocks. *Nature*. 2010; 464(7290):890.
- 805 **Nagy M**, Vásárhelyi G, Pettit B, Roberts-Mariani I, Vicsek T, Biro D. Context-dependent hierarchies in pigeons. *Proceedings of the*  
806 *National Academy of Sciences*. 2013; 110(32):13049–13054.
- 807 **Nath T**, Mathis A, Chen AC, Patel A, Bethge M, Mathis MW. Using DeepLabCut for 3D markerless pose estimation across species  
808 and behaviors. *Nature protocols*. 2019; .
- 809 **Newell A**, Yang K, Deng J. Stacked Hourglass Networks for Human Pose Estimation. In: *Computer Vision - ECCV 2016 -*  
810 *14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part VIII*; 2016. p. 483–499. doi:  
811 10.1007/978-3-319-46484-8\_29.
- 812 **Van den Oord A**, Kalchbrenner N, Espeholt L, Vinyals O, Graves A, et al. Conditional image generation with pixelcnn decoders.  
813 In: *Advances in neural information processing systems*; 2016. p. 4790–4798.
- 814 **Oord Avd**, Dieleman S, Zen H, Simonyan K, Vinyals O, Graves A, Kalchbrenner N, Senior A, Kavukcuoglu K. Wavenet: A generative  
815 model for raw audio. arXiv preprint arXiv:160903499. 2016; .
- 816 **Pereira TD**, Aldarondo DE, Willmore L, Kislin M, Wang SSH, Murthy M, Shaevitz JW. Fast animal pose estimation using deep  
817 neural networks. *Nature methods*. 2019; 16(1):117.
- 818 **Pérez-Escudero A**, Vicente-Page J, Hinz RC, Arganda S, De Polavieja GG. idTracker: tracking individuals in a group by automatic  
819 identification of unmarked animals. *Nature methods*. 2014; 11(7):743.
- 820 **Pratt LY**. Discriminability-based transfer between neural networks. In: *Advances in neural information processing systems*; 1993.  
821 p. 204–211.
- 822 **Prechelt L**. Automatic early stopping using cross validation: quantifying the criteria. *Neural Networks*. 1998; 11(4):761–767.
- 823 **Ran FA**, Hsu PD, Wright J, Agarwala V, Scott DA, Zhang F. Genome engineering using the CRISPR-Cas9 system. *Nature protocols*.  
824 2013; 8(11):2281.
- 825 **Ren S**, He K, Girshick R, Sun J. Faster r-cnn: Towards real-time object detection with region proposal networks. In: *Advances in*  
826 *neural information processing systems*; 2015. p. 91–99.
- 827 **Robbins H**, Monro S. A stochastic approximation method. *The annals of mathematical statistics*. 1951; p. 400–407.
- 828 **Romero-Ferrero F**, Bergomi MG, Hinz RC, Heras FJ, de Polavieja GG. idtracker. ai: tracking all individuals in small or large  
829 collectives of unmarked animals. *Nature methods*. 2019; 16(2):179.
- 830 **Ronneberger O**, Fischer P, Brox T. U-net: Convolutional networks for biomedical image segmentation. In: *International*  
831 *Conference on Medical image computing and computer-assisted intervention* Springer; 2015. p. 234–241.
- 832 **Rosenthal SB**, Twomey CR, Hartnett AT, Wu HS, Couzin ID. Revealing the hidden networks of interaction in mobile animal groups  
833 allows prediction of complex behavioral contagion. *Proceedings of the National Academy of Sciences*. 2015; 112(15):4690–  
834 4695.



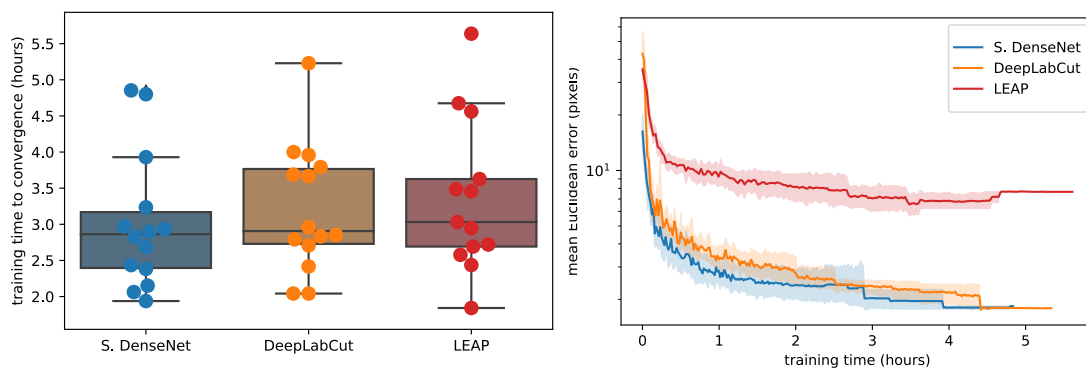
- 835 **Roy AG**, Conjeti S, Navab N, Wachinger C. Bayesian QuickNAT: Model Uncertainty in Deep Whole-Brain Segmentation for  
836 Structure-wise Quality Control. *CoRR*. 2018; abs/1811.09800. <http://arxiv.org/abs/1811.09800>.
- 837 **Sabour S**, Frosst N, Hinton GE. Dynamic routing between capsules. In: *Advances in neural information processing systems*; 2017.  
838 p. 3856–3866.
- 839 **Sandler M**, Howard A, Zhu M, Zhmoginov A, Chen LC. Mobilenetv2: Inverted residuals and linear bottlenecks. In: *Proceedings of*  
840 *the IEEE Conference on Computer Vision and Pattern Recognition*; 2018. p. 4510–4520.
- 841 **Schiffman R**, Drones flying high as new tool for field biologists. *American Association for the Advancement of Science*; 2014.
- 842 **Seethapathi N**, Wang S, Saluja R, Blohm G, Kording KP. Movement science needs different pose tracking algorithms. *arXiv*  
843 preprint arXiv:190710226. 2019; .
- 844 **Stephens GJ**, de Mesquita MB, Ryu WS, Bialek W. Emergence of long timescales and stereotyped behaviors in *Caenorhabditis*  
845 *elegans*. *Proceedings of the National Academy of Sciences*. 2011; 108(18):7286–7289.
- 846 **Stowers JR**, Hofbauer M, Bastien R, Griessner J, Higgins P, Farooqui S, Fischer RM, Nowikovsky K, Haubensak W, Couzin ID, et al.  
847 Virtual reality for freely moving animals. *Nature methods*. 2017; 14(10):995.
- 848 **Strandburg-Peshkin A**, Farine DR, Couzin ID, Crofoot MC. Shared decision-making drives collective movement in wild baboons.  
849 *Science*. 2015; 348(6241):1358–1361.
- 850 **Strandburg-Peshkin A**, Farine DR, Crofoot MC, Couzin ID. Habitat and social factors shape individual decisions and emergent  
851 group structure during baboon collective movement. *Elife*. 2017; 6:e19505.
- 852 **Strandburg-Peshkin A**, Twomey CR, Bode NW, Kao AB, Katz Y, Ioannou CC, Rosenthal SB, Torney CJ, Wu HS, Levin SA, et al.  
853 Visual sensory networks and effective information transfer in animal groups. *Current Biology*. 2013; 23(17):R709–R711.
- 854 **Todd JG**, Kain JS, de Bivort BL. Systematic exploration of unsupervised methods for mapping behavior. *Physical biology*. 2017;  
855 14(1):015002.
- 856 **Tran D**, Hoffman MW, Moore D, Suter C, Vasudevan S, Radul A. Simple, distributed, and accelerated probabilistic programming.  
857 In: *Advances in Neural Information Processing Systems*; 2018. p. 7609–7620.
- 858 **Uhlmann V**, Ramdya P, Delgado-Gonzalo R, Benton R, Unser M. FlyLimbTracker: An active contour based approach for leg  
859 segment tracking in unmarked, freely behaving *Drosophila*. *PLoS One*. 2017; 12(4):e0173433.
- 860 **Valentin J**, Keskin C, Pidlypenskyi P, Makadia A, Sud A, Bouaziz S. TensorFlow Graphics: Computer Graphics Meets Deep  
861 Learning. In: ; 2019. .
- 862 **Versace E**, Caffini M, Werkhoven Z, de Bivort BL. Individual, but not population asymmetries, are modulated by social  
863 environment and genotype in *Drosophila melanogaster*. *bioRxiv*. 2019; p. 694901.
- 864 **Weigert M**, Schmidt U, Boothe T, Müller A, Dibrov A, Jain A, Wilhelm B, Schmidt D, Broaddus C, Culley S, et al. Content-aware  
865 image restoration: pushing the limits of fluorescence microscopy. *Nature methods*. 2018; 15(12):1090.
- 866 **Werkhoven Z**, Rohrsen C, Qin C, Brems B, de Bivort B. MARGO (Massively Automated Real-time GUI for Object-tracking), a  
867 platform for high-throughput ethology. *BioRxiv*. 2019; p. 593046.
- 868 **Wild B**, Sixt L, Landgraf T. Automatic localization and decoding of honeybee markers using deep convolutional neural networks.  
869 *CoRR*. 2018; abs/1802.04557. <http://arxiv.org/abs/1802.04557>.
- 870 **Wiltchko AB**, Johnson MJ, Iurilli G, Peterson RE, Katon JM, Pashkovski SL, Abaira VE, Adams RP, Datta SR. Mapping sub-second  
871 structure in mouse behavior. *Neuron*. 2015; 88(6):1121–1135.
- 872 **Zuffi S**, Kanazawa A, Jacobs DW, Black MJ. 3D menagerie: Modeling the 3D shape and pose of animals. In: *Proceedings of the*  
873 *IEEE Conference on Computer Vision and Pattern Recognition*; 2017. p. 6365–6373.



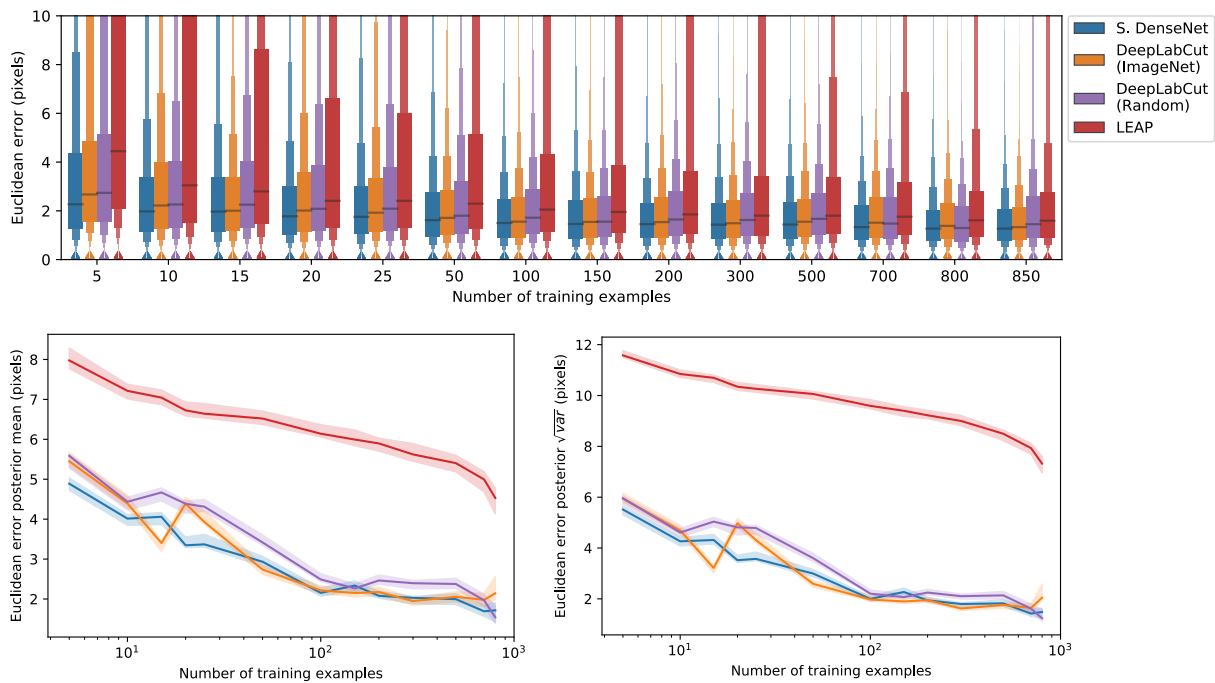
**Appendix 0 Figure 5.** Our subpixel maxima algorithm increases speed without decreasing accuracy. Prediction accuracy on the fly dataset is maintained across downsampling configurations (a). Letter-value plots (a-top) show the raw error distributions for each configuration. Visualizations of the credible intervals (99% highest-density region) of the posterior distributions for the mean and variance (a-bottom) illustrate statistical differences between the error distributions, where using subpixel maxima decreases both the mean and variance of the error distribution. Inference speed is fast and can be run in real-time on single images (batch size = 1) at ~30-110Hz or offline (batch size = 100) upwards of 1000Hz (b). Plots show the inference speeds for our Stacked DenseNet model across downsampling configurations as well as the other models we tested for each of our datasets.



**Appendix 0 Figure 6.** Predicting the multi-scale geometry of the posture graph reduces error. Letter-value plots (**top**) show the raw error distributions for each experiment. Visualizations of the posterior distributions for the mean and variance (**bottom**) show statistical differences between the error distributions. Predicting the posture graph decreases both the mean and variance of the error distribution.



**Appendix 0 Figure 7.** Training time required for our Stacked DenseNet model, the DeepLabCut model (*Mathis et al., 2018*), and the LEAP model (*Pereira et al., 2019*) ( $n=15$  per model) using our zebra dataset. Boxplots and swarm plots (**left**) show the total training time to convergence ( $<0.001$  improvement in validation loss for 50 epochs). Line plots (**right**) illustrate the Euclidean error of the validation set during training, where error bars show bootstrapped ( $n=1000$ ) 99% confidence intervals of the mean. Fully training models to convergence requires only a few hours of optimization (**left**) with reasonable accuracy reached after only 1 hour (**right**) for our Stacked DenseNet model.



**Appendix 0 Figure 8.** A comparison of prediction accuracy with different numbers of training examples from our zebra dataset. The error distributions shown as letter-value plots (**top**) illustrate the Euclidean error for the remainder of the dataset not used for training—with a total of 900 labeled examples in the dataset. Line plots (**bottom**) show posterior credible intervals (99% highest-density region) for the mean and variance of the error distributions. We tested our Stacked DenseNet model; the DeepLabCut model (*Mathis et al., 2018*) with transfer learning—i.e., with weights pretrained on ImageNet (*Deng et al., 2009*); the same model without transfer learning—i.e., with randomly-initialized weights; and the LEAP model (*Pereira et al., 2019*). Our Stacked DenseNet model achieves high accuracy using few training examples without the use the transfer learning. Using pretrained weights does slightly decrease overall prediction error for the DeepLabCut model (*Mathis et al., 2018*), but the effect size is relatively small.

## 874 Appendix 1

### 875 Convolutional neural networks (CNNs)

876 *Artificial neural networks* like CNNs are complex, non-linear regression models that "learn" a hierar-  
877 chically-organized set of parameters from real-world data via optimization. These machine learning  
878 models are now commonplace in science and industry and have proven to be surprisingly effective for  
879 a large number of applications where more conventional statistical models have failed (**LeCun et al.,**  
880 **2015**). For computer vision tasks, CNN parameters typically take the form of two-dimensional convo-  
881 lutional filters that are optimized to detect spatial features needed to model relationships between  
882 high-dimensional image data and some related variable(s) of interest, such as locations in space—e.g.  
883 posture keypoints—or semantic labels (**Long et al., 2015; Badrinarayanan et al., 2015**).

884 Once a training set is generated (Appendix 2), a CNN model must be selected and optimized to  
885 perform the prediction task. CNNs are incredibly flexible with regard to how models are specified and  
886 trained, which is both an advantage and a disadvantage. This flexibility means models can be adapted  
887 to almost any computer vision task, but it also means the number of possible model architectures and  
888 optimization schemes is very large. This can make selecting an architecture and specifying hyperpa-  
889 rameters a challenging process. However, most research on pose estimation has converged on a set of  
890 models that generally work well for this task (Appendix 3).

891 After selecting an architecture, the parameters of the model are set to an initial value and then  
892 iteratively updated to minimize some objective function, or *loss function*, that describes the difference  
893 between the model's predictive distribution and the true distribution of the data—in other words,  
894 the likelihood of the model's output is maximized. These parameter updates are performed using a  
895 modified version of the gradient descent algorithm (**Cauchy 1847**) known as *mini-batch stochastic gradient*  
896 *descent*—often referred to as simply *stochastic gradient descent* or *SGD* (**Robbins and Monro, 1951; Kiefer**  
897 **et al., 1952**). SGD iteratively optimizes the model parameters using small randomly-selected subsamples,  
898 or *batches*, of training data. Using SGD allows the model to be trained on extremely large datasets  
899 in an iterative "online" fashion without the need to load the entire dataset into memory. The model  
900 parameters are updated with each batch by adjusting the parameter values in a direction that minimizes  
901 the error—where one round of training on the full dataset is commonly referred to as an *epoch*. The  
902 original SGD algorithm requires careful selection and tuning of hyperparameters to successfully optimize  
903 a model, but modern versions of the algorithm, such as *ADAM* (**Kingma and Ba, 2014**), automatically  
904 tune these hyperparameters, which makes optimization more straightforward.

905 The model parameters are optimized until they reach a convergence criterion, which is some measure  
906 of performance that indicates the model has reached a good location in parameter space. The most  
907 commonly used convergence criterion is a measure of predictive accuracy—often the loss function  
908 used for optimization—on a held-out *validation set*—a subsample of the training data not used for  
909 optimization—that evaluates the model's ability to generalize to new "out-of-sample" data. The model is  
910 typically evaluated at the end of each training epoch to assess performance on the validation set. Once  
911 performance on the validation set stops improving, training is usually stopped to prevent the model  
912 from overfitting to the training set—a technique known as *early stopping* (**Prechelt, 1998**).



## 913 Appendix 2

### 914 Collecting training data

915 Depending on the variability of the data, CNNs usually require thousands or tens of thousands of  
916 manually-annotated examples in order to reach human-level accuracy. However, in laboratory settings,  
917 sources of image variation like lighting and spatial scale can be more easily controlled, which minimizes  
918 the number of training examples needed to achieve accurate predictions.

919 This need for a large training set can be further reduced in a number of ways. Two commonly used  
920 methods include (1) *transfer learning*—using a model with parameters that are pre-trained on a larger  
921 set of images, such as the ImageNet database (**Deng et al., 2009**), containing diverse features (**Pratt,**  
922 **1993; Insafutdinov et al., 2016; Mathis et al., 2018**)— and (2) *augmentation*— artificially increasing data  
923 variance by applying spatial and noise transformations such as flipping (mirroring), rotating, scaling, and  
924 adding different forms of noise or artificial occlusions. Both of these methods act as useful forms of  
925 *regularization*—incorporating a prior distribution—that allows the model to generalize well to new data  
926 even when the training set is small. Transfer learning incorporates prior information that images from  
927 the full dataset should contain statistical features similar to other images of the natural world, while  
928 augmentation incorporates prior knowledge that animals are bilaterally symmetric, can vary in their  
929 body size, position, and orientation, and that noise and occlusions sometimes occur.

930 **Pereira et al. (2019)** introduced two especially clever solutions for collecting an adequate training  
931 set. First, they cluster unannotated images based on pixel variance and uniformly sample images from  
932 each cluster, which reduces correlation between training examples and ensures the training data are  
933 representative of the entire distribution of possible images. Second, they use *active learning* where  
934 a CNN is trained on a small number of annotated examples and is then used to initialize keypoint  
935 locations for a larger set of unannotated data. These pre-initialized data are then manually corrected  
936 by the annotator, the model is retrained, and the unannotated data are re-initialized. The annotator  
937 applies this process iteratively as the training set grows larger until they are providing only minor  
938 adjustments to the pre-initialized data. This “human-in-the-loop”-style annotation expedites the process  
939 of generating an adequately large training set by reducing the cognitive load on the annotator—where  
940 the pose estimation model serves as a “cognitive partner”. Such a strategy also allows the annotator to  
941 automatically select new training examples based on the performance of the current iteration—where  
942 low-confidence predictions indicate examples that should be annotated for maximum improvement  
943 (Figure 1).

944 Of course, annotating image data requires software made for this purpose. **Pereira et al. (2019)**  
945 provide a custom annotation GUI written in MATLAB specifically designed for annotating posture using  
946 an active learning strategy. **Mathis et al. (2018)** recently added a Python-based GUI in an updated  
947 version of their software—including active learning and image sampling methods (see **Nath et al. 2019**).  
948 Our framework also includes a Python-based GUI for annotating data with similar features to **Mathis**  
949 **et al. (2018)** and **Pereira et al. (2019)**.

## 950 Appendix 3

### 951 Fully-convolutional regression

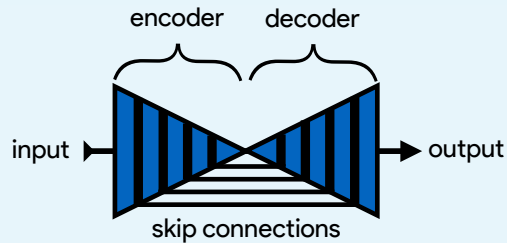
952 For the task of pose estimation, a CNN is optimized to predict the locations of postural keypoints in an  
953 image. One approach is to use a CNN to directly predict the numerical value of each keypoint coordinate  
954 as an output. However, making predictions in this way removes real-world constraints on the model's  
955 predictive distribution by destroying spatial relationships within images, which negates many of the  
956 advantages of using CNNs in the first place.

957 CNNs are particularly good at transforming one image to produce another related image, or set  
958 of images, while preserving spatial relationships and allowing for translation-invariant predictions—a  
959 configuration known as a *fully-convolutional neural network* or *F-CNN* (Long et al., 2015). Therefore,  
960 instead of directly regressing images to coordinate values, a popular solution (Newell et al., 2016;  
961 Insafutdinov et al., 2016; Mathis et al., 2018; Pereira et al., 2019) is to optimize a F-CNN that transforms  
962 images to predict a stack of output images known as *confidence maps*—one for each keypoint. Each  
963 confidence map in the output volume contains a single, two-dimensional, symmetric Gaussian indicating  
964 the location of each joint, and the scalar value of the peak indicates the confidence score of the  
965 prediction—typically a value between 0 and 1. The confidence maps are then processed to produce the  
966 coordinates of each keypoint.

967 In the case of *multiple pose estimation* where an image contains many individuals, the global geom-  
968 etry of the posture graph is also predicted by training the model to produce *part affinity fields* (Cao  
969 et al., 2017)—directional vector fields drawn between joints in the posture graph—or *pairwise terms*  
970 (Insafutdinov et al., 2016)—vector fields of the conditional distributions between posture keypoints  
971 (e.g.  $p(\text{foot}|\text{head})$ ). This allows multiple posture graphs to be disentangled from the image using graph  
972 partitioning as the vector fields indicate the probability of the connection between joints (see Cao et al.  
973 2017 for details).

975

## Box 1. Encoder-decoder models



976

977

978

**Box 1 Figure 1.** An illustration of the basic encoder-decoder design. The encoder converts the input images into spatial features, and the decoder transforms spatial features to the desired output.

979

980

981

982

983

984

985

986

987

A popular type of F-CNN (Appendix 3) for solving posture regression problems is known as an *encoder-decoder* model (Figure 1), which first gained popularity for the task of *semantic segmentation*—a supervised computer vision problem where each pixel in an image is classified into a one of several labeled categories like “dog”, “tree”, or “road” (Long et al., 2015). This model is designed to repeatedly convolve and downsample input images in the bottom-up *encoder* step and then convolve and upsample the encoder’s output in the top-down *decoder* step to produce the final output. Repeatedly applying convolutions and non-linear functions, or *activations*, to the input images transforms pixel values into higher-order spatial features, while downsampling and upsampling respectively increases and decreases the scale and complexity of these features.

988

989

990

991

992

993

994

995

996

997

998

999

1000

1001

1002

1003

1004

1005

1006

1007

1008

1009

1010

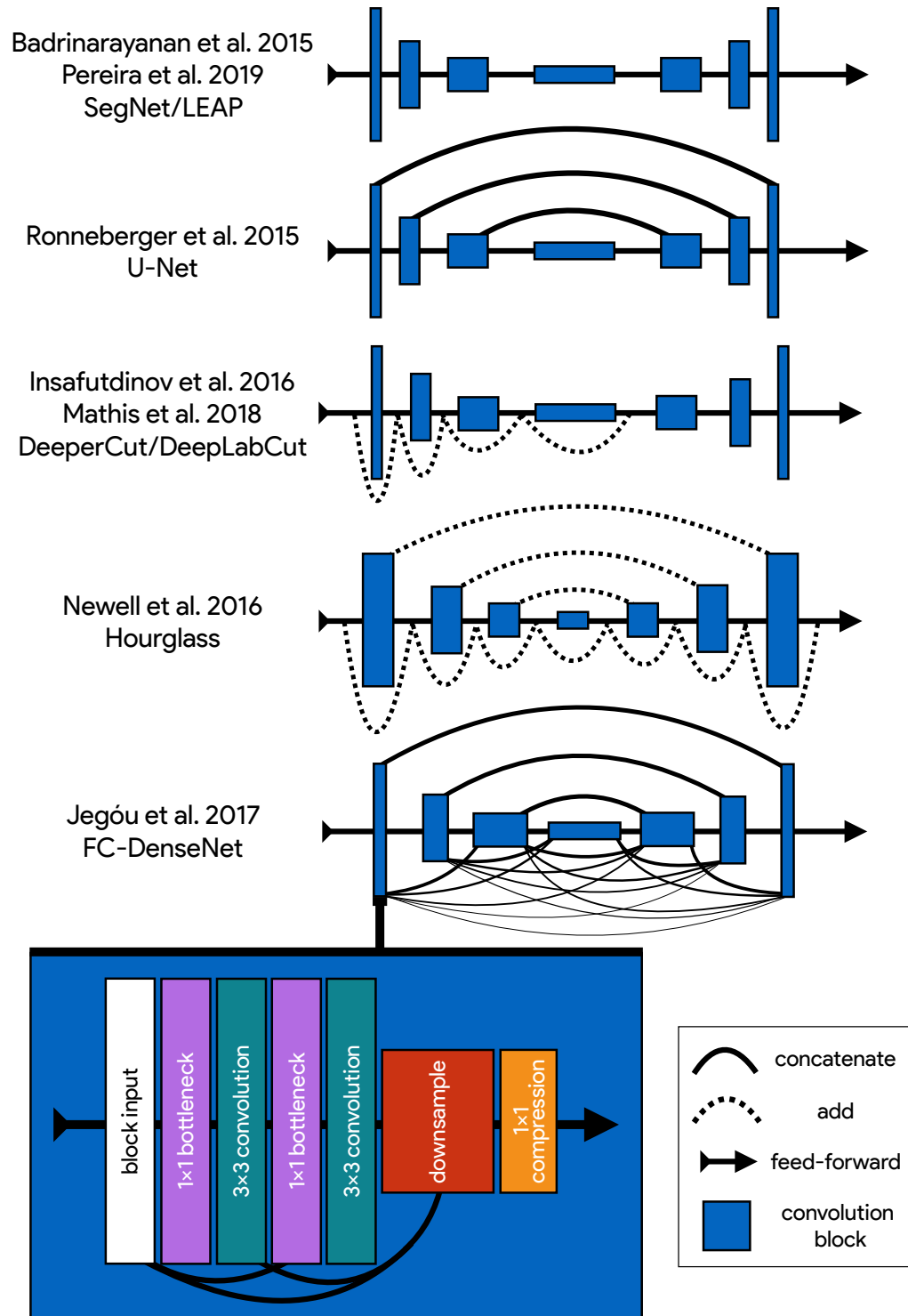
1011

1012

1013

**Badrinarayanan et al. (2015)** were the first to popularize a form of this model—known as *SegNet*—for semantic segmentation. However, this basic design is inherently limited because the decoder relies solely on the downsampled output from the encoder, which restricts the features used for predictions to those with the largest spatial scale and highest complexity. For example, a very deep network might learn a complex spatial pattern for predicting “grass” or “trees”, but because it cannot directly access information from the earliest layers of the network, it cannot use the simplest features that plants are green and brown. Subsequent work by **Ronneberger et al. (2015)** improved on these problems with the addition of *residual* or *skip connections* between the encoder and decoder, where feature maps from encoder layers are concatenated to those decoder layers with the same spatial scale. This set of connections then allows the optimizer, rather than the user, to select the most relevant spatial scale(s) for making predictions.

**Jégou et al. (2017)** are the latest to advance the encoder-decoder paradigm. These researchers introduced a fully-convolutional version of **Huang et al.’s (2017a) DenseNet** architecture known as a *fully-convolutional DenseNet*, or *FC-DenseNet*. FC-DenseNet’s key improvement is an elaborate set of feed-forward residual connections where the input to each convolutional layer is a concatenated stack of feature maps from all previous layers. This densely-connected design was motivated by the insight that many state-of-the-art models learn a large proportion of redundant features. Most CNNs are not designed so that the final output layers can access all feature maps in the network simultaneously, and this limitation causes these networks to “forget” and “relearn” important features as the input images are transformed to produce the output. In the case of the incredibly popular ResNet-101 (**He et al., 2016**) nearly 40% of the features can be classified as redundant (**Ayinde and Zurada, 2018**). A densely-connected architecture has the advantages of reduced feature redundancy, increased feature reuse, enhanced feature propagation from early layers to later layers, and subsequently, a substantial reduction in the number of parameters needed to achieve state-of-the-art results (**Huang et al., 2017a**). Recent work has also shown that DenseNet’s elaborate residual connections have the pleasant side-effect of convexifying the loss landscape during optimization (**Li et al., 2018**), which allows for faster optimization and increases the likelihood of reaching a good optimum.



**Appendix 3 Figure 1.** An illustration showing the progression of encoder-decoder architectures from the literature—ordered by performance from top to bottom (see Appendix 3 Box 1 for further details). Most advances in performance have come from adding connections between layers in the network, culminating in FC-DenseNet from Jégou et al. (2017). Lines in each illustration indicate connections between convolutional blocks with the thickness of the line indicating the magnitude of information flow between layers in the network. The size of each convolution block indicates the relative number of feature maps (width) and spatial scale (height). The callout for FC-DenseNet (Jégou et al. 2017; **bottom-left**) shows the elaborate set of skip connections within each densely-connected convolutional block as well as our additions of bottleneck and compression layers (described by Huang et al. 2017a) to increase efficiency (Appendix 7)

## 1014 Appendix 4

### 1015 The state of the art for individual pose estimation

1016 Many of the current state-of-the-art models for individual posture estimation are based on the design  
1017 from *Newell et al. (2016)* (e.g. *Ke et al. 2018*, *Chen et al. 2017*; also see benchmark results from *An-*  
1018 *driluka et al. 2014*), but employ various modifications that increase complexity to improve performance.  
1019 *Newell et al. (2016)* employ what they call a *Stacked Hourglass* network (Appendix 3 Figure 1), which consists  
1020 of a series of multi-scale encoder-decoder *hourglass* modules connected together in a feed-forward  
1021 configuration (Figure 2). The main novelties these researchers introduce include (1) stacking multiple  
1022 hourglass networks together for repeated top-down-bottom-up inference, (2) using convolutional blocks  
1023 based on the ResNet architecture (*He et al., 2016*) with residual connections between the input and  
1024 output of each block, and (3) using residual connections between the encoder and decoder (similar to  
1025 *Ronneberger et al. 2015*) with residual blocks in between. *Newell et al. (2016)* also apply a technique  
1026 known as *intermediate supervision* (Figure 2) where the loss function used for model training is applied to  
1027 the output of each hourglass as a way of improving optimization across the model's many layers. Recent  
1028 work by *Jégou et al. (2017)* has further improved on this encoder-decoder design (see Appendix 3 Box 1  
1029 and Appendix 3 Figure 1), but to the best of our knowledge, the model introduced by *Jégou et al. (2017)*  
1030 has not been previously applied to pose estimation.



## 1031 Appendix 5

### 1032 Overparameterization and the limitations of LEAP

1033 Overparameterization is a key limitation for many pose estimation methods, and addressing this  
1034 problem is critical for high-performance applications. *Pereira et al. (2019)* approached this problem by  
1035 designing their LEAP model after the model from *Badrinarayanan et al. (2015)*, which is a straightforward  
1036 encoder-decoder design (Appendix 3 Figure 1; Appendix 3 Box 1). They benchmarked their model on  
1037 posture estimation tasks for laboratory animals and compared performance with the more-complex  
1038 Stacked Hourglass model from *Newell et al. (2016)*. They found their smaller, simplified model achieved  
1039 equal or better median accuracy with dramatic improvements in inference speed up to 185 Hz. However,  
1040 *Pereira et al. (2019)* first rotationally and translationally aligned each image to improve performance,  
1041 and their reported inference speeds do not include this computationally expensive preprocessing step.  
1042 Additionally, rotationally and translationally aligning images is not always possible when the background  
1043 is complex or highly-variable—such as in field settings—or the study animal has a non-rigid body. This  
1044 limitation makes the LEAP model (*Pereira et al., 2019*) unsuitable in many cases. While their approach is  
1045 simple and effective for a multitude of experimental setups, the LEAP model (*Pereira et al., 2019*) is also  
1046 implicitly limited in the same ways as *Badrinarayanan et al.*'s SegNet model (see Appendix 3 Box 1 for  
1047 details). The LEAP model cannot make predictions using multiple spatial scales and is not robust to data  
1048 variance such as rotations (*Pereira et al., 2019*).

## 1049 Appendix 6

### 1050 Linear model fitting with Stan

1051 We estimated the joint posterior  $p(\theta_\mu, \theta_\phi | X, y)$  for each model using the No-U-Turn Sampler (NUTS;  
1052 **Hoffman and Gelman 2014**), a self-tuning variant of the Hamiltonian Monte Carlo (HMC) algorithm  
1053 (**Duane et al., 1987**), implemented in Stan (**Carpenter et al., 2017**). We drew HMC samples using 4  
1054 independent Markov chains consisting of 1,000 warm-up iterations and 1,000 sampling iterations for  
1055 a total of 4,000 sampling iterations. To speed up sampling, we randomly subsampled 20% of the data  
1056 from each replicate when fitting each linear model, and we fit each model 5 times to ensure the results  
1057 were consistent. All models converged without any signs of pathological behavior. We performed a  
1058 posterior predictive check by visually inspecting predictive samples to assess model fit. For our priors  
1059 we chose relatively uninformative distributions  $\theta_\mu \sim \text{Cauchy}(0, 5)$  and  $\theta_\phi \sim \text{Cauchy}(0, 10)$ , but we found  
1060 that the choice of prior generally did not have an effect on the final result due to the large amount of  
1061 data used to fit each model.

## 1062 Appendix 7

### 1063 Stacked DenseNet

1064 Our Stacked DenseNet model consists of an initial 7×7 convolutional layer with stride 2, to efficiently  
1065 downsample the input resolution—following *Newell et al. (2016)*—followed by a stack of densely-  
1066 connected hourglass networks with intermediate supervision (Appendix 4) applied at the output of  
1067 each network. We also include hyperparameters for the bottleneck and compression layers described  
1068 by *Huang et al. (2017a)* to make the model as efficient as possible. These consist of applying a 1×1  
1069 convolution to inexpensively compress the number of feature maps before each 3×3 convolution as well  
1070 as when downsampling and upsampling (see *Huang et al. 2017a* and Appendix 3 Figure 1 for details).

### 1071 Model hyperparameters

1072 For our Stacked Hourglass model we used a block size of 64 filters (64 filters per 3×3 convolution) with a  
1073 bottleneck factor of 2 ( $64/2 = 32$  filters per 1×1 bottleneck block). For our Stacked DenseNet model we  
1074 used a growth rate of 48 (48 filters per 3×3 convolution), a bottleneck factor of 1 ( $1 \times \text{growth rate} = 48$   
1075 filters per 1×1 bottleneck block), and a compression factor of 0.5 (feature maps compressed with 1×1  
1076 convolution to  $0.5m$  when upsampling and downsampling, where  $m$  is the number of feature maps). For  
1077 our Stacked DenseNet model we also replaced the typical configuration of batch normalization and ReLU  
1078 activations (*Goodfellow et al., 2016*) with the more recently-developed self-normalizing SELU activation  
1079 function (*Klambauer et al., 2017*), as we found this modification increased inference speed. For the  
1080 LEAP model (*Pereira et al., 2019*) we used a 1× resolution output with integer-based global maxima  
1081 because we wanted to compare our more complex models with this model in the original configuration  
1082 described by *Pereira et al. (2019)*. The LEAP model could be modified to output smaller confidence  
1083 maps and increase inference speed, but because there is no obvious "best" way to alter the model to  
1084 achieve this, we forgo any modification. Additionally, applying our subpixel maxima algorithm at high  
1085 resolution reduces inference speed compared to integer-based maxima, so this would bias our speed  
1086 comparisons.

### 1087 Our implementation of the DeepLabCut model

1088 Because the DeepLabCut model from *Mathis et al. (2018)* was not implemented in Keras (a requirement  
1089 for our pose estimation framework), we re-implemented it. Implementing this model directly in our  
1090 framework is important to ensure model training and data augmentation are identical when making  
1091 comparisons between models. As a consequence, our version of this model does not exactly match the  
1092 description in the paper but is identical except for the output. Rather than using the location refinement  
1093 maps described by *Insafutdinov et al. (2016)* and post-processing confidence maps on the CPU, our  
1094 version of the DeepLabCut model (*Mathis et al., 2018*) has an additional transposed convolutional layer  
1095 to upsample the output to  $\frac{1}{4} \times$  resolution and uses our subpixel maxima algorithm.

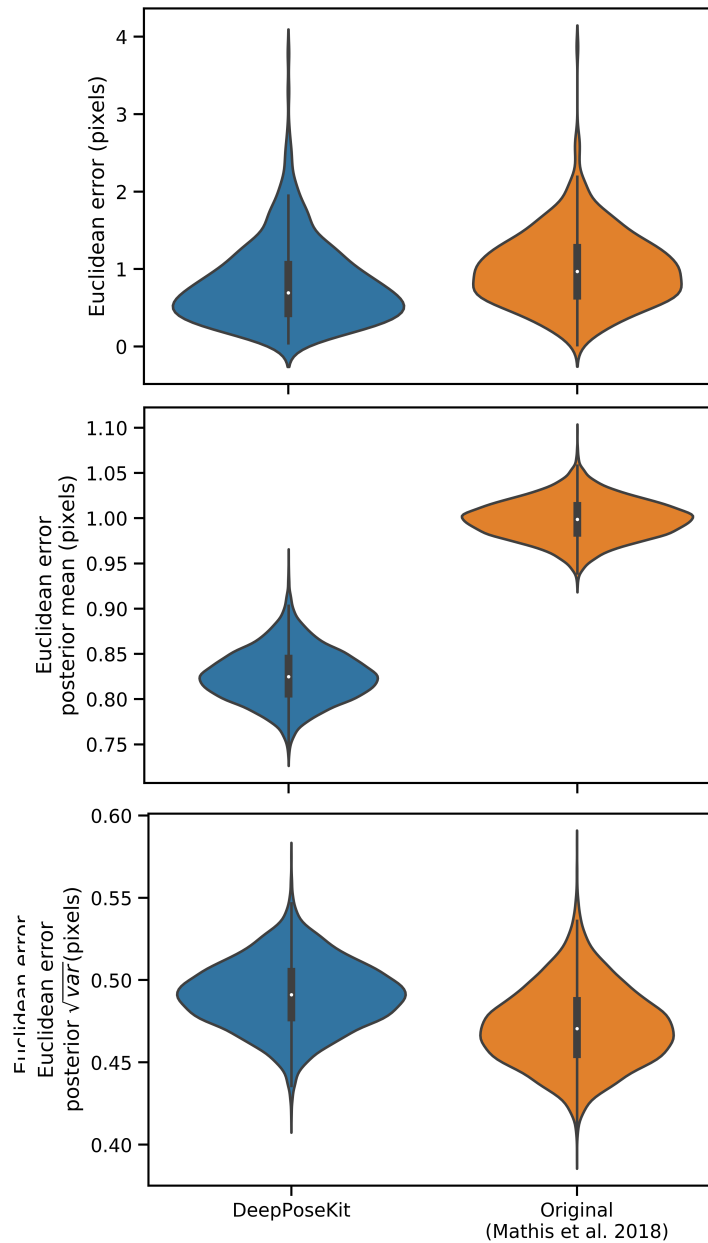
1096 To demonstrate that our implementation of the DeepLabCut model matches the performance de-  
1097 scribed by *Mathis et al. (2018)*, we compared prediction accuracy between the two frameworks using the  
1098 odor-trail mouse dataset provided by *Mathis et al. (2018)* (downloaded from <https://github.com/AlexEMG/DeepLabCut>).  
1099 This dataset consists of 116 images of a freely-moving individual mouse labeled with four keypoints  
1100 describing the location of the snout, ears, and the base of the tail. See *Mathis et al. (2018)* for further  
1101 details on this dataset. We trained both models using 95% training and 5% validation data and applied  
1102 data augmentations for both frameworks using the data augmentation procedure described by *Nath*  
1103 *et al. (2019)*. We tried to match these data augmentations as best as possible in DeepPoseKit; however,  
1104 rather than cropping images as described by *Nath et al. (2019)*, we randomly translated the images  
1105 independently along the horizontal and vertical axis by drawing from a uniform distribution in the range  
1106 [-100%, +100%]—where percentages are relative to the size of each axis. Translating the images in this  
1107 way should serve the same purpose as cropping them.

We trained the original DeepLabCut model (*Mathis et al., 2018*) using the default settings and  
recommendations from *Nath et al. (2019)* for 1 million training iterations. See *Mathis et al. (2018)*;  
*Nath et al. (2019)* for further details on the data augmentation and training routine for the original  
implementation of the DeepLabCut model (*Mathis et al., 2018*). For our re-implementation of the  
DeepLabCut model (*Mathis et al., 2018*) we trained the model with the same batch size and optimization  
scheme described in the "Model training" section. We then calculated the the prediction accuracy on the

1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132

full data set. We repeated this procedure five times for each model and fit a Bayesian linear model to a randomly selected subset of the evaluation data to compare the results statistically (see Appendix 6).

These results demonstrate that our re-implementation of and modification to the DeepLabCut model (*Mathis et al., 2018*) have little effect on prediction accuracy (Appendix 7 Figure 1). We also provide qualitative comparisons of these results in Appendix 7 Figure 1-Figure supplement 1 and Appendix 7 Figure 1-video 1. For these qualitative comparisons, we also added an additional rotational augmentation (drawing from a uniform distribution in the range [-180°, +180°]) when training our implementation of the DeepLabCut model (*Mathis et al., 2018*) as we noticed this improved generalization to the video for situations where the mouse rotated its body axis. To the best of our knowledge, rotational augmentations are not currently available when using the software from *Mathis et al. (2018)*; *Nath et al. (2019)*, which demonstrates the flexibility of the data augmentation pipeline (*Jung, 2018*) for DeepPoseKit. The inference speed for the odor-trail mouse dataset using our implementation of the DeepLabCut model (*Mathis et al., 2018*) is ~49Hz with a batch size of 64 (offline speeds) and ~35Hz with a batch size of 1 (real-time speeds) at full resolution 640×480, which matches well with results from *Mathis and Warren (2018)* of ~47Hz and ~32Hz respectively. This suggests our modifications did not affect the speed of the model and that our speed comparisons are also reasonable. Because the training routine could be changed for any underlying model—including the new models we present in this paper—this factor is not relevant when making comparisons as long as training is identical for all models being compared, which we ensure when performing our comparisons.



**Appendix 7 Figure 1.** Prediction errors for the odor-trail mouse dataset from *Mathis et al. (2018)* using the original implementation of the DeepLabCut model (*Mathis et al., 2018; Nath et al., 2019*) and our modified version of this model implemented in DeepPoseKit. Mean prediction error is slightly lower for the DeepPoseKit implementation, but there is no discernible difference in variance. These results indicate that the models achieve nearly identical prediction accuracy despite modification. We also provide qualitative comparisons of these results in Appendix 7 Figure 1-Figure supplements 1 and 2, and Appendix 7 Figure 1-video 1.

**Figure 1-Figure supplement 1.** Plots of the predicted output for Appendix 7 Figure 1-video 1 comparing our implementation of the DeepLabCut model (*Mathis et al., 2018*) in DeepPoseKit vs. the original implementation from *Mathis et al. (2018); Nath et al. (2019)*. Note the many fast jumps in position for the original version from *Mathis et al. (2018)*, which indicates prediction errors.

**Figure 1-Figure supplement 2.** Plots of the temporal derivatives of the predicted output for Appendix 7 Figure 1-video 1 comparing our implementation of the DeepLabCut model (*Mathis et al., 2018*) in DeepPoseKit vs. the original implementation from *Mathis et al. (2018); Nath et al. (2019)*. Note the many fast jumps in position for the original version from *Mathis et al. (2018)*, which indicates prediction errors.

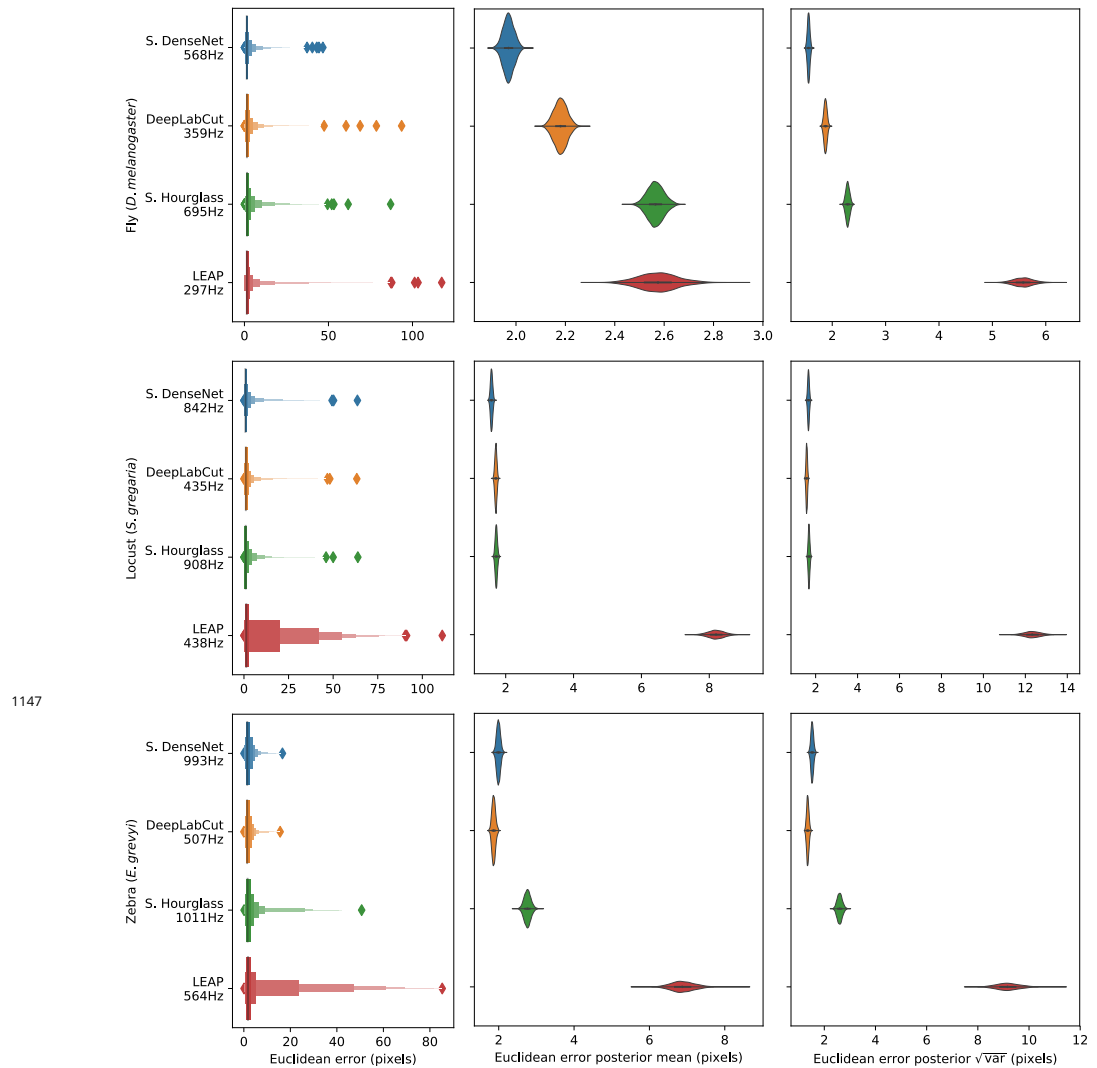
**Figure 1-video 1.** A video comparison of the tracking output of our implementation of the DeepLabCut model (*Mathis et al., 2018*) in DeepPoseKit vs. the original implementation from *Mathis et al. (2018); Nath et al. (2019)*.



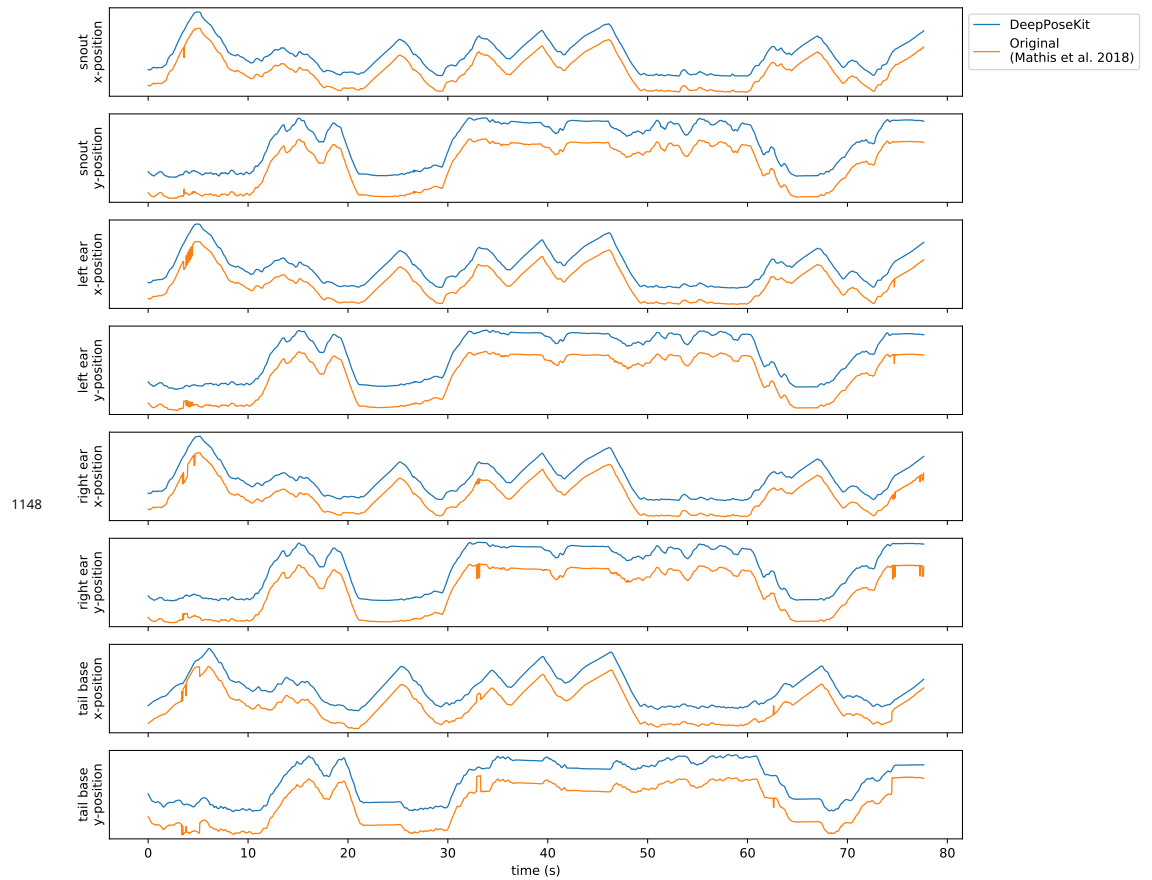
## 1133 Appendix 8

### 1134 **Depthwise-separable convolutions for memory-limited applications**

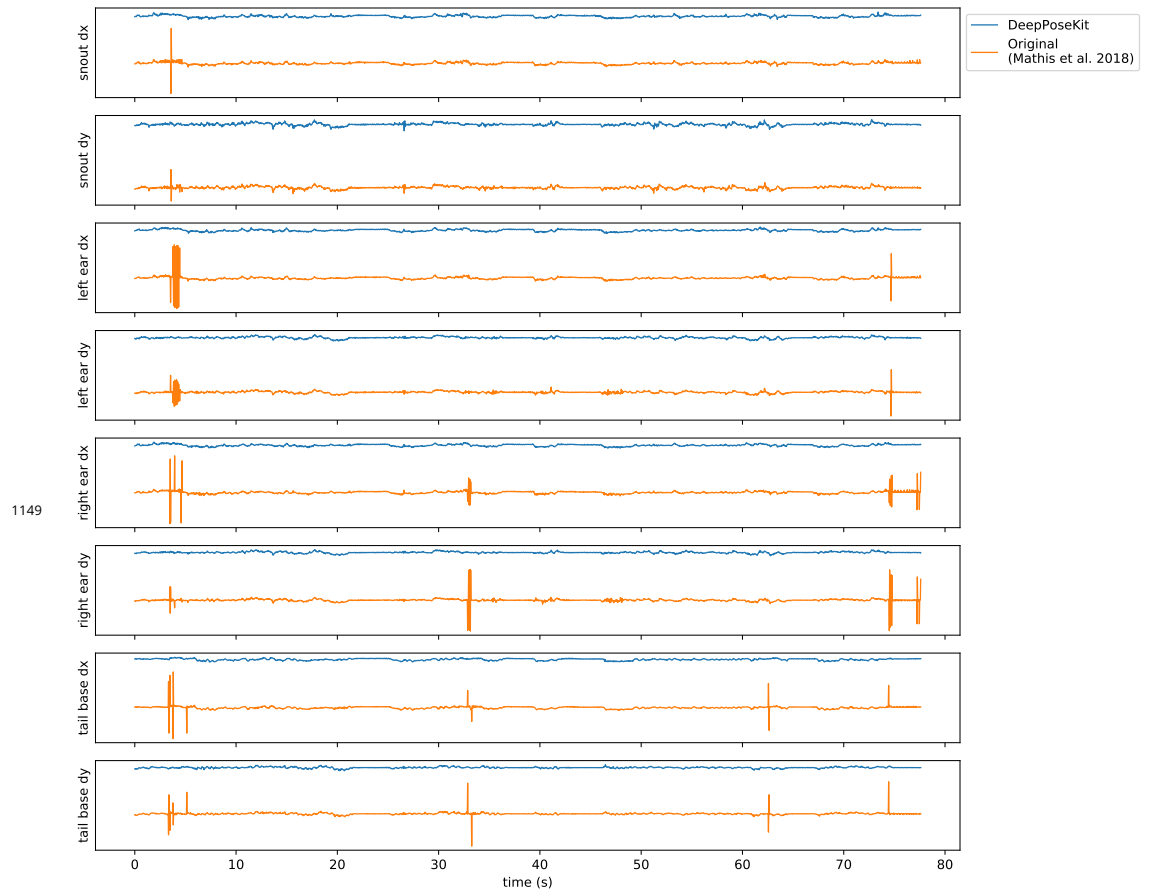
1135 In an effort to maximize model efficiency, we also experimented with replacing 3×3 convolutions in  
1136 our model implementations with 3×3 depthwise-separable convolutions —first introduced by *Chollet*  
1137 (*2017*) and now commonly used in fast, efficient “mobile” CNNs (e.g. *Sandler et al. 2018*). In theory this  
1138 modification should both reduce the memory footprint of the model and increase inference speed.  
1139 However we found that, while this does drastically decrease the memory footprint of our already  
1140 memory-efficient models, it slightly decreases accuracy and does not improve inference speed, so we  
1141 opt for a full 3×3 convolution instead. We suspect that this discrepancy between theory and application is  
1142 due to inefficient implementations of depthwise-separable convolutions in many popular deep learning  
1143 frameworks, which will hopefully improve in the near future. At the moment we include this option as  
1144 a hyperparameter for our Stacked DenseNet model, but we recommend using depthwise-separable  
1145 convolutions only for applications that require a small memory footprint such as training on a lower-end  
1146 GPU with limited memory or running inference on a mobile device.



**Figure 4-Figure supplement 1.** Euclidean error distributions for each model across our three datasets. Letter-value plots (left) show the raw error distributions for each model. Violinplots of the posterior distributions for the mean and variance (right) show statistical differences between the error distributions. Overall the LEAP model (Pereira et al., 2019) was the worst performer on every dataset in terms of both mean and variance. Our Stacked Densenet model was the best performer for the fly dataset, while our Stacked DenseNet model and the DeepLabCut model (Mathis et al., 2018) both performed equally well on the locust and zebra datasets. The posteriors for the DeepLabCut model (Mathis et al., 2018) and our Stacked DenseNet model are highly overlapping for these datasets, which suggests they are not statistically discernible from one another. Our Stacked Hourglass model (Newell et al., 2016) performed equally to the DeepLabCut model (Mathis et al., 2018) and our Stacked DenseNet model for the locust dataset but performed slightly worse for the fly and zebra datasets.



**Figure 1-Figure supplement 1.** Plots of the predicted output for Appendix 7 Figure 1-video 1 comparing our implementation of the DeepLabCut model (*Mathis et al., 2018*) in DeepPoseKit vs. the original implementation from *Mathis et al. (2018)*; *Nath et al. (2019)*. Note the many fast jumps in position for the original version from *Mathis et al. (2018)*, which indicates prediction errors.



**Figure 1-Figure supplement 2.** Plots of the temporal derivatives of the predicted output for Appendix 7 Figure 1-video 1 comparing our implementation of the DeepLabCut model (*Mathis et al., 2018*) in DeepPoseKit vs. the original implementation from *Mathis et al. (2018)*; *Nath et al. (2019)*. Note the many fast jumps in position for the original version from *Mathis et al. (2018)*, which indicates prediction errors.