# Automated Location Invariant Animal Detection In Camera Trap Images Using Publicly Available Data Sources

**Andrew Shepley[1] | Greg Falzon[2] | Paul Meek[3,4] | Paul Kwan[5]**

[1]School of Science and Technology, University of New England; Armidale, NSW, Australia

[2]College of Science and Engineering , Flinders University; Adelaide, SA, Australia

[3]Vertebrate Pest Research Unit, NSW Department of Primary Industries, PO Box 530, Coffs Harbour, NSW, Australia

[4]School of Environmental and Rural Science, University of New England, Armidale, NSW, Australia

[5]School of IT and Engineering, Melbourne Institute of Technology, Australia

**Correspondence:** Andrew Shepley: asheple2@une.edu.au

## Abstract

1. A time-consuming challenge faced by camera trap practitioners is the extraction of meaningful data from images to inform ecological management. An increasingly popular solution is automated image classification software. However, most software solutions are not sufficiently robust to be deployed on a large scale due to lack of location invariance

23    when transferring models between sites. This prevents optimal use of ecological data and

24    results in significant expenditure of time and resources to annotate and retrain deep

25    learning models.

26    2.  In this study, we aimed to (a) assess the value of publicly available image datasets in the

27        training of deep learning models for camera trap object detection focusing on images

28        obtained from FlickR and iNaturalist (FiN), (b) develop a method to be used by ecologists to

29        train location invariant image processing object detection models and (c) explore the use of

30        small subsets of camera trap images in the optimization of FiN training.

31    3.  We collected and annotated 3 datasets of images of the following classes; striped hyena,

32        rhinoceros and pig, from the image sharing websites, and used transfer learning to train 3

33        object detection models in the task of animal detection. We compared the performance of

34        these models to the performance of 3 models trained on the Wildlife Conservation Society

35        and Camera CATalogue datasets, when tested on out of sample Snapshot Serengeti datasets.

36        Furthermore, we explored optimization of the FiN trained models via infusion of small

37        subsets of camera trap images to increase robustness for challenging detection cases.

38    4.  In all experiments, the mean Average Precision (mAP) of the FiN trained models was

39        significantly higher (82.33-88.59%) than that achieved by the models trained only on

40        camera trap datasets (38.5-66.74%).  The infusion of camera trap images into FiN training

41        further improved mAP, with increases ranging from 1.78-32.08%.

42    5.  Ecology researchers can use FiN images for training deep learning object detection

43        solutions for camera trap image processing to develop location invariant, robust, out-of-the-

44        box software. This would allow AI technologies to be deployed on a large scale in ecological

45        applications. Datasets and code related to this study are open source and available on this

46        repository: https://github.com/ashep29/infusion

47

48

## 1. Introduction

Automated survey methods such as camera trapping and passive acoustic monitoring are widely used in ecological research (Rovero and Zimmermann 2016, Sugai, Silva et al. 2018, Gibb, Browning et al. 2019). These methods provide invaluable insight into a plethora of ecological information including species occurrence, activity patterns and behavior (O'Connell, Nichols et al. 2011). However, they often result in the collection of large quantities of data, which must be processed, requiring a significant commitment of time and resources for manual or supervised classification (Swinnen, Reijniers et al. 2014, Young, Rode-Margono et al. 2018). Reducing the processing time and resources necessary for traditional data analysis such as manual analysis and citizen science (Swanson, Kosmala et al. 2015, Nguyen, Maclagan et al. 2017) has prompted increasing research into the adoption of Artificial Intelligence (AI) software in automated data classification (Falzon, Meek et al. 2014, Norouzzadeh, Nguyen et al. 2018, Willi, Pitman et al. 2018).

Object detector and image classifier software (models) have already been adopted to some extent in the processing of camera trap images (Yu, Jiangping et al. 2013, Gomez Villa, Salazar et al. 2016, Norouzzadeh, Nguyen et al. 2018, Willi, Pitman et al. 2018, Tabak, Norouzzadeh et al. 2019, Falzon, Lawson et al. 2020). These tools rely on data-driven deep learning to identify complex patterns which can be used for classification without feature engineering as described by (Miao, Gaynor et al. 2019). However, most solutions presented thus far have shown limited transferability to image data outside the domain of the training data (Beery, Van Horn et al. 2018, Willi, Pitman et al. 2018). This results in the need to develop models specific to each domain, however this process is time and resource intensive, requiring repeated collection and manual annotation of camera trap data, and

73    computationally expensive training of deep neural networks (Falzon, Lawson et al. 2020).

74    Thus, there is a clear need to develop location invariant object detectors, which are deep

75    learning models that can be transferred from one location to another, achieving acceptable

76    results without having to be retrained. Such out-of-the-box solutions are attractive due to

77    their potential for extensive application, particularly in circumstances where the

78    development of domain or study-specific models is prohibitively expensive or otherwise

79    unattainable.

80

81    Achieving location invariance requires training data to be characterized by high intra-

82    dataset variability. This is because neural networks learn patterns in data, meaning low

83    intra-dataset variability can result in learning of domain specific features such as camera

84    angle, lighting, and vegetation, reducing location invariance (Torralba and Sinha 2003, Miao,

85    Gaynor et al. 2019, Singh, Lindshield et al. 2020). Therefore, camera trap images must be

86    obtained from many sources to be able to train effective object detectors and classifiers.

87    However, the process of collecting camera trap images from an extensive network of

88    cameras from many domains is time and resource intensive and may be unfeasible for

89    smaller scale studies or those focusing on rare or elusive species. Even when researchers

90    have access to camera trap network, collecting enough images for training object detectors

91    can prove difficult. (Maurice 2019) deployed 15 cameras for 2 months resulting in the

92    collection of only 41 images of the pangolin (the target species), a number which would be

93    insufficient for effective neural network training (Shahinfar, Meek et al. 2020). Other factors

94    which limit the accessibility and availability of camera trap images include the reticence of

95    researchers to share existing camera trap data, or lack of data for novel species studies.

96

97    These limitations in data accessibility and availability limit the adoption of automated AI

98   solutions in ecological camera trap image processing (Schneider, Taylor et al. 2018). Thus,

99   alternative data sources must be identified and evaluated to assist in the development of

100  object detectors capable of being deployed in any domain, at any location, achieving

101  acceptable results regardless of camera trap image availability. Possible solutions include

102  publicly available sources of animal imagery, such as FlickR (flickr.com) and iNaturalist

103  (inaturalist.org). FlickR is a consumer photo sharing website, hosting approximately 10

104  billion images, shared by over 90 million monthly users. It is characterized by high intra-

105  dataset variability, high accessibility and a wide range of species types in highly varying

106  contexts, with minimal unintentional bias, as images are not collected for a specific purpose

107  (Everingham, Van Gool et al. 2010). It is arguably the most extensively used source of image

108  data in object detection benchmark datasets, including ImageNet (Deng, Dong et al. 2009),

109  MS COCO (Lin, Maire et al. 2014), the Open Images Dataset (Kuznetsova, Rom et al. 2020)

110  and PASCAL VOC (Everingham, Van Gool et al. 2010). iNaturalist contains over 45 million

111  observations of biodiversity data including both flora and fauna. Labelling of images on

112  iNaturalist may be more accurate than FlickR due to its purpose as a biodiversity data

113  sharing website and it does contain more camera trap images than FlickR. Other potential

114  image sources include Pinterest (www.pinterest.com), Imgur (www.imgur.com), pixabay

115  (www.pixabay.com) and 500px (www.web.500px.com). These image sources are highly

116  beneficial in training general, location invariant neural networks as they exhibit an

117  extensive range of contextual features, not necessarily present in camera trap imagery.

118

119  Despite their benefits as out-of-the-box solutions, universal or general object detectors

120  usually fail to achieve the high accuracy attainable by domain-specific object detectors

121  (Rebuffi, Bilen et al. 2017, Wang, Cai et al. 2019). Due to the need to achieve high accuracy

122  object detection and classification in ecological research, it may therefore be necessary to

123       optimize location invariant models for domain-specific studies. This is particularly relevant

124       when processing camera trap imagery characterized by features which differ strongly from

125       non-camera trap data, including infrared imagery, poor quality illumination and blurry

126       images.

127

128       Therefore the aims of this study are twofold:

      i)      To evaluate the use of publicly available image sources, in the development of location invariant camera trap object detectors.

      ii)      To develop an optimization strategy dubbed 'infusion' to improve the performance of location invariant object detectors in domain-specific applications.

129

130       In this study, we will demonstrate our proposed approach on three single class applications.

131       The rare species Striped Hyena (*Hyaena hyaena*) was chosen due to the sparsity of camera

132       trap training data, and the difficulty in discriminating between the striped hyena and the

133       more common spotted hyena. Furthermore, other studies have highlighted it as a species of

134       particular interest due to the difficulty they faced in detecting its presence in camera trap

135       images, for example, (Willi, Pitman et al. 2018) failed to detect any of the 27 striped hyenas

136       present in their test dataset. Next, the iconic and critically endangered Rhinoceros

137       (*Rhinocerotidae*) was also chosen, due to the high research interest in monitoring its

138       prevalence and changes in populations. Finally, the pest family *Suidae* (pigs, boars and hogs)

139       was included due to the significant role it plays across global ecosystems and its host status

140       for a range of diseases such as Swine Fever, which are a major threat to agricultural

141       industries.

## 2. Related Work

### a. Traditional Methods: Manual Analysis and Citizen Science

The majority of camera trap image processing is achieved by manual analysis conducted by ecologists, or via citizen science. Manual analysis involves the use of software programs to manually tag animals in images/capture events. Each image sequence or capture event is treated as a detection, and the ecologist must manually select a tag reflecting the identity of the animal. Once tagging is complete, a verification process is undertaken to identify and correct mistaken classifications. These tagged images can then be interrogated according to the purpose of the study, using tools such as R scripts, or specially developed GUI programs. Manual analysis of images is a significant resource demand on ecologists and research teams, requiring large expenditures in time and resources, hindering effective biodiversity management.

This time-consuming task may also be undertaken by citizen scientists, who are volunteers that contribute to scientific enquiry by collecting or processing image data (Nguyen, Maclagan et al. 2017). Large citizen science-based programs such as Zooniverse (www.zooniverse.org) enable the effective classification of millions of camera trap images (Jones, Allen et al. 2018). Citizen science projects have many benefits for researchers including customization of projects and annotation requirements in accordance with the aims of projects. However, the effectiveness of citizen science in rapidly processing large volumes of image data with sufficient accuracy is limited (Meek and Zimmerman 2016), causing large delays between the data collection and interpretation stages, which may be detrimental to ecological management (Fox, Bourn et al. 2019). Furthermore, the need to upload significant

167        amounts of data onto publicly accessible websites may pose privacy risks (Sagarra,

168        Gutiérrez-Roig et al. 2015) or poaching concerns and undermine the protection of

169        rare or endangered species by revealing their geographical location and behavioral

170        habits to poachers (Falzon, Lawson et al. 2020).

171

172        **b. Automated Image Processing Using Deep Learning**

173        Due to the shortcomings of traditional methods, research has centered primarily on

174        integration of automated image processing within camera trap research (Meek,

175        Fleming et al. 2014, Meek, Ballard et al. 2015, Fegraus and MacCarthy 2016, Willi,

176        Pitman et al. 2018, Young, Rode-Margono et al. 2018). To achieve this, neural

177        networks such as Deep Convolutional Neural Networks (DCNNs) are trained on

178        large amounts of annotated image data (thousands to millions of images) to

179        recognize discriminative features belonging to target classes (Zhao, Zheng et al.

180        2019). Handcrafted features specified by researchers are not used, instead the

181        features are 'learned' via updating of weights during training. When the DCNN is

182        confident in the presence of an object in an image, it maps bounding boxes,

183        segmentation masks, or classification labels to the image or object (Ren, He et al.

184        2015).  If a DCNN is very deep, consisting of many layers, it will have many trainable

185        parameters (usually millions) which gives rise to the need for large annotated image

186        datasets used in training these parameters from scratch. This is necessary for the

187        network to learn complex features (Samala, Chan et al. 2016). Although DCNNs can

188        be used to classify data with high accuracy, their usability can be limited by

189        insufficient training data which may lead to overfitting (memorization of training

190        data), and consequently, inability of the model to generalize to new data (Zhao

191        2017).

192

193     Early attempts at automated camera trap classification and object detection tasks

194     using neural networks were dependent on significant amounts of pre-processing

195     (Yu, Jiangping et al. 2013) and resulted in relatively poor accuracy (Swinnen,

196     Reijniers et al. 2014, Chen, Han et al. 2015). However, most modern solutions use

197     minimal pre-processing, or automate pre-processing (Giraldo Zuluaga, Salazar et al.

198     2017). Accuracy and recall attained by deep learning solutions is also increasing

199     significantly, as large annotated datasets become available and progress is achieved

200     in training methods, such as the adoption of transfer learning (Gomez Villa, Salazar

201     et al. 2016, Willi, Pitman et al. 2018). Transfer learning involves the repurposing of

202     learned features for another task (Yosinski, Clune et al. 2014). This allows general

203     features learned on a large, highly varied dataset such as ImageNet (Deng, Dong et

204     al. 2009) which contains 3.2 million images, or Snapshot Serengeti (Swanson,

205     Kosmala et al. 2015), which contains 7.3 million images to be transferred to a

206     smaller, similar dataset containing only hundreds to thousands of images. Transfer

207     learning has been shown to improve accuracy and the ability to generalize as well as

208     reducing training time and the quantity of data needed (Khan, Hon et al. 2019). Its

209     effectiveness in ecological camera trap applications has been established by

210     (Norouzzadeh, Nguyen et al. 2017) and (Willi, Pitman et al. 2018).

211

212     **c. Image Classification vs. Object Detection**

213     The majority of camera trap image processing solutions achieve image classification

214     rather than object detection (Gomez Villa, Salazar et al. 2016, Nguyen, Maclagan et

215     al. 2017, Norouzzadeh, Nguyen et al. 2017, Willi, Pitman et al. 2018, Miao, Gaynor et

216     al. 2019, Tabak, Norouzzadeh et al. 2019). Image classification is a process by which

217        a whole image is labeled as containing a given object, for example, if a pig is featured

218        in an image, it will be labelled 'pig. However, image classification is limited in

219        situations where an image contains more than one species, e.g. a pig and a

220        wildebeest (Schneider, Taylor et al. 2018). Object localization and counting is also

221        not effectively achieved by image classification and models tend to struggle to

222        distinguish between an empty frame and a small background object (Yousif, Yuan et

223        al. 2019). In contrast, object detection is the process of locating and identifying one

224        or more objects in an image. The model plots bounding boxes of varying

225        classification confidence and association class labels, around each object in an image

226        (see Figure 1 for comparison). It is more useful than image classification because it

227        allows more information to be extracted from the images, such as the number of

228        objects in an image, as well as information about reproduction, distribution,

229        quantification and comparison of behavior across individual animals within a

230        species group based on factors such as age and gender (Schneider, Taylor et al.

231        2018).

232

233        Another major benefit of object detection is the reduced impact of background and

234        environmental features on object classification. Unlike image classifiers, which learn

235        patterns in the entire image, object detectors only learn patterns within the

236        constraints of the bounding boxes, and actively negative sample on the image

237        background (area not included in the bounding boxes) (Wang, Hu et al. 2019, Zhao,

238        Zheng et al. 2019). This enables object detectors  to better generalize to new

239        domains, thus facilitating location invariance. Despite these benefits, object

240        detection necessitates a significantly higher expenditure of time and resources, due

241        to the need to annotate all training images with bounding boxes and labels.

242    Consequently, most studies achieve image classification rather than object

243    detection. In contrast, due to the major benefits provided by object detectors for

244    automated camera trap image processing , this study focuses on object detection

245    rather than image classification. For a more detailed overview of available image

246    classification methods, refer to Appendix S1.

247

248    Several studies have achieved object detection in the context of camera trap image

249    processing, however none have achieved location invariance, with testing using

250    restricted to in-sample datasets. (Yousif, Yuan et al. 2019) employed sequence-level

251    background subtraction using handcrafted Histogram of Oriented Gradient (HOG)

252    (Dalal and Triggs 2005) features to localize moving objects in camera trap images.

253    This study did not aim to identify animal species, instead simply distinguished

254    between humans and animals, and eliminated empty frames. Although it achieved

255    high accuracy in this task, its application was not extended beyond eastern North

256    America.

257

258    A novel ecological image processing software solution for use on a laptop by field

259    ecologists and wildlife managers was developed by (Falzon, Lawson et al. 2020). It

260    provides object detection and localization as well as species classification and object

261    counting capabilities via training of YOLOv2 DarkNet-19 (Redmon and Farhadi

262    2016) Deep Convolutional Neural Networks (DCNN) on both daytime and infrared

263    imagery. It boasts fast processing speeds and acceptable accuracy, achieved on a

264    local machine, within a dedicated on-demand application. Tailored models can be

265    applied to trap sites in Australia, New Zealand, North America, Serengeti and the

266    USA. However, optimal performance is only achieved when models are trained and

267      developed for a specific environment, camera trap imaging configuration and

268      species cohort. Thus, it suffers from lack of location invariance and robustness, as its

269      accuracy and recall decrease significantly when it is used outside the scope of the

270      environments on which it was trained.

271

272      (Schneider, Taylor et al. 2018) addressed the problem of object detection in camera

273      trap images, with the aim of identifying, quantifying and localizing animal species.

274      They used transfer learning to train a YOLOv2 model, achieving recall of 93% and

275      accuracy of 80.4% on the Reconyx (www.reconyx.com) and Snapshot Serengeti

276      (Swanson, Kosmala et al. 2015) datasets. The Reconyx dataset contained 946 images

277      of 20 species, while the Snapshot Serengeti dataset contained 4,097 images of 48

278      species. They also trained a Faster R-CNN model (Ren, He et al. 2015) achieving

279      76.7% recall and 72.2% accuracy. They used a model pretrained on the MS COCO

280      dataset (Lin, Maire et al. 2014) to initialize transfer learning. However, the

281      robustness of the model was not evaluated on out of sample images, which is

282      camera trap imagery obtained from traps and geographical locations not included in

283      the training data. It also suffered from class imbalance with lower accuracy and

284      recall for classes with fewer instances. Our research indicates this limitation can be

285      overcome by sourcing images from publicly available data sources.

286

287      **d. Improving Location Invariance via Dataset Construction**

288      The suboptimal performance and inability of neural networks to generalize to

289      contexts beyond the domain of the training data is a strong area of research interest.

290      As early as 2008, studies in contextual object detection examined the consequences

291      of 'unintentional regularities' in datasets resulting in object detectors learning

292          associations between objects and their backgrounds, inhibiting their ability to

293          detect objects out of context (Hoiem, Efros et al. 2008, Sudderth, Torralba et al.

294          2008). (Everingham, Van Gool et al. 2010) noted that classifiers tend to learn the

295          context of an object rather than model the appearance of the object. Thus, when the

296          object is dissociated with its context, the classifier fails to detect it due to extensive

297          use of image composition and context, resulting in a significant drop in

298          performance. These findings were confirmed by (Miao, Gaynor et al. 2019) in an

299          ecological context via the use of GRAD-CAM technology applied to models trained

300          solely on camera trap images, illustrating the tendency of neural networks to learn

301          background features as elements of an object if image background and context is not

302          highly varied. It is therefore essential to broaden the context of animal imagery to

303          extend beyond a restricted range of camera traps to ensure robustness and location

304          and context invariance.

305

306          This phenomena of contextual association was also found by (Everingham, Van Gool

307          et al. 2010) to be particularly prevalent in neural networks trained on images taken

308          by researchers for a specific purpose. Consistencies within datasets, such as camera

309          trap images collected within the context of a specific project, create an inner dataset

310          bias, which results in the development of models less capable of generalization to

311          other camera trap contexts. On this basis, we postulate that collection of camera

312          trap images for neural network training mimics collection of images under

313          laboratory or controlled conditions, whereby features such as lighting, camera

314          angle, distance of objects from the camera, and background features are consistent

315          across many images, thus encouraging contextual association. This is supported by

316          (Willi, Pitman et al. 2018) who noted that their models, trained on camera trap

317          images, would need to be retrained for use out of sample in other camera traps

318          which did not form part of the training set. In contrast, networks trained on data

319          sourced from consumer photo sharing websites such as FlickR are more capable of

320          generalization (Torralba and Efros 2011) due to the inherently high intra-dataset

321          variability and reduced likelihood of inner dataset bias.

322

## 3. Datasets and Annotation

324         The datasets used in this study were collated using images from FlickR and iNaturalist. We

325         also used camera trap image datasets obtained from www.lila.science including Snapshot

326         Serengeti (SS), Wildlife Conservation Society (WCS) Camera Traps, as well as other sites

327         specified in more detail below. All datasets, annotations, and the algorithms used for dataset

328         collection and processing, as well as auto-annotation of images are available here:

329         https://github.com/ashep29/infusion.

330

### a. FlickR and iNaturalist

332          We developed and used a Python script to download images from FlickR using the

333          FlickR API. This allowed us to download images with multiple keywords at once.

334          The keywords used are shown in Table 1. We downloaded a maximum of 200

335          images per keyword, to maximize the variety of search results. Our datasets were

336          restricted to Creative Commons images. We also developed a Python script to

337          download images from iNaturalist using a csv file containing URLs of relevant

338          observations downloaded from inaturalist.org.

339

340

**Table 1:** *Keyword searches used to download images from FlickR and iNaturalist.*

*Scientific names tended to return more accurately labelled images.*

| Rhinocerotidae | Hyaena hyaena | Suidae |
|---|---|---|
| *diceros AND bicornis* | *striped AND hyena* | *Phacochoerus AND africanus* |
| *ceratotherium AND simum* | *Hyaena AND* | *Sus AND scrofa* |
| *dicerorhinus AND* | *hyaena* | *sanglier* |
| *sumatrensis* | | *warthog OR warthogs* |
| *white AND rhinoceros* | | *wild AND pig OR boar OR hog* |
| *rhinoceros* | | *feral AND pig OR boar OR hog* |

Duplicates and near duplicates were removed using a Structural Similarity Index (SSIM) (Zhou, Bovik et al. 2004) clustering algorithm we developed (see Appendix S4). We deleted all images with a similarity score above 0.8, where a score of 1.0 represents a 100% similarity between 2 images. Near duplicates are images with strong visual similarity, containing only small distortions, slight variations and occlusions (Everingham, Van Gool et al. 2010). Interestingly, the datasets downloaded from FlickR and iNaturalist were mutually exclusive, with not one image present on one site, being also present on the other. Although this does not mean that images obtained from FlickR will not be available via iNaturalist, it does suggest that users of FlickR may often not be users of iNaturalist. Details about the final datasets are shown in Table 2. Subsamples of the final datasets are illustrated by Figure 1.

**Table 2:** *Final number of images obtained from FlickR and iNaturalist for both the single class and multi-class experiments, after duplicate removal and cleaning. Datasets are referred to hereon according to their source, abbreviated as FiN (FlickR-iNaturalist) and class name.*

| Dataset Name | Class | FlickR | iNaturalist | Total Images |
|---:|:---|:---|:---|:---|
| FiN_rhino | Rhino | 784 | 881 | 1665 |
| FiN_striped_hyena | Striped hyena | 401 | 71 | 472 |
| FiN_pig | Pig | 606 | 0 | 606 |



**Figure 1**: *Subsamples of the FiN datasets. Top to bottom: striped hyena, rhinoceros, and pig. Images of were highly varied, and included both color/daytime and infrared images, as well as a large range of contexts and distances from the camera.*

## b. Camera Trap Datasets

We obtained all camera trap data of rhinoceros and striped hyena from lila.science

370      using a Python script we developed, which we have made available on our GitHub

371      repository. We scoured all images of striped hyena and rhinoceros from both WCS

372      Camera Traps (*WCS_striped_hyena* and *WCS_rhino*) and Snapshot Serengeti

373      (*SS_striped_hyena* and *SS_rhino*) datasets (Swanson, Kosmala et al. 2015). We used

374      the same script to obtain our *EU_pig* and *NA_pig* datasets from the Missouri Camera

375      Traps (Zhang, He et al. 2016) and North American Camera Trap Images (Tabak,

376      Norouzzadeh et al. 2018) datasets respectively, also from lila.science. A summary of

377      all camera trap datasets is provided in Table 3.

378

379      **Table 3**: *Summary of the characteristics of the camera trap datasets used in this*

380      *study. The term 'quality' refers to characteristics such as blurriness, pixilation,*

381      *illumination etc. A poor-quality dataset will contain many images that are over or*

382      *underexposed, blurriness caused by poor focus, or other features which make it harder*

383      *to distinguish the identity of a target class and distort or damage key features. A visual*

384      *subsample of these datasets is provided (see Figure 2).*

| Dataset | Source | Location | Size | Characteristics |
|---------|--------|----------|------|-----------------|
| *WCS_striped_hyena* | Wildlife Conservation Society | Multiple | 582 | Moderate quality<br>Night and day |
| *SS_striped_hyena* | Snapshot Serengeti | Tanzania | 478 | Moderate quality<br>Infrared and day<br>Includes partials |
| *WCS_rhino* | Wildlife Conservation Society | Multiple | 333 | Low quality<br>Mostly infrared<br>Many partials |
| *SS_rhino* | Snapshot Serengeti | Tanzania | 153 | Moderate quality<br>Daytime<br>Many partials |
| *AU_pig* | Custom | NSW, Australia | 589 | Low quality<br>Mostly infrared<br>High occlusion<br>High density |
| *SS_pig* | Snapshot Serengeti | Tanzania | 574 | Moderate quality<br>Mostly daytime |

| CC_pig | Camera CATalogue | South Africa | 559 | Moderate quality Partials Low density |
|--------|------------------|--------------|-----|----------------------------------------|
| NA_pig | North America Camera Trap Images | United States | 514 | High quality |
| EU_pig | Missouri Camera Traps | Europe | 501 | Difficult High occlusion |

The *SS_pig* dataset is a subset of the Snapshot Serengeti dataset, and *CC_pig* is a

subset of the Camera CATalogue project conducted by Panthera

(www.panthera.org). Both are available from the Data Repository for the University

of Minnesota, used by (Willi, Pitman et al. 2018) and released under a CC0 1.0

Universal Public Domain Dedication license. The Australian pig dataset (AU_pig) is a

custom dataset, obtained during feral pig trapping and control operations. More

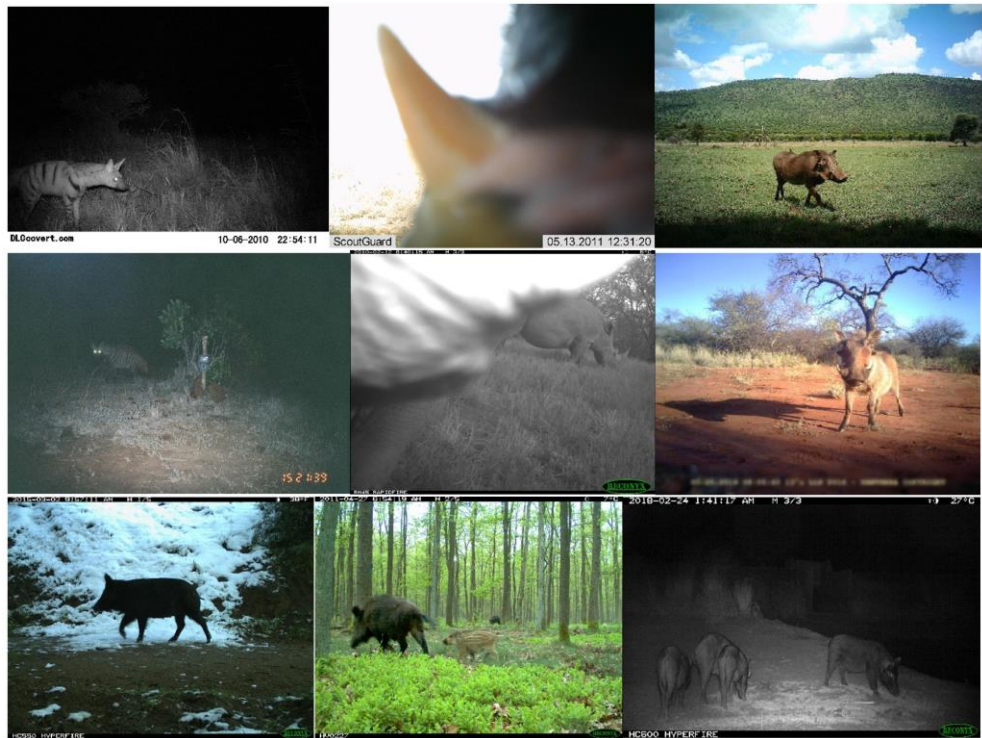information about each dataset is provided in Table 3, and a subset is shown in

Figure 2.

396 **Figure 2**: *Subsamples of the camera trap datasets. Top row: SS datasets, left to right;*

397 *striped hyena, rhino, and pig. Middle row: WCS datasets, left to right; striped hyena,*

398 *rhino, and pig. Bottom row: left; pig from NA_pig, middle; pigs from EU_pig and right;*

399 *pigs from AU_pig.*

400

401 Each image in the final datasets were annotated with bounding boxes and corresponding

402 class labels. Bounding box annotation involves the positioning of an axis aligned box

403 surrounding an object. We used an auto-annotator tool we developed to roughly annotate

404 all the images. We then edited any suboptimal bounding boxes using the graphical

405 annotation tool labelImg (Tzutalin 2015)[i] to ensure all objects were correctly annotated.

406 Annotations were saved in PASCAL VOC format.

407

408 ## 4. Training and Evaluation Methodology

409 In this study, we conducted two major experiments. Firstly, we compared the performance

410 of models trained on FlickR-iNaturalist (FiN) datasets only to those trained only on camera

411 trap data using evaluation on out of sample test sets. Next, we optimized the FiN models by

412 infusing small subsets of camera trap imagery into the FiN training set, evaluating

413 performance on out of sample test sets. Details about the model architecture and training

414 parameters are provided in Appendix S3. Additional information on transfer learning is also

415 provided. The experiments outlined in this section were also verified on a multi-class

416 application documented in Appendix S5.

417

418 ### a. Comparison between FiN and Camera Trap Data in Developing Location

419 ### Invariant Object Detectors

420 To evaluate the potential for publicly available data from FlickR and iNaturalist to

| 421 | be used in the development of location invariant object detectors for camera trap |
| 422 | image processing, we trained Keras-RetinaNet (Lin, Goyal et al. 2018) models on FiN |
| 423 | datasets, and compared their performance to that of RetinaNet models trained on |
| 424 | camera trap data when tested on out of sample camera trap images. |
| 425 | |
| 426 | We trained 3 single-class RetinaNet models on FiN datasets. These models are |
| 427 | referred to as *FiN_Classname*, e.g. *FiN_rhino* refers to a rhino detector trained on FiN |
| 428 | data. We also trained 2 single class (rhino and striped hyena) RetinaNet models |
| 429 | using the *WCS_striped_hyena* and *WCS_rhino* datasets, as well as 4 pig detectors, on |
| 430 | the *AU_pig*, *CC_pig*, *NA_pig* and *EU_pig* datasets. All models are named based on the |
| 431 | source of their training data. Note, we were able to train 4 pig models due to greater |
| 432 | availability of data when compared with rare species such as rhino and striped |
| 433 | hyena. |
| 434 | |
| 435 | The datasets were randomly split into training and validation sets, with 90% of |
| 436 | images reserved for training, and 10% used for validation. Each training set was |
| 437 | supplemented with 800 explicit negative samples to improve discrimination |
| 438 | between target species and non-target species or background. A detailed |
| 439 | breakdown of the training and validation splits as well as the out of sample test set |
| 440 | is provided in Table 4. |
| 441 | |
| 442 | **Table 4**: *Data distribution for models trained on datasets obtained from* |
| 443 | *FlickR/iNaturalist, abbreviated as FiN (FlickR-iNaturalist), and models trained using* |
| 444 | *camera trap images alone abbreviated as follows; WCS (Wildlife Conservation* |
| 445 | *Society), AU (Australia), NA (North America), CC (Camera CATalogue) and EU* |

446 *(Europe). All models were tested on out of sample images obtained from Snapshot*

447 *Serengeti.*

448

449

| Models | Training set (90%) | Validation set (10%) | Out of Sample Test set (SS) |
|---|---|---|---|
| *FiN_striped_hyena* | 425 | 47 | 478 |
| *WCS_striped_hyena* | 524 | 58 | |
| *FiN_rhino* | 1499 | 166 | 153 |
| *WCS_rhino* | 300 | 33 | |
| *FiN_pig* | 545 | 61 | 574 |
| *AU_pig* | 530 | 59 | |
| *CC_pig* | 503 | 56 | |
| *NA_pig* | 463 | 51 | |
| *EU_pig* | 451 | 50 | |

450

451

452

453

454

455

456 All models were tested using out of sample images from the Snapshot Serengeti (SS)

457 datasets, i.e. *SS_striped_hyena*, *SS_rhino* and *SS_pig*. Each test set was supplemented

458 with 200 negative samples to prevent biased evaluation of false positives. These

459 negative samples were derived from the Snapshot Serenget, and consisted of empty

460 images, or images of non-target species.  For more information relating to the

461 negative sampling data collection process, refer to Appendix S2.

462

463 **b. Infusion: Optimization of Location Invariant Models Using Camera Trap**

464 **Imagery**

465 Next, we conducted experiments to evaluate an optimization process that would

466 allow ecologists to improve object detection performance with minimal infusion of

467 camera trap images into the FiN training set. Infusion is the process of

468 supplementing the training set with a small subset of camera trap images, to

469 improve robustness to the particularities of camera trap data, such as infrared, high

470        occlusion, blurriness etc. Infusion was conducted both out of sample and in-sample.

471        Out of sample results are presented in this manuscript. For in-sample results, refer

472        to Appendix S6.

473

474        Due to the large number of highly similar images present within camera trap

475        datasets, the infusion subsets were not randomly selected. Instead, our SSIM

476        algorithm was used to retain only images with low SSIM scores, with the aim of

477        maximizing intra-dataset variability. The SSIM algorithm allowed us to randomly

478        select one frame from each cluster of images (usually one capture event, or different

479        capture events with very similar properties).

480



481

482        **Figure 3**: *Graphical illustration of image clustering using an SSIM algorithm. The test*

483        *image represented by 1.0 is compared with every other image. Highly dissimilar*

484        *images have low SSIM scores (<0.4).*

485

486          Our research indicates that image pairs with an SSIM value above 0.4 have

487          sufficiently high similarity to be clustered. For example, Figure 3 illustrates the

488          output of the SSIM algorithm graphically, clearly showing the three clusters formed

489          by visually similar images, the image denoted by the arrow (the test image) is

490          compared to each other image, with values closest to 1 indicating high similarity

491          with the test image, This method allows researchers to compile highly varied

492          datasets automatically, minimizing the need for extensive time-consuming image

493          sorting and annotation.

494

495          Out of sample infusion was conducted by training 4 additional models for each

496          species, with incremental infusion of the SSIM sorted camera trap images from the

497          WCS and CC datasets into the FiN training data. These images were added in

498          increments of 5% from 5-20%, as shown by Table 5. For example, the *FiN_rhino*

499          dataset comprised of 1665 images. To achieve 5% infusion, 83 images from the

500          *WCS_rhino* dataset were added to the *FiN_rhino* dataset. 90% of these images were

501          retained for training, with 10% reserved for monitoring training via the validation

502          set. This process was repeated for all percentiles and species shown in Table 5.

503

504    **Table 5**: *Incremental infusion of camera trap images into FiN training. An additional*

505    *800 negative samples were included in the training set. Models are named according*

506    *to the class name and infusion percentile. Note the infusion images are trap images.*

507    *The infusion training set is made up of FiN + infusion images. The validation set is FiN*

508    *validation + infusion images.*

| Class | Model name | Infusion Source | Nº infusion images | Infusion training set | Infusion Validation set |
|-------|-----------|-----------------|--------------------|-----------------------|-------------------------|
| Hyaena | hyaena_inf_05 | WCS_hyena | 24 | 446 | 50 |
|        | hyaena_inf_10 |           | 47 | 467 | 52 |
|        | hyaena_inf_15 |           | 71 | 489 | 54 |
|        | hyaena_inf_20 |           | 94 | 509 | 57 |
| Rhino | rhino_inf_05 | WCS_rhino | 83 | 1573 | 175 |
|       | rhino_inf_10 |           | 167 | 1649 | 183 |
|       | rhino_inf_15 |           | 250 | 1723 | 192 |
|       | rhino_inf_20 |           | 333 | 1798 | 200 |
| Pig | pig_inf_05 | CC_pig | 30 | 572 | 64 |
|     | pig_inf_10 |        | 61 | 600 | 67 |
|     | pig_inf_15 |        | 91 | 627 | 70 |
|     | pig_inf_20 |        | 121 | 654 | 73 |

509

510    The models were then tested on the out of sample Snapshot Serengeti test sets

511    presented in Section 4(a). Both the training and test sets were supplemented with

512    negative samples as described in Section 4(a).

513

514    **c. Model Evaluation**

515    To evaluate the performance of our models, mean Average Precision (mAP) results

516    will be provided. mAP is calculated as documented in the PASCAL VOC benchmark

517    (Everingham, Van Gool et al. 2010). A high mAP indicates that the model is detecting

518    the majority of objects with high accuracy, and minimal retention of false positives.

519    Accuracy is measured using Intersection over Union (IoU), which is a measure of the

520        overlap between the detection box and the ground truth bounding box.

521

## 5. Results

### a. Comparison between FiN and Camera Trap Data in Developing Location Invariant Object Detectors

The results of training on FiN data compared with training on camera trap data are presented in Figure 4. All results were collected on the out of sample Snapshot Serengeti test sets. The models trained on FiN datasets achieved mAP results ranging between 82.33% and 88.59%, while the models trained on camera trap data achieved mAP results ranging from 38.5% to 66.74%. In all cases, the FiN models outperformed the models trained on camera trap images.



**Figure 4:** *Comparison of the mAP results achieved by the models trained on FiN data, and those trained on camera trap datasets. In all cases, the FiN models outperformed*

535    *the camera trap models.*

536

537    The *FiN_pig* model achieved a mAP of 88.59% when tested on the out of sample

538    *SS_pig* dataset. This was far superior to the *CC_pig* model, which was trained on

539    camera trap images of warthogs from the Camera CATalogue  (CC) dataset,

540    achieving a mAP of only 53.87%. Although both the *CC_pig* dataset and the *SS_pig*

541    dataset contained the same subspecies (*Phacochoerus africanus*), the *CC_pig* model

542    did not generalize well to the *SS_pig* test set. This may be because the *SS_pig* dataset

543    was characterized by more variation in background, greater variation in the

544    distance of pigs from the camera and greater contrast. Notably, the worst

545    performing pig model was trained on data from Australia (*AU_pig*). This is very

546    likely due to the large number of low quality infrared images present in the training

547    data, which encouraged the model to return a high rate of false positives, and the

548    large disparity between contextual features such as vegetation and species type (the

549    Australia subspecies was *Sus scrofa*, while the SS subspecies was *Phacochoerus*

550    *africanus*).

551

552    In comparison, the significantly greater intra-dataset variability present in the FiN

553    datasets allowed for better model generalization when compared to the models

554    trained only on single location camera trap data. This trend was observed across all

555    classes, with the *FiN_striped_hyena* and *FiN_rhino* models significantly

556    outperforming the *WCS_striped_hyena* and *WCS_rhino* models.

557

558

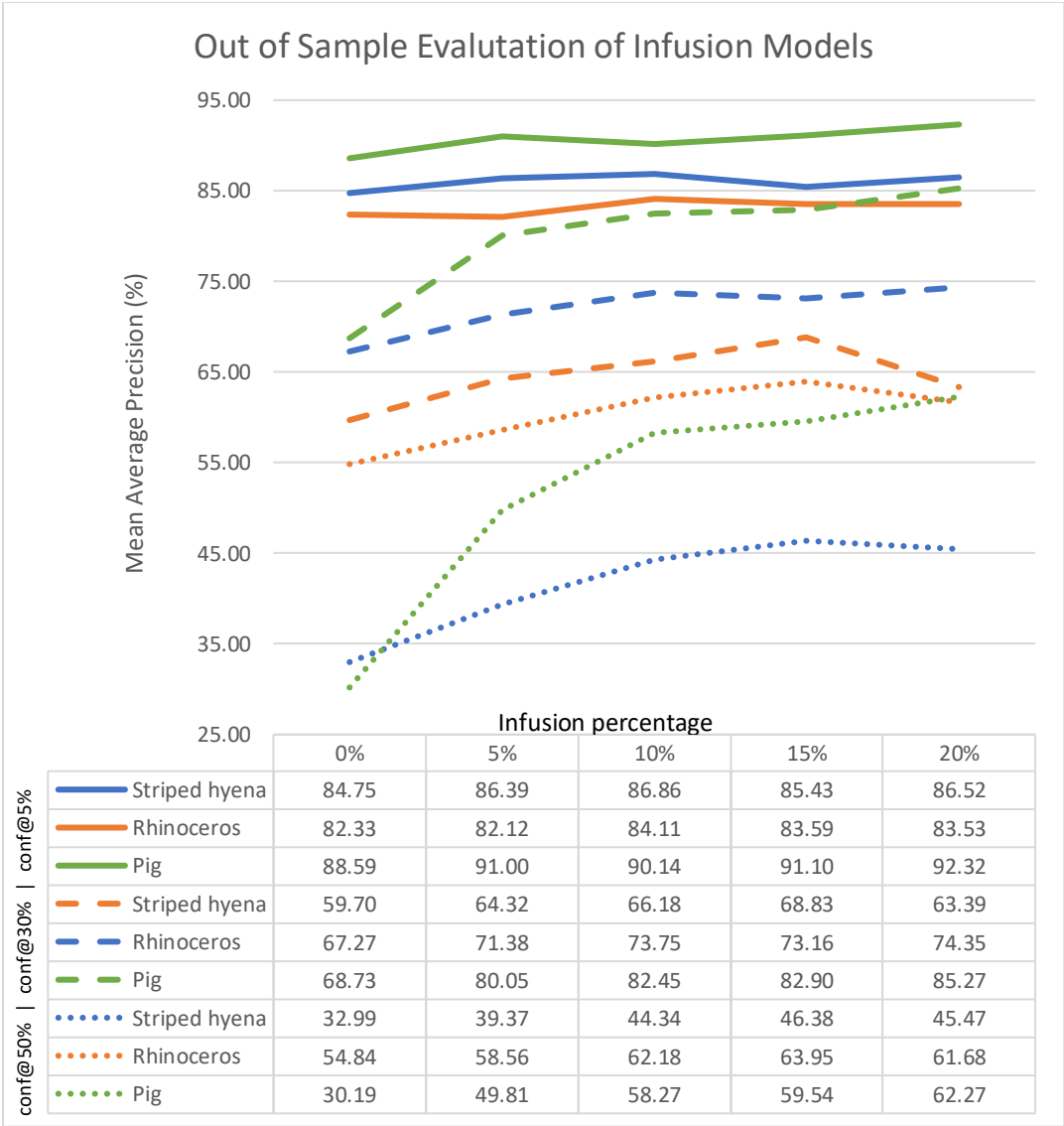| 559 | **b. Infusion: Optimization of Location Invariant Models Using Camera Trap** |
|---|---|
| 560 | **Imagery** |
| 561 | The results presented in the previous section indicate that the models trained on |
| 562 | FiN datasets can be used to effectively process images collected at any camera trap |
| 563 | site with an acceptable level of location invariance. However, camera trap images |
| 564 | possess particular characteristics which differentiate them from FiN images. In |
| 565 | difficult cases, the mAP achieved by FiN models may not be sufficiently high for |
| 566 | practical purposes, particularly when higher confidence thresholds are used, for |
| 567 | example, for a given study, the confidence threshold may be set to 50%, meaning all |
| 568 | detections with a classification score lower than 50% would be ignored. Thus, we |
| 569 | present the results of our infusion optimization experiments, illustrated by Figure 5. |
| 570 | In all cases, infusion resulted in an increase in mAP when evaluated on out of |
| 571 | sample images. |
| 572 | |

## Out of Sample Evalutation of Infusion Models

**Mean Average Precision (%)** vs **Infusion percentage**

|  |  | 0% | 5% | 10% | 15% | 20% |
|---|---|---|---|---|---|---|
| conf@5% | Striped hyena | 84.75 | 86.39 | 86.86 | 85.43 | 86.52 |
| | Rhinoceros | 82.33 | 82.12 | 84.11 | 83.59 | 83.53 |
| | Pig | 88.59 | 91.00 | 90.14 | 91.10 | 92.32 |
| conf@30% | Striped hyena | 59.70 | 64.32 | 66.18 | 68.83 | 63.39 |
| | Rhinoceros | 67.27 | 71.38 | 73.75 | 73.16 | 74.35 |
| | Pig | 68.73 | 80.05 | 82.45 | 82.90 | 85.27 |
| conf@50% | Striped hyena | 32.99 | 39.37 | 44.34 | 46.38 | 45.47 |
| | Rhinoceros | 54.84 | 58.56 | 62.18 | 63.95 | 61.68 |
| | Pig | 30.19 | 49.81 | 58.27 | 59.54 | 62.27 |

**Figure 5:** *Results of the infusion experiments on the out of sample SS test set. Infusion resulted in improvement across all models, particularly when evaluated at higher confidence thresholds. Infusion of 5% significantly improves performance, however optimum performance occurs at 10-15%, with the mAP results plateauing beyond 15%.*

At a confidence threshold of 5% (the standard threshold for mAP measurement (Lin, Goyal et al. 2018)), out of sample infusion did not result in a pronounced

| 582 | improvement, with gains in mAP results ranging from 1.78-3.73%. However, in |
| 583 | practical deployment, a confidence threshold of 5% would rarely be used, with |
| 584 | ecologists favoring higher thresholds to ensure confident classification of species. It |
| 585 | is at these higher thresholds that the benefits of infusion are best demonstrated. For |
| 586 | example, at a confidence threshold of 30%, the mAP improved by 7.08-16.54%, |
| 587 | while at a confidence threshold of 50% it improved by 9.11- 32.08%. It is well |
| 588 | established that increasing the confidence threshold decreases recall (the number of |
| 589 | true positives retained in the final output), and consequently decreases mAP (Willi, |
| 590 | Pitman et al. 2018). Note, we did not conduct evaluations of the models at |
| 591 | confidence thresholds above 50% because almost all detections with scores above |
| 592 | 50% were true positives, which meant increasing the threshold simply removed |
| 593 | true positives. Selecting a confidence threshold for a given application is highly |
| 594 | dependent on the quality of training data, extent of negative sampling and the model |
| 595 | used. The supplementation of FiN training with out of sample camera trap imagery |
| 596 | is therefore highly beneficial as it allows more true positives to be retained, because |
| 597 | the overall confidence of correctly detected objects is improved. This is a result of |
| 598 | the improved robustness to the particularities of camera trap imagery. |
| 599 | |
| 600 | The results presented in Figure 5 indicate that the addition of a small percentage of |
| 601 | camera trap images into the FiN training dataset can significantly improve |
| 602 | performance. In most cases, the greatest improvement occurred with infusion of |
| 603 | 5%, with performance continuing to improve as infusion was increased to 15%. As |
| 604 | infusion was increased beyond 15%, performance plateaued, or decreased, with |
| 605 | only 4 out of 9 results improving beyond 15%. |
| 606 | |

## 6. Discussion

We investigated the use of FiN images as an alternative to camera trap images in the task of DCNN training for location invariant camera trap image processing tasks, on three case studies, namely striped hyena, rhinoceros and pig. Specifically, we established the greater transferability of the FiN trained models when compared to models trained on camera trap datasets, and their high usability as location invariant object detectors. We then demonstrated how such models can be optimized via out of sample infusion, which was shown to increase the confidence of detections, allowing more true positives to be retained at higher confidence thresholds.

Our results show that FiN training significantly improves model robustness and location invariance. Particularly, it provides ecologists with a practical, cost effective, out of the box solution, capable of detecting animals even in the most challenging camera trap environments. We not only established that FiN data alone can be used to achieve good results, but these models can be improved with minimal infusion of camera trap data to improve robustness to the particularities of camera trap imagery. This suggests that ecologists can train object detectors using FiN imagery, and if camera trap data is available for their target species, use it to infuse the FiN training data. This model can then be used to process out of sample images from any camera trap, achieving a sufficiently high mAP to be deployed in most applications.

Furthermore, in circumstances where model performance is still considered suboptimal, they may then infuse the model with in-sample camera trap images, for further optimization. Although in-sample infusion makes the model more location variant, it does provide a means by which ecologists can train powerful models capable of achieving results

632     in the 90th percentile, with very few training images, as demonstrated by the results of in-

633     sample infusion presented in Appendix S6. As demonstrated by various studies in

634     automated camera trap image processing, achieving robust object detectors via training

635     solely on camera trap images usually requires thousands to millions of images

636     (Norouzzadeh, Nguyen et al. 2017, Willi, Pitman et al. 2018, Tabak, Norouzzadeh et al.

637     2019). In-sample infusion overcomes this requirement by leveraging off the robustness of

638     the FiN model, and the strong availability of FiN imagery to allow ecologists to train high

639     accuracy optimized deep leaning models with very few camera trap images, significantly

640     reducing the time and resources necessary to develop automated deep leaning object

641     detectors.

642

643     In light of the growing number of camera trap based projects undertaken by ecologists, this

644     research provides an invaluable method by which researchers can process extensive image

645     data regardless of the location from which the images were obtained, and the particularities

646     of the camera trap site or species. This method has been proven on several species,

647     including rare species, for which camera trap data for training models is often sparse. As

648     illustrated by (Willi, Pitman et al. 2018), the lack of camera trap data for rare species poses

649     significant problems when training multi-class object detectors, as the large class imbalance

650     between common species and rare species causes object detectors to misclassify species, by

651     over enthusiastically classifying species based on how common they are in the dataset

652     rather than via their features. This was observed by (Willi, Pitman et al. 2018) who noted

653     that insufficient images of the rare striped hyena in their dataset resulted in their model

654     achieving a mAP of 0% on this class. We have specifically addressed this problem by

655     proposing the use of FiN images of striped hyena to rectify limitations in data availability.

656

657 The use of FlickR as the principal training data also rectifies another major problem faced

658 by researchers. Studies have indicated that deep learning models have a tendency to return

659 overly confident predictions (Willi, Pitman et al. 2018) when trained on camera trap data

660 and deployed in-sample. This is due to the high consistency in image quality, lighting,

661 camera angle and geographical and vegetation features in camera trap data. Furthermore,

662 many trap images feature obscured or poor quality imagery of animals which if used in the

663 training set, may cause the network to make unrealistically optimistic predictions, by

664 attributing 100% confidence to visual features which may not display sufficiently distinct

665 characteristics present solely in the target class. In contrast, the higher resolution of FiN

666 images and large variations between images forces the model to reduce the confidence

667 attributed to poor quality or obscured animals. Their greater robustness allows them to be

668 deployed out of sample, further minimizing this problem.

669

670 One potential benefit in using FiN imagery for training image processing models is the high

671 availability of already annotated animal images. Because FlickR is a major source of images

672 used in datasets such as ImageNet and MS COCO, many animal classes have already been

673 annotated with bounding boxes, which are freely available for downloading. Using the

674 method proposed in the paper would therefore significantly reduce the time and resource

675 expenditure necessary for model development, by leveraging off the work already

676 completed by the broader object detection community. We were unable to use annotated

677 FlickR images from ImageNet as it was under maintenance, however it may prove to be a

678 valuable resource in the development of future models. This study was limited to the

679 evaluation of FlickR and iNaturalist images, and did not evaluate alternative images sources

680 mentioned in Section 1.

681

682   This research did not investigate the application of the FiN and infusion training method

683   using alternative object detectors such as YOLO (Redmon and Farhadi 2016), and Faster R-

684   CNN (Ren, He et al. 2015). Applying the findings of this study to these architectures may be

685   beneficial. YOLO is a faster, more efficient object detector, which may be more suited to

686   video processing, while Faster RCNN generally achieves higher accuracies, but is slower.

687   RetinaNet was chosen as it achieves a good balance between the computational efficiency of

688   YOLO and the accuracy of Faster-RCNN, which made it an appropriate choice for the difficult

689   task of camera trap image processing. In this study, we have only demonstrated location

690   invariance using RetinaNet. Although it goes beyond the scope of this study, it would be

691   interesting to ascertain whether changes in model architecture would influence the

692   robustness of location invariance models. Another possible area of research could be the

693   application of this method to object segmentation-based image processing. Object

694   segmentation builds upon the benefits of object detection by excluding background

695   features. This limits the influence of contextual features on model performance, thus

696   improving model accuracy and overall performance, however it is likely that they would

697   encounter the same modelling bias faced by bounding box-based object detection models.

698

699   One limitation of this study is that it only evaluates the models in terms of the Snapshot

700   Serengeti dataset. We could only evaluate on one dataset for the classes 'striped hyena' and

701   'rhinoceros' due to lack of data availability. To maintain consistency, we also only presented

702   results for the class 'pig' on Snapshot Serengeti in this manuscript. However, to verify the

703   usability of this method at any location and for any dataset, we present more extensive

704   results in Appendix S7 for the class pig, for which we had more data available, thus showing

705   location invariance across 4 extra test locations.

706

707 Finally, the proposed method may be extended to other image modalities. For example, it

708 could be extended to drone imagery (Kellenberger, Volpi et al. 2017, Xu, Wang et al. 2020).

709 Drone images are often captured from an aerial perspective, meaning they would contain

710 quite different features to those present available on FlickR. Applying our findings to object

711 detection in the context of drone imagery would be interesting, particularly with infusion of

712 a small subset of drone images to boost performance and allow better generalization to the

713 particularities of drone imagery. This would determine how transferable FiN images are to

714 new modalities. It could also be extended to other applications such as underwater animal

715 imagery (Dawkins, Sherrill et al. 2017, Christensen, Mogensen et al. 2018), surveillance

716 footage (Raghunandan, Mohana et al. 2018), and thermal camera imagery (Rodin, Lima et al.

717 2018, Bondi, Jain et al. 2020). This may present opportunities to rectify image shortages, or

718 problems with low intra-dataset variability, particularly in novel studies.

719

720

721 **7. Conclusion**

722 This study successfully demonstrated the use of FiN datasets in training location invariant

723 deep learning object detection models in the task of camera trap image processing. It also

724 evaluated an optimization process dubbed infusion, to improve robustness to the

725 particularities of camera trap imagery. Results presented across three single class models

726 on out of sample test sets indicate the aims of this study have been achieved. However, our

727 approach is limited by its inability to achieve high precision out of sample object detection,

728 which is still best achieved via in-sample training or infusion. Furthermore, this method was

729 not evaluated on alternative object detection frameworks and did not provide findings on

730 an extensive multi-class dataset. Nevertheless, this study provides a promising pathway to

731 develop robust, location invariant models using publicly accessible data sources.

Furthermore, development of these models will facilitate the widespread deployment of AI

in ecological management. The findings of this study could also be extended beyond camera

trapping to other object detection tasks and image modalities such as drone imagery.

Furthermore, the methodology of using transfer learning and publicly available datasets

characterized by high intra-dataset variability and minimal unintentional bias to train

location and context invariant AI-based data processing software could be extended beyond

images to other forms of data.

## 8. Acknowledgements

## 9. Author Contributions

Andrew Shepley, Greg Falzon and Paul Kwan conceived the ideas and designed

methodology; Andrew Shepley, Paul Meek and Greg Falzon collected the data; Andrew

Shepley and Greg Falzon analysed the data; Andrew Shepley led the writing of the

manuscript. All authors contributed critically to the drafts and gave final approval for

757     publication.

758

## 10.      Data Accessibility

760     *Image and Annotation datasets:* All image datasets and corresponding annotations

761     will be made available via Dryad upon acceptance.

762     *Code and scripts:* All code and scripts will be made available via Dryad upon

763     acceptance.

764

765

766

## 11.      Reference List

768 Beery, S., G. Van Horn and P. Perona (2018). Recognition in Terra Incognita, Cham, Springer

769 International Publishing.

770 Bondi, E., R. Jain, P. Aggrawal, S. Anand, R. Hannaford, A. Kapoor, J. Piavis, S. Shah, L. Joppa, B.

771 Dilkina and M. Tambe (2020). BIRDSAI: A Dataset for Detection and Tracking in Aerial Thermal

772 Infrared Videos. 2020 IEEE Winter Conference on Applications of Computer Vision (WACV).

773 Chen, G., T. Han, Z. He, R. Kays and T. Forrester (2015). "Deep convolutional neural network based

774 species recognition for wild animal monitoring." 2014 IEEE International Conference on Image

775 Processing, ICIP 2014: 858-862.

776 Christensen, J. H., L. V. Mogensen, R. Galeazzi and J. C. Andersen (2018). Detection, Localization and

777 Classification of Fish and Fish Species in Poor Conditions using Convolutional Neural Networks.

778 2018 IEEE/OES Autonomous Underwater Vehicle Workshop (AUV).

779 Dalal, N. and B. Triggs (2005). "Histograms of Oriented Gradients for Human Detection." IEEE

780 Conference on Computer Vision and Pattern Recognition (CVPR 2005) **2**.

781  Dawkins, M., L. Sherrill, K. Fieldhouse, A. Hoogs, B. Richards, D. Zhang, L. Prasad, K. Williams, N.

782  Lauffenburger and G. Wang (2017). An Open-Source Platform for Underwater Image and Video

783  Analytics. 2017 IEEE Winter Conference on Applications of Computer Vision (WACV).

784  Deng, J., W. Dong, R. Socher, L.-J. Li, K. Li and F. F. Li (2009). "ImageNet: a Large-Scale Hierarchical

785  Image Database." IEEE Conference on Computer Vision and Pattern Recognition: 248-255.

786  Everingham, M., L. Van Gool, C. Williams, J. Winn and A. Zisserman (2010). "The Pascal Visual Object

787  Classes (VOC) challenge." International Journal of Computer Vision **88**: 303-338.

788  Falzon, G., C. Lawson, K.-W. Cheung, K. Vernes, G. A. Ballard, P. J. S. Fleming, A. S. Glen, H. Milne, A.

789  Mather-Zardain and P. D. Meek (2020). "ClassifyMe: A Field-Scouting Software for the Identification

790  of Wildlife in Camera Trap Images." Animals **10**(1): 58.

791  Falzon, G., P. D. Meek and K. Vernes (2014). Computer Assisted identification of small Australian

792  mammals in camera trap imagery. Camera Trapping: Wildlife Management and Research. Paul

793  Meek, Peter Fleming, Guy Ballard et al. Melbourne, Australia, CSIRO Publishing**:** 299-306.

794  Fegraus, E. H. and J. MacCarthy (2016). Camera Trap Data Management and Interoperability.

795  Camera Trapping for Wildlife Research. F. R. a. F. Zimmerman. Exeter UK, Pelagic Publishing**:** 33-42.

796  Fox, R., N. Bourn, E. Dennis, R. Heafield, I. Maclean and R. Wilson (2019). "Opinions of citizen

797  scientists on open access to UK butterfly and moth occurrence data." Biodiversity and Conservation.

798  Gibb, R., E. Browning, P. Glover-Kapfer and K. E. Jones (2019). "Emerging opportunities and

799  challenges for passive acoustics in ecological assessment and monitoring." Methods in Ecology and

800  Evolution **10**(2): 169-185.

801  Giraldo Zuluaga, J., A. Salazar, A. Gomez Villa and A. Diaz-Pulido (2017). "Automatic Recognition of

802  Mammal Genera on Camera-Trap Images using Multi-Layer Robust Principal Component Analysis

803  and Mixture Neural Networks."

804    Gomez Villa, A., A. Salazar and J. Vargas-Bonilla (2016). "Towards Automatic Wild Animal

805    Monitoring: Identification of Animal Species in Camera-trap Images using Very Deep Convolutional

806    Neural Networks." Ecological Informatics **41**.

807    Hoiem, D., A. Efros and M. Hebert (2008). "Putting Objects in Perspective." International Journal of

808    Computer Vision **80**: 3-15.

809    Jones, F., C. Allen, C. Arteta, J. Arthur, C. Black, L. Emmerson, R. Freeman, G. Hines, C. Lintott, Z.

810    Macháčková, G. Miller, R. Simpson, C. Southwell, H. Torsey, A. Zisserman and T. Hart (2018). "Time-

811    lapse imagery and volunteer classifications from the Zooniverse Penguin Watch project." Scientific

812    Data **5**: 180124.

813    Kellenberger, B., M. Volpi and D. Tuia (2017). Fast animal detection in UAV images using

814    convolutional neural networks. 2017 IEEE International Geoscience and Remote Sensing

815    Symposium (IGARSS).

816    Khan, N., M. Hon and N. Abraham (2019). "Transfer Learning with intelligent training data selection

817    for prediction of Alzheimer's Disease."

818    Kuznetsova, A., H. Rom, N. Alldrin, J. Uijlings, I. Krasin, J. Pont-Tuset, S. Kamali, S. Popov, M. Malloci,

819    A. Kolesnikov, T. Duerig and V. Ferrari (2020). "The Open Images Dataset V4." International Journal

820    of Computer Vision **128**(7): 1956-1981.

821    Lin, T.-Y., P. Goyal, R. Girshick, K. He and P. Dollar (2018). "Focal Loss for Dense Object Detection."

822    IEEE Transactions on Pattern Analysis and Machine Intelligence **PP**: 1-1.

823    Lin, T.-Y., M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár and C. Zitnick (2014).

824    "Microsoft COCO: Common Objects in Context." **8693**.

825    Maurice, M. E. (2019). A Survey on the Status of Pangolins By Camera Trapping in Deng-Deng

826    National Park, Eastern Region, Cameroon, Ommega Internationals.

827    Meek, P. D., G.-A. Ballard and P. J. S. Fleming (2015). "The pitfalls of wildlife camera trapping as a

828    survey tool in Australia." Australian Mammalogy **37**(1): 13-22.

829    Meek, P. D., P. Fleming, A. G. Ballard, P. B. Banks, A. W. Claridge, S. McMahon, J. Sanderson and D. E.

830    Swann (2014). Putting contemporary camera trapping in focus. Camera Trapping in Wildlife

831    Research and Management. B. PD Meek, A. G., Banks, P. B., Claridge, A. W., Fleming, P. J. S.,

832    Sanderson, J. G., and Swann, D. Melbourne, Victoria, CSIRP Publishing**:** 349-356.

833    Meek, P. D. and F. Zimmerman (2016). Camera Traps and Public Engagement. Camera Trapping for

834    Wildlife Research. F. a. Z. Rovero, F. Exeter, Pelagic Publishing UK**:** 219-231.

835    Miao, Z., K. Gaynor, J. Wang, Z. Liu, O. Muellerklein, M. S. Norouzzadeh, A. McInturff, R. Bowie, R.

836    Nathan, S. Yu and W. Getz (2019). "Insights and approaches using deep learning to classify wildlife."

837    Scientific Reports **9**.

838    Nguyen, H., S. Maclagan, T. Nguyen, T. Nguyen, P. Flemons, K. Andrews, E. Ritchie and D. Phung

839    (2017). "Animal Recognition and Identification with Deep Convolutional Neural Networks for

840    Automated Wildlife Monitoring."

841    Norouzzadeh, M. S., A. Nguyen, M. Kosmala, A. Swanson, C. Packer and J. Clune (2017).

842    "Automatically identifying wild animals in camera trap images with deep learning." Proceedings of

843    the National Academy of Sciences **115**.

844    Norouzzadeh, M. S., A. Nguyen, M. Kosmala, A. Swanson, M. S. Palmer, C. Packer and J. Clune (2018).

845    "Automatically identifying, counting, and describing wild animals in camera-trap images with deep

846    learning." Proceedings of the National Academy of Sciences **115**(25): E5716.

847    O'Connell, A., J. D. Nichols and K. U. Karanth (2011). Camera traps in animal ecology: Methods and

848    analyses.

849    Raghunandan, A., Mohana, P. Raghav and H. V. R. Aradhya (2018). Object Detection Algorithms for

850    Video Surveillance Applications. 2018 International Conference on Communication and Signal

851    Processing (ICCSP).

852    Rebuffi, S.-A., H. Bilen and A. Vedaldi (2017). "Learning multiple visual domains with residual

853    adapters." 506--516.

854    Redmon, J. and A. Farhadi (2016). "YOLO9000: Better, Faster, Stronger."

855    Ren, S., K. He, R. Girshick and J. Sun (2015). "Faster R-CNN: Towards Real-Time Object Detection

856    with Region Proposal Networks." IEEE Transactions on Pattern Analysis and Machine Intelligence

857    **39**.

858    Ren, S., K. He, R. Girshick, X. Zhang and J. Sun (2015). "Object Detection Networks on Convolutional

859    Feature Maps." IEEE Transactions on Pattern Analysis and Machine Intelligence **39**.

860    Rodin, C. D., L. N. d. Lima, F. A. d. A. Andrade, D. B. Haddad, T. A. Johansen and R. Storvold (2018).

861    Object Classification in Thermal Images using Convolutional Neural Networks for Search and

862    Rescue Missions with Unmanned Aerial Systems. 2018 International Joint Conference on Neural

863    Networks (IJCNN).

864    Rovero, F. and F. Zimmermann (2016). Camera Trapping for Wildlife Research, Pelagic Publishing.

865    Sagarra, O., M. Gutiérrez-Roig, I. Bonhoure and J. Perelló (2015). "Citizen Science Practices for

866    Computational Social Science Research: The Conceptualization of Pop-Up Experiments." Frontiers

867    in Physics **3**.

868    Samala, R., H.-P. Chan, L. Hadjiiski, M. Helvie, J. Wei and K. Cha (2016). "Mass detection in digital

869    breast tomosynthesis: Deep convolutional neural network with transfer learning from

870    mammography." Medical Physics **43**: 6654.

871    Schneider, S., G. Taylor and S. Kremer (2018). "Deep Learning Object Detection Methods for

872    Ecological Camera Trap Data." 321-328.

873    Shahinfar, S., P. Meek and G. Falzon (2020). " "How many images do I need?" Understanding how

874    sample size per class affects deep learning model performance metrics for balanced designs in

875    autonomous wildlife monitoring." Ecological Informatics **57**: 101085.

876    Singh, P., S. M. Lindshield, F. Zhu and A. R. Reibman (2020). Animal Localization in Camera-Trap

877    Images with Complex Backgrounds. 2020 IEEE Southwest Symposium on Image Analysis and

878    Interpretation (SSIAI).

879    Sudderth, E., A. Torralba, W. Freeman and A. Willsky (2008). "Describing Visual Scenes Using

880    Transformed Objects and Parts." <u>International Journal of Computer Vision</u> **77**: 291-330.

881    Sugai, L., T. Silva, J. Ribeiro Jr and D. Llusia (2018). "Terrestrial Passive Acoustic Monitoring: Review

882    and Perspectives." <u>BioScience</u> **69**.

883    Swanson, A., M. Kosmala, C. Lintott, R. Simpson, A. Smith and C. Packer (2015). "Snapshot Serengeti,

884    high-frequency annotated camera trap images of 40 mammalian species in an African savanna."

885    <u>Scientific Data</u> **2**: 150026.

886    Swinnen, K., J. Reijniers, M. Breno and H. Leirs (2014). "A Novel Method to Reduce Time Investment

887    When Processing Videos from Camera Trap Studies." <u>PloS one</u> **9**: e98881.

888    Tabak, M., M. S. Norouzzadeh, S. Sweeney, K. Vercauteren, N. Snow, J. Halseth, P. Salvo, J. Lewis, M.

889    White, B. Teton, R. Boughton, B. Wight, E. Newkirk, E. Odell, R. Brook, A. Moeller, E. Mandeville, J.

890    Clune, R. Miller and P. Schlichting (2019). "Machine learning to classify animal species in camera

891    trap images: Applications in ecology." <u>Methods in Ecology and Evolution</u> **10**: 585-590.

892    Tabak, M., M. S. Norouzzadeh, D. Wolfson, S. Sweeney, K. Vercauteren, N. Snow, J. Halseth, P. Salvo, J.

893    Lewis, M. White, B. Teton, J. Beasley, P. Schlichting, R. Boughton, B. Wight, E. Newkirk, J. Ivan, E.

894    Odell, R. Brook and R. Miller (2018). "Machine learning to classify animal species in camera trap

895    images: applications in ecology."

896    Torralba, A. and A. Efros (2011). "Unbiased look at dataset bias." <u>Proceedings of the IEEE Computer

897    Society Conference on Computer Vision and Pattern Recognition</u>: 1521 - 1528.

898    Torralba, A. and P. Sinha (2003). "Contextual Priming for Object Detection." <u>International Journal of

899    Computer Vision</u> **53**.

900    Wang, X., Z. Cai, D. Gao and N. Vasconcelos (2019). <u>Towards Universal Object Detection by Domain

901    Attention</u>. 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR).

902    Wang, X., X. Hu, C. Chen, Z. Fan and S. Peng (2019). <u>Improving Object Detection with Consistent

903    Negative Sample Mining</u>. Neural Information Processing, Cham, Springer International Publishing.

904 Willi, M., R. Pitman, A. Cardoso, C. Locke, A. Swanson, A. Boyer, M. Veldthuis and L. Fortson (2018).

905 "Identifying Animal Species in Camera Trap Images using Deep Learning and Citizen Science."

906 Methods in Ecology and Evolution **10**.

907 Xu, B., W. Wang, G. Falzon, P. Kwan, L. Guo, Z. Sun and C. Li (2020). "Livestock classification and

908 counting in quadcopter aerial images using Mask R-CNN." International Journal of Remote Sensing

909 **41**(21): 8121-8142.

910 Yosinski, J., J. Clune, Y. Bengio and H. Lipson (2014). "How transferable are features in deep neural

911 networks?": 3320-3328.

912 Young, S., J. Rode-Margono and R. Amin (2018). "Software to facilitate and streamline camera trap

913 data management: A review." Ecology and Evolution **8**(19): 9947-9957.

914 Yousif, H., J. Yuan, R. Kays and Z. He (2019). "Animal Scanner: Software for classifying humans,

915 animals, and empty frames in camera trap images." Ecology and Evolution **9**.

916 Yu, X., W. Jiangping, R. Kays, P. Jansen, T. Wang and T. Huang (2013). "Automated identification of

917 animal species in camera trap images." EURASIP Journal on Image and Video Processing **1**.

918 Zhang, Z., Z. He, G. Cao and C. Wenming (2016). "Animal Detection From Highly Cluttered Natural

919 Scenes Using Spatiotemporal Object Region Proposals and Patch Verification." IEEE Transactions on

920 Multimedia **18**: 1-1.

921 Zhao, W. (2017). "Research on the deep learning of the small sample data based on transfer

922 learning." AIP Conference Proceedings **1864**: 020018.

923 Zhao, Z.-Q., P. Zheng, S.-T. Xu and X. Wu (2019). "Object Detection With Deep Learning: A Review."

924 IEEE Transactions on Neural Networks and Learning Systems **PP**: 1-21.

925 Zhou, W., A. C. Bovik, H. R. Sheikh and E. P. Simoncelli (2004). "Image quality assessment: from

926 error visibility to structural similarity." IEEE Transactions on Image Processing **13**(4): 600-612.

927

[i] https://github.com/tzutalin/labelImg