



HAL
open science

Automatic detection and classification of honey bee comb cells using deep learning

Thiago S Alves, M Alice Pinto, Paulo Ventura, Cátia J Neves, David Biron, Arnaldo C Junior, Pedro L. de Paula Filho, Pedro J Rodrigues

► **To cite this version:**

Thiago S Alves, M Alice Pinto, Paulo Ventura, Cátia J Neves, David Biron, et al.. Automatic detection and classification of honey bee comb cells using deep learning. *Computers and Electronics in Agriculture*, 2020, 170, pp.105244. 10.1016/j.compag.2020.105244 . hal-03017457

HAL Id: hal-03017457

<https://hal.science/hal-03017457>

Submitted on 20 Nov 2020

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial - NoDerivatives 4.0 International License

AUTOMATIC DETECTION AND CLASSIFICATION OF HONEY BEE COMB CELLS USING DEEP LEARNING

Thiago S. Alves^{1,2*}, M. Alice Pinto³, Paulo Ventura³, Cátia J. Neves³, David G. Biron⁴, Arnaldo C. Junior², Pedro L. De Paula Filho², Pedro J. Rodrigues^{1*}

ABSTRACT: In a scenario of worldwide honey bee decline, assessing colony strength is becoming increasingly important for sustainable beekeeping. Temporal counts of number of comb cells with brood and food reserves offers researchers data for multiple applications, such as modelling colony dynamics, and beekeepers information on colony strength, an indicator of colony health and honey yield. Counting cells manually in comb images is labour intensive, tedious, and prone to error. Herein, we developed a free software, named *DeepBee*©, capable of automatically detecting cells in comb images and classifying their contents into seven classes. By distinguishing cells occupied by eggs, larvae, capped brood, pollen, nectar, honey, and other, *DeepBee*© allows an unprecedented level of accuracy in cell classification. Using Circle Hough Transform and the semantic segmentation technique, we obtained a cell detection rate of 98.7%, which is 16.2% higher than the best result found in the literature. For classification of comb cells, we trained and evaluated thirteen different convolutional neural network (CNN) architectures, including: DenseNet (121, 169 and 201); InceptionResNetV2; InceptionV3;

¹ Research Center in Digitalization and Intelligent Robotics (CeDRI), Instituto Politécnico de Bragança, Campus de Santa Apolónia, 5300-253 Bragança, Portugal

² Department of Computer Science, Federal Technological University of Paraná (UTFPR), Medianeira, Paraná, Brazil

³ Centro de Investigação de Montanha (CIMO), Instituto Politécnico de Bragança, Campus de Santa Apolónia, 5300-253, Portugal

⁴ Laboratoire Microorganismes: Génome et Environnement, UMR CNRS 6023, Université Clermont-Auvergne, Campus Universitaire des Cézeaux, France

* Corresponding author: alvest@alunos.utfpr.edu.br, Rua Mário Fontana, 1331, Toledo, Paraná CEP: 85910-190, Brazil.

E-mail addresses: apinto@ipb.pt (M. A. Pinto), paulo.j.c.ventura@gmail.com (P. Ventura), catia.jose7@ipb.pt (C. J. Neves), biron@mpl.ird.fr (D. G. Biron), arnaldoc@utfpr.edu.br (A. Candido Junior), pedrol@utfpr.edu.br (P. L. De Paula Filho), pjsr@ipb.pt (P. J. S. Rodrigues).

22 MobileNet; MobileNetV2; NasNet; NasNetMobile; ResNet50; VGG (16 and 19) and
23 Xception. MobileNet revealed to be the best compromise between training cost, with
24 ~9s for processing all cells in a comb image, and accuracy, with an F1-Score of
25 94.3%. We show the technical details to build a complete pipeline for classifying and
26 counting comb cells and we made the CNN models, source code, and datasets
27 publicly available. With this effort, we hope to have expanded the frontier of
28 apicultural precision analysis by providing a tool with high performance and source
29 codes to foster improvement by third parties ([https://github.com/AvsThiago/DeepBee-](https://github.com/AvsThiago/DeepBee-source)
30 [source](https://github.com/AvsThiago/DeepBee-source)).

31 **KEYWORDS:** cell classification; *Apis mellifera* L.; semantic segmentation; machine
32 learning; deep learning; *DeepBee* software.

33

34 1 INTRODUCTION

35 In a scenario of worldwide honey bee (*Apis mellifera* L.) decline, assessing
36 colony strength is becoming increasingly important as it can assist apiary
37 management strategies and provide valuable information for research purposes.
38 Counts of comb cells with brood and food reserves offer beekeepers information on
39 colony nutritional status, colony health status, queen quality, honey yield, etc. The
40 same data collected across time can be used by researchers in multiple applications,
41 including assessment of queen genotypes, assessment of new treatments against
42 parasites and pathogens, modelling colony dynamics in response to abiotic
43 (pesticides) or biotic (parasites, pathogens, predators) stressors, among others.

44 Delaplane et al. (2013) reviewed the methods for assessing colony strength,
45 among which is the Liebefeld method. This method has been included in the
46 HEALTHY-B toolbox compiled by EFSA (European Food Safety Authority) for

47 harmonising data collection on the health status of honey bees in Europe (EFSA
48 AHAW Panel, 2016). The Liebefeld method is based on direct observations of comb
49 frames in the apiary. Estimates are made with the help of an eight-quadrant grid
50 placed in front of each frame. With the grid in place, the observer estimates the
51 surface occupied by the targeted class (e.g., honey or capped brood) in every
52 quadrant. This method is prone to error, time-consuming and produces data with a
53 high level of subjectivity. Additionally, at least two well-trained observers are required
54 for gathering the data and for obtaining better estimates through averaging the two
55 subjective evaluations (Delaplane et al., 2013). In this context, semi-automatic or
56 automatic methods offer a better alternative for assessing colony strength as
57 subjectivity is eliminated.

58 Herein we developed a free software, DeepBee©, that can be readily
59 employed by honey bee researchers and apiculturists to assess colony strength
60 through analysis of comb images. DeepBee© is capable of automatically detecting
61 cells in comb images and classifying their contents with an unprecedented
62 discriminating power and level of accuracy. DeepBee© evaluates a set of comb
63 images at a high speed, allows edition of the automatic predictions, if needed, and
64 produces a spreadsheet file for downstream analysis. While developing this tool, we
65 further contributed to a wider community including machine learning, image
66 processing and other software developers (i) by providing a method in the
67 segmentation process to automatically readjust the scale of the images driven by the
68 size of the cells, (ii) by testing different neural network architectures related to
69 performance and quality of results, (iii) by providing datasets that can be employed
70 by others when testing new methods, and (iv) by making available the source codes
71 enabling the reproducibility of our results.

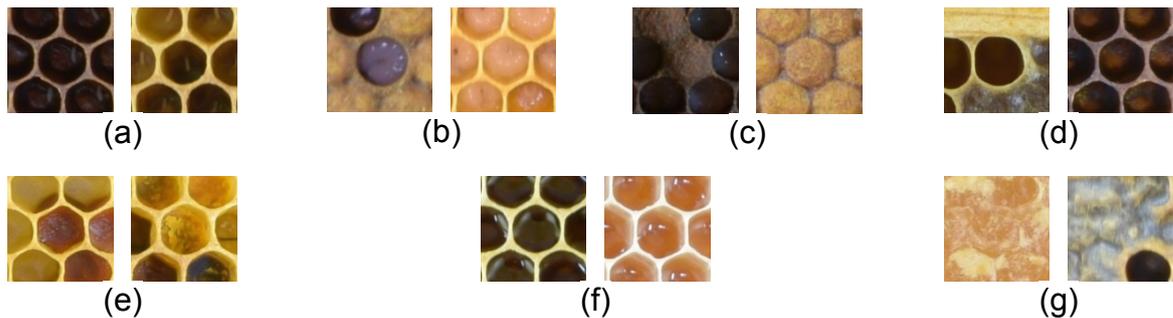
72 1.1 Motivation

73 Several semi-automatic methods for assessing colony strength using digital
74 images of comb frames have been proposed in the last years (reviewed below).
75 When compared with manual methods, *post-hoc* analysis of comb images reduces
76 the time of information collection, provides more accurate data, and assures
77 reproducibility of results even with different users. Furthermore, the images
78 themselves are permanent records of the data, representing an important step
79 towards more accountable and objective assessment of colony strength, as no
80 record is usually available after the commonly used visual estimation of combs.

81 Developing tools for analysis of comb images requires pre-defining the cell
82 classes that will be targeted. The higher the number of cell classes to be
83 distinguished in a comb image, the greater the complexity of the classification model.
84 During a colony lifetime, comb cells may be (i) momentarily empty, (ii) occupied by
85 the honey bee in its different immature stages (egg, larva, pupa), or (iii) filled with
86 food resources (pollen, nectar, honey) required for colony development and
87 maintenance (Fig. 1). To reflect the high level of cell-content diversity, at least seven
88 different classes should be pre-defined when developing models for cell
89 classification. In addition to class diversity, there might be a wide array of colours and
90 textures within each class, making cell classification a challenging endeavour.

91 Previous works (reviewed below) have addressed this challenge by
92 developing tools for assessing only the number of capped brood cells (Fig. 1c), a
93 task that is greatly facilitated by the striking visual differences between capped brood
94 and the remaining cells. Here, we developed a tool capable of assigning cell contents
95 to seven different classes, which represents an unprecedented level of accuracy in
96 classification of comb images.

97



98 **Fig. 1:** Comb cell classes considered in this study: (a) egg, (b) larva (uncapped
 99 brood), (c) pupa (capped brood), (d) other (e.g. empty), (e) pollen, (f) nectar, (g)
 100 honey (bee-processed nectar).

101

102 1.2 Related works

103 Early works developed semi-automatic tools for assessing colony strength by
 104 measuring the comb area occupied by brood (Emsen, 2006; Yoshiyama et al., 2011).
 105 Emsen (2006) performed the segmentation using mainly the selection tools of the
 106 software *Adobe Photoshop*[®] CS2. Yoshiyama et al. (2011) developed an approach
 107 similar to that of Emsen (2006) by using a tool similar to *Adobe Photoshop*[®] to create
 108 a semi-supervised segmentation. The novelty was the plugin called *LarvaeArea*
 109 created for the image processing software *ImageJ*
 110 (<https://imagej.nih.gov/ij/index.html>). With this plugin, the user is able to open the
 111 previously segmented image and calculate automatically the area occupied by both
 112 capped and uncapped cells.

113 Cornelissen et al. (2009) further advanced comb image assessment by
 114 developing a semi-automatic method that counted the number of capped brood cells,
 115 instead of measuring occupied area. This method was on average 23s slower than
 116 the Liebefeld method, as it required human intervention during segmentation.
 117 However, the estimates of capped cells were more accurate with the semi-automatic

118 method (correlation with the actual number of cells = 0.99) than with the Liebefeld
119 method (correlation with the actual number of cells = 0.91).

120 One of the first digital methods capable of detecting and counting individual
121 cells was developed by Liew et al. (2010). The authors used pre-processing methods
122 to highlight the edges and applied the Circle Hough Transform (CHT) to detect the
123 cells. They obtained a detector with a cell detection rate of 82.6%.

124 More recently, Rodrigues et al. (2016) developed a method for automatic
125 detection and counting of capped brood cells using circular convolution. The circular
126 mask has the same size of a comb cell, it stops in each cell position, and it calculates
127 the contrast between the pixels of the cell edge and its interior. From this contrast, it
128 is possible to know whether there is a cell, and according to established thresholds, it
129 is possible to know whether the cell is capped or uncapped.

130 Meanwhile, various software packages have been developed for comb
131 assessment. In a presentation of the *HoneybeeComplete*, Wang & Brewer (2013)
132 showed the ability of this commercial software to correctly classify capped brood cells
133 97.4% of the time, with the rate increasing up to 99.5% when the user pre-selects the
134 search area. The developers did not provide methodological details on software
135 development, only reporting the use of an unspecified set of pattern algorithms.

136 The commercial *HiveAnalyzer* software, developed by Höferlin et al. (2013),
137 represented an important step forward in comb assessment by classifying cells other
138 than capped brood. The developers categorised the comb cells into seven classes by
139 using a cascade of classifiers based on linear Support Vector Machines. The
140 accuracy of the classifier was 94%, as measured on a subset of cells classified with
141 high confidence.

142 More recently, Colin et al. (2018) developed the software *CombCount* to
143 assess capped brood and capped honey. Although the software is able to detect both
144 classes, it requires a user to distinguish the contents using selection tools.

145 To the best of our knowledge, none of the studies published so far has
146 employed Neural Networks such as CNNs in honey bee comb image analysis.
147 Furthermore, most of the available methods are limited to detecting capped brood,
148 with only two of them being capable of assessing food resources (Colin et al., 2018;
149 Höferlin et al., 2013). As such, there is an excellent opportunity for innovation using
150 new methods to address a major challenge in honey bee research, which is
151 assessing brood and food resources in the hive in a time- and cost-effective manner
152 and with a high degree of accuracy.

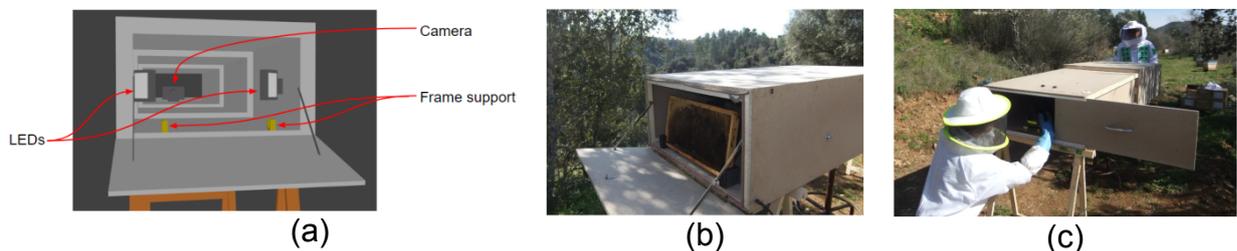
153

154 **1.3 Goal**

155 The goal of this study is threefold: (i) to develop a pipeline capable of detecting
156 all cells in a comb image, (ii) to reliably classify the cells into seven different classes,
157 and (iii) to encapsulate this pipeline in a free software. In accomplishing this goal, the
158 following research questions will be addressed: (i) Is it possible to develop an image
159 processing method to detect cells in comb images, even when the edges are hard to
160 identify? (ii) Is it possible to develop computational models capable of reliably
161 classifying the contents of comb images using Deep Learning? (iii) What are the
162 implementation details to achieve the best functional performance? (iv) Which neural
163 models provide the best results for our problem among many Deep Learning
164 architectures available? Finally, (v) How does the proposed approach compare to the
165 related published works?

166 **2 IMAGES CAPTURE AND ANALYSIS**

167 To assure image capture standardization, we developed a wooden tunnel
 168 sealed for external light and with optimized dimensions (Fig. 2). This tunnel had a
 169 retractable architecture for easy transportation, having a length of 247cm, when fully
 170 opened, and 92cm, when retracted. The comb frame and the camera were placed at
 171 the opposite sides of the tunnel. The comb frame was positioned in two holders (Fig.
 172 2b). The holders had an angle of 11° (to make up for to the $9\text{-}13^\circ$ natural inclination
 173 of comb cells) for a better image capturing of the interior of the comb cells (see the
 174 3D model file in <https://github.com/AvsThiago/DeepBee-source>). Close to the comb
 175 frame (40cm from the top), there was pair of Light-Emitting Diode (LED) sources, with
 176 7 Watts of power. The LEDs were turned to the walls at 45° to provide homogeneous
 177 light conditions and to avoid shadows during image acquisition. The camera was
 178 fixed with a screw on the opposite side of the tunnel (Fig. 2c). Further details about
 179 the tunnel features are shown in Appendix B.



180 **Fig. 2:** (a) Details of the interior of the tunnel. (b) Tunnel installed in an apiary
 181 showing the comb frame placed on holders. (c) Researcher adjusting the camera
 182 before shooting.

183

184 In this study we used a digital camera Nikon D3300, with lens AF-S DX VR
 185 Zoom-Nikon ED 55-200mm F4-5.6G, and the following settings: aperture - 10; ISO -
 186 100; shutter speed - 1/60; autofocus - on; flash - no; compression – JPEG; white
 187 balance - on. During image capturing, the tunnel was closed on both sides and the
 188 camera was activated by an external trigger. The images captured had a resolution
 189 of 24MPixels (6000x4000px). Using these settings, 1,102 comb frames were

190 photographed on both sides making a total of 2,204 images. The image dataset is
191 available online at <https://cloud.ipb.pt/d/aa29c989ab1944aaa222/?p=/DS-COMB-PT>.

192

193 **3 SCALE INVARIANT DETECTION AND FALSE DETECTION REMOVAL**

194 After obtaining the 2,204 images, we searched for a method capable of
195 reliably detecting individual comb cells, a task that precedes cell classification.
196 Below, the different approaches and steps followed in this study are described.

197

198 **3.1 Circle Hough Transform**

199 Duda and Hart (1972) developed the classical Hough Transform method
200 currently available. This method was originally developed for detecting lines in
201 images. Later on, it was discovered that it could also be used for identifying arbitrary
202 shapes, such as circles and ellipses, even when they were partially occluded.

203 The Circle Hough Transform (CHT) method uses a voting process to calculate the
204 probability that a set of pixels form a circle. There are several implementations of the
205 method. Herein, we used the implementation contained in the OpenCV v.4.0 library
206 (<https://github.com/opencv/opencv/releases/tag/4.0.0>). The parameters expected by
207 this method are: a grayscale image, size of an internal accumulator that will store
208 intermediate results, minimum distance between two detections centre, threshold to
209 be applied to the internal Canny operator, number of votes that a circle must have in
210 the accumulator to be set as true, minimum circle radius, and maximum circle radius.

211

212 **3.2 Using CHT to detect cells**

213 Prior to detecting comb cells using CHT, we applied a pre-processing method
214 to normalise illumination, remove noise, and enhance cell edges. The pre-processing

215 pipeline was developed using empirical tests, as follows: (i) extract only the red
216 channel from the image; (ii) apply a Contrast Limited Adaptive Histogram
217 Equalization (CLAHE) (Zuiderveld, 1994) with 8x8 tiles and clip limit of 2.0; (iii)
218 remove the noise keeping the edges by using the bilateral filter (Tomasi & Manduchi,
219 1998) with diameter 5, sigma colour 50, and sigma space 50.

220 Having the pre-processing method defined, we started looking for the best
221 parameters combination for CHT. Typically, comb cells have well-behaved diameters
222 and distances and do not overlap. These features greatly facilitated cell detection by
223 the CHT method. On the other hand, finding a combination of parameters to make
224 CHT capable of detecting uncapped and capped brood and honey cells revealed to
225 be a challenging task to be done by guessing and checking. Therefore, we used a
226 grid search-based algorithm to find the optimal combination. The best result obtained
227 from our images was: internal accumulator size – 3, minimum distance – 51, Canny
228 threshold – 100, minimum number of votes – 25, minimum radius – 31, and
229 maximum radius - 37.

230 With this combination of parameters, we were able to successfully detect all
231 different types of cells. However, setting fixed values for the minimum distance,
232 minimum radius and maximum radius makes detection less generalist, requiring that
233 images are acquired using a setup like that described above. Therefore, to
234 generalize the method, we developed a scale-invariant detection method, which has
235 two main stages. First, the average cell size and the mean distance between them
236 are sought. Second, all cells are detected using the other parameters discovered by
237 the grid search. The detailed operation involved the following steps (the distances
238 are in pixels):

239 I. **Detect cells with radius belonging to different ranges.** In this step, we only
 240 considered detections by the CHT method with high levels of confidence (Fig. 3a).
 241 The fixed parameters were: internal accumulator - 2, Canny - 145, minimum
 242 number of votes (confidence) – 55, and minimum distance - 12. We also iterated
 243 a loop with i ranging from 5 to 50 and step 5. At each iteration, the CHT method
 244 was executed with the parameters $minRadius = i+1$ and $maxRadius = i+5$. After
 245 running this method, a list with the number of detections made for each radius i
 246 was returned (Fig. 3b).

247 II. **Find the most frequent radius.** In this step, we selected the most frequent
 248 radius from the detections made on the given image with different radius i . Fig. 3b
 249 shows how many radius-based circles were found in the image. Most of the
 250 circles were detected with the radius 18, since it is the most frequent cell size for
 251 this image.

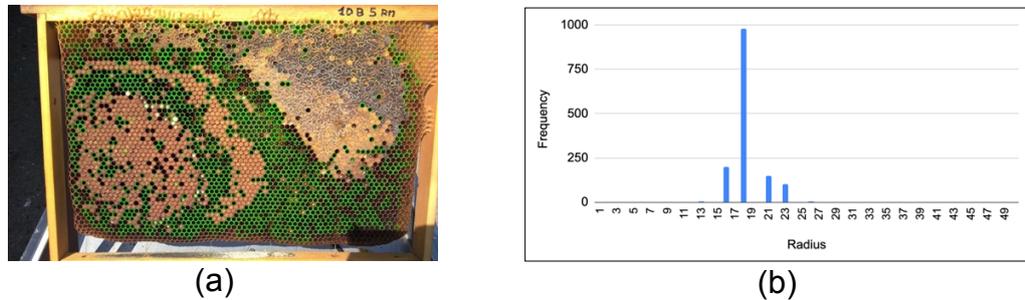
252 III. **Define the minimum distance between two detections.** Due to different cell
 253 sizes and imperfections made by honey bees during comb construction, in this
 254 step, we had to choose values for the minimum distance parameter smaller than
 255 $2 \times radius$. After analysing the average distance between two cells of numerous
 256 frame images of different sizes, we found that usually the minimum distance fits
 257 the equation 1, where r stands for radius.

$$258 \quad minDist(r) = 1.65r - 3, \quad (Equation 1)$$

259 IV. **Find the parameters minRadius and maxRadius for a given radius.** To deal
 260 with the natural variation of comb cell size and with images taken from different
 261 distances, in this step, we created a range based on the average cell radius. The
 262 range is defined by equation 2.

$$263 \quad range(r) = r \pm \begin{cases} 0.1r, & \text{if } 0.1r > 1 \\ 1, & \text{if } 0.1r \leq 1 \end{cases}, \quad (Equation 2)$$

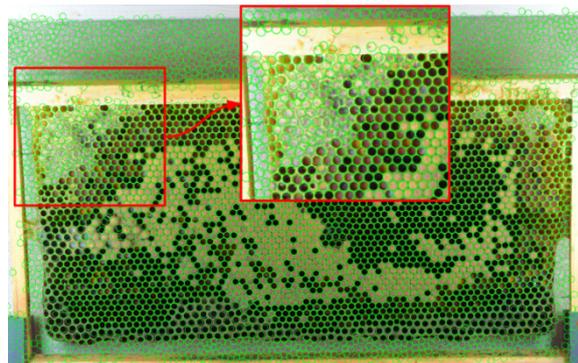
264 V. **Perform CHT with the obtained parameters.** After obtaining the parameters
 265 needed, we processed the images again with the CHT method, but this time using
 266 the remaining parameters found by the grid search (accumulator size, Canny
 267 threshold, minimum number of votes).



268 **Fig. 3:** (a) Detection of cells with high confidence by the CHT method. (b) Number of
 269 cells detected for radius size.

270

271 Using parameters that accept more detections as true increases the power of
 272 detecting different types of cells, even those with fuzzy edges like honey cells (Fig.
 273 3). However, by reducing the threshold for detecting all true cells there is a risk of
 274 false detections (Fig 4). To alleviate this problem, we developed a method based on
 275 CNNs, which is described in the ensuing section.



276

277 **Fig. 4:** Comb cells, with different contents (honey, pollen, capped brood), detected by
 278 our approach. Detected cells, including false detections outside of the comb, are
 279 marked by a green hexagon. The close-up square shows honey cells (on the left half)
 280 with fuzzy edges.

281

282 **3.3 Removing false detections using semantic segmentation**

283 Semantic image segmentation, also called pixel-level classification, is the task
284 of clustering parts of an image together, which belong to the same object class
285 (Thoma, 2016). We used this technique for detecting comb cells in the image and
286 from this segmentation remove false cells falling outside the comb area (Fig. 4). To
287 that end, we used a CNN encoder-decoder architecture based on U-Net
288 (Ronneberger et al., 2015).

289

290 3.3.1 Dataset creation for comb segmentation

291 We created the annotations using the *Quick Selection Tool* from the software
292 *Adobe Photoshop® CS6*. To define the classes, we painted white the comb area and
293 black the remaining area. We labelled 61 comb images (Alves et al., 2019), which
294 were selected to represent a high diversity of cell content (e.g. honey, pollen, capped
295 and uncapped brood) and age (the older the comb the darker it gets). Fig. 5
296 illustrates annotations made on those images. The annotations were split in three
297 sets: training (85%); validation (10%), and testing (5%).

298 Using the strategy proposed in Ronneberger et al. (2015), we divided the input
299 images and the labels in tiles, as shown in Fig. 6. Prior to transforming each image in
300 tiles, a mirrored border with size 184px (top-bottom) by 148px (left-right) was added
301 to create more space for the tiles, and the images were resized to a constant size of
302 1000x1500px. During tiles extraction, overlaps between tiles were taken into account
303 to reduce border problems in the reassembly phase. Using these processes, we
304 transformed the 61 images into 7,137 tiles.

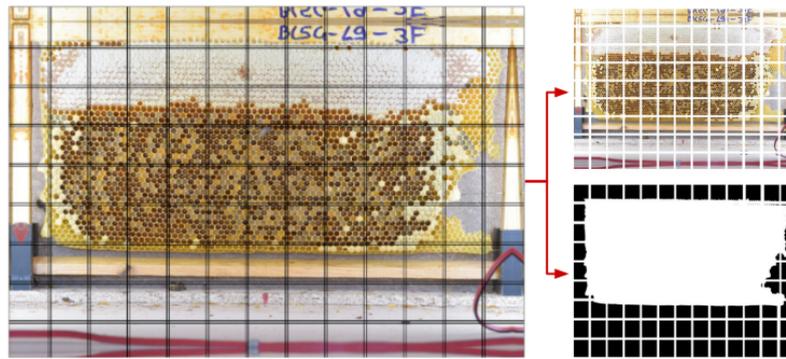


(a)



(b)

305 **Fig. 5:** Samples from the dataset created for semantic segmentation: (a) original
 306 image; (b) label.



307

308 **Fig. 6:** Tiles created from the original images and respective labels.

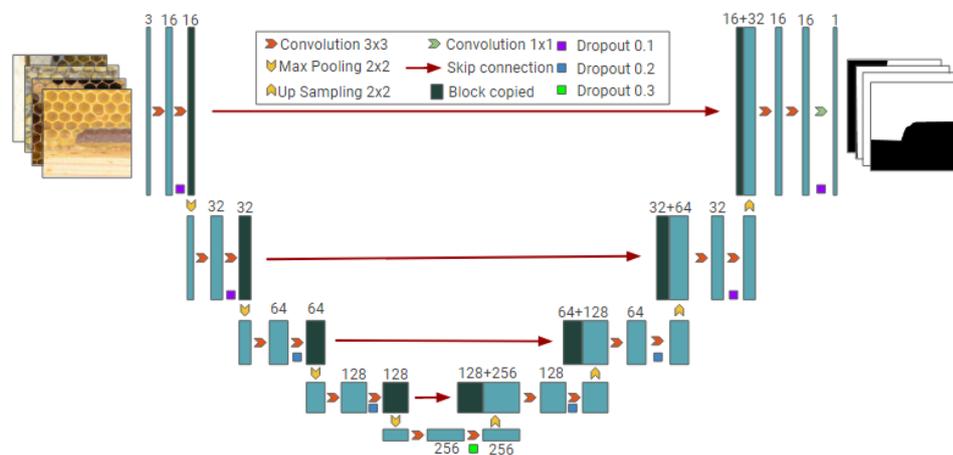
309

310 3.3.2 Semantic segmentation architecture and training policy

311 The architecture had a depth of 5 convolutions with 3×3 filters and layers of
 312 maximum pooling with 2×2 filters and a stride of 2, as proposed by Ronneberger et
 313 al. (2015). Our modifications to the original model, made to obtain the best results in
 314 our semantic segmentation dataset, were tuned using a trial-error approach and the
 315 following settings: input image with 128×128 resolution, dropout (Srivastava et al.,
 316 2014) ranging from 0.1 to 0.3 between the convolution layers, use of Exponential
 317 Linear Units (ELU) activation function (Clevert et al., 2015), and use of 16 filters
 318 (channels) in the first layer, doubling the amount at each inner level and returning 16
 319 filters in the penultimate layer with the last layer only having two dimensions. The
 320 CNN architecture is shown in Fig. 7.

321 We normalised the input images by dividing each pixel by 255. As we only had
 322 two regions to be classified, we used the binary cross-entropy as loss function. The
 323 network weights were initialised using the He Normal initialisation (He et al., 2015).
 324 For the output layer, we chose a sigmoid function, so we could easily transform the
 325 output in a binary image applying a threshold, where values < 0.5 became zero and
 326 the remaining became one. The best training was carried out using the Adam

327 optimiser (Kingma & Ba, 2014) with parameters $\beta_1 = 0.9$ and $\beta_2 = 0.999$. The
 328 Learning Rate (LR) was 10^{-3} , which was preserved during 50 epochs. Due to the
 329 Early Stopping method used here, the training may last less than 50 epochs. This
 330 method can halt the training if a chosen metric does not improve after a pre-defined
 331 number of epochs (6, in this study). The architecture was built using the framework
 332 Keras 2.1.4 with TensorFlow 1.4 as backend. The configurations of the computer
 333 used for training were two GPUs: NVIDIA GeForce GTX 1080Ti and NVIDIA
 334 GeForce GTX 1070; RAM: 16GB; CPU: Intel® Core™ i7-7700K CPU @ 4.20GHz×8;
 335 operating system: Ubuntu 17.10. All tests were performed in this computer.



336 **Fig. 7:** CNN architecture based on U-Net to handle comb segmentation.

337

338 3.3.3 Post-processing of comb segmentation

339 A post-processing step was undertaken so that the segmented image could be
 340 used to minimize false detections. To have a binary output after the CNN
 341 computation, we applied a threshold of 0.5 and the tiles were reassembled.
 342 Subsequently, we found the largest contour using an OpenCV method for finding
 343 contours and draw it filled using a method to draw polygons (Fig. 8). Filtering using
 344 the largest contour helped removing false segmentations inside and outside the
 345 comb, as shown in Section 5.1.2.



346
347 **Fig. 8:** Process developed to reduce the number of false segmented areas.

348

349 **3.4 Semantic segmentation experiments**

350 In the first experiment, we assessed the quality of detections on an
351 independent set of 10 comb images, which were downloaded from the Internet. This
352 set was submitted to the scale-invariant cell detection algorithm and the false
353 detections removal method.

354 In the second experiment, we measured the cells detection rate of our
355 algorithm and compared with that proposed by Liew et al. (2010). Following the
356 methodology of Liew et al. (2010), we selected 10 images from our dataset to be
357 analysed by our detection algorithm. Subsequently, we evaluated manually the false
358 positive (FP) and false negative (FN) detections. From these annotations, we
359 collected the following metrics: number of cells identified by humans; number of cells
360 detected by the algorithm; true positive (TP) detections; true cells detected correctly;
361 FP detections on inexistent cells; and FN cells undetected by the algorithm. Finally,
362 as in Liew et al. (2010), we calculated for each image the cells detection rate using
363 equation 3. This metric is based on the total number of cells automatically detected,
364 excluding the falsely identified cells, divided by the manual count.

365

$$366 \text{ CellDetectionRate} = \frac{\text{DetectionCount} - \text{FP}}{\text{ManualCount}} \times 100\% , \quad (\text{Equation 3})$$

367

368 **4. CELLS CLASSIFICATION**

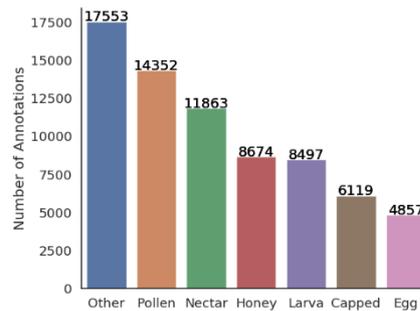
369 The cell classification was carried out using CNNs. This supervised approach
370 gained momentum after the release of Krizhevsky et al. (2012) work. At the time, this
371 work was considered the state-of-the-art in the ImageNet Large Scale Visual
372 Recognition Challenge, where the goal was to accurately classify more than a million
373 images into 1000 distinct classes. After this seminal work, important advances have
374 been made in CNN architectures, with major consequential breakthroughs in various
375 fields of study such as agriculture (Kamilaris & Prenafeta-Boldú, 2018). A key feature
376 for training a CNN architecture is the massive amount of data needed. Following, we
377 show how we gathered the cells images for our dataset.

378

379 **4.1 Dataset gathering for cells classification**

380 The dataset created for classification should contain cells representing all
381 different classes in different comb images. In this study, the annotations were made
382 by an experienced beekeeper, assuring the high quality required for developing the
383 models. A piece of software was developed to facilitate the beekeeper's work. The
384 software allowed choosing a label corresponding to each class (nectar, honey,
385 pollen, egg, larva, capped brood, and other) and, at the same time, pointing the
386 centre of the cells, adjusting the contrast, brightness and gamma of the images. A
387 total of 71,915 cells were annotated on 1,202 comb images, with an average of 25
388 cells per image. The number of annotations by class is shown in Fig. 9.

389 We divided the annotations in three sets for the training as follows: 80% of the
390 original dataset for training, 20% of the training set for validation, and 20% of the
391 original dataset for testing. Then, we selected 15 additional comb images and
392 annotated all the cells. We made a total of 39,533 new annotations and added them
393 to the test set.

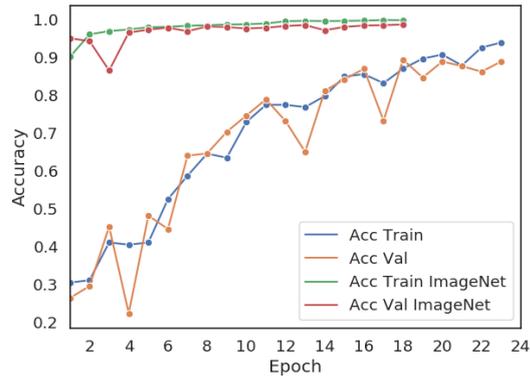


394
395 **Fig. 9:** Number of annotations per class.

396

397 **4.2 Transfer learning for cells classification**

398 Using the transfer learning technique, it is possible to transfer the weights of
 399 feature extraction layers (e.g. convolutions) from a trained model over a dataset to
 400 another model that will be trained in a new dataset (Oquab et al., 2014). Because the
 401 new model received kernels already trained to recognise generic features, like lines
 402 and curves, it will be easier for the model to generalise a new dataset being
 403 unnecessary to learn the filters from scratch. We made a sanity check and trained an
 404 architecture in our dataset with and without the transfer learning before applying this
 405 approach to the next experiments. In this experiment, we used the same policy in
 406 both trainings. We transferred the weights from a pre-trained model on the ImageNet
 407 dataset (Deng et al., 2009). We used the architecture InceptionV3 (Szegedy et al.,
 408 2015) and, as shown in Fig. 10, the model converged faster, in a lower number of
 409 epochs and with a higher accuracy with the transfer learning than with training from
 410 scratch.



411
 412 **Fig. 10:** Comparison between models trained in our dataset from scratch and using
 413 pre-trained weights from ImageNet.

414

415 4.3 Finding the best region of interest (ROI) size to crop the cells

416 Before we defined the CNN architecture for the classifier, we needed to find at
 417 which window size around the cells the image should be cropped. If we had cropped
 418 only the interior of the cells, our classifier could have had difficulty in distinguishing
 419 between capped brood and honey cells, as these classes typically exhibit a similar
 420 texture (Fig. 11). To select the best input size, we created and tested 16 datasets, all
 421 of them with the same annotations but with different ROI size. We defined each
 422 dataset with 10% of the annotations from the main dataset. The sizes (pixels) of the
 423 squared crops were 40, 50, 60, 80, 100, 120, 140, 160, 180, 200, 220, 240, 280, 298,
 424 400 and 500.

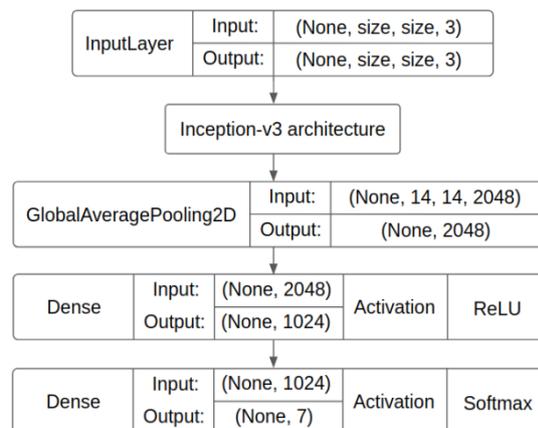
425



426 **Fig. 11:** Comparison between (a) honey and (b) capped brood classes. In addition to
 427 the cells interior (a1, b1), the cells neighbourhood (a2, b2) was taken into
 428 consideration to better define the two classes.

429

430 The training framework for the tests was the Keras version 2.2 and the
 431 InceptionV3 architecture (Szegedy et al., 2015). For the feature extraction layers, we
 432 used weights pre-trained on the ImageNet dataset using the transfer learning
 433 technique. The architecture and trained weights were provided by Keras library. This
 434 architecture does not allow inputs lower than 139×139px. Image datasets lower than
 435 139×139px were resized to the minimum input size. We added three layers at the
 436 end of the architecture to create the specific learning on the classifier. The first one
 437 was a flattening layer applied to the network output. Afterwards, we included two
 438 Dense layers (fully-connected) with the last one having 7 neurons and a Softmax
 439 activation function to represent our classes in a linear probabilistic domain. Details
 440 about the architecture are shown in Fig. 12.



441 **Fig. 12:** Developed architecture based on InceptionV3.
 442

443

444 The classifiers were compiled using the Categorical Cross Entropy loss
 445 function. We chose Adam with default parameters $\beta_1 = 0.9$ and $\beta_2 = 0.999$ as the
 446 optimiser. The training started with the LR at 10^{-3} and using the technique Reduce
 447 Learning Rate on Plateau (He et al., 2015a). We defined that the LR would be halved
 448 after 3 epochs without improvement in the Loss metric, being the minimum value
 449 10^{-6} . We established 50 for the maximum number of epochs. Again, we used the
 450 Early Stopping technique with the maximum number of epochs, without
 451 improvement, set at 5. We saved the model with the lowest validation Loss in each

452 dataset. The best results were provided by a window size of 224x224px. We opted
453 for 224x224px rather than 220x220px, for the reasons pointed out in Section 5.2.1.

454

455 **4.4 Tests with different CNN architectures**

456 After obtaining the best window size, as described previously, we trained
457 different CNN architectures with the input size 224x224px to identify which one
458 produced the best results on our dataset. We trained 13 distinct architectures
459 selected by their superior performance in image classification competitions in the
460 past years. These architectures included DenseNet 121, DenseNet 169,
461 DenseNet201 (Huang et al., 2016), InceptionResNetV2 (Szegedy et al., 2016),
462 InceptionV3 (Szegedy et al., 2015), MobileNet (Howard et al., 2017), MobileNetV2
463 (Sandler et al., 2018), NasNet; NasNetMobile (Zoph et al., 2017), ResNet50 (He et
464 al., 2015a), VGG 16, VGG 19 (Simonyan & Zisserman, 2014), and Xception.

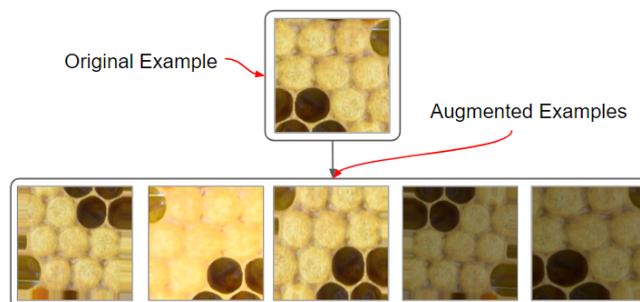
465 Each architecture was trained using all training and validation images from the
466 core dataset. The images were extracted with the cell centralised 224x224px. Before
467 starting training each model, we made some modification in the CNN architectures.
468 We added a set of fully-connected layers, as shown in Fig. 12. Prior to processing by
469 the model, the images were normalised by subtracting the ImageNet Mean Image
470 (103.939, 116.779, 123.68). The weights were transferred from previous ImageNet
471 trainings. The training was performed with batches of 40 images. We defined the
472 initial LR at 10^{-3} and used the Early Stopping and Reduce LR on Plateau, as in
473 Section 4.3. For comparing the model, we extracted some information and metrics,
474 including total architecture parameters (weights), time to process an image batch,
475 training time, accuracy, loss, precision, recall and F1-score.

476

477 4.4 Data augmentation

478 It is not always possible to obtain large datasets for training CNNs, either due
 479 to difficulties in gathering images with the object or in affording human resources to
 480 annotate the datasets. One way to enlarge the working dataset is through the Data
 481 Augmentation (DA) technique (Krizhevsky et al., 2012). Using DA, it is possible to
 482 create virtual examples from a set of images. Different transformations with random
 483 values are applied to these images and new ones are generated. Examples of
 484 transformations include changes in brightness, contrast, translates, rotations, zoom,
 485 and perspective.

486 Based on the models with the best metrics discovered in the experiment
 487 described in Section 4.3, we made a new training with additional data. We generated
 488 the new images using DA. Flips, brightness changes, rotations, shift, and zoom were
 489 applied in the newly created images. As a result, we generated a dataset of 250,000
 490 images evenly spread across the 7 different classes. These images were used in the
 491 training set. Validation and test sets were kept with the original images. We
 492 compared the resultant models using the metrics described in Section 4.3. Fig. 13
 493 shows several examples generated using DA.



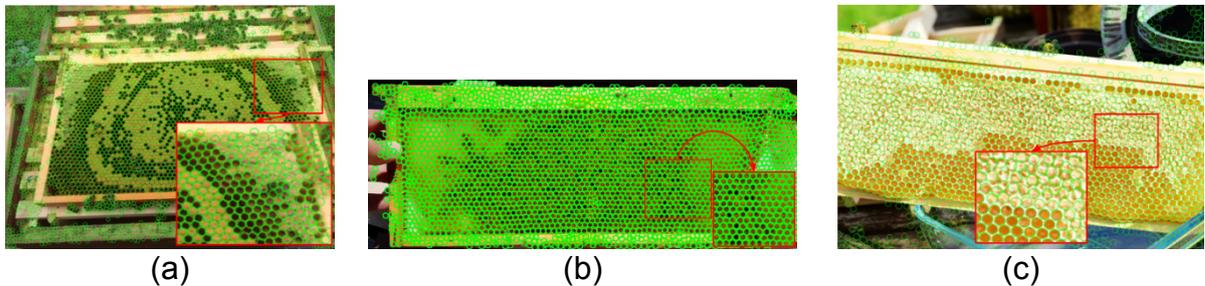
494 **Fig. 13:** Augmented examples generated from an original image.
 495
 496

497 5 RESULTS AND DISCUSSION

498 5.1 Cells detection and false detection removal

499 5.1.1 Cells detection

500 In the first cell detection experiment, we assessed the performance of the
501 developed algorithm in a set of independent images downloaded from the internet.
502 After selecting images that were captured under a wide range of conditions, we
503 generated the results, some of which are illustrated in Fig. 14.



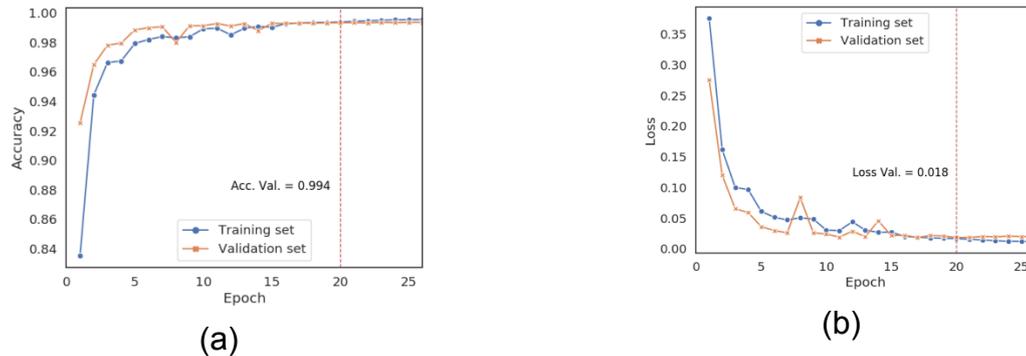
504 **Fig. 14:** Cells detected with (a) radius 18 (leahybeekeeping.com), (b) radius 8
505 (mudsongs.org), and (c) radius 48 (beekeepercenter.com).
506

507 The algorithm successfully detected most comb cells of the selected images
508 (Fig. 14). These images represented a wide range of hive frame type, cell size, cell
509 content (e.g. honey, capped and uncapped brood), and even varying illumination,
510 texture and resolution, suggesting that the algorithm developed in this study is
511 robust.

512

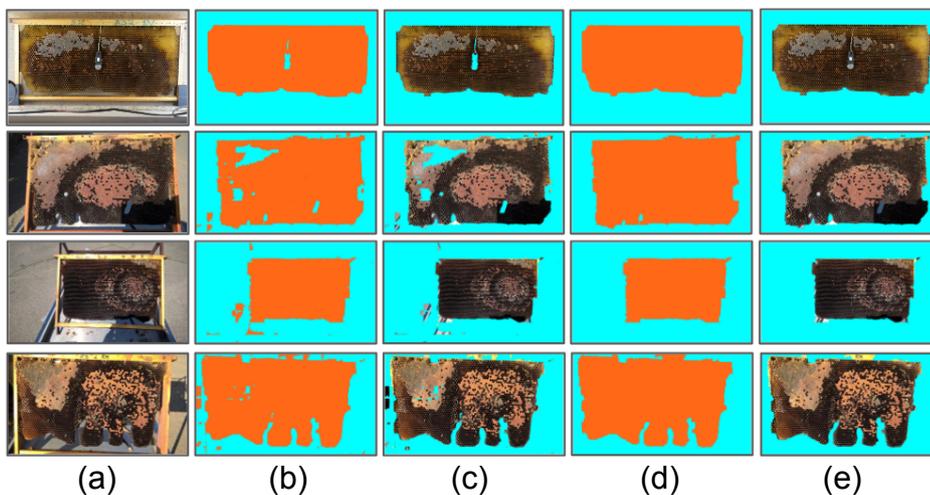
513 5.1.2 False detection removal

514 The training of the CNN for segmentation was carried out with 23 epochs in
515 3.45 min on the computational architecture referred in section 3.3. Fig. 15 shows the
516 evolution of accuracy and loss metrics in the training and validation sets along the
517 epochs. The model loss improved quickly before the 10th epoch; after that, only small
518 improvements were achieved, even with the constant learning rate.



519 **Fig. 15:** Evolution of (a) accuracy and (b) loss during training of the comb semantic
 520 segmentation model. The vertical dashed line over the 20th epoch represents the
 521 best calculated results (lowest loss).
 522

523 Using a CNN to segment the comb proved to be a robust solution, as it
 524 delivered great results even on the independent set of images. Fig. 16 shows some
 525 examples of segmentations performed in the independent and in our image sets. The
 526 downside is that CNN may have difficulty in segmenting when the comb cells have
 527 not yet been developed (there is only wax foundation) and when honey cells are very
 528 bright. The decision of selecting only the largest polygon and fill it to create a mask
 529 contributed to removing false segmentations inside and outside the comb area, as
 530 illustrated in Fig. 16. This approach has a poorer performance on combs that are
 531 broken (see bottom image of Fig. 16a) or have objects on front (see the thermohygro
 532 sensor on the top image of Fig. 16a), for example. However, these comb defects are
 533 unusual.



534
 535

536 **Fig. 16:** (a) Original image, (b) segmentation mask provided by the CNN without
 537 post-processing, (c) segmentation mask applied to the original image, (d) largest
 538 contour used as a mask, (e) largest contour applied to the original image.

539

540 5.1.3 Comparative analysis

541 To measure the quality of the detections using the segmentation, we
 542 performed several tests, similar to those of Liew et al. (2010), on 10 selected images.

543 The results of these tests are shown in Table 1. False cell-related or noise-related
 544 detections were negligible in our tests, as found by Liew et al. (2010). A factor that
 545 had a negative impact on the results of Liew et al. (2010) work was the low contrast
 546 in some cells. As detailed on Section 3.2, we dealt with this problem by applying the
 547 CLAHE filter to our images before detecting the cells.

548

549 **Table 1:** Comparison between cells detected automatically and cells detected
 550 automatically and manually corrected.

Image name	Manual count	Automatic count	TP	TP (%)	FP	FP (%)	FN	FN (%)	CDR (%)
DSC_1940.JPG	3024	2949	2944	97.35	5	0.17	80	2.65	97.35
DSC_1992.JPG	2795	2742	2735	97.85	7	0.26	60	2.15	97.85
DSC_2832.JPG	2869	2833	2794	97.39	39	1.38	75	2.61	97.39
DSC_2839.JPG	3082	3062	3041	98.67	21	0.69	41	1.33	98.67
DSC_2864.JPG	2961	2982	2948	99.56	34	1.14	13	0.44	99.56
DSC_2951.JPG	2910	2889	2857	98.18	32	1.11	53	1.82	98.18
DSC_2443.JPG	2077	2088	2075	99.90	13	0.62	2	0.10	99.90
DSC_3475.JPG	2875	2876	2852	99.20	24	0.83	23	0.80	99.20
DSC_4326.JPG	3061	3092	3054	99.77	38	1.23	7	0.23	99.77
DSC_4496.JPG	3072	3056	3044	99.09	12	0.39	28	0.91	99.09

551 CDR – Cell detection rate; FP False Positive; FN False Negative; TP True Positive.

552

553 As in Liew et al. (2010), we calculated the cell detection rate for each image
 554 using Equation 3. The resulting detection rates varied between 97.35% and 99.9%,
 555 with an average of 98.7% (Table 2). These rates were substantially higher than those
 556 reported by Liew et al. (2010). However, the method of Liew et al. (2010) produced a
 557 lower number of FP (5.1) than our method (22.5). This difference can be explained

558 by the fact that Liew et al. (2010) used more stringent parameters for CHT, allowing
 559 only detections with a high level of confidence. On the other hand, when we
 560 examined the number of FN, our approach produced substantially better results (38.2
 561 compared with 274.7; Table 2). Overall, the method developed in this study revealed
 562 to be well balanced regarding FP and FN and these metrics had a small impact on
 563 cell detection rate (98.7%).

564

565 **Table 2:** Performance metrics for the detection methods developed in this study and
 566 in Liew et al. (2010).

Method	Min FP	Max FP	Avg FP	Min FN	Max FN	Avg FN	Avg CDR (%)
Liew et al.	1	11	5.1	139	530	274.7	82.5
Ours	5	39	22.5	2	80	38.2	98.7

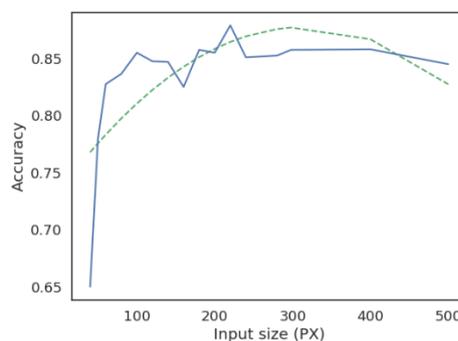
567 CDR – Cell detection rate; FP False Positive; FN False Negative.

568

569 5.2 Cells Classification

570 5.2.1 Finding the optimal input image size

571 After training 16 different classification models, we obtained a plot relating ROI
 572 size with cell classification accuracy (Fig 17). The trending line shows a steady
 573 increase in the accuracy as the ROI size increases up to 300px; after that point the
 574 quality decreases. Accordingly, input images with sizes between 200px and 300px
 575 should be preferred as they tend to produce better results.



576

577 **Fig. 17:** Testing accuracy according to the ROI size. The dashed green line
 578 represents the trend.

579

580 We chose 224×224px as the default input size for all the following tests. We
581 based this choice on the computational cost. Additionally, although the best result in
582 the test set was the one trained with 220×220px images, it was not possible to use
583 this input size because architectures like MobileNetV2 only had pre-trained weights
584 in the ImageNet dataset for sizes 128, 160, 192, or 224px (<https://bit.ly/2DhnLRb>). To
585 avoid the hassle of comparing models trained with varying input sizes, we opted for
586 using only 224×224px across all tests.

587

588 5.2.2 Comparing different CNN architectures

589 During the training of different architectures with the 224×224px input size, we
590 faced some difficulties. We noticed that models VGG 16 and 19 were unable to
591 converge over the used dataset (10% of our original dataset), even after trying
592 different LRs, loss functions and optimisers. Therefore, we decided to remove these
593 two models from the tests. Another model that caused problems during the training
594 stage was NasNet (Large). This model suffers from a known bug (see
595 <https://github.com/keras-team/keras/issues/8711#issuecomment-354585187>) and
596 there are bypasses for it, but we decided to discard it due to the large amount of time
597 required for retraining.

598 Table 3 presents some metrics for the computational performance of 11
599 models. The MobileNet model produced the best results regarding the number of
600 epochs required for reaching convergence and the average time per epoch.
601 MobileNet weights number was also among the lowest, only behind of its second
602 version.

603 To better understand the performance of the models, we also analysed the
604 loading time and the time to predict a batch of 100 images (Fig. 17). These additional

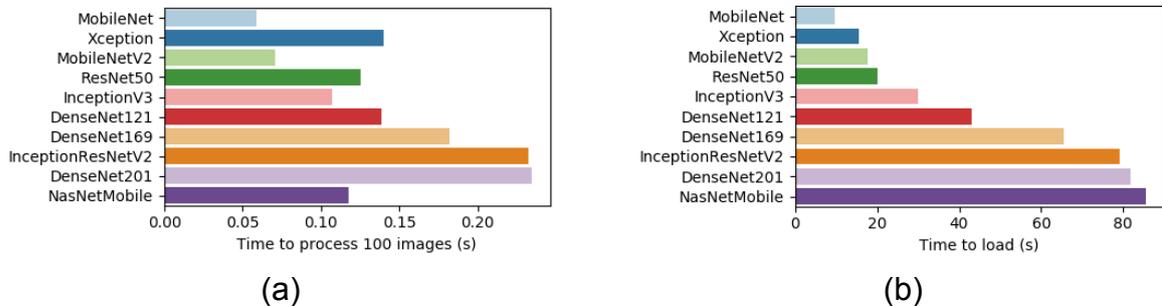
605 analyses were important to predict how the models behave, regarding time-efficiency
 606 performance, when used after training. Once again, MobileNet showed the best
 607 performance when compared with the other models. Even though MobileNet had
 608 fewer parameters than most models (Table 3), it was still the fast one to be loaded
 609 into memory.

610

611 **Table 3** Comparison among models regarding the training time and weights number.

Model name	Epochs to converge	Average time/epoch (min)	Number of weights
DenseNet121	16	362.77	8,094,279
DenseNet169	22	450.16	14,355,015
DenseNet201	21	577.14	20,296,263
InceptionResNetV2	19	606.31	55,917,799
InceptionV3	18	253.63	23,908,135
MobileNet	11	211.58	4,285,639
MobileNetV2	31	235.98	3,576,903
NASNet	28	2332.04	89,053,785
NasNetMobile	21	393.81	5,359,259
ResNet50	25	343.45	25,693,063
Xception	14	503.87	22,966,831

612



613 **Fig. 17:** (a) Comparison among models regarding (a) time to process 100 images of
 614 224×224px and (b) time to be loaded in memory.

615

616

To evaluate the functional quality of the classifications, we first compared the
 617 loss and accuracy of each model in their best epochs (Table 4). The model
 618 ResNet50 showed a higher capacity to predict training examples, but performed
 619 worse on the other sets, probably due to overfitting. DenseNet201 exhibited the best
 620 accuracy in the validation and test sets, but its predictions were made with less

621 confidence when compared with MobileNet, which had the lowest loss score in the
 622 validation and test sets.

623

624 **Table 4:** Comparison of loss and accuracy between models in different sets.

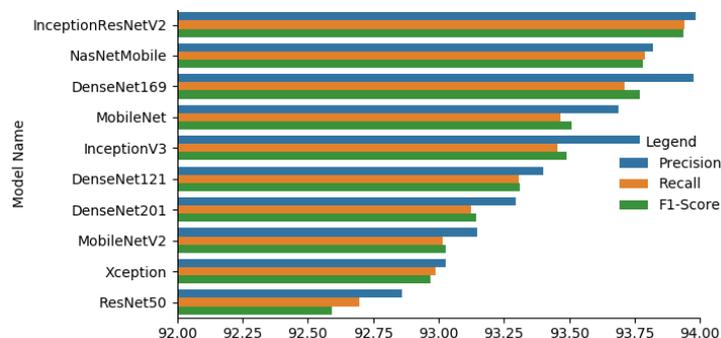
Model name	Loss train	Acc train	Loss val	Acc val	Loss test	Acc test
DenseNet121	0.00818	99.75%	0.05213	98.56%	0.25716	93.71%
DenseNet169	0.00159	99.95%	0.06365	98.58%	0.37087	93.12%
DenseNet201	0.00115	99.97%	0.05990	98.66%	0.31397	93.94%
InceptionResNetV2	0.00425	99.89%	0.05986	98.55%	0.29882	93.45%
InceptionV3	0.00415	99.90%	0.05594	98.58%	0.27237	93.47%
MobileNet	0.01563	99.57%	0.05106	98.48%	0.23944	93.31%
MobileNetV2	0.00942	99.69%	0.06468	98.57%	0.37828	93.02%
NasNetMobile	0.00162	99.94%	0.07417	98.56%	0.37836	93.79%
ResNet50	0.00033	99.99%	0.08845	98.44%	0.39329	92.99%
Xception	0.01173	99.73%	0.06574	98.54%	0.36011	92.70%

625

626 Sometimes accuracy can be biased by the majority class. To better
 627 understand this bias, we calculated precision, recall and F1-score for each model on
 628 the test set. As shown in Fig. 18, InceptionResNetV2 had the best performance,
 629 exhibiting a good balance between accuracy and recall. DenseNet201 had the best
 630 accuracy in the test set, yet it was positioned in seventh place for the F1-Score
 631 metric. This result suggests that DenseNet201 suffered by overfitting and favoured
 632 the majority class.

633

634



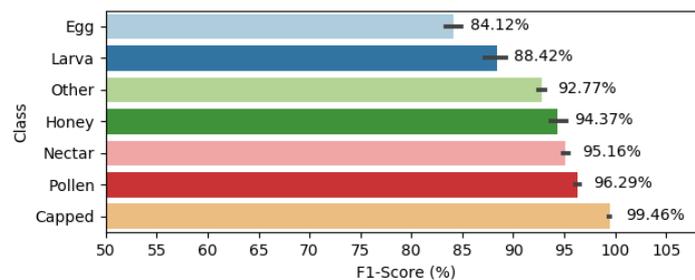
635

636 **Fig. 18:** Precision, recall and F1-score calculated using different models. The models
 637 are sorted by F1-Score.

638

639 In the per class analysis, we assessed how well each model classified the
 640 comb cell contents into the seven established classes. As shown in Fig. 19, the egg
 641 class exhibited the highest proportion of incorrect predictions (15.88%) whereas the
 642 capped brood class was very close to 100% correct predictions, based on F1-Score.
 643 This is an expected result as capped brood cells are simpler to classify due to their
 644 striking visual differences when compared with the remaining classes (Fig. 9). The
 645 egg class may have suffered from being the minority class. Eggs are small objects
 646 placed at the bottom of the cells and are easily confused with light reflections.
 647 Moreover, establishing thresholds for egg/empty and egg/young larva cells is a
 648 challenging endeavour. This issue will be further addressed in Section 5.2.5.

649



650

651 **Fig. 19:** Average F1-score per class.

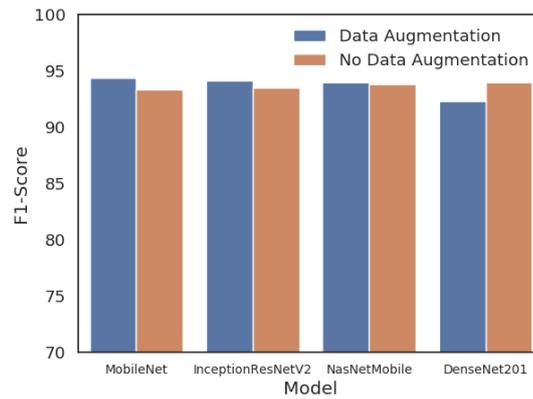
652

653

5.2.3 Data augmentation for cells classification

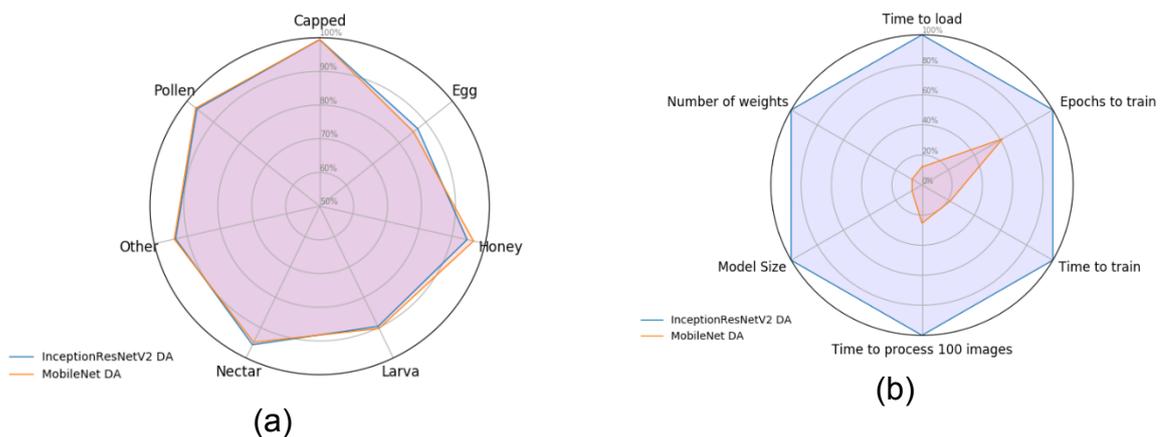
654

655 We chose four models to be trained using data augmentation (DA) for the
 656 following reasons: InceptionResNetV2 and NasNetMobile for having the best F1-
 657 score, MobileNet for having the lowest loss in the validation and test sets, and
 658 DenseNet201 for having the best accuracy in the validation and test sets. Fig. 20
 659 presents the F1-score calculated for the selected models trained on the dataset
 660 without and with DA. Except for DenseNet201, all other models showed a higher F1-
 661 score after training using DA. MobileNet outperformed more complex models, with an
 F1-score over 94%.



662
663 **Fig. 20:** Comparison of models trained with and without data augmentation.
664

665 Given that InceptionResNetV2 and MobileNet exhibited the highest F1-score
666 after training with DA, we decided to employ again the F1-Score to compare the
667 resources required to train and use these models for each class. As shown in Fig. 21,
668 while the difference between the two models in the quality of results per class is
669 modest, MobileNet revealed to be superior for computational resources across all
670 performance metrics. The superior performance of MobileNet can be attributed to the
671 lower number of trainable parameters. Hence, the model complexity is reduced,
672 regarding the number of training examples, counteracting overfitting.



673 **Fig. 21:** Comparison between the models MobileNet and InceptionResNetV2 using
674 (a) F1-score by class and (b) resources by DA model.
675

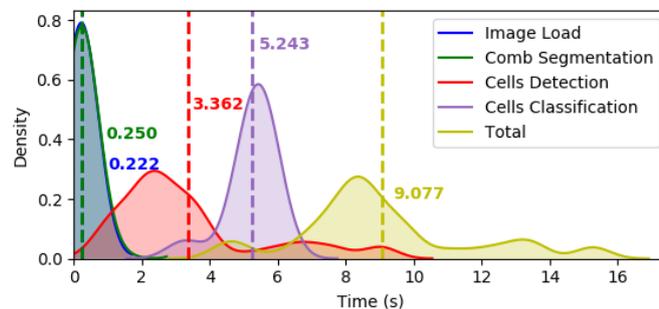
676 5.2.4 Comparison with methods from the literature

677 Here we compared our method with those reported in the literature, although
 678 this endeavor may not always be fair for three main reasons: (i) we classified a wider
 679 array of cell types, (ii) we do not have the same dataset, and (iii) some works
 680 employed different classifiers.

681 Cornelissen et al. (2009) compared their semi-automatic method of counting
 682 capped brood cells in comb images with the Liebefeld method. While annotations
 683 with the Liebefeld method took 26s per frame, the semi-automatic approach took 19s
 684 for image capturing plus 30s for image processing. This semi-automatic method
 685 consists of manual segmentation of the capped brood area followed by automatic
 686 count of cell number.

687 Fig. 22 presents the time distribution required by each phase of our cell
 688 detection and classification approach. The results were obtained from processing all
 689 61 images of the segmentation dataset using the scaled invariant detection algorithm
 690 and the MobileNet model trained with DA. The time required to fully process an
 691 image varied between ~4 and ~16s, with an average of 9.07s. Considering only the
 692 average value, the time to photograph a frame and process the image was 28.07s
 693 using our setup, which was about 2s slower than the Liebefeld method for capped
 694 brood cells.

695



696
 697
 698

Fig. 22: Time distribution to detect and classify all cells in a comb image.

699 Cornelissen et al. (2009) reported a correlation of 99.37% between the actual
700 and the predicted number of cells, which was substantially higher than the 90.85%
701 obtained with the Liebefeld method. Our approach correctly detected 98.7% of the
702 cells. Using CNNs, we obtained an F1-Score of 99.47% and 99.77% for the capped
703 brood class with the MobileNet-DA model and the InceptionResNetV2-DA model,
704 respectively.

705 Our approach of cell detection and classification overcomes some important
706 challenges pointed out by Colin et al. (2018). By using a CNN model, we were able to
707 distinguish capped honey from capped brood. Furthermore, by using a grid search
708 for finding good parameters for the CHT and the semantic segmentation, we
709 dismissed the user interaction for detecting the cells.

710 Rodrigues et al. (2016) obtained a precision of 99.04% and a recall of 97.2%
711 for the capped brood class. In our analyses, using the MobileNet DA model, we were
712 able to improve those metrics up to 99.47% and 99.41%, respectively.

713 Wang and Brewer (2013) reported a 97.4% hit rate with the
714 *HoneybeeComplete* commercial software developed for counting capped brood cells.
715 This value increased to 99.5% when the search area was delineated by the user in
716 the comb image. Herein, the MobileNet DA architecture produced a value for the
717 capped brood class very close (99.47%) to that obtained by the *HoneybeeComplete*
718 software, but without human assistance.

719 The commercial software *HiveAnalyzer* (Höferlin et al. 2013), which is able to
720 classify detected comb cells into seven classes, achieved 94.3% accuracy on cells
721 that were classified with high confidence (78% of 20,000 analysed cells). In this
722 study, we achieved 94.31% accuracy (very close to the 94.3% F1-Score value) using
723 the MobileNet model DA on 100% of the test set (53,914 analysed cells). When we

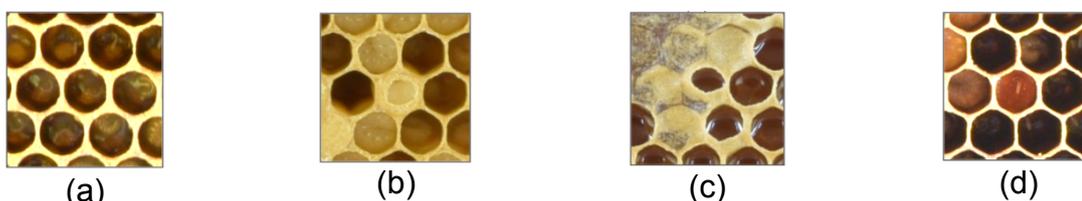
724 selected predictions with confidence >99.6% (corresponding to 42,410 cells and
725 78.66% of the dataset), we obtained 99.35% accuracy, a substantially higher rate
726 than that of Höferlin et al. (2013).

727

728 5.2.5 Further analyses on datasets creation and cells classification

729 The task of correctly classifying all the cells in a comb image is not trivial
730 because of the wide range of colours, shapes, and textures of cell contents and wax
731 types typically found in a hive. While working with the datasets, we realized that
732 comb classification is further challenging due to the impact of some factors on the
733 results quality. One such factor was related with cell contents. Cells with multiple
734 contents (e.g. pollen and egg) or cells with contents in a transition stage, such as
735 from larva to pupa (capped cell) or from egg to larva, revealed to be problematic (Fig.
736 23).

737



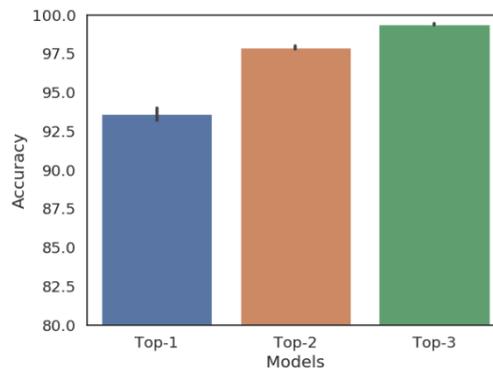
738 **Fig. 23:** (a) Transition from egg to young larva; (b) transition from old larva to
739 moulting, when cells will be capped; (c) transition from nectar to honey, when cells
740 will be capped; (d) central cell containing pollen and an egg.

741

742 The co-occurrence of different cell contents makes evaluation of the classifier
743 less precise, as there may be cases where it hits one of the classes but the
744 alternative class has been defined as the ground truth of the image. This problem
745 has been handled in competitions of image classification using Top-n accuracy
746 (Krizhevsky et al., 2012). With this methodology, the model earns credit for correctly
747 classifying the image in its Top N guesses. We evaluated our model using Top-2
748 accuracy, which corrected the cell content if the correct class was between the two

749 more likely predictions. By reprocessing the test set with the Top-2 accuracy method,
 750 the quality of our results improved 5% on average (Fig. 24).

751

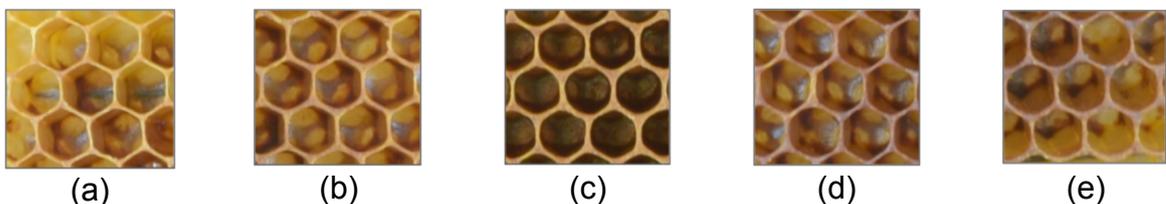


752

753 **Fig. 24:** Accuracies obtained for the models Top-1, Top-2 and Top-3.

754

755 Another factor affecting the quality of the results is related with the positions
 756 most annotated by the beekeeper. We noticed that there was an inverse relationship
 757 between the areas most frequently labelled by the beekeeper and the areas where
 758 most incorrect predictions occur. Due to the camera-optical behaviour, cells in
 759 different regions of the comb may display different areas of their interior, as illustrated
 760 in Fig. 25. This effect could impact cell classification if certain regions of the images
 761 were favoured during the annotation process. One way to alleviate this problem is to
 762 place the camera far from the frame and use a uniform light to diminish shadows
 763 during image capture.



764

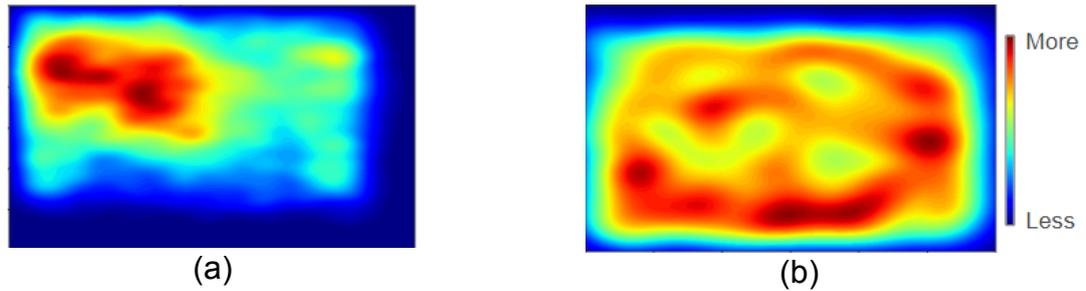
765 **Fig. 25:** Different cells interior captured in different regions due to lens effects (a)
 766 Upper left cell; (b) Upper right cell; (c) Central cell; (d) Lower left cell; (e) Lower right
 767 cell.

767

768 To test this effect, we assessed the distribution of the annotations across
 769 comb regions. To that end, using the annotations of the main dataset, we generated
 770 a heatmap plot displaying the areas of the comb that were preferentially annotated.

771 As shown in Fig. 26a, annotations were more concentrated in the upper left area of
 772 the comb, suggesting that models trained in this dataset would have a better
 773 classifying performance in that region. Next, we predicted all cells that were
 774 homogeneously annotated in the test set and generated a new heatmap showing the
 775 location of most of the wrong predictions (Fig. 26b).

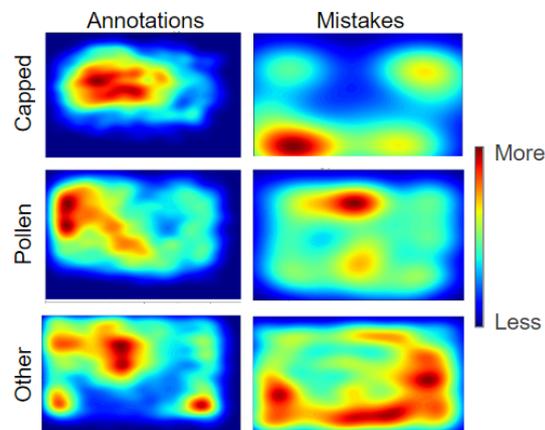
776



777 **Fig. 26:** Distribution of annotations in the comb. Comparison between (a) most
 778 annotated areas and (b) with more errors.

779

780 According to both heatmaps, incorrect predictions occurred mostly in the
 781 lower-right regions of the comb, and this pattern was inversely related with the
 782 regions where more annotations were made. This spatial pattern of annotations
 783 becomes more striking in the heatmaps generated by class (Fig. 27). Altogether,
 784 these results suggest that training a good classifier requires not only a large number
 785 of annotations but also a homogeneous distribution across the comb. Only then the
 786 annotations can inform the model, during the training, about the different angles that
 787 cells can present and help in the generalisation.



788

789 **Fig. 27:** Comparison among the most annotated areas and with more errors by class.

790

791 **5.3 DeepBee© software**

792

793

794 With all the methods developed and presented herein we built a software that
 795 we named *DeepBee©* (Fig. 28). This software allows the user to automatically
 796 process a batch of comb images. After processing the images, the user can view the
 797 results, change prediction labels, if needed, add and remove new detections, and
 798 export all results for further analysis into a spreadsheet like excel. *DeepBee©* is
 799 freely available at <https://avsthiago.github.io/DeepBee/>.



799

800 **Fig. 28:** *DeepBee©* software developed for the interaction of the users with the
 801 predictions.

802

803 **6 FINAL REMARKS**

804 In this study we developed a free software, *DeepBee©*, capable of
 805 automatically detecting and classifying comb cells. We demonstrated how we found a
 806 pre-processing pipeline able to enhance cells edges, filtering colour channels and
 807 equalizing small image regions using the CLAHE method. We demonstrated how we
 808 found parameters for the Circle Hough Transform that enables the method to detect
 809 cells in a comb even when it is difficult to visually distinguish the edges. We
 810 demonstrated that by applying the semantic segmentation technique it is possible to

811 remove false detections that may occur on the background. Although we obtained a
812 cell detection rate of 98.7%, we believe that the false positive rate may decrease by
813 training the semantic segmentation model with an input larger than 128×128px.

814 After we trained over thirty CNN models with different training techniques,
815 such as transfer learning and data augmentation, and comparing them using different
816 perspectives, we recommend MobileNet. While InceptionResNetV2 showed the best
817 results in our dataset, the time performance of MobileNet was superior, due to 93%
818 fewer weights. Using MobileNet, we achieved 94.3% of correctness with the metric
819 F1-score weighted over the seven classes. We believe this rate can be further
820 improved using annotations more evenly spread across comb images. The model
821 learned some human biases during the training and became better in classifying cells
822 in some comb regions in detriment of others.

823 To the best of our knowledge, the cell detection rate and the cell classification
824 accuracy of our model outperformed similar works reported in the literature. Future
825 work will focus on development of a service that enables users to process images
826 remotely. Using this web service, even devices with less power, such as
827 smartphones, will be able to run *DeepBee*©. To deal with low resolution images from
828 smartphones, we intend to create one composition of many images taken near to the
829 comb frame. With this web service, it will be possible to use detection corrections of
830 the users to improve future results by retraining the classifier.

831

832 **ACKNOWLEDGEMENTS**

833 This research was developed in the framework of the project “BeeHope - Honeybee
834 conservation centers in Western Europe: an innovative strategy using sustainable
835 beekeeping to reduce honeybee decline”, funded through the 2013-2014

836 BiodivERsA/FACCE-JPI Joint call for research proposals, with the national funders
837 FCT (Portugal), CNRS (France), and MEC (Spain).

838

839 REFERENCES

840 Alves, T., Pinto, M. A., Candido Junior, A., De Paula Filho, P. L., Rodrigues, P. J. S.,
841 Ventura, P., & Neves, C. (2019). DS-COMB-SEG-BEEHOPE, Mendeley Data,
842 v1 <http://dx.doi.org/10.17632/db35fj73x5.1>

843 Clevert, D.-A., Unterthiner, T., & Hochreiter, S. (2015). Fast and Accurate Deep
844 Network Learning by Exponential Linear Units (ELUs). *CoRR*, *abs/1511.0*.
845 Retrieved from <http://arxiv.org/abs/1511.07289>

846 Colin, T., Bruce, J., Meikle, W. G., & Barron, A. B. (2018). The development of honey
847 bee colonies assessed using a new semi-automated brood counting method:
848 Combcount. *PLoS ONE*, *13*(10), 1–14.
849 <https://doi.org/10.1371/journal.pone.0205816>

850 Cornelissen, B., Schmid, S., Henning, J., & Der, J. Van. (2009). Estimating colony
851 size using digital photography. *Proceedings of 41st International Apicultural*
852 *Congress*, 48.

853 Delaplane, K. S., Steen, J. Van Der, & Guzman-novoa, E. (2013). Standard methods
854 for estimating strength parameters of *Apis mellifera* colonies Métodos estándar
855 para estimar parámetros sobre la fortaleza de las colonias de *Apis mellifera*.
856 *Journal of Apicultural Research*, *52*(1), 1–12.
857 <https://doi.org/10.3896/IBRA.1.52.1.03>

858 Deng, J. D. J., Dong, W. D. W., Socher, R., Li, L.-J. L. L.-J., Li, K. L. K., & Fei-Fei, L.
859 F.-F. L. (2009). ImageNet: A Large-Scale Hierarchical Image Database. *2009*
860 *IEEE Conference on Computer Vision and Pattern Recognition*, 2–9.
861 <https://doi.org/10.1109/CVPR.2009.5206848>

862 Duda, R. O., & Hart, P. E. (1972). Use of the Hough Transformation to Detect Lines
863 and Curves in Pictures. *Commun. ACM*, *15*(1), 11–15.
864 <https://doi.org/10.1145/361237.361242>

865 EFSA AHAW Panel (EFSA Panel on Animal Health and Welfare), 2016. Scientific
866 opinion on assessing the health status of managed honeybee colonies
867 (HEALTHY- B): a toolbox to facilitate harmonised data collection. *EFSA Journal*
868 *2016*;14(10):4578, 241 pp. doi:10.2903/j.efsa.2016.4578

869

870 Emsen, B. (2006). Semi-automated measuring capped brood areas of honey bee
871 colonies. *Journal of Animal and Veterinary Advances*.

872 He, K., Zhang, X., Ren, S., & Sun, J. (2015a). Deep Residual Learning for Image
873 Recognition. *CoRR*, *abs/1512.03385*. Retrieved from
874 <http://arxiv.org/abs/1512.03385>

875 He, K., Zhang, X., Ren, S., & Sun, J. (2015b). Delving Deep into Rectifiers:

- 876 Surpassing Human-Level Performance on ImageNet Classification. *CoRR*,
877 *abs/1502.0*. Retrieved from <http://arxiv.org/abs/1502.01852>
- 878 Höferlin, B., Höferlin, M., Kleinhenz, M., & Bargaen, H. (2013). Automatic analysis of
879 apis mellifera comb photos and brood development. In *Association of Institutes*
880 *for Bee Research Report of the 60 th Seminar in Würzburg* (Vol. 44, p. 19).
881 Apidologie. Retrieved from
882 [https://www.springer.com/cda/content/document/cda_downloaddocument/AGIB -](https://www.springer.com/cda/content/document/cda_downloaddocument/AGIB-Abstracts2013_Final_Final.pdf?SGWID=0-0-45-1417002-p174076256)
883 [Abstracts 2013_Final_Final.pdf?SGWID=0-0-45-1417002-p174076256](https://www.springer.com/cda/content/document/cda_downloaddocument/AGIB-Abstracts2013_Final_Final.pdf?SGWID=0-0-45-1417002-p174076256)
- 884 Howard, A. G., Zhu, M., Chen, B., Kalenichenko, D., Wang, W., Weyand, T., ...
885 Adam, H. (2017). MobileNets: Efficient Convolutional Neural Networks for Mobile
886 Vision Applications. [https://doi.org/10.1016/S1507-1367\(10\)60022-3](https://doi.org/10.1016/S1507-1367(10)60022-3)
- 887 Huang, G., Liu, Z., & Weinberger, K. Q. (2016). Densely Connected Convolutional
888 Networks. *CoRR*, *abs/1608.06993*. Retrieved from
889 <http://arxiv.org/abs/1608.06993>
- 890 Kamilaris, A., & Prenafeta-Boldú, F. X. (2018). Deep learning in agriculture: A survey.
891 *Computers and Electronics in Agriculture*, *147*, 70–90.
892 <https://doi.org/https://doi.org/10.1016/j.compag.2018.02.016>
- 893 Kingma, D. P., & Ba, J. (2014). Adam: {A} Method for Stochastic Optimization. *CoRR*,
894 *abs/1412.6*. Retrieved from <http://arxiv.org/abs/1412.6980>
- 895 Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). ImageNet Classification with
896 Deep Convolutional Neural Networks. In F. Pereira, C. J. C. Burges, L. Bottou, &
897 K. Q. Weinberger (Eds.), *Advances in Neural Information Processing Systems*
898 *25* (pp. 1097–1105). Curran Associates, Inc. Retrieved from
899 [http://papers.nips.cc/paper/4824-imagenet-classification-with-deep-](http://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf)
900 [convolutional-neural-networks.pdf](http://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf)
- 901 Liew, L. H., Lee, B. Y., & Chan, M. (2010). Cell detection for bee comb images using
902 Circular hough transformation. *CSSR 2010 - 2010 International Conference on*
903 *Science and Social Research*, (C SSR), 191–195.
904 <https://doi.org/10.1109/CSSR.2010.5773764>
- 905 Oquab, M., Bottou, L., Laptev, I., & Sivic, J. (2014). Learning and Transferring Mid-
906 level Image Representations Using Convolutional Neural Networks. In *2014*
907 *IEEE Conference on Computer Vision and Pattern Recognition* (pp. 1717–1724).
908 <https://doi.org/10.1109/CVPR.2014.222>
- 909 Rodrigues, P., Neves, C., & Pinto, M. A. (2016). Geometric contrast feature for
910 automatic visual counting of honey bee brood capped cells. *EURBEE 2016: 7th*
911 *European Conference of Apidology*, *7*. Retrieved from
912 <http://hdl.handle.net/10198/17318>
- 913 Ronneberger, O., Fischer, P., & Brox, T. (2015). U-Net: Convolutional Networks for
914 Biomedical Image Segmentation. *CoRR*, *abs/1505.0*. Retrieved from
915 <http://arxiv.org/abs/1505.04597>
- 916 Sandler, M., Howard, A. G., Zhu, M., Zhmoginov, A., & Chen, L.-C. (2018). Inverted
917 Residuals and Linear Bottlenecks: Mobile Networks for Classification, Detection
918 and Segmentation. *CoRR*, *abs/1801.04381*. Retrieved from
919 <http://arxiv.org/abs/1801.04381>

- 920 Simonyan, K., & Zisserman, A. (2014). Very Deep Convolutional Networks for Large-
921 Scale Image Recognition. *CoRR*, *abs/1409.1556*. Retrieved from
922 <http://arxiv.org/abs/1409.1556>
- 923 Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., & Salakhutdinov, R. (2014).
924 Dropout: A Simple Way to Prevent Neural Networks from Overfitting. *Journal of*
925 *Machine Learning Research*, *15*, 1929–1958. Retrieved from
926 <http://jmlr.org/papers/v15/srivastava14a.html>
- 927 Szegedy, C., Ioffe, S., & Vanhoucke, V. (2016). Inception-v4, Inception-ResNet and
928 the Impact of Residual Connections on Learning. *CoRR*, *abs/1602.07261*.
929 Retrieved from <http://arxiv.org/abs/1602.07261>
- 930 Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., & Wojna, Z. (2015). Rethinking the
931 Inception Architecture for Computer Vision. *CoRR*, *abs/1512.00567*. Retrieved
932 from <http://arxiv.org/abs/1512.00567>
- 933 Thoma, M. (2016). A Survey of Semantic Segmentation. *CoRR*, *abs/1602.0*.
934 Retrieved from <http://arxiv.org/abs/1602.06541>
- 935 Tomasi, C., & Manduchi, R. (1998). Bilateral Filtering for Gray and Color Images. In
936 *Proceedings of the Sixth International Conference on Computer Vision* (p. 839--
937). Washington, DC, USA: IEEE Computer Society. Retrieved from
938 <http://dl.acm.org/citation.cfm?id=938978.939190>
- 939 Wang, M., & Brewer, L. (2013). New Computer Methods for Honeybee Colony
940 Assessments. *8th SETAC Europe Special Science Symposium*. Retrieved from
941 [http://sesss08.setac.eu/embed/sesss08/Larry_Brewer_-](http://sesss08.setac.eu/embed/sesss08/Larry_Brewer_-_New_Computer_Methods_for_Honeybee_Colony_Assessments.pdf)
942 [_New_Computer_Methods_for_Honeybee_Colony_Assessments.pdf](http://sesss08.setac.eu/embed/sesss08/Larry_Brewer_-_New_Computer_Methods_for_Honeybee_Colony_Assessments.pdf)
- 943 Yoshiyama, M., Kimura, K., Saitoh, K., & Iwata, H. (2011). Measuring colony
944 development in honey bees by simple digital image analysis. *Journal of*
945 *Apicultural Research*, *50*(2), 170–172. <https://doi.org/10.3896/IBRA.1.50.2.10>
- 946 Zoph, B., Vasudevan, V., Shlens, J., & Le, Q. V. (2017). Learning Transferable
947 Architectures for Scalable Image Recognition. *CoRR*, *abs/1707.07012*.
948 Retrieved from <http://arxiv.org/abs/1707.07012>
- 949 Zuiderveld, K. (1994). Graphics Gems IV. In P. S. Heckbert (Ed.) (pp. 474–485). San
950 Diego, CA, USA: Academic Press Professional, Inc. Retrieved from
951 <http://dl.acm.org/citation.cfm?id=180895.180940>
- 952

Declaration of interests

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

The authors declare the following financial interests/personal relationships which may be considered as potential competing interests: