1 **'Tailception': using neural networks for assessing tail lesions on pictures**

2 **of pig carcasses**

3 J. Brünger[1], S. Dippel[2a], R. Koch[1], C. Veit[2b]

4 *[1] Multimedia Information Processing Group, Computer Science Institute, University of*

5 *Kiel, Hermann-Rodewald-Str. 3, 24118 Kiel, Germany, jobr@informatik.uni-kiel.de,*

6 *rk@informatik.uni-kiel.de*

7 *[2] Institute of Animal Welfare and Animal Husbandry, Friedrich-Loeffler-Institut,*

8 *Dörnbergstr. 25/27, 29223 Celle, Germany, Sabine.Dippel@fli.de,*

9 *christina.maria.veit@nmbu.no*

10 *[a] Corresponding author: Sabine.Dippel@fli.de, phone: +49 5141 3846-200*

11 *[b] Present address: Department of Food Safety and Infection Biology, Faculty of*
12 *Veterinary Medicine, Norwegian University of Life Sciences, Ullevålsveien 72, 0454*
13 *Oslo, Norway*
14

15

16 Johannes Brünger, Sabine Dippel, Christina Veit: authors contributed equally to the

17 research and article preparation

18

19 **Short title**: Assessing tail lesions from pig carcase pictures

20

21 **Abstract**

22 Tail lesions caused by tail biting are a widespread welfare issue in pig husbandry.

23 Determining their prevalence currently involves labour intensive, subjective scoring

24 methods. Increased societal interest in tail lesions requires fast, reliable and cheap

25 systems for assessing tail status. In the present study, we aimed to test the reliability

26 of neural networks for assessing tail pictures from carcasses against trained human

27 observers. Three trained observers scored tail lesions from automatically recorded

pictures of 13 124 pigs. Nearly all pigs had been tail docked. Tail lesions were classified using a 4-point score (0 = no lesion, to 3 = severe lesion). In addition, total tail loss was recorded. Agreement between observers was tested prior and during the assessment in a total of seven inter-observer tests with 80 pictures each. We calculated agreement between observer pairs as exact agreement (%) and prevalence-adjusted bias-adjusted kappa (PABAK; value 1 = optimal agreement). Out of the 13 124 scored pictures, we used 80% for training and 20% for validating our neural networks. As the position of the tail in the pictures varied (high, low, left, right), we first trained a part detection network to find the tail in the picture and select a rectangular part of the picture which includes the tail. We then trained a classification network to categorise tail lesion severity using pictures scored by human observers whereby the classification network only analysed the selected picture parts. Median exact agreement between the three observers was 80% for tail lesions and 94% for tail loss. Median PABAK for tail lesions and loss were 0.75 and 0.87, respectively. The agreement between classification by the neural network and human observers was 74% for tail lesions and 95% for tail loss. In other words, the agreement between the networks and human observers were very similar to the agreement between human observers. The main reason for disagreement between observers and thereby higher variation in network training material were picture quality issues. Therefore, we expect even better results for neural network application to tail lesions if training is based on high quality pictures. Very reliable and repeatable tail lesion assessment from pictures would allow automated tail classification of all pigs slaughtered, which is something that some animal welfare labels would like to do.

53

**Implications**

Lesions caused by tail biting are a big welfare problem in pig production. Pigs reaching slaughter without tail lesions could be rewarded with premium payments, but this requires reliable lesion assessment in large numbers of pigs. We showed, that neural networks could help in automating this assessment.

59

**Introduction**

Tail biting is a widespread welfare problem in pig husbandry during which pigs manipulate the tails of their group mates with their mouth. This results in tail lesions of different degrees of severity, from superficial skin lesions over deep lesions to completely bitten-off tails (Taylor *et al.*, 2010). Tail biting is influenced by multiple risk factors which makes it difficult to prevent (EFSA, 2007). Cutting off the tails shortly after birth (tail docking) greatly reduces the risk of tail biting occurring later in life (EFSA, 2007). In the EU, tail docking is only allowed in exceptional cases (Council Directive 2008/120/EC; Council of the European Union, 2008) but nevertheless frequently applied. This discrepancy has led to a formal complaint to the European Commission (Marzocchi, 2014), which in turn caused increased public awareness and political pressure. As a result, animal welfare labels started to include tail status as a criterion ("Für mehr Tierschutz", Germany; "Beter leven", The Netherlands; "Bedre dyrevelfærd", Denmark) and programmes were launched which pay a premium for non-docked pigs (e.g. "Ringelschwanzprämie" by the German state Lower Saxony: 16.50 € per slaughter pig with not docked, not injured tail; ML Niedersachsen, 2015). The status of a pig tail is therefore now economically relevant. At the same time, large numbers of pig tails have to be evaluated. Thus, there is a need for fast, reliable, valid

78  and cheap systems to assess tail status. Currently, tail status for the German

79  "Ringelschwanzprämie", for example, is assessed by more or less trained observers

80  who travel to farms and walk through pens where they score pig tails (oral information;

81  S. Dippel). Assessing pigs on multiple farms in short periods of time requires

82  considerable resources in terms of time and money for travelling and assessment, with

83  added biosecurity risk through entering pens. Furthermore, tail lesion scoring by

84  multiple observers has a strong subjective component (Mullan *et al.*, 2011). Tails can

85  also be scored with minimal logistical input at slaughterhouses, where pig carcasses

86  are already inspected for signs of disease or severe injury. Studies have investigated

87  possible integration of tail lesion scoring in this inspection but found significant

88  influences of e.g. inspector work shift (Teixeira *et al.*, 2016).

89  Neural networks could be a low-cost, objective and indefatigable alternative to human

90  observers. Their development distinctly improved automated object recognition in

91  images (Russakovsky *et al.*, 2015) and they have already been used for e.g.

92  classification of hams (Muñoz *et al.*, 2015). Some attempts have been made at

93  developing automated assessment of lesions on slaughter carcasses using various

94  forms of algorithms. To our knowledge, the only published systems are a system for

95  assessing footpad dermatitis in broilers (Vanderhasselt *et al.*, 2013) and a system for

96  recording presence or absence of tail and ear lesions in pigs (Blömke and Kemper,

97  2017). However, many research and industry institutions are still struggling with the

98  reliability of their systems, which are mostly based on linear algorithms. The aim of the

99  present study was to test the reliability of neural networks for assessing tail lesions

100 from pictures of pig carcasses.

## Methods

*Tail pictures*

Tail pictures were taken of all pigs slaughtered on six days between March 27 and May 12, 2017 in one abattoir in North-Western Germany. Two synchronized RGB (red-green-blue) cameras automatically photographed tails from two dorsal angles after scalding and dehairing (cameras: UI-5480RE-C-HQ rev.2, lenses: IDS 25 HB Tamron Focal Length 12 mm, casing: Videotec Type NXM; all by IDS Imaging Development Systems, Obersulm, Germany). The two angles were stitched together in a single picture per pig. Lighting was provided by standard fluorescent lamps (tubes) with luminous colour 840 (cold white). Four double tubes were installed at the height of the carcase transport rails and provided light from above at distances of approximately 1 and 2 m from the carcase. One additional double tube was installed at the height of the back of the carcase with a distance of 2.8 m in order to reduce shadows from the top lights.

A total of 100,000 pictures were taken during the six days, out of which approximately 90% showed tails without lesions. As the aim was to determine agreement across all lesion severities, which may be influenced by unequal severity prevalences (Kottner *et al.*, 2011), we deleted a random sample of pictures without lesions in order to equalize the distribution of severity classes. For this, the most experienced observer screened pictures by recording hour (pictures from the same hour had been saved in one folder) in order to estimate the respective severity prevalences. She then first deleted all blurry pictures and then deleted every second picture without lesions until roughly similar proportions of pictures with different lesion severities were left in each folder. All pictures left were used for human observer training or training and testing the neural networks, respectively.

5

126 The first author made a software tool for picture scoring utilizing the OpenCV-library

127 (OpenCV team, 2018) which allowed observers to look at an image and directly enter

128 the scores. In addition, observers marked the position of the anal drill hole with a

129 mouse click in both angles. Pictures were brightened using IrfanView$^{©}$ (version 4.44)

130 and assessed on screens calibrated with dccw.exe (Windows®).

131 *Human assessment of tail pictures*

132 After training inter-human agreement, a total of 13 124 pictures scored by three

133 human observers was used for training and testing the neural networks.

134 *Scoring key*

135 We scored tail lesions on a scale from 0 to 3 and tail losses as presence (1) or absence

136 (0) of total tail loss (Figure 1). Discolouration at the tail base was not taken into account

137 because in direct observations it seemed to be associated with brushing during

138 scalding rather than with biting. Different degrees of partial tail loss could not be

139 assessed because of tail docking.

*Observers and training*

Pictures were scored by three observers in order to distribute the workload. Observers were chosen based on availability and previous experience with scoring tail lesions. One observer had experience in scoring tails on pictures from carcasses, one observer had experience in scoring tails on live pigs and one observer was naïve regarding scoring of pig tails. Observers trained by discussing and scoring tail pictures and tested their agreement at regular intervals using 80 unknown pictures for each test. The pictures were preselected by the last author (who led observer training and tests) to make sure, each test batch contained several pictures for each of the scores. We calculated agreement between observer pairs as exact agreement (%) and prevalence-adjusted bias-adjusted kappa (**PABAK** $= [(k*p)-1]/(k-1)$ where k = number of categories and p = proportion of matchings). PABAK values > 0.6 to 0.8 were regarded as satisfactory to good agreement and values > 0.8 as very good agreement (Fleiss *et al.*, 2003). Before picture assessment started, five inter-observer tests were required until satisfactory agreement was achieved. Two inter-observer tests were performed during the assessment to monitor potential drifts.

*Neural network assessment of tail pictures*

There are different approaches regarding the respective proportions of training and validation pictures. Many image datasets supplied for developing visual recognition systems use 95% training and 5% validation pictures (Russakovsky *et al.*, 2015). However, if the pictures (or mathematical outcome parameter) are highly variable such as the appearance of tail lesions in pictures, somewhat larger validation data sets in the range of 20% are recommended (Dohoo *et al.*, 2012) and used (e.g. image

163  datasets CIFAR-10 and CIFAR-100[1] or MNIST[2]). This is why out of the 13 124 scored

164  pictures we used 10 499 (80%) for training and 2 625 (20%) for subsequent validation

165  of the networks (Table 1).

166  *Localization of the tail region*

167  In order to train a classification network properly, it is important to use only relevant

168  picture sections as input. As the position of the tail varied from picture to picture, we

169  first trained a part detection network to locate the relevant region in each picture before

170  it was handed to the classification network. The part detection network (Figure 2) was

171  based on the idea from Bulat and Tzimiropoulos (2016) and realized using a fully

172  convolutional residual layer (ResNet)-50 backbone (He *et al.*, 2016). To preserve the

173  local information of the input data, we extracted, scaled up and added the feature maps

174  after the 7th (8-fold downsampling), 13th (16-fold downsampling) and 16th (32-fold

175  downsampling) building block of the ResNet before applying the pixelwise sigmoid-

176  loss. We initialized the network with pretrained Imagenet weights (Russakovsky *et al.*,

177  2015) and fine-tuned it for 30 epochs with the Adam-optimizer (Kingma and Ba, 2015)

178  at a learning-rate set to 0.0001. In order to subjectively verify that the network used the

179  tail-region to identify the injury patterns we used the *Image-Specific Class Saliency*

180  *Visualisation* from Simonyan *et al.* (2014).

181  *Classification of tail lesion and tail loss*

182  The part detection network predicted the location of the anal drill hole, which was then

183  used to position the region-of-interest window. The original pictures were scaled down,

184  so that the selection window for each angle covered 320 x 256 px. The two windows

185  for the two angles joined together resulted in the input of 320 x 512 px for the classifier

[1] https://www.cs.toronto.edu/~kriz/cifar.html
[2] yann.lecun.com/exdb/mnist/

186 network. For tail lesion classification, we used a modification of the standard Inception-
187 ResNet-v2 classifier network by Szegedy *et al.* (2017) for predicting the four tail lesion
188 scores in our dataset. To compensate for the large imbalance between scores, we
189 used sub-/oversampling until 4 000 training pictures were available for each score.
190 This meant that pictures from lesion score 2 and 3 were duplicated many times (Table
191 1). During training, the pictures were augmented online by rotating the two picture-
192 halves randomly (± 10 degrees) before cutting the region of interest and by applying
193 picture manipulations like adaptive noise, brightness-changes and blurring to the final
194 input-pictures. Again, we initialized the network with pretrained Imagenet weights and
195 fine-tuned it for 30 epochs with the Adam-optimizer (learning-rate set to 0.00001). We
196 used a categorical-crossentropy loss on the final four-classes-softmax activation. Due
197 to the pre-trained weights, the network started to overfit quickly so we applied early-
198 stopping. The tail loss classification was done on the same pre-processed input
199 pictures and the same classification network architecture, but with binary-crossentropy
200 loss on a single sigmoid activated decision-neuron.

201 **Results**

202 *Agreement between human observers*

203 For lesions, exact agreement between observer pairs ranged from 65 to 88% with 50%
204 of agreement values between 71 to 84% (first (Q25) to third (Q75) quartile; median =
205 80%; Figure 3). PABAK for lesions ranged from 0.56 to 0.84 with 50% of values
206 between 0.64 and 0.80 (median = 0.75). For tail loss, exact agreement ranged from 85
207 to 98% (Q25 to Q75: 90 to 95%, median = 94%) and PABAK ranged from 0.70 to 0.95
208 (Q25 to Q75: 0.80 to 0.90, median = 0.87).

209 *Agreement between neural network and human assessment*

210 The trained tail lesion classification network yielded an agreement of 74% with the

211 human observer scores, while agreement for tail losses was 95%. For tail lesions,

212 normalized values on the confusion matrix diagonal ranged from 0.59 to 0.85 with

213 uncertainty occurring on both sides of the diagonal (Figure 4).

214 The classification network mostly used information from the correct region for

215 classification (*Image-Specific Class Saliency Visualisation*; Figure 5). In pictures with

216 many optical structures in non-tail regions, especially reddish-coloured structures, the

217 network used more non-relevant pixels for its decision. Misclassifications were often

218 associated with shadows or overlapping structures (Figure 6).

219 **Discussion**

220 In the present study, human observers evaluated pictures of pig carcasses regarding

221 tail lesions and tail losses. The scored pictures were used to train and test neural

222 networks. Agreement between network and observer scores were similar to agreement

223 between human observers.

224 Agreement between human observers was acceptable in most tests for lesions and

225 good in most tests for tail loss, but fluctuated over time for both parameters. This was

226 mostly dependent on the prevalence of blurry pictures or lesions or losses on the

227 border between two categories in the test pictures. Even though lighting had been

228 optimised as much as possible, all pictures were more or less blurred due to high speed

229 of the carcasses on the line. In addition, most carcasses had discolourations and marks

230 from the scalding and dehairing process. The latter were also present on some tails

231 and thus interfered with assessment of low severity lesions. Overall, the greatest

232 difficulty was, where to distinguish between two lesion severity categories, i.e. "is this

233 still score 0 or already score 1". The issue remained despite training, due to the great

234   variation regarding colour and size along continuous gradients. This problem of

235   categorising continuous characteristics has been described before. In a study where

236   three observers scored 80 pictures and videos of sheep feet regarding lesions on a 5-

237   point scale (Foddai *et al.*, 2012), the width of the categories varied significantly

238   between observers, and categories also overlapped within observers. Similar results

239   were found for scoring lameness in sheep on an ordinal versus visual analogue

240   (continuous) scale (Vieira *et al.*, 2015). Therefore, assessment of lesions on a

241   continuous scale might be recommendable for reducing variation in training data by

242   improving agreement between observers who annotate training pictures.

243   In tasks of supervised learning like the one presented here, neural networks can only

244   be as good as the data they are trained with. This is why the disagreement between

245   human observers in our study is reflected in the uncertainty in the confusion matrix of

246   the tail lesion network. Using averaged annotations from several trained observers

247   (Muñoz *et al.*, 2015) could additionally improve training material quality. However,

248   neural networks also require large datasets in order to be trained on complex

249   parameters, such as tail lesions. Several observers re-scoring the same pictures

250   considerably increases labour input. Therefore, calculations on trade-off between large

251   numbers of training pictures annotated with greater variability by single observers

252   versus fewer pictures with average annotations with less variability should be made.

253   Nevertheless, improving human agreement is the necessary first step towards better

254   network assessment results. Based on our study, high quality pictures are a

255   prerequisite for good agreement. In addition, using continuous scales rather than

256   categorical scores might help to raise agreement for lesions above 90%.

257   Overall, the neural network assessment results in our study are very promising

258   because the agreement between network and human observers was similar to the

259   agreement between human observers. So far, only few studies investigated automatic

260   computerised injury assessment on carcasses. Vanderhasselt *et al.* (2013) tested a

261   system for assessing footpad dermatitis in broiler chickens. The maximum correlation

262   between scores assigned by humans and the automated system was 0.77. However,

263   even though there is less spatial variation regarding the position of broiler footpads

264   compared to pig tails on a line, the system found the relevant areas only in 86 of 197

265   recorded chickens (44%). Blömke and Kemper (2017) achieved much better results

266   with a system for automated assessment of presence or absence of ear and tail injuries

267   in pigs. Their system found the relevant areas in an average of 95% of pictures.

268   Sensitivity and specificity for detecting lesions were > 70% and > 94%, respectively,

269   for tail as well as for ear lesions (2 634 to 2 684 pigs). Only presence or absence of

270   lesions were assessed. Neither the threshold for lesion detection nor the algorithms for

271   picture analysis were reported yet.

272   **Conclusions**

273   Neural networks can assess tail lesions in pictures from slaughter pigs with a

274   reliability comparable to human observers. If supervised learning is used, high quality

275   training material (i.e. pictures) is necessary for achieving good network results. In

276   order to be able to generalise such complex parameters like tail lesions, neural

277   networks require large numbers of training pictures with equal representation of

278   different severities. Using continuous lesion severity scales instead of predefined

279   categorical scores might help to make the system more repeatable and versatile. In

280   sum, neural network analysis of tail pictures poses a promising technique which

281   might allow all pigs in a welfare label to be scored for tail lesions with little labour

282   input.

283 **Acknowledgments**

293 **Declaration of interest**

297 **Ethics committee**

298 Pig-related data in this study were collected without causing harm to the animals for

299 the purpose of the study. All experimental work was conducted in accordance with

300 relevant national legislation and approval by an ethics committee was not required.

301 **Software and data repository resources**

302 Data are available from the authors upon reasonable request.

303 **List of references**

304 Blömke L and Kemper N 2017. Automated assessment of animal welfare indicators in pigs at
305 slaughter. In Proceedings of the 12th International Symposium on the Epidemiology and
306 Control of Biological, Chemical and Physical Hazards in Pigs and Pork (SAFEPORK), 21-24
307 August 2017, Foz do Iguacu, Brasil, Foz do Iguacu, Brasil, pp. 241-244.
308 Bulat A and Tzimiropoulos G 2016. Human Pose Estimation via Convolutional Part Heatmap
309 Regression. In Proceedings of the 14th European Conference on Computer Vision (ECCV),
310 11–14 October 2016, Amsterdam, The Netherlands, pp. 717-732.

311     Council of the European Union 2008. Council Directive 2008/120/EC of 18 December 2008
312     laying down minimum standards for the protection of pigs. Official Journal of the European
313     Union L47, 5-13.
314     Dohoo I, Martin W and Stryhn H 2012. Methods in epidemiologic research. VER Inc.,
315     Charlottetown, Prince Edward Island, Canada.
316     EFSA 2007. The risks associated with tail biting in pigs and possible means to reduce the
317     need for tail docking considering the different housing and husbandry systems. The EFSA
318     Journal 611, 1-13.
319     Fleiss JL, Levin B and Paik MC 2003. The measurement of interrater agreement. In
320     Statistical methods for rates and proportions (ed. JL Fleiss, B Levin and MC Paik), pp. 598-
321     626, Wiley Interscience, Hoboken, NJ, United States of America.
322     Foddai A, Green LE, Mason SA and Kaler J 2012. Evaluating observer agreement of scoring
323     systems for foot integrity and footrot lesions in sheep. BMC Veterinary Research 8, 65.
324     He K, Zhang X, Ren S and Sun J 2016. Deep residual learning for image recognition. In
325     Proceedings of the 29th IEEE Conference on Computer Vision and Pattern Recognition
326     (CVPR), 26 June – 1 July 2016, Las Vegas, Nevada, United States of America, pp. 770-778.
327     Kingma DP and Ba JL 2015. Adam: A Method for Stochastic Optimization. In Proceedings of
328     the 3rd International Conference for Learning Representations, 7 - 9 May 2015, San Diego,
329     CA, United States of America, p. abs/1412.6980.
330     Kottner J, Audigé L, Brorson S, Donner A, Gajewski BJ, Hróbjartsson A, Roberts C, Shoukri
331     M and Streiner DL 2011. Guidelines for Reporting Reliability and Agreement Studies
332     (GRRAS) were proposed. Journal of Clinical Epidemiology 64, 96-106.
333     Marzocchi O 2014. Routine tail-docking of pigs. European Union, Brussels, Belgium.
334     ML Niedersachsen 2015. Agrarminister Meyer: Ringelschwanzprämie startet mit 16,50 Euro
335     [19.06.2015]. Retrieved on 06.02.2018 from
336     http://www.ml.niedersachsen.de/service/pressemitteilungen/agrarminister-meyer-
337     ringelschwanzpraemie-startet-mit-1650-euro-134624.html
338     Mullan S, Edwards SA, Butterworth A, Whay HR and Main DCJ 2011. Inter-observer
339     reliability testing of pig welfare outcome measures proposed for inclusion within farm
340     assurance schemes. The Veterinary Journal 190, e100-e109.
341     Muñoz I, Rubio-Celorio M, Garcia-Gil N, Guàrdia MD and Fulladosa E 2015. Computer
342     image analysis as a tool for classifying marbling: A case study in dry-cured ham. Journal of
343     Food Engineering 166, 148-155.
344     OpenCV team 2018. Open Source Computer Vision Library 3.3. Retrieved on 06.02.2018
345     from https://opencv.org/
346     Russakovsky O, Deng J, Su H, Krause J, Satheesh S, Ma S, Huang Z, Karpathy A, Khosla
347     A, Bernstein M, Berg AC and Fei-Fei L 2015. ImageNet Large Scale Visual Recognition
348     Challenge. International Journal of Computer Vision 115, 211-252.
349     Simonyan K, Vedaldi A and Zisserman A 2014. Deep Inside Convolutional Networks:
350     Visualising Image Classification Models and Saliency Maps. arXiv:1312.6034.
351     Szegedy C, Ioffe S, Vanhoucke V and Alemi AA 2017. Inception-v4, Inception-ResNet and
352     the Impact of Residual Connections on Learning. In Proceedings of the Thirty-First AAAI
353     Conference on Artificial Intelligence (AAAI-17), 04.-09.02.2017, San Francisco, California
354     USA, pp. 4278-4284.
355     Taylor NR, Main DCJ, Mendl M and Edwards SA 2010. Tail-biting: A new perspective. The
356     Veterinary Journal 186, 137-147.
357     Teixeira DL, Harley S, Hanlon A, O'Connell NE, More SJ, Manzanilla EG and Boyle LA 2016.
358     Study on the association between tail lesion score, cold carcass weight, and viscera
359     condemnations in slaughter pigs. Frontiers in Veterinary Science 3, 24.
360     Vanderhasselt RF, Sprenger M, Duchateau L and Tuyttens FAM 2013. Automated
361     assessment of footpad dermatitis in broiler chickens at the slaughter-line: Evaluation and
362     correspondence with human expert scores. Poultry science 92, 12-18.
363     Vieira A, Oliveira MD, Nunes T and Stilwell G 2015. Making the case for developing
364     alternative lameness scoring systems for dairy goats. Applied Animal Behaviour Science
365     171, 94-100.
366

367 **Tables**

368 Table 1: Number of pig carcase pictures scored by human observers and used for

369 training and validating neural networks. Numbers are given for each score assigned

370 by human observers for tail lesion and tail loss, respectively (Figure 1). Tail loss  was

371 only scored as present or absent. Out of the 13 124 scored pictures, 80% were used

372 for training and 20% for subsequent validation of the networks. n.a. = not applicable.

| Score | Tail lesions | | Tail losses | |
|---|---|---|---|---|
| | Training | Validation | Training | Validation |
| 0 | 6052 | 1460 | 9469 | 2359 |
| 1 | 3905 | 1041 | 1030 | 266 |
| 2 | 457 | 108 | n.a. | n.a. |
| 3 | 85 | 16 | n.a. | n.a |

373
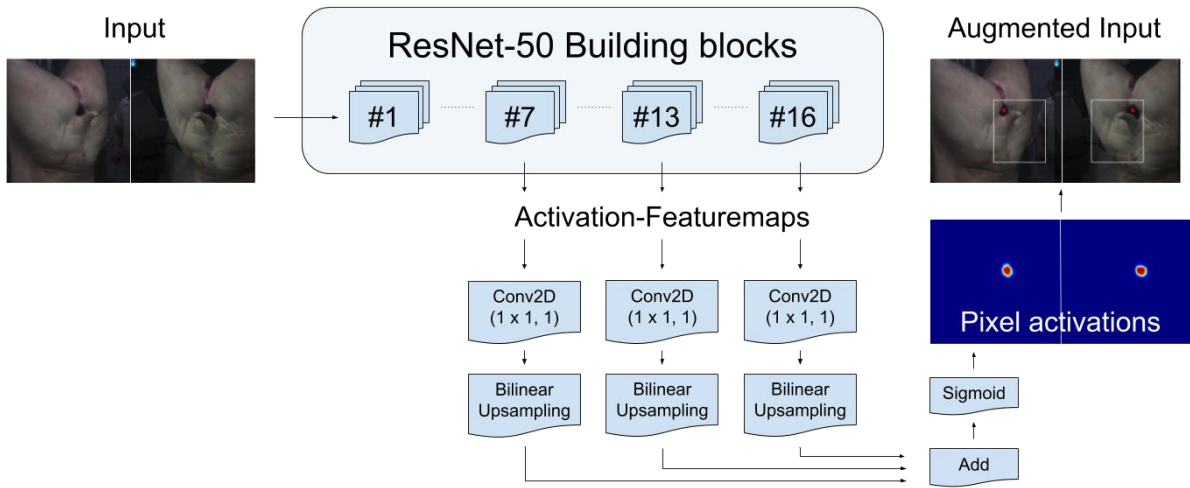
374 **List of figure captions**

375 Figure 1: Scoring key used for assessing tail lesions and total tail loss on pictures from

376 pig carcasses. Tail lesions and losses were scored independently of each other.

377 "Lesion" was defined as broken skin. The tail loss 1 picture shows the longest

378 remaining "stump" which was still considered as tail loss (longer stumps would be

379 classified as tail loss 0). Centimetres given are subjective estimates from a picture.

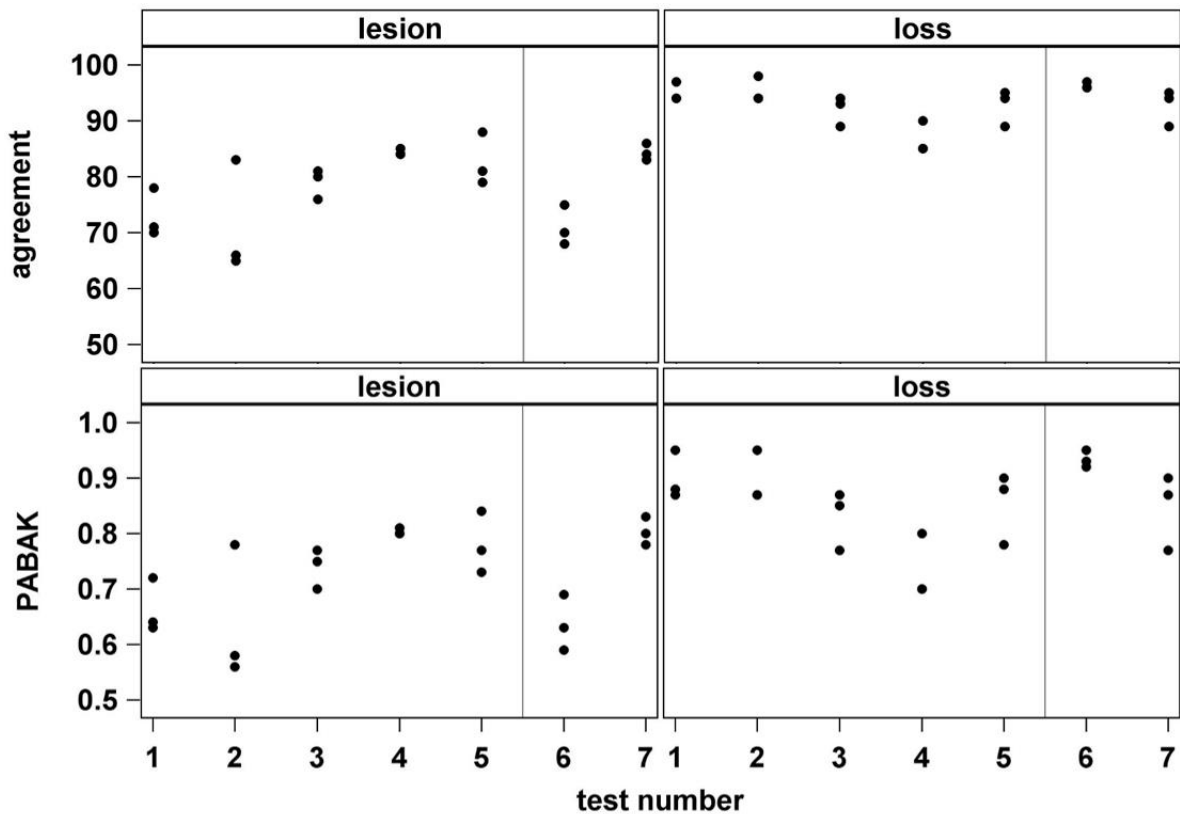| Score | Tail lesion | | Tail loss | |
|---|---|---|---|---|
| 0 | No visible lesion or reddish / violet / brownish discoloration the size of a pinhead. Skin looks intact |  | No loss or partial loss with more than a "stump" left (> 3 cm) |  |
| 1 | Lesion < tail diameter at respective location, with or without loss of tail substance |  | Total loss: only a "stump" protruding from tail base (≤ 3 cm) |  |
| 2 | Lesion ≥ tail diameter at respective location, with or without loss of tail substance |  | n.a. | |
| 3 | Tail tip with irregular outline (abrasion and / or elevations) in combination with dark reddish / brownish / blackish discoloration (necrosis) |  | n.a. | |

380

381

382

383 Figure 2: Architecture of a part detection network used for locating tails in pictures of

384 pig carcases. The network learns to activate pixels in the specified areas which can

385 then be used for positioning the region-of-interest windows for cutting out the relevant

386 picture section (tail) for subsequent classification.
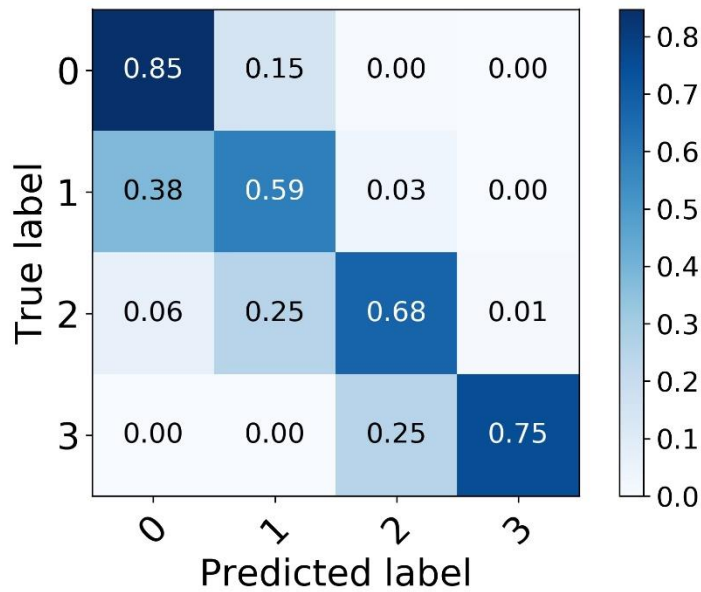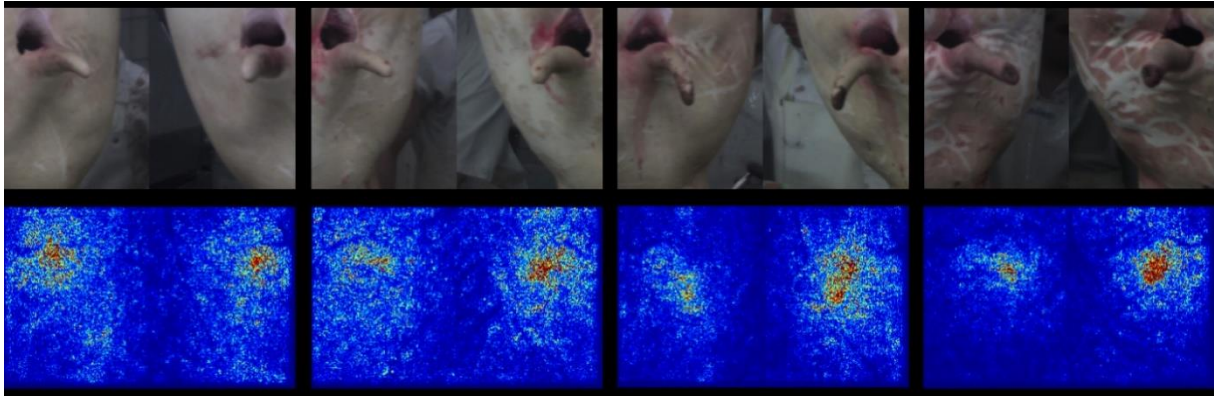
387



388

389 Figure 3: Results of inter-observer agreement tests of three human observers scoring

390    tail lesions or tail loss, respectively, from pig carcase pictures. Each dot represents

391    the exact agreement (%) or prevalence-adjusted bias-adjusted kappa (PABAK; range

392    0 to 1), respectively, for one observer-pair during one test (consecutive test number

393    on X-axis; n = 80 pictures per test). Grey vertical line = start of data collection.
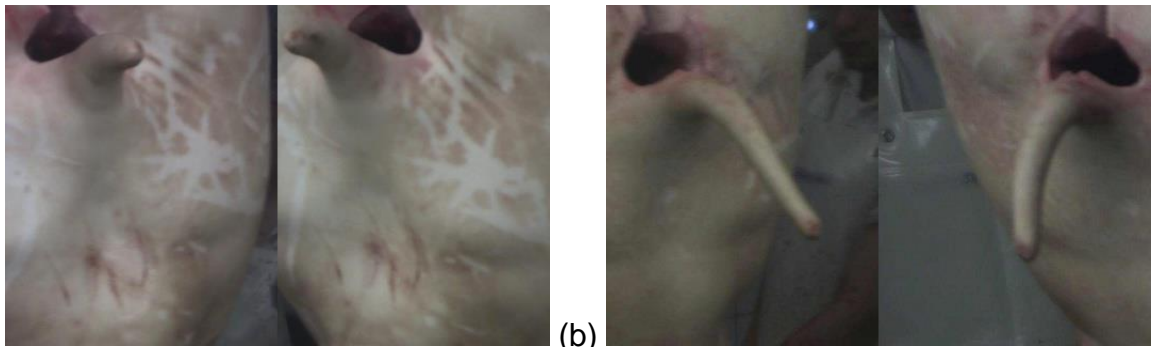
394



395

396    Figure 4: Normalized confusion matrix for the predictions of the tail lesion

397    classification network based on 13 124 pig tail pictures annotated by human

398    observers. True label = tail lesion severity score assigned by humans, Predicted label

399    = score predicted by neural network. The colouring indicates the normalised

400    distribution of numbers of pictures per cell.

401

402

403 Figure 5: Example pictures of slaughter pig tails from the verification of the tail lesion

404 severity classification network (top row). From left to right, pictures represent tail

405 lesion scores 0, 1, 2 and 3, respectively (Figure 1). The bottom row shows the

406 respective gradient-map made by the network, in which warmer colours indicate a

407 larger influence of the respective pixel on the final classification result.

408



409 (a)          (b)



410 (c)

411 Figure 6: Three examples for misclassification of pig tail lesion severity scores by the

412 network. Pictures (a) and (b) were assigned lesion score 1 by a human and lesion

19

413    score 0 by the network, picture (c) was assigned lesion score 3 by a human and

414    score 2 by the network.