



This document is a postprint version of an article published in Food Control© Elsevier after peer review. To access the final edited and published work see <https://doi.org/10.1016/j.foodcont.2019.06.019>

Document downloaded from:



1 **Computer image analysis for intramuscular fat segmentation in dry-cured ham**  
2 **slices using convolutional neural networks.**

3 \*I. Muñoz, P. Gou, E. Fulladosa,

4  
5 IRTA. Food Technology. Finca Camps i Armet 17121 Monells, Girona, Spain

6 \*Corresponding author: Israel Muñoz. israel.munoz@irta.es

7 **ABSTRACT**

8 Determination of intramuscular fat (IMF) content in dry cured meats is critical because it  
9 affects the sensory quality and consumer's acceptability. Recently, deep learning has  
10 become one of the most promising techniques in machine learning for image analysis.  
11 However, few applications in food products are found in the literature. This study presents  
12 the application of deep learning for the detection of intramuscular fat (IMF) in images of  
13 slices of dry cured ham. 8 convolutional neural networks (CNNs) have been studied and  
14 compared using segmented images (252 for training, 61 for validation and 62 for testing).  
15 The performance was compared to other simple CNNs. CNNs were able to segment IMF  
16 with an overall pixel accuracy of 0.99 and a recall and precision rates for fat near 0.82  
17 and 0.84, respectively, using a limited number of training images. However, performance  
18 is affected by the quality of the ground truth due to the difficulty of labelling correctly  
19 pixels.

20 **Keywords: Convolutional neural network, deep learning, intramuscular fat, image**  
21 **analysis, dry-cured ham**

22 **1. INTRODUCTION**

23 The amount of visible fat in dry-cured ham and distribution of fat streaks, affects  
24 palatability and consumers acceptability. Marbling is used as a visual cue by consumers  
25 to judge dry-cured ham quality. Although high IMF content is closely related to positive  
26 emotional responses during consumption of dry-cured ham (Lorido, Pizarro, Estévez &  
27 Ventanas, 2019), consumers prefer to purchase ham with moderate amounts of IMF,  
28 linked to positive nutritional and flavour characteristics (Morales, Guerrero, Aguiar,  
29 Guàrdia & Gou, 2013). This is a challenge for the industry, since the amount of IMF,  
30 even within the same breed, can widely vary. For the industry, it is of interest to

31 characterize online the IMF of slices of dry cured ham. This will allow the companies to  
32 segment the market, and offer products tailored to the consumers' needs.

33 Computer image analysis (CIA) is a reliable alternative for fast and non-destructive  
34 assessment of food characteristics such as colour, freshness, textural properties and other  
35 quality aspects. [Some applications include](#) the determination of marbling scores in pork  
36 meat ([Liu, Ngadi, Prasher & Gariépy, 2012](#)), the assessment of fish quality and freshness  
37 ([Dutta, Issac, Minhas & Sarkar, 2016](#)) and the quality [assessment](#) of pizza ([Sun &](#)  
38 [Brosnan, 2003](#)), cheese ([Caccamo et al., 2004](#)) and bread ([Srivastava, Vaddadi &](#)  
39 [Sadistap, 2015](#)). CIA has also been applied to grading of fruits and vegetables ([Blasco,](#)  
40 [Munera, Aleixos, Cubero & Molto, 2017](#)).

41 IMF detection using CIA is challenging because IMF cannot be easily characterized. For  
42 this reason, simple segmentation approaches are not useful and more sophisticated  
43 techniques are needed. For example, a segmentation-based approach was reported by  
44 [Jackman, Sun and Allen \(2009\)](#), which used K-means clustering to segment images of  
45 beef *Longissimus dorsi* muscle into background, lean muscle, and intramuscular fat areas.  
46 Results showed that IMF pixels were underestimated by 12.4% with respect to ground  
47 truth images. One of the most usual techniques for IMF detection is line detection  
48 algorithms. [Faucitano, Huff, Teuscher, Gariépy and Wegner \(2005\)](#) evaluated marbling  
49 by enhancing the colour contrast of pork meat samples using chemical pre-treatments and  
50 line detection algorithms. The authors did not check the accuracy of this approach. [Liu,](#)  
51 [Milan, Shen and Reid \(2012\)](#) and [Huang, Liu, Ngadi and Gariépy \(2013\)](#) used a line  
52 detection algorithm for determining a marbling score of pork loins and pork chops,  
53 respectively. [Qiao, Ngadi, Wang, Gariépy and Prasher \(2007\)](#) studied the potential of  
54 hyperspectral imaging techniques to assess pork quality and marbling levels using a  
55 hyperspectral imaging system and artificial neural networks. Both authors focused on the  
56 ability of these algorithms to predict marbling scores.

57 Recently, [Lohumi et al. \(2016\)](#) applied hyperspectral imaging for the characterization of  
58 intramuscular fat in beef. Several methods were evaluated and the accuracy ranged from  
59 91% to 96%. [Velázquez, Cruz-Tirado, Siche and Quevedo \(2017\)](#) segmented fat and  
60 classified the degree of marbling in beef from hyperspectral images using decision trees.  
61 Decision trees were able to reach an accuracy of 99.92% for the classification of lean and  
62 fat pixels during the construction of the tree (training). [Liu, Ngadi, Prasher and Gariépy](#)

63 (2018) segmented fat by automatically estimating the threshold between the lean and fat  
64 tissues. No information on accuracy was given.

65 In dry-cured ham, segmentation of IMF is more complex. The variation of dryness and  
66 colour across the slice, the presence of phosphates and tyrosine crystals and, in some  
67 cases, of nitrification rings make image segmentation more difficult. [Cernadas, Dur and  
68 Antequera \(2002\)](#), by using a multi-scale line detection framework for the recognition of  
69 fat streaks in the image, correctly classified 90% of the fat streaks with an acceptable rate  
70 of false positives. [Widiyanto et al. \(2013\)](#) segmented correctly IMF and lean in slices of  
71 dry-cured ham using fuzzy c-means and bias field estimation, obtaining a dice similarity  
72 coefficient of 0.94 for lean and 0.88 for IMF. [Muñoz, Rubio-Celorio, Garcia-Gil, Guardia  
73 and Fulladosa \(2015\)](#) and [Santos Garcés, Muñoz, Gou, Garcia-Gil and Fulladosa \(2014\)](#)  
74 used gradient-based techniques, such as discrete Fourier transform (DFT), but not  
75 evaluated the accuracy of the IMF estimation. However, new approaches for image  
76 analyses have been developed in the previous decade, which allow researchers to develop  
77 powerful algorithms for complex tasks. One of these new tools is deep learning  
78 ([Goodfellow, Bengio, Courville & Bengio, 2016](#)), in particular, deep convolutional  
79 networks. A convolutional neural network (also known as CNN or ConvNet) is a type of  
80 neural network used for deep learning in image applications. CNNs are used in a wide  
81 range of applications in image analysis. For example, object recognition ([He, Zhang, Ren  
82 & Sun 2016](#)), image classification ([Krizhevsky, Sutskever & Hinton 2012](#)) or image  
83 segmentation ([Badrinarayanan, Kendall & Cipolla, 2017](#)). The main advantage of this  
84 technique is that eliminates the need for hand-engineered filter design, as those are  
85 learned by the CNN itself. In the last few years, this technique has been applied to an  
86 increasing number of problems. In many cases, the performance of CNN has  
87 outperformed conventional CIA algorithms, becoming the state of the art solutions for  
88 many real applications. One of the applications is pixelwise classification, also known as  
89 semantic segmentation, which aims at assigning labels to pixels in an image ([Long,  
90 Shelhamer & Darrell, 2015](#); [Lin, Milan, Shen & Reid, 2017](#)). This is one the approaches  
91 that can be used to segment IMF in images.

92 Deep learning is a promising and very powerful tool to solve computer image problems.  
93 However, there are still very few applications of deep learning in the food sector.  
94 Recently, deep learning techniques have been applied to evaluate automatically the  
95 quality of fresh-cut lettuce ([Cavallo, Cefola, Pace, Logrieco & Attolico, 2018](#)), to assess

96 nutrient concentrations of commercially prepared pureed food (Pfisterer, Amelard, Chung  
97 & Wong, 2018), or to automate the segmentation of the skeleton of pigs using CT images  
98 (Kvam, Gangsei, Kongsro & Schistad-Solberg, 2018) or the detection of salmon muscle  
99 gaping (Xu & Sun, 2018). Other applications in food include food localization and  
100 recognition in images (Bolaños & Radeva, 2016).

101 This study aims at the segmentation of IMF in slices of dry-cured ham using deep  
102 learning. This problem has already been addressed using conventional image analysis  
103 techniques.

## 104 **2. MATERIAL AND METHODS**

### 105 **2.1. Sampling**

106 Ham slices were sampled from 190 dry-cured hams as it was described in Muñoz, Rubio-  
107 Celorio, Garcia-Gil, Guardia and Fulladosa (2015).

### 108 **2.2. Image acquisition**

109 Images were acquired with the photographic system depicted in Fig. 1. The exact  
110 methodology was described in Muñoz, Rubio-Celorio, Garcia-Gil, Guardia and Fulladosa  
111 (2015).

### 112 **2.3 Ground Truth**

113 Two regions of interest (ROIs) (both sides of a 1 cm thick slices) corresponding to the  
114 Biceps femoris (BF) muscles were manually selected from each image (Fig. 2). BF  
115 muscle was chosen because it is the biggest and the most representative muscle in dry  
116 cured ham slices. Besides, together with Semitendinosus (ST) muscle, it may have an  
117 considerable amount of intramuscular fat, which is also correlated to the ST muscle (the  
118 fattiest muscle). 375 ROIs were evaluated and 5 ROIs were discarded from the study due  
119 to defects on the surface (such as cuts and phosphate crystals) which made them  
120 unsuitable for the CIA. Patches of 64x64 pixels (one patch per image) were automatically  
121 extracted from the ROIs with three channels of information corresponding to R,G,B  
122 colour channels. All patches were treated as independent samples, as IMF distribution  
123 and colour of IMF and lean showed big differences in patches obtained from both sides  
124 of the same ham slice.

125 Next, reference images of correctly segmented IMF (Ground Truth) were obtained from  
126 these patches similarly to the methodology described in Muñoz, Rubio-Celorio, Garcia-

127 Gil, Guardia & Fulladosa, 2015) and (Santos Garcés, Muñoz, Gou, Garcia-Gil &  
128 Fulladosa, 2014). For each ROI, IMF was segmented using edge detection based on the  
129 discrete Fourier transform (DFT) (Rangayyan, 2004). DFT followed by a gaussian high  
130 pass filter with a cut-off frequency of 250 was applied to each image. After filtering, the  
131 images were transformed back using the Inverse Discrete Fourier Transform (IDFT). The  
132 real component of the transformed matrix was used for further processing. Pixels with  
133 values equal or below a threshold value were labelled as IMF. This threshold value was  
134 set manually. An expert in the field of food technology, trained for the sensory evaluation  
135 of foods and specially for dry-cured ham visual evaluation was responsible for adjusting  
136 the threshold values following the guidelines established for dry-cured ham by Claret,  
137 Guerrero, Guàrdia, Garcia-Gil and Arnau (2009). After this, most of the IMF was  
138 correctly segmented. However, several thresholding operations in combination with  
139 different logical operators (AND, OR, NOT) were applied to the image (combining the  
140 IDFT transform image and the RGB image) for the segmentation of still incorrectly  
141 segmented pixels. This work was also carried out by the trained food technologist and the  
142 threshold values adjusted accordingly. In some cases, even for a trained expert, it was  
143 difficult to decide whether a pixel should be labelled as fat or lean, in particular for small  
144 fat streaks and contour pixels due to the wide range of RGB values for fat and lean,  
145 structure of fat, etc

146 Small size patches (3x64x64 pixels) were used in order to have same size samples for  
147 training (BF muscles are different in shape and in the number of pixels) and speed up  
148 learning.

#### 149 **2.4 Convolutional neural network architecture**

150 The convolutional neural network (CNN) was trained to classify pixels into two different  
151 classes (class 0: lean, class 1: fat) using pixelwise classification (semantic segmentation).  
152 Ground Truths for images were determined as described in section 2.3 and were used as  
153 labelled images during training of the CNNs. In the CNN architecture used in this work,  
154 four of the most common types of operation in a CNN were used: convolution, non-linear,  
155 pooling and upsampling layers.

156 In Convolution layers, a filter (also known as kernel) performs the convolution operation  
157 over a matrix (images). Convolution can be thought as a sliding window function applied  
158 to a matrix. The number of parameters to be learned in these filters is equal to the number  
159 of elements of these filters (depth x height x width) (Fig. 3a). In this work (as in others

160 studies in the field), when referring to filters only the height and the width is given,  
161 whereas the depth can be obtained from the depth of the input matrix (image).

162 Non-linear layers are usually placed right after convolution layers. Non-linear layers  
163 perform a non-linear operation on the matrices resulting from the convolution operation,  
164 similar to the sigmoid function. The most common function in CNN is the rectifier linear  
165 function (ReLU) (Fig. 3a).

166 Pooling layer reduces the size of the image, also known as downsampling. Among the  
167 existing pooling layers, average and max pooling are the most common ones. Pooling  
168 layer consists of a sliding window function that moves over the matrix and takes the  
169 largest value in the window. The matrix is partitioned into several non-overlapping  
170 regions where the operation associated with pooling is applied. Pooling reduces the size  
171 of the matrix. In Fig. 3b, the max pooling is applied.

172 Upsampling layers can be considered as a kind of reverse convolution (Fig. 3b),  
173 sometimes denoted as deconvolution. Upsampling resizes an input matrix to the desired  
174 size by upsampling and interpolation (e.g. bilinear interpolation). In CNN, it can be used  
175 to resize the output of a CNN to the original size after convolutional and pooling  
176 operations have reduced the size of the original image.

177 Fig. 4 depicts the basic architecture used in this work, in the case of using 512 filters in  
178 the first layer. This architecture is based on the work by Long, Shelhamer and Darrell  
179 (2015), in which information from different layers of the CNN are combined to make  
180 predictions, and the VGG net (Simonyan & Zisserman, 2015), in which the number of  
181 filters increases with the depth of the network. Prior to the final selection of the CNN  
182 architecture used in this study, several parameters were evaluated, namely number of  
183 convolutional layers (1-4), kernel size (3x3,5x5,7x7) and number of filters.

184 In this architecture, a RGB patch (3x64x64 pixels) is convolved by 512 3x3 convolution  
185 filters (and depth 3, as the image has three channels: R, G and B) and zero padding is  
186 applied to ensure that after convolution the height and width of the image remains the  
187 same. Zero padding consists in adding “0” around the border of the matrix of data. For  
188 3x3 convolution filters, a zero padding of size 1 must be applied to ensure that the size  
189 does not change after convolving. Therefore, convolution, including padding, transforms  
190 the input image into 512 64x64 matrices. After convolution, the rectifier function (ReLU)  
191 is applied to each element of the obtained matrices and next, the 2x2 max pooling is

192 applied. A 2x2 pooling reduces the size of the matrix by a factor of 2 (i.e. from 64x64 to  
193 32x32), but it does not change the number of matrices (512). The whole structure  
194 (network layer) (Conv-ReLU-pool) is repeated 3 more times. At the end of the process,  
195 there are 4096 8x8 matrices. After each max pooling the number of matrices is increased  
196 by a factor 2 at the next convolution operation in order to keep the complexity of the  
197 network (Simonyan & Zisserman, 2015). After each pooling operation, an upsampling  
198 operation is applied, using bilinear interpolation, to obtain two matrices (two classes) with  
199 the original size (64x64 pixels). Additionally, an upsampling operation is applied to the  
200 matrices obtained at ReLU1. All 64x64 pixels obtained by upsampling at different layers  
201 are finally added together (2x64x64 pixels). According to Long, Shelhamer & Darrell  
202 (2015) the combination of information from different layers is equivalent to combine  
203 coarse, high level information with fine low layer information. This integration of  
204 information allows the network to predict finer details. Next, the output of the network  
205 (2x64x64 matrices) is passed through a softmax classifier. The output after the softmax  
206 classifier is a probability map having the same size as the input image (64x64) with each  
207 pixel having two values, the probability of belonging to class 0 (lean) and class 1 (fat).  
208 The class with the highest probability value is selected as the segmented class. The  
209 performance of this CNNs architecture is compared to other more simple CNNs  
210 architectures in which all upsampling operations are removed from the architecture with  
211 the exception of the last upsampling operation previous to the softmax classifier  
212 (Upsample 5 in Fig. 4).

213 In this study, different parameters of this architecture were studied, namely, the depth of  
214 the network (from 1 to 4 Conv-ReLU-pool layers) and the number of filters at Conv1  
215 (128 and 512). In the results and discussion section, network architectures will be denoted  
216 as 2\_128, first figure indicates the number of Conv-ReLU-pool structures (network  
217 layers) and the second figure indicates the number of filters at the first convolutional layer  
218 (Table 1). In total, 8 different combinations of depth of the network (1-4) and number of  
219 filters were studied (128, 256). According to this notation, Fig. 4 depicts a 4\_512  
220 architecture (4 Conv-ReLU-pool layers and 512 filters in layer 1). The same notation is  
221 used for the simple networks with the difference that only the last upsampling is included  
222 in the network (Upsample5 in Fig. 4). In this case, the number of filters at Conv 1 is 128  
223 and 512 and the depth of the network from 1 to 3. A subscript (s) has been added to denote  
224 a simple network (Table 1).



## 225 **2.5 Software and Hardware**

226 Matlab 2008b and its image processing toolbox (The MathWorks, Inc., United States)  
227 were used to select the ROI and segment IMF in images using the procedure described in  
228 the previous section.

229 Caffe was used to create, train, validate and test the CNN architecture. Caffe is a deep  
230 learnig framework and stands for Convolutional Architecture for Fast Feature Embedding  
231 (Jia et al., 2014). Results were processed using Python 2.7.13. The following parameters  
232 were used in this study in Caffe: Batch size 32, the base learning rate 1e-4, the momentum  
233 0.9, the weight decay 0.05. The learning rate policy was “inv” (learning rate decay over  
234 time) and the parameters for this policy were gamma 0.01 and power 0.5.

235 Caffe tries to minimize the multinominal logistic loss (also known as cross-entropy  
236 classification loss) and it was used to compute the error classification during training and  
237 optimization. Stochastic gradient descent was used for the optimization of the network.  
238 Each network was trained for 50,000 iterations. In Caffe the term iteration is used instead  
239 of epoch, for this reason the term iteration is used across the text. The equivalency  
240 between epoch and iteration is as follows:

241 
$$\text{Epoch\_index} = \text{floor}(\text{iteration\_index} \times \text{batch\_size}) / (\text{number of training data samples})$$

242 Convolution layers: Weights were initialized using “xavier” method. Bias were of type  
243 “constant” which initialises biases to zero. A learning rate multiplier of 1 was selected for  
244 the weights and a multiplier of 2 for the biases. Kernel sizes of 3x3 were used in this study  
245 and zero padding was of size 1. The stride was 1.

246 Upsampling layers: Upsampling layers used the “bilinear” method. For upsampling from  
247 64x64, 32x32, 16x16 and 8x8 to 64x64 a kernel size of 1,4,8,16, a stride of 1,4,8,16 and  
248 a zero padding of size 0,1,2,4 were used, respectively.

249 Prior to the tests several parameters of the network were tested and adjusted: base learning  
250 rate, batch number, momentum, weight decay, gamma and power. Once, these parameters  
251 were determined, all networks structures were trained using the same values

252 Training, validation and testing of CNNs was performed on a Z820 workstation with 512  
253 GB of RAM and 16 cores Intel Xeon ES-2687W at 3.10 GHz

## 254 **2.6 Testing**

255 2/3 of patches were randomly selected for training (252), 1/6 for validation (61) and 1/6  
 256 for testing (62) by assigning a random number to each image patch. Patches were assigned  
 257 to each group based on the value of the random number. This means that for the training  
 258 set a total of 1,032,192 pixels (252 images x 64 rows x 64 columns) were available for  
 259 training.

260 The following metrics were used to evaluate the performance on the test set:

261  $tp$ ,  $tn$ ,  $fp$ ,  $fn$  denote true positive, true negative, false positive and false negative  
 262 respectively. Positive class denotes fat, negative class denotes lean.

263 1) Overall pixel accuracy: percentage of pixels correctly predicted (fat and lean  
 264 pixels)

$$265 \quad Accuracy = \frac{t_p + t_n}{t_p + t_n + f_p + f_n}$$

266 2) Fat recall rate: rate of pixels correctly predicted as fat into the total number of  
 267 pixels labelled as fat.

$$268 \quad Fat\ recall\ rate = \frac{t_p}{t_p + f_n}$$

269 3) Fat precision rate: rate of pixels correctly predicted as fat into the total number of  
 270 pixels predicted as fat.

$$271 \quad Fat\ precision\ rate = \frac{t_p}{t_p + f_p}$$

272 4) Rate of false negatives near the areas predicted as fat (FnPFc rate): rate of false  
 273 negatives that are correctly predicted as fat ( $t'_p$ ) after applying a 3x3 dilation  
 274 operation on the pixels predicted as fat by the CNN.

$$275 \quad FnPFc\ rate = \frac{t'_p}{f_n}$$

276 5) Rate of false positives near the areas labelled as fat (FpLFc rate): rate of false  
 277 positives that are correctly predicted as non-fat ( $t'_p$ ) after applying a 3x3 dilation  
 278 operation on the pixels labelled as fat (ground truth).

$$279 \quad FpLFc\ rate = \frac{t'_p}{f_p}$$

280 At evaluation, special attention was given to segmentation of fat in the images. FnPFc  
 281 and FpLFc rates attempt to incorporate the uncertainty of manual classification  
 282 (classification of contour pixels) during the preparation of the ground truth images.

283 Dilation operations are used in computer vision for expanding the shapes contained in an  
284 image. The size of the expansion depends on the size of the operation (3x3 in this case)  
285 or the number of times the operation is applied (1 in this case).The application of a dilate  
286 operation on the pixels predicted or labelled as fat may incorporate this uncertainty into  
287 the evaluation of performance. The results presented for the different network  
288 architectures correspond to the iteration with the best overall pixel accuracy of the test set  
289 for the last 5.000 training iterations. Then, the learned parameters of the network were  
290 used to evaluate the test set. Training data was recorded every 500 iterations.

291 Moreover, the average time needed for the segmentation of images of the test set was also  
292 recorded.

293 After training, validation and testing, four representative images from the test set were  
294 segmented and analysed using the worst and the best performing (overall pixel accuracy)  
295 architectures and the segmentation was compared for the two architectures.

296 Results presented in the result and discussion section lack any statistical significance as  
297 the training, validation and testing was done only once, because of the long training times  
298 of the architectures studied in this investigation (4 months). This is quite common in many  
299 works in the field of CNN (i.e. Long, Shelhamer & Darrell, 2015; Ronneberger, Fischer  
300 & Brox, 2016; Lin, Milan, Shen & Reid, 2017 ) and many times comparisons are based  
301 on one single training (on training, validation and test sets) due to this limitation.

### 302 **3. RESULTS AND DISCUSSION**

303 [Fig. 5](#) shows the change in the multinomial logistic loss with the number of iterations for  
304 the training and test sets of the CNNs 3\_128 and 3\_512. These two CNNs were chosen  
305 as an example to study the learning of the network. Logistic Loss decreased more rapidly  
306 for CNN 3\_128 than for CNN 3\_512 due to the lower number of learnable parameters  
307 (1,478,914 vs. 23,610,368) of the network. After approximately 4000 and 6000 iterations  
308 for CNNs 3\_128 and 3\_512 respectively, the loss for the training and test set tended to  
309 decrease very slowly, even though the loss for the training set decreased more rapidly  
310 than for the test set. After around 25.000 iterations, the logistic losses barely changed.  
311 One of the reasons for this result is the learning rate decay and the convergence of the  
312 learning of the network. Overfitting was not observed, as the logistic loss for the test set  
313 did not increase with the number of iterations. However, training the networks for a larger  
314 number of iterations might have increased the logistic loss for the test set (not tested).

315 Logistic loss of the test set was lower for CNN 3\_512 than for CNN 3\_128. The larger  
316 number of learnable parameters of CNN 3\_512 may have captured better the complexity  
317 of the segmentation for this task. However, a 16 fold increase in the number of learnable  
318 parameters only brought about a small improvement in the performance of the network.  
319 For CNN 3\_128 logistic loss for the test set was only slightly lower than that of the  
320 training set. This result can be surprising, but it is not uncommon for small size sets of  
321 test data (62 images) as chance during random selection of training, validation and test  
322 sets may produce this result. The difference between the loss for the training and test set  
323 decreased with the number of iterations.

324 Table 2 shows the results for the simple CNNs and the CNN architectures developed for  
325 this work. Simple CNN architectures performed worse (performance, lower recall and  
326 precision rates for fat segmentation) than those architectures specifically conceived for  
327 this work with the same number of filters in the first layer. Results also showed that  
328 performance in simple CNN increased with the number of filters in the first layer, but  
329 decreased with the number of layers. This latter result is not clearly observed for the  
330 complex CNNs presented in Table 2. However, it seems that 1\_128 and 1\_512, performed  
331 worse than other architectures with more layers. This result can be specially observed for  
332 1\_512 vs 3\_512 and 4\_512, even though it cannot be considered conclusive due to the  
333 lack of statistical significance. As expected processing time was much lower for the  
334 simple architectures.

335 For the architectures conceived for this study, as the number of filters and/or the number  
336 of layers increase, the number of parameters to be adjusted increases (Table 2) and thus,  
337 the CNN is expected to fit more accurately the training set, improving overall pixel  
338 accuracy of the training set. However, in our study, the overall pixel accuracy was very  
339 high (0.988) in the most simple CNN architecture (1\_128) and increased to 0.991 in the  
340 most complex CNN architectures (3\_512 and 4\_512). These values are similar to those  
341 obtained in other works. During training, [Velázquez, Cruz-Tirado, Siche & Quevedo](#)  
342 [\(2017\)](#) obtained an accuracy of 0.9992 using decision trees for the segmentation of IMF  
343 in beef.

344 In general, increasing the number of filters and layers allows capturing better the  
345 complexity of the problem, but the overall pixel accuracy of the test set can decrease due  
346 to the well-known problem of overfitting ([Hawkins, 2004](#)), which is originated in models  
347 with more terms or more complicated approaches than necessary. In our study, the overall

348 pixel accuracy of the test set hardly increased with the number of filters, from 0.988 for  
349 CNNs with 1\_128 filters to 0.989 for CNNs with 3\_512 filters and 4\_512 and it did not  
350 change with the number of layers. The CNN with the highest overall pixel accuracy was  
351 3\_512. The overall accuracy tended to increase slightly with the number of learnable  
352 parameter. No drop in performance was observed with the increase of learnable  
353 parameters. Therefore, overfitting was not observed for this data and the studied  
354 architectures.

355 The overall pixel accuracy was highly influenced by the lean tissue segmentation of the  
356 CNNs, due to the large ratio of pixels corresponding to lean tissue in the images. The  
357 precision and recall rates of fat were also studied, as they give more accurate information  
358 than overall pixel accuracy on the performance of the CNNs for the segmentation of fat.  
359 For a similar overall pixel accuracy and for a given CNN, the fat recall and precision rates  
360 are related to each other, as an increase in one of them usually results in a decrease in the  
361 other one. In general, the fat recall rates were higher for CNN x\_512 than for CNN x\_128,  
362 whereas the precision rate was similar in both cases (around 0.84). These results were  
363 similar to other works found in the literature, even though metrics were not fully  
364 comparable. [Jackman, Sun and Allen \(2009\)](#) underestimated the number of marbling  
365 pixels (12.4% not classified as IMF) for beef. No information was given on misclassified  
366 lean pixels. For dry-cured ham, [Cernadas, Dur and Antequera \(2002\)](#) classified correctly  
367 90% of the fat streaks with an acceptable rate of false positives, whereas [Widiyanto et al.  
368 \(2013\)](#) using a slightly different metric for accuracy (dice similarity coefficient) obtained  
369 0.94 and 0.88 for the ham and IMF regions, respectively. For CNN 3\_512 the dice  
370 similarity coefficient was calculated and similar values were obtained, 0.99 and 0.83 for  
371 lean and IMF regions, respectively.

372 FnPFc and FpLFC rates showed that for x\_128 and x\_512, around 35-40% of the  
373 misclassified pixels were found near the contours of the fat patches in the images. Similar  
374 rates were observed for the simple CNNs 1\_128\_s and 1\_512\_s. One of the reasons for  
375 these results are the difficulties faced by the trained expert during the preparation of the  
376 ground truth images. This amounts to using noisy data during training. Another possible  
377 reason was the lack of enough samples for training, due to the wide range of possible  
378 RGB values for the fat and lean, structures of the fat, etc. For 2\_128\_s, 3\_128\_s, 2\_512\_s  
379 and 3\_512\_s, the FpLFC rate was much higher than the FnPFc rate. This may indicate that  
380 these CNNs was overestimating the contours of the fat patches.

381 Processing time increased with the number of filters in the first layer and the depth of the  
382 CNN. In general, an increase in performance resulted in an increase in the processing  
383 time. High processing times (i.e 410 ms for CNN 4\_512) might be a problem for the  
384 segmentation of images in real-time applications. As expected, processing time increased  
385 with the number of filters and the depth of the network. For example, for CNN 2\_128  
386 average processing time was 20 ms, whereas for CNN 2\_512 was 58 ms. This represents  
387 an increase of the processing time by a factor of almost 3, for an increase by a factor of 2  
388 in the number of filters in the first layer. CNN 1\_512\_s and CNN 1\_128 had a similar  
389 performance (fat precision and recall rates) on the test set. However, processing time was  
390 lower for CNN 1\_128 (19 ms vs. 9 ms). This fact should be further investigated.

391 The CNN 3\_512 and CNN 1\_128 was selected (the best and worst performing  
392 architectures) to evaluate the segmentation of images from the test set. In general, the best  
393 performing CNN (3\_512) was able to segment correctly raw images (Fig. 6a). Some small  
394 divergences can be observed in the contours of the fat regions between the segmented  
395 image using the CNN and the ground truth. In some particular images, some areas were  
396 not correctly segmented (Figs. 6b, 6c and 7). The reason for these divergences have  
397 already been discussed above. In some other cases, the convolutional network segmented  
398 fat patches that were not correctly selected during the preparation of the ground truth  
399 images (Fig. 6c).

400 Results for CNN 3\_512 and CNN 1\_128 showed that in 49 images out of 62 images of  
401 the evaluation set, the CNN 3\_512 had equal or higher overall pixel accuracy than CNN  
402 1\_128. In those images where CNN 1\_128 performed better, the differences in pixel  
403 accuracy were very small. However, in some cases CNN 3\_512 was able to segment fat  
404 much better than CNN 1\_128. For example, in Fig. 7, both cases did not segment correctly  
405 some of the fat pixels. However, CNN 3\_512 was able to segment IMF better than CNN  
406 1\_128. Probably, due to the larger number of filters and layers, the CNN 3\_512 was able  
407 to store more information on fat detection from the training samples. However, the small  
408 number of samples in the training set would rather explain the poor performance of the  
409 CNN in this case for both architectures.

410 This study was applied to 3x64x64 patches obtained from images. However, using  
411 different strategies, the algorithm could be applied to segment larger images. For  
412 example, using an overlap-tile strategy (Ronneberger, Fischer & Brox, 2015).

413 The good results obtained for the detection of intramuscular fat in sliced dry-cured ham  
414 suggests that this methodology can be of interest for the dry-cured ham industry and might  
415 be used to develop systems for food quality analysis in other food products. One of the  
416 advantages of this machine learning technique is that no specialized knowledge and skills  
417 in computer vision are required. However, some challenges must be addressed. Image  
418 processing with CNNs might be too slow for real-time image segmentation in industrial  
419 processes, especially for CNNs with many filters and layers. Moreover, training samples  
420 must be collected and labelled before training the CNN. Although, in food elaboration  
421 processes (i.e dry cured ham), training examples are available in large quantities,  
422 preparation of ground truth images can be time consuming and may require the expertise  
423 of food technologists. Although CNNs provides state-of-the-art performance in many  
424 computer vision applications, other algorithms should be also evaluated (Support Vector  
425 Machines, Decision Trees, etc) as long image processing times might be a problem for  
426 real-time applications in industry.

427 Detection of intramuscular fat is the first step to efficiently quantify intramuscular fat  
428 content. Deep learning algorithms in combination with information obtained using other  
429 non-destructive technologies (Fulladosa, Gou & Muñoz, 2016; Fulladosa, Rubio-Celorio,  
430 Skytte, Muñoz & Picouet 2017; Fulladosa et al., 2018; Garrido-Novell, Garrido-Varo,  
431 Perez-Marin, Guerrero-Ginel & Kim, 2015; Gou et al., 2013) might help to find a  
432 nutritional label specific for each sliced ham pack and thus encourage consumers to adopt  
433 healthier eating habits and/or buy products according to their needs and/or preferences.

434 Detection of colour defects, for example, due to oxidation, could be performed (i.e. using  
435 deep learning) simultaneously with the IMF segmentation. With these systems,  
436 companies could discard and/or redirect to other process the defective products. Besides,  
437 a good detection of IMF in images may also provide a tool to improve prediction precision  
438 in other technologies. Prediction error of salt and water contents using computed  
439 tomography can be reduced with a good detection and quantification of fat content  
440 (Santos Garcés, Muñoz, Gou, Garcia-Gil & Fulladosa, 2014), leading to models that  
441 improve the optimization of the dry-cured ham elaboration process.

## 442 **CONCLUSIONS**

443 Results show that deep learning is able to segment correctly IMF in dry cured ham by just  
444 using training samples in combination with CNN. CNN attained a similar performance to

445 that of conventional image analysis algorithms, reducing development time, at the cost of  
446 requiring greater computing resources.

447 The increase in the complexity of the network helps to improve the performance, but up  
448 to a certain level, as the network may end up overfitting and processing time of images  
449 may increase considerably. One of the challenges is the need to obtain good training data  
450 for training the CNN, due to the difficulty in classifying pixels correctly and objectively  
451 even by trained experts. CNN opens new possibilities to solve complex detection  
452 problems in the food industry without the need of developing complex algorithms,  
453 facilitating the deployment of these technologies in the food industry.

#### 454 **ACKNOWLEDGMENTS**

455 This work was partially supported by INIA (grant number RTA2013-00030-CO3-01) and  
456 CERCA programme from Generalitat de Catalunya.

#### 457 **BIBLIOGRAPHY**

458 Badrinarayanan, V., Kendall, A., & Cipolla, R. (2017). Segnet: A deep convolutional  
459 encoder-decoder architecture for image segmentation. *IEEE transactions on pattern  
460 analysis and machine intelligence*, 39(12), 2481-2495.

461 Blasco, J., Munera, S., Aleixos, N., Cubero, S., Molto, E. (2017). Machine vision-based  
462 measurement systems for fruit and vegetable quality control in postharvest. In advances  
463 in Biochemical Engineering/Biotechnology book series, 161, 71-91.

464 Bolaños, M., & Radeva, P. (2016). Simultaneous Food Localization and Recognition. In  
465 23rd International Conference on Pattern Recognition (ICPR) (pp. 3140–3145). Cancun,  
466 Mexico.

467 Caccamo, M., Melilli, C., Barbano, D.M., Portelli, G., Marino, G., & Licitra, G. (2004).  
468 Measurement of gas holes and mechanical openness in cheese by image analysis. *Journal  
469 of Dairy Science*, 87(3), 739-748.

470 Cavallo, D.P., Cefola, M., Pace, B., Logrieco, A.F., & Attolico, G. (2018). Non-  
471 destructive automatic quality evaluation of fresh-cut iceberg lettuce through packaging  
472 material, *Journal of Food Engineering*, 223, 46-52.

473 Cernadas, E., Dur, M. L., & Antequera, T. (2002). Recognizing marbling in dry-cured  
474 Iberian ham by multiscale analysis. *Pattern Recognition Letters*, 23, 1311–1321.



475 Claret, A., Guerrero, L., Guàrdia, M.D., Garcia-Gil, N. & Arnau, J. (2009). Desarrollo de  
476 escalas de referencia para determinados atributos sensoriales del jamón curado de cerdo  
477 blanco. V Congreso Mundial del jamón, Aracena, Huelva (Spain)

478 Dutta, M.K., Issac, A., Minhas, N. & Sarkar, B. (2016). Image processing based method  
479 to assess fish quality and freshness. *Journal of Food Engineering*, 177, 50-58.

480 Faucitano, L., Huff, P., Teuscher, F., Garipey, C. & Wegner, J. (2005). Application of  
481 computer image analysis to measure pork marbling characteristics. *Meat Science*, 69,  
482 537-543.

483 Fulladosa, E., Austrich, A., Muñoz, I., Guerrero, L., Benedito, J., Lorenzo, J. M. & Gou,  
484 P. (2018). Texture characterization of dry-cured ham using multi energy X-ray analysis.  
485 *Food Control*, 89 46-53.

486 Fulladosa, E., Gou, P. & Muñoz, I. (2016). Effect of dry-cured ham composition on X-  
487 ray multi energy spectra. *Food Control*, 70 41-47.

488 Fulladosa, E., Rubio-Celorio, M., Skytte, J. L., Muñoz, I. & Picouet, P. (2017). Laser-  
489 light backscattering response to water content and proteolysis in dry-cured ham. *Food*  
490 *Control*, 77 235-242.

491 Garrido-Novell, C., Garrido-Varo, A., Perez-Marin, D., Guerrero-Ginel, J. E. & Kim, M.  
492 (2015). Quantification and spatial characterization of moisture and NaCl content of  
493 Iberian dry-cured ham slices using NIR hyperspectral imaging. *Journal of Food*  
494 *Engineering*, 153 117-123.

495 Goodfellow, I., Bengio, Y., Courville, A., & Bengio, Y. (2016). *Deep learning* (Vol. 1).  
496 Cambridge: MIT press.

497 Gou, P., Santos-Garcés, E., Hoy, M., Wold, J. P., Liland, K. H. & Fulladosa, E. (2013).  
498 Feasibility of NIR interactance hyperspectral imaging for on-line measurement of crude  
499 composition in vacuum packed dry-cured ham slices. *Meat Science*, 95 (2), 250-255.

500 Hawkins, D.M. (2004). The problem of overfitting. *Journal of Chemical Information and*  
501 *Computer Sciences*, 44(1), 1-12.

502 He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image  
503 recognition. In *Proceedings of the IEEE conference on computer vision and pattern*  
504 *recognition* (pp. 770-778).

505 Huang, H., Liu, L., Ngadi, M.O., & Gariépy, C. (2013). Prediction of pork marbling  
506 scores using pattern analysis techniques. *Food Control*, 31, 224-229.

507 Jackman, P., Sun, D.-W., & Allen, P. (2009). Automatic segmentation of beef  
508 longissimus dorsi muscle and marbling by an adaptable algorithm. *Meat Science*, 83(2),  
509 187-194.

510 Jia, Y., Shelhamer, E., Donahue, J., Karayev, S., Long, J., Girshick, R., Guadarrama, S.,  
511 & Darrell, T. (2014). Caffe: Convolutional Architecture for Fast Feature Embedding. In  
512 Proceedings of the 22nd ACM International conference on multimedia (pp. 675-678).  
513 Orlando, USA.

514 Kvam, J., Gangsei, L.E., Kongsro, J., & Schistad-Solberg, A.H. (2018). The use of deep  
515 learning to automate the segmentarion of the skeleton from CT volume pigs. *Translational*  
516 *Animal Science*, 2(3), 324-335.

517 Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). Imagenet classification with deep  
518 convolutional neural networks. In Proceedings of the Advances in neural information  
519 processing systems (pp. 1097-1105), Lake Tahoe, USA.

520 Lin, G., Milan A., Shen, C., & Reid, I. (2017). RefineNet Multi-path refinement networks  
521 for high-resolution semantic segmentation. In Proceedings of the IEEE Conference on  
522 Computer Vision and Pattern Recognition (CVPR), Milan, Italy.

523 Liu, L., Ngadi, M.O., Prasher, S.O., & Gariépy, C. (2012). Objective determination of  
524 pork marbling scores using the wide line detector. *Journal of Food Engineering*, 110(3),  
525 497-504.

526 Liu J.-H., Sun, X., Young, J.M., Bachmeier, L.A., & Newman, D.J. (2018). Predicting  
527 pork loin intramuscular fat using computer vision system. *Meat Science*, 143, 18-23.

528 Lohumi, S., Lee, S., Lee, H., Kim, M.S., Lee, W.H., & Cho, B.-K. (2016). Application  
529 of hyperspectral imaging for characterization of intramuscular fat distribution in beef.  
530 *Infrared Physics & Technology*, 74, 1-10.

531 Long, J., Shelhamer, E., & Darrell, T. (2015). Fully Convolutional Networks for Semantic  
532 Segmentation. In Proceedings of the IEEE Conference on Computer Vision and Pattern  
533 Recognition (CVPR) (pp 3431-3440), Boston, USA.

534 Llorido, L., Pizarro, E., Estévez, M., & Ventanas, S. (2019). Emotional responses to the  
535 consumption of dry-cured hams by Spanish consumers: A temporal approach. *Meat*  
536 *Science*, 149, 126-133.

537 Morales, R., Guerrero, L., Aguiar, A.P.S, Guàrdia, M.D., & Gou, P. (2013). Factors  
538 affecting dry-cured ham acceptability. *Meat Science*, 95(3), 652-657.

539 Muñoz, I., Rubio-Celorio, M., Garcia-Gil, N., Guardia, M.D., & Fulladosa, E. (2015).  
540 Computer image analysis as a tool for classifying marbling: A case study in dry-cured  
541 ham. *Journal of Food Engineering*, 166, 148-155.

542 Pfisterer, K.J., Amelard, R., Chung, A.G., & Wong, A. (2018). A new take on measuring  
543 relative nutritional density: The feasibility of using a deep neural network to assess  
544 commercially-prepared puréed food concentrations. *Journal of Food Engineering*, 223,  
545 220-235.

546 Qiao, Jun , Ngadi, M.O., Wang, N., Gariépy, C., & Prasher, S.O. (2007). Pork quality and  
547 marbling level assessment using a hyperspectral imaging system, *Journal of Food*  
548 *Engineering*, 83(1), 10-16.

549 Rangayyan, M.R. (2014). *Biomedical image analysis*. CRC Press.

550 Ronneberger, O., Fischer, P., & Brox, T. (2015). U-Net: Convolutional networks for  
551 biomedical image segmentation. *Medical Image Computing and Computer-Assisted*  
552 *Intervention (MICCAI)*, 9351, 234-241.

553 Santos-Garcés, E., Muñoz, I., Gou, P., Garcia-Gil, N., & Fulladosa, E. (2014) Including  
554 estimated intramuscular fat content from computed tomography images improves  
555 prediction accuracy of dry-cured ham composition. *Meat Science*, 96(1), 943-947.

556 Simonyan, K., & Zisserman, A. (2015). Very deep convolutional networks for large-scale  
557 image recognition. In *Proceedings of the 3th International Conference on Learning*  
558 *Representations (ICLR)* (pp. 1-14). San Diego, USA.

559 Srivastava, S., Vaddadi, S. & Sadistap, S. (2015). Quality assessment of commercial  
560 bread samples based on breadcrumb features and freshness analysis using and ultrasonic  
561 machine vision (UVS) system. *Journal of Food Measurement and Characterization*, 9(4),  
562 525-540.

563 Sun, D.W., & Brosnan, T. (2003). Pizza quality evaluation using computer visión- part 1.  
564 Pizza base and sauce spread. *Journal of Food Engineering*, 57(1), 81-89.

- 565 Velazquez, L., Cruz-Tirado, J.P., Siche, R., & Quevedo, R. (2017). An application based  
566 on the decision tree to classify the marbling of beef by hyperspectral imaging. *Meat*  
567 *Science*, 133, 43-50.
- 568 Widiyanto, S., Cufí, X., Rubio, M., Muñoz, I., Fulladosa, E., & Martí, R. (2013).  
569 Automatic intra muscular fat analysis on dry-cured ham slices. In *Proceedings of ibPRIA*,  
570 873-880.
- 571 Xu, J.-L., Sun, D.W. (2018). Computer vision detection of salmon muscle gaping using  
572 convolutional neural network features. *Food Analytical Methods*, 11(1), 34-47.  
573

574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588

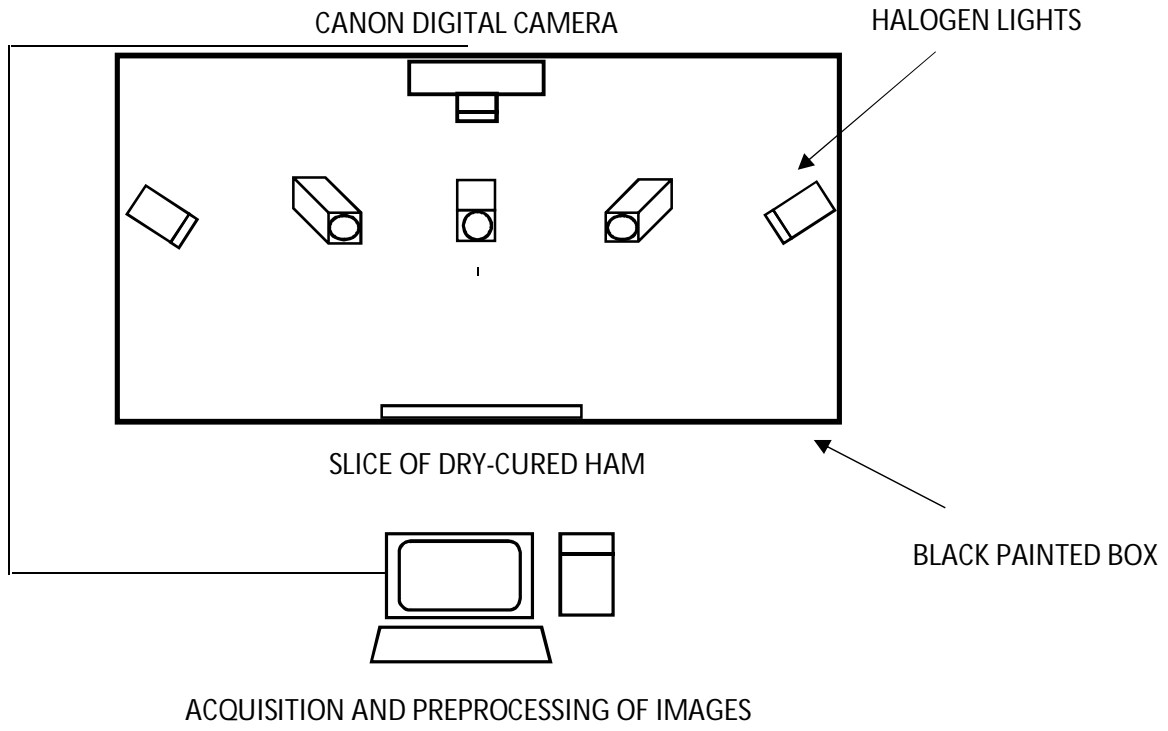
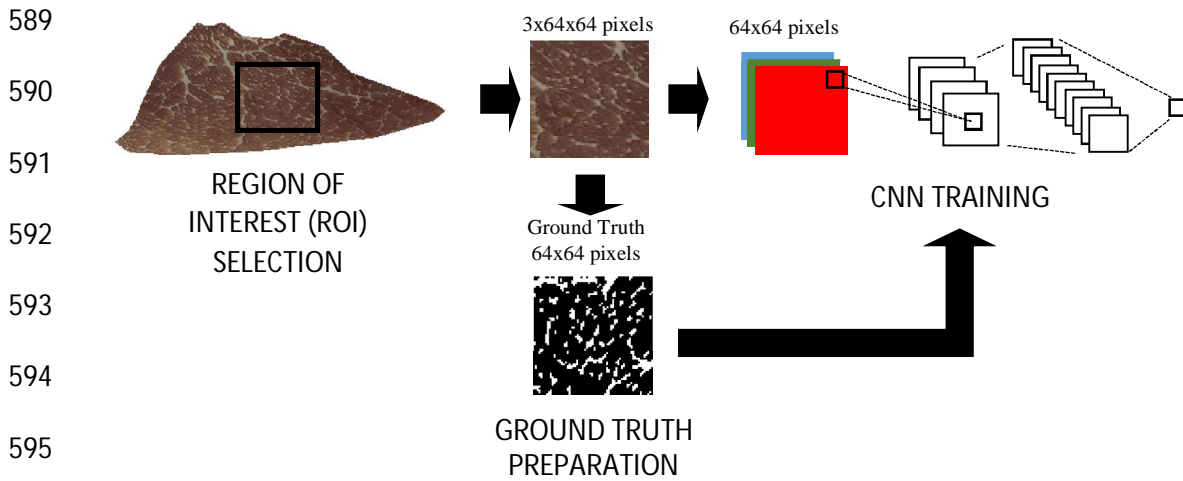


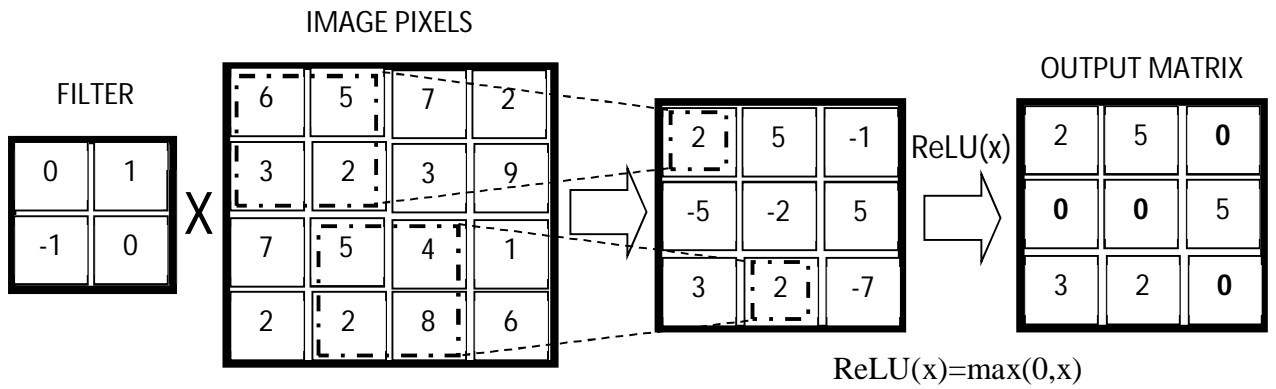
Figure 1. Overview of the image acquisition system.



596 Figure 2. Overview of the learning scheme for fat segmentation using a convolutional  
 597 neural network (CNN)  
 598

599  
 600  
 601  
 602  
 603  
 604  
 605  
 606  
 607  
 608  
 609  
 610  
 611  
 612  
 613  
 614  
 615  
 616  
 617  
 618  
 619  
 620  
 621  
 622  
 623  
 624  
 625  
 626

a)



b)

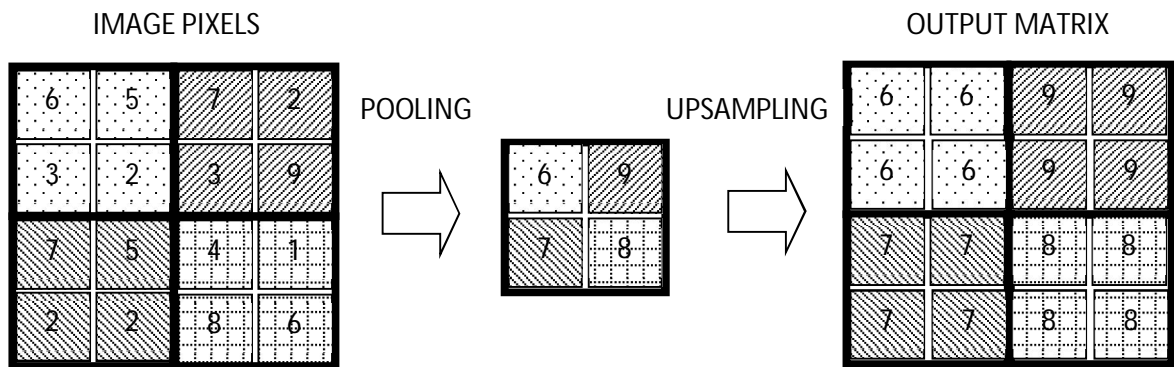
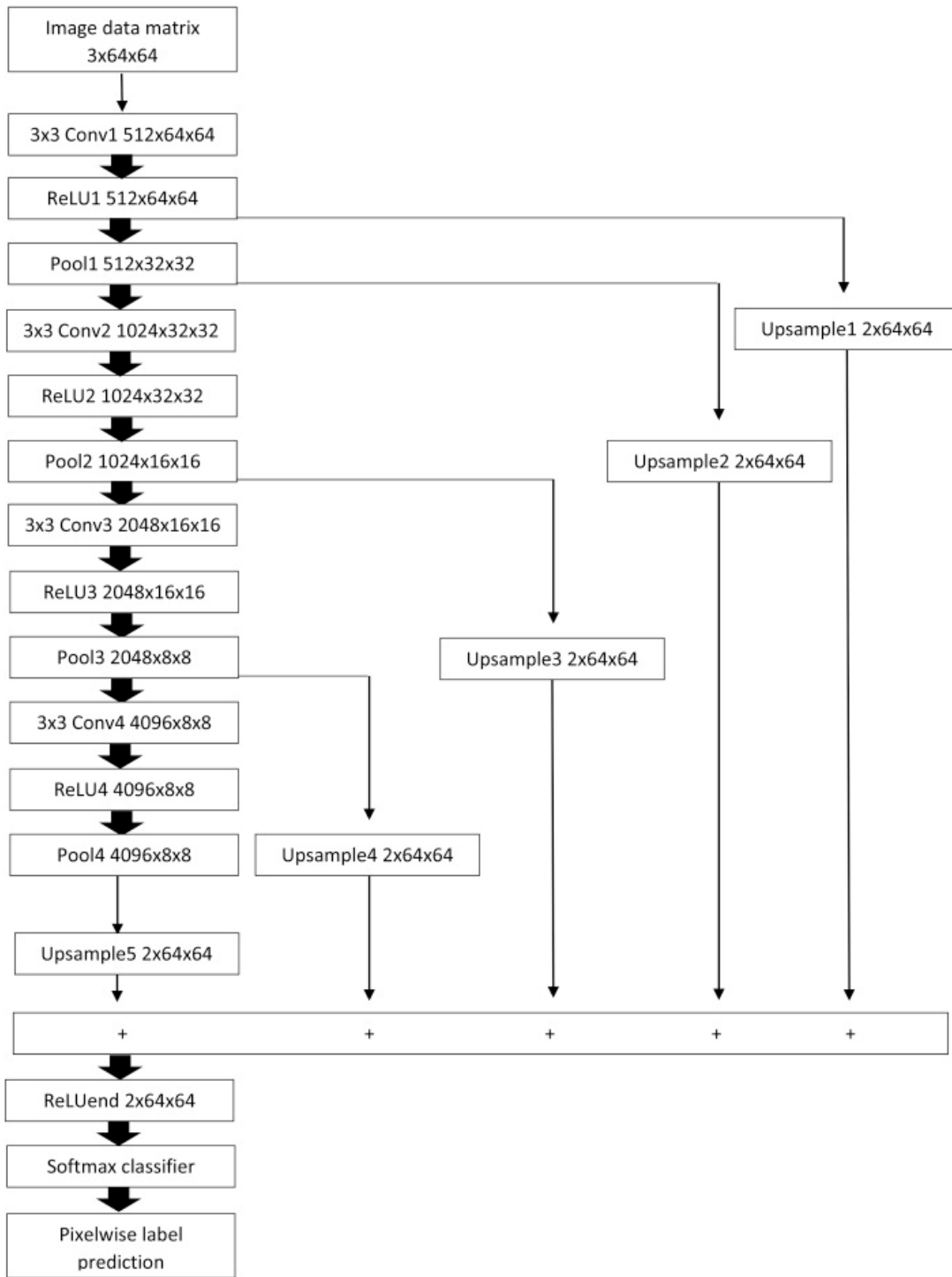


Figure 3. Main types of layers in CNNs: a) a convolution operation with a filter of 2x2 pixels and depth 1 followed by a Rectifier Linear Unit (ReLU) activation function; b) A 2x2 pixels max pooling layer followed by a nearest neighbour upsampling layer from a 2x2 to a 4x4 pixels matrix.



627

628

629

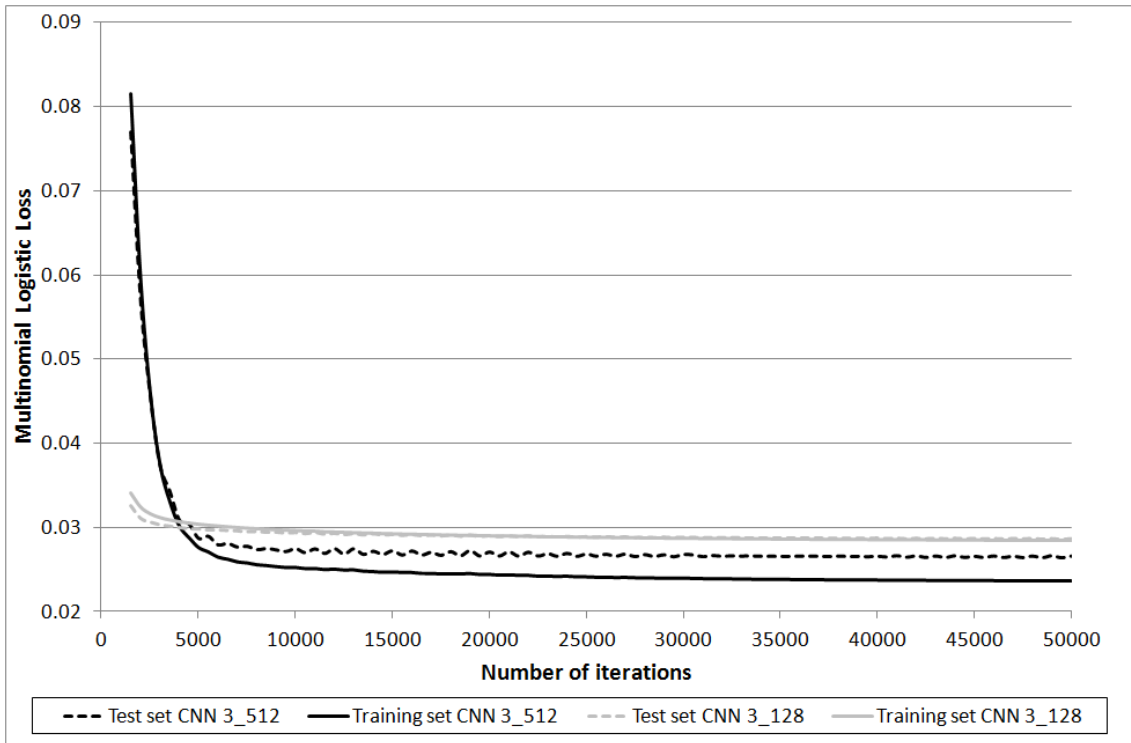
630

631

632

Figure 4. Architecture of the convolutional neural network architecture used in this work with 512 filters in the first layer and four layers. Convx, ReLUx, Poolx, Upsamplex indicate convolutional, rectified linear unit, pooling and upsampling operations respectively.



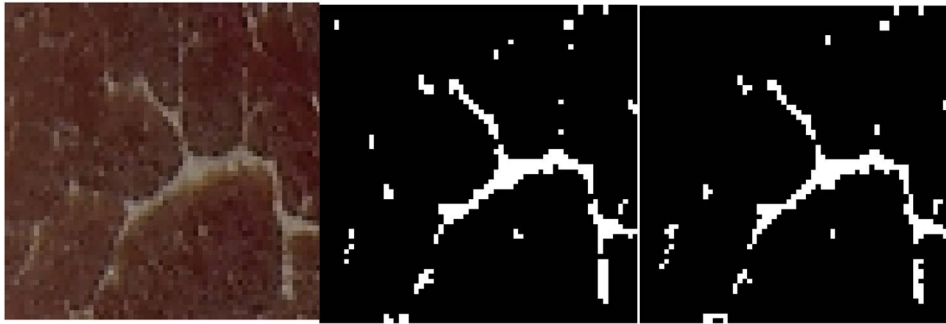


633

634 Figure 5. Multinomial Logistic Loss vs number of iterations (from 1000 to 50000  
 635 iterations) for the training and test sets of CNN 3\_128 and CNN 3\_512.  
 636

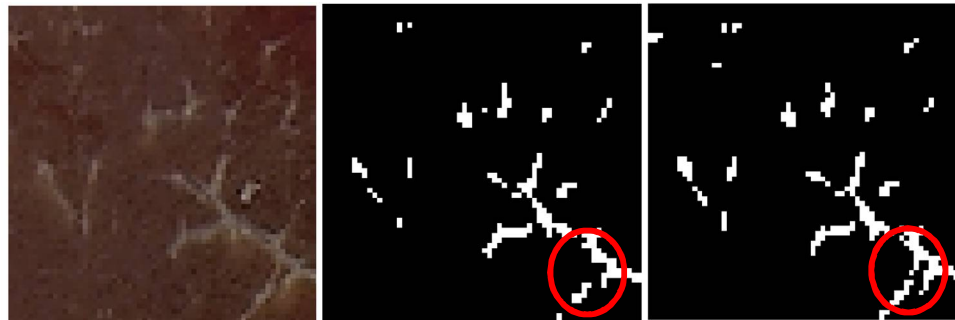
637

638 a)



639

640 b)



641

642 c)



643

644

I

II

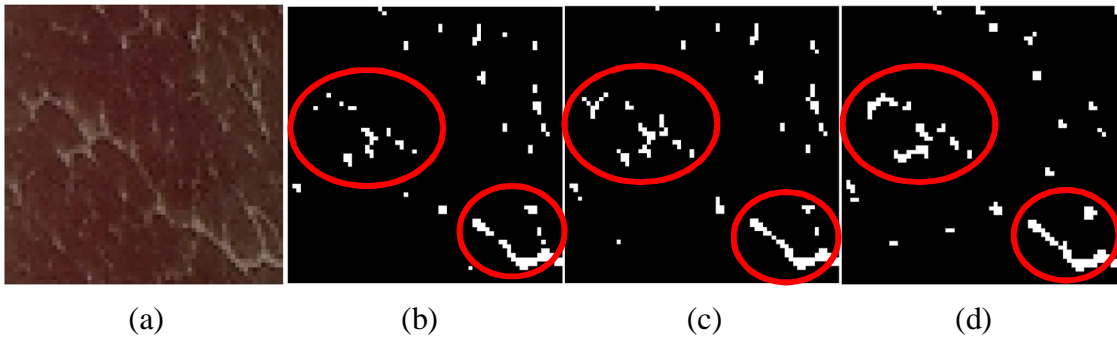
III

645 Figure 6. Images of slices of dry cured ham segmented with CNN 3\_512. Raw image (I),  
646 segmented image with CNN 3\_512 (II) and ground truth image (III). Red circle denotes  
647 areas with poor segmentation (b) and possible misclassified ground truth pixels (c)

648

649

650



654 Figure 7. Segmentation of an image of a slice of dry cured ham. Raw image (a),  
655 segmented image by CNN 1\_128 (b), segmented image by CNN 3\_512 (c) and ground  
656 truth image (d). Red circles denotes areas with poor segmentation.  
657

Architecture name	Kernel size	Number of layers	Number of filters in layers 1/2/3/4	Upsampling included	Number of learnable parameters
1_128_s	3x3	1	128	2	3,586
2_128_s	3x3	2	128	3	298,754
3_128_s	3x3	3	128	4	1,478,914
1_512_s	3x3	1	512	2	7,170
2_512_s	3x3	2	512	3	4,733,954
3_512_s	3x3	3	512	4	23,610,368
1_128	3x3	1	128	1,2	3,586
2_128	3x3	2	128/256	1,2,3	298,754
3_128	3x3	3	128/256/512	1,2,3,4	1,478,914
4_128	3x3	4	128/256/512/1024	1,2,3,4,5	6,198,530
1_512	3x3	1	512	1,2	7,170
2_512	3x3	2	512/1024	1,2,3	4,733,954
3_512	3x3	3	512/1024/2048	1,2,3,4	23,610,368
4_512	3x3	4	512/1024/2048/4096	1,2,3,4,5	99,111,938

659 Table 1. Description of the parameters of several CNN architectures.

Architecture	Overall pixel accuracy (training set)	Overall pixel accuracy (test set)	Fat recall rate (test set)	Fat precision rate (test set)	Rate of false negatives in the predicted fat contour (test set)	Rate of false positives in the labelled fat contour (test set)	Processing time per image (ms)
1_128_s	0.986	0.987	0.741	0.834	0.360	0.409	10
2_128_s	0.981	0.982	0.562	0.807	0.234	0.613	15
3_128_s	0.97	0.971	0.180	0.683	0.066	0.572	20
1_512_s	0.988	0.988	0.770	0.834	0.388	0.455	19
2_512_s	0.985	0.985	0.668	0.820	0.305	0.551	55
3_512_s	0.975	0.975	0.312	0.742	0.127	0.623	101
1_128	0.988	0.988	0.778	0.840	0.376	0.347	9
2_128	0.988	0.988	0.776	0.846	0.393	0.360	16
3_128	0.989	0.989	0.785	0.842	0.377	0.395	22
4_128	0.989	0.989	0.790	0.843	0.405	0.371	42
1_512	0.989	0.989	0.793	0.847	0.396	0.377	22
2_512	0.99	0.989	0.81	0.841	0.415	0.391	58
3_512	0.991	0.989	0.816	0.84	0.412	0.399	114
4_512	0.991	0.989	0.803	0.846	0.395	0.394	410

Table 2. Performance results of the studied CNN architectures

661

662