047

# Graph convolutional neural networks with uncertainty modelling applied to edge detection in mammograms

Anonymous Full Paper Submission 41

### **OD1** Abstract

This paper addresses the use of a graph convolutional 002 network for delineation of structures and estimation 003 of landmarks in mammograms, a critical step in the 004 evaluation of image quality in breast cancer screen-005 006 ing. In this context, the ability to estimate the uncertainty of the predicted positions is crucial. In 007 the current work we focus on the pectoral muscle, 008 where the variability in muscle visibility across im-009 ages introduces significant uncertainty. Our main 010 contribution is a novel modification of a deep graph 011 convolutional network (GCN) that not only locates 012 key points along the muscle boundary but also pro-013 vides uncertainty estimates, which are useful for 014 selecting images that must be evaluated by a hu-015 man. We introduce a novel approach to estimate 016 both aleatoric and epistemic uncertainties using a 017 GCN framework. Aleatoric uncertainty captures 018 variability in ground truth due to annotator dif-019 ferences, while epistemic uncertainty accounts for 020 the model's inherent limitations. Our method was 021 tested on in-house annotated mammograms and the 022 external InBreast dataset, demonstrating compara-023 ble accuracy to human annotators and robustness 024 025 in the presence of domain shifts. The uncertainty estimates were found to be highly accurate, confirm-026 ing their potential for identifying cases that require 027 human review. 028

### 029 1 Introduction

Breast cancer is the most common form of cancer in
females, worldwide, and mammographic screening
is an effective way of detecting cancers at an early
stage.

In mammography screening it is crucial to main-034 tain high image quality to ensure the best possible 035 visualization of breast tissue and the identification 036 of potential breast cancer indicators. This includes 037 technical quality in the form of image sharpness and 038 contrast as well as the proper positioning, which 039 ensures that the relevant parts of the breast are 040 depicted. The quality of the images significantly 041 impacts the rates of patient recall and the detec-042 tion of cancer through screening, thus affecting the 043 accuracy and reliability of the screening process. 044 Therefore, mammograms routinely undergo quality 045 assessment by radiographers, which would benefit 046



Figure 1. Sample mammograms where the pectoral muscle is clearly depicted (left) and blurry (right).

from automatic analysis.

An important part of the quality assessment con-048 cerns the depiction of the pectoral muscle in the 049 mediolateral oblique (MLO) view, relating to its 050 size, its shape, and its orientation (Waade et al. [1]). 051 The muscle should be located in the top left (right) 052 corner of a standard X-ray mammogram in MLO-053 view. The complexity of the task varies considerably, 054 as shown in Figure 1, where the muscle is clearly 055 visible as a light area in the left image, while the 056 right image shows a more difficult case where the 057 lower part of the muscle is blurry. 058

Previous studies have used different methods for 059 segmenting the muscle, using various proprietary 060 versions of convolutional neural nets (CNNs). Ma 061 et al. [2] apply a method similar to a U-Net and 062 concludes it works better than traditional imaging 063 techniques. Brahim et al. [3] use a CNN model 064 heatmaps generated by a GradCam method. PeM-065 Net (Yu et al. [4]) uses a InceptionResNetv2 back-066 bone and a complex up-sampling scheme to generate 067 pixel masks. Guo et al. [5] use a U-Net to segment 068 the muscle. Yang et al. [6] uses deep learning in the 069 way of a modified U-Net model to segment pectoral 070 muscle volume from computed tomography images. 071

While these methods may provide state-of-the-art 072 segmentation masks, their handling of uncertainty 073 leaves something to be desired, as they focus only on 074 average error (over different datasets and different 075 criteria). They do not provide explicit probability 076 distributions or confidence intervals for their location 077 estimates. We argue that the application to image 078 quality control would benefit from such estimates, 079 so that a human could perform the evaluations incritical cases where the predicted locations are likelyto be substantially off target.

Graph convolutional neural nets (GCNs) have 083 been introduced as an alternative to heatmap based 084 methods for landmark detection to better capture 085 relationships between landmarks. The aim is to 086 enable the model to learn relations between point 087 positions and thereby improve on the relative po-088 sitioning of the landmarks. In the present study, 089 we use the GCN approach to identify points along 090 the border of the muscle. Since the muscle is lo-091 cated in a corner, the task of identifying the border 092 that extends from the top edge to the side edge 093 is equivalent to segmenting it. However, landmark 094 based modelling extends directly to locating of sin-095 gular key points that are used for quality assessment 096 of mammogram, such as the nipple and the infra-097 mammary fold. While these are not included in the 098 present study, this is a natural extension of the land-099 mark approach, which can also benefit from learning 100 the geometric relations among the points. While 101 heatmap based methods can also be used for these tasks, the location of singular points is handled more 103 directly in landmark based methods. 104

We adopt a modification of the "Deep Adaptive 105 Graph" (DAG) framework, which has previously 106 been used for identifying key points in faces (Li et 107 al. [7]). Our contribution is an extension of the DAG 108 framework that provides an explicit probabilistic 109 model for the location of the key points. We train 110 the model to output estimates of the aleatoric (truly 111 random) uncertainties together with the key point 112 coordinates. 113

To our knowledge, no attempts to explicit uncer-114 tainty modelling in GCNs have been published to 115 date, and such modelling has been called for: "To 116 date, the vast majority of existing works do not 117 take into account the uncertainty of a GCN regard-118 ing its prediction, which is alarming especially in 119 high-stake scenarios." (Kang et al. [8]). We use the 120 Laplace distribution, which has been used as a prior 121 for Bayesian uncertainty estimation in a general con-122 text (Kendall and Gal [9]), but to our knowledge 123 not in a GCN setting. 124

In addition we estimate the epistemic (model related) uncertainties through the variations within a model ensemble, and combine the aleatoric and epistemic uncertainties to create accurate confidence intervals.

#### $_{130}$ 2 Method

We include only an informal description of the "Deep Adaptive Graph" (DAG) framework (Li et al. [7]), which we build on, and refer to this source for technical details. The model's basic task is to locate a vector of points  $v \in \mathbb{R}^{n \times 2}$  in an image as close as possible to the ground truth of correct locations 136  $v^*$ . It uses a high-resolution convolutional net (HR-137 Net) to generate a feature map for the given image. 138 These features are fed into a GCN together with 139 the geometry of a current location estimate v and 140 give a vector  $\Delta v$  as output. The process starts with 141 an initial point vector  $v_0$ , which is updated itera-142 tively:  $v_{i+1} = v_i + \Delta v_i$ , which is intended to move 143 toward  $v^*$ . The algorithm described by (Li et al. [7]) 144 also included a so-called global step prior to the  $\Delta v$ 145 updates, but this was not included in the present 146 application. We let T be the number of local steps 147 and for convenience define  $\mu = v_T$ . 148

#### 2.1 Aleatoric uncertainty estimation 149

By aleatoric uncertainty we mean uncertainty about 150 the ground truth that is "truly random" in the 151 sense that it cannot be eliminated by any amount 152 of training data or any kind of model (Hüllermeier 153 and Waegeman [10]). It represents the randomness 154 among different annotators - or even the same an-155 notator at different times - in how they place the 156 markings. In the right image of Figure 1, e.g. the 157 delineation of the lower part of the muscle is likely 158 to vary substantially in this way. 159

Earlier research on aleatoric uncertainty in GCN 160 models focus on labelling problems, where the task 161 is to assign properties to the nodes in a graph 162 (Vashishth et al. [11]). A recent overview of uncer-163 tainty in GCN models is given in (Wang et al. [12]), 164 which gives a taxonomy of types and sources for un-165 certainties, and ways to estimate them. Our work is 166 different in the way that we develop an architecture 167 that allows the model to estimate 2-dimensional lo-168 cations in parallel with uncertainty estimates for the 169 same locations. We accomplish this by extending the 170 method of Li et al with the inclusion of a separate 171 GCN module that estimates aleatoric uncertainties 172 through parameterized random distributions. It is 173 structurally equivalent to the GCN module that 174 computes  $\Delta v$  and is applied after the last iteration 175 T. It takes as input the HR-net features evaluated 176 at the locations  $\mu$  together with the geometric fea-177 tures of  $\mu$ . The output is denoted by  $\log(b) \in \mathbb{R}^{n \times 2}$ , 178 where  $b_x^i$  and  $b_y^i$  represent the uncertainties for  $\mu_x^i$ 179  $\mu_{y}^{i}$ . 180

Rather than the traditional deep learning ap-181 proach of defining a loss function which measures 182 the distance from the desired outcome, we view the 183 entire model as a parameterized statistical model 184 of the training data. For each point in  $v^{*i}$ , we 185 assume that  $v_x^{*i}$  has the Laplace (double exponen-186 tial) distribution with median  $\mu_x^i$  scale parameter 187  $b_x^i$ , which we write  $v_x^{*i} \sim \mathcal{L}(\mu_x^i, b_x^i)$  and similarly  $v_y^{*i} \sim \mathcal{L}(\mu_y^i, b_y^i)$ . We assumed all components to be 188 189

238

248



Figure 2. Computational pipeline.

<sup>190</sup> independent, which gives the following likelihood:

191 
$$L = \prod_{k}^{M} \prod_{i}^{n} \prod_{j \in \{x,y\}} \frac{1}{2b_{j}^{k,i}} \exp\left(-\frac{|v_{j}^{*k,i} - \mu_{j}^{k,i}|}{b_{j}^{k,i}}\right)$$

Here, k runs over the M training images, i runs over the n key points and j runs over the x and y dimension.

Our approach is to optimize the model parameters with respect to this likelihood function, but as usual in statistical modelling, we minimize  $-\log(L)$ instead, which has the same optimum and more favorable numerical properties:

200 
$$-\log \mathbf{L} = \sum_{k}^{M} \sum_{i}^{n} \sum_{j \in \{x, y\}} \log(2b_{j}^{k, i}) + \frac{1}{b_{j}^{k, i}} |v_{j}^{*k, i} - \mu_{j}^{k, i}|$$

Minimizing this expression amounts to training the model end-to-end to estimate locations  $\mu_x^i, \mu_y^i$  and uncertainties  $b_{x}^i, b_y^i$  for a given input image. To our knowledge, this approach to estimating aleatoric uncertainties for GCN models is novel.

We also considered using normal distributions in-206 stead of Laplace, but chose the latter because it is 207 the natural generalization of the L1-error that was 208 used by (Li et al. [7]). It has the property that it pe-209 nalizes deviations linearly rather than quadratically, 210 which makes it more robust. Normal distributions 211 might place too much emphasize on the cases where 212 the model has trouble reproducing the ground truth. 213

#### 214 2.1.1 Model pipeline

The model pipeline is shown in Figure 2. The HR-215 Net (blue color) computes localized features, which 216 are fed into the GCN models. GCN-1 (yellow) reads 217 the HR-Net features at the initial points, combined 218 with their geometry, and moves the points (hope-219 fully) toward the correct locations. This process 220 is repeated T times. Then GCN-2 (red) takes the 221 HR-Net features of these final locations as input, 222 together with their geometry, and outputs the esti-223 mated x- and y- uncertainties, illustrated with red 224 ellipses. The entire model is trained end-to-end to 225 optimize the Laplace log-likelihood of the ground 226 227 truth data.

#### 2.2 Epistemic uncertainty estimation 228

Our approach for estimating epistemic (model re-229 lated) uncertainty (Hüllermeier and Waegeman [10]) 230 is more standard and straight forward. We train a set 231 of models with cross-validation, and use them as an 232 ensemble when evaluating the test sets (Dutschmann 233 et al. [13]). For each x- and y-value of each output 234 point for a given image we estimate the standard 235 deviation among the ensemble model outputs, and 236 treat this as the epistemic uncertainty. 237

#### 2.3 Combined uncertainty

For the aleatoric uncertainty, we computed the av-239 erage  $\log(b)$  tensors over the ensemble and used the 240 exponential of this average as our b-values. We 241 computed the component-wise total variance V242 by adding the ensemble variance estimate to the 243 Laplace variance  $2b^2$ . We then inverted the Laplace 244 variance function to get a modified  $b = \sqrt{V/2}$ , which 245 was used to calculate confidence intervals according 246 to the Laplace distribution. 247

### 3 Data sets

For training data only in-house annotation was used 249 (details to be included after the anonymous review 250 is finished). The images were sampled randomly 251 from a set of mammograms from screening. All 252 annotations were made by the first author, who has 253 no formal background in radiology or radiography. 254 The number of in-house annotated images was 545. 255

From the same source as the in-house training data 256 we sampled a non-overlapping set of mammograms 257 to be annotated by two radiographers. A total of 94 258 images were annotated by both. 259

In order to evaluate our models' generalizability, 260 we also tested it on the external dataset InBreast 261 (Moreira et al. [14]), which had 200 annotated images. 263

### 4 Experimental setup

264

The annotations of the pectoral muscle in the 265 datasets were represented as a list of points along 266 the border of the muscle. The number of points var-267 ied among the data sets and also among the images 268 in each set. To facilitate the subsequent use of key 269 points, the annotations were standardized to n = 10270 equidistant points along the annotated path, where 271 the first one was on the upper edge and the last one 272 was on the vertical edge. 273

The HR-net was set up with a depth of 32, while 274 the GCN modules that compute the  $\Delta v$  and b ten-275 sors had 6 layers and 256 filters. The number of 276 coordinate iterations T was set to 3. In line with Li 277 et al. [7], We use the average of the ground truth 278 279 locations  $v^*$  over the training set as starting values 280  $v_0$ .

The models were trained with the ADAM opti-281 mizer with a learning rate of 0.0001 and a batch size 282 of 4 over 200 epochs, which was enough for over-283 fitting the models. The cross validation used 5 folds, 284 where each model was trained on 4 of them. The 285 last one was used for monitoring the log-likelihood 286 and parameters that gave the highest value on the 287 validation fold was saved. This procedure gives in-288 flated performance on the validation fold, but our 289 purpose was only to create a model ensemble to 290 be used on the test sets, not to cross-validate the 291 models' performance on the training set. 292

Only a minimum of image preprocessing was performed. Right-side images were flipped, so that all images had the breast on the left side, with the pectoral muscle in the top left corner, if present. The images were resized to 512 x 512 pixels and the pixel values were re-scaled to [0, 1]. No data augmentation was used.

### 300 5 Results

#### 301 5.1 Predictive performance

We measure the location errors by the average absolute difference between the predicted coordinates and the ground truth:

305 
$$\frac{1}{2nM^{\text{test}}} \sum_{k}^{M^{\text{test}}} \sum_{i}^{n} \sum_{j \in \{x,y\}} |v_j^{*k,i} - \mu_j^{k,i}|$$

The location values  $v^*$  and  $\mu$  were scaled to the range [0, 1], so the error estimates can be interpreted as fractions of the height and width of the images.

Recall that the internal test set had a ground truth 309 annotation by two radiographers, and  $v^*$  above was 310 defined as the average of these. The model's average 311 error on this data set was 0.0054, while the average 312 error compared to each of the radiographers indi-313 vidually gave 0.0061 and 0.0062, respectively. For 314 comparison, the average between-radiographer error 315 was 0.0057, so the model performance is essentially 316 on par with our human experts. 317

For the INBreast data set, the average error was higher, as would be expected due to the domain shift, but still quite acceptable: 0.0120.

#### **5.2** Uncertainty estimates

The crucial property of uncertainty estimates is that they accurately model the empirical errors, so that the model can identify the cases where its location predictions should not be trusted. To test this we standardize the model errors by divided them with the predicted standard deviation. Under the model assumptions, these standardized errors should follow



Figure 3. Histogram of standardized error distribution for the internal test set.



Figure 4. Comparison of the empirical and theoretical cumulative distribution for the internal test set.

the standard Laplace distribution. Figure 3 gives 329 the histogram of these for the internal test set, to-330 gether with the probability density function for the 331 standard Laplace. Figure 4 shows a plot of the corre-332 sponding empirical cumulative distribution together 333 with the theoretical one. These plots show a very 334 strong correspondence. The maximum difference 335 between the empirical and theoretical cumulative 336 distributions is 0.0459, which confirms a very good 337 match. 338

Figure 5 illustrates the model output on the two 339 mammograms shown in the introduction, where we 340 have zoomed in on the muscle. The red dots show the 341 average annotation points of the two radiographers. 342 The white and orange ellipses show the epistemic 343 and aleatoric 99% confidence areas, while the vellow 344 ellipses give the combined 99% confidence areas. As 345 expected, the model predictions are more accurate 346 for the left image with clearly visible muscle, and the 347 small uncertainty ellipses confirm that the model is 348 more certain about these. In the right image, we see 349 that the model is more certain (and accurate) for 350



**Figure 5.** Sample mammograms with ground truth and uncertainty ellipses. The center of the ellipses are the predicted locations (not depicted)

the top points and more uncertain about the lower ones, which fits the visual impression. We also see that the epistemic uncertainties are smaller than the the aleatoric ones, which is desirable.

In Figure 6 the ellipses are converted to uncertainty bands (yellow dotted lines), while the red lines show the interpolated muscle annotations.

Figures 7 and 8 give the histogram and cumu-358 lative distribution plots for the external test set, 359 which also shows an acceptable match with a maxi-360 mum difference of 0.1248. We see that the empirical 361 distribution is shifted to the left, compared to the 362 standard Laplace. This is likely due to a slightly 363 different annotation practice, in that the annotators 364 of the external images may have included a larger 365 part of the blurry areas. This is confirmed by visual 366 inspection of the images, as illustrated in 9 where we 367 have problems seeing the lower part of the annotated 368 muscle. Here we have included the border predicted 369 by our model in yellow, and we see a perfect match 370 in the upper, more visible part of the muscle. 371

### 372 6 Discussion

We have successfully trained a GCN model to predict 373 the border points of pectoral muscles, while simul-374 taneously estimating the error distribution of these 375 estimates, by interpreting the model outputs as the 376 parameters of a Laplace distributions. The model's 377 uncertainty estimates were remarkably accurate, in 378 that the standardized errors followed the standard 379 Laplace distribution almost perfectly. This means 380 that the uncertainty estimates were almost perfectly 381 calibrated out-of-the-box. A practical implication 382 is that if we define a desired confidence level of, say 383 95%, the confidence intervals would in fact cover the 384



Figure 6. Sample mammograms with ground truth and uncertainty bands.



Figure 7. Histogram of standardized error distribution for the external test set.

ground truth in 95% of the cases. This will be very 385 useful for later applications of the model, since it 386 can reliably identify images that should be evaluated 387 by a human expert. The model's point prediction 388 performance on the test set was also convincing, 389 with an L1-error similar to the difference between 390 the two radiographers, which might be considered 391 a lower bound on the possible performance. One 392 might expect weaker performance because the train-393 ing set was annotated by a non-professional, but 394 this apparently made little difference. The model 395 performed reasonably well on the external data set, 396 despite substantial domain shift. We also suspect 397 that the main reason for the weaker results may be 398 a different annotation practice for uncertain cases. 399

Using the negative log-likelihood as a "loss function" may be unfamiliar to some, since it does not represent model errors directly, and is not even bounded downwards. This is not a problem, however, 403



Figure 8. Comparison of the empirical and theoretical cumulative distribution for the external test set.

as long as it points the gradients in the direction 404 that improves the model. The use of early stopping 405 through monitoring of the log-likelihood on the vali-406 dation folds was necessary, as over-fitting produces 407 models that radically underestimate the aleatoric 408 uncertainty, even though the predicted locations 409 might still be good. Therefore, monitoring of the 410 full log-likelihood was preferable to monitoring of 411 the location errors. 412

## 413 7 Conclusion

In this work, we presented a novel GCN model for de-414 tecting pectoral muscle boundaries in mammograms, 415 416 with integrated uncertainty estimation. By modeling both aleatoric and epistemic uncertainties, we were 417 able to produce accurate predictions of key bound-418 ary points, alongside uncertainty estimates that can 419 help identify cases requiring human review. Our 420 approach achieved results on par with human anno-421 tators and demonstrated robustness across domain 422 shifts, as shown in tests on the InBreast dataset. 423

The proposed uncertainty-aware framework has 424 potential applications in clinical workflows, where 425 the ability to flag uncertain cases could assist radiog-426 raphers in prioritizing manual reviews. Future work 427 will explore the integration of expert-annotated data 428 to further refine the model and improve generaliz-429 ability to diverse datasets, as well as investigate its 430 practical use in a clinical setting. 431

### 432 References

G. Waade, A. Skyrud Danielsen, Å. Holen,
M. Larsen, B. Hanestad, N.-M. Hopland, V.
Kalcheva, and S. Hofvind. "Assessment of
breast positioning criteria in mammographic
screening: Agreement between artificial intelligence software and radiographers". In:



Figure 9. Example image from the INBreast test set, where the annotated boundary may seem large.

Journal of Medical Screening 28 (Mar. 9, 439 2021), p. 096914132199871. DOI: 10.1177/ 440 0969141321998718. 441

- X. Ma, J. Wei, C. Zhou, M. A. Helvie, H.-P. 442 Chan, L. M. Hadjiiski, and Y. Lu. "Automated 443 pectoral muscle identification on MLO-view 444 mammograms: Comparison of deep neural net- 445 work to conventional computer vision". In: 446 *Medical Physics* 46.5 (May 2019), pp. 2103– 447 2114. ISSN: 2473-4209. DOI: 10.1002/mp. 448 13451. 449
- M. Brahim, K. Westerkamp, L. Hempel, R. 450 [3] Lehmann, D. Hempel, and P. Philipp. "Auto- 451 mated Assessment of Breast Positioning Qual- 452 ity in Screening Mammography". In: Can-453 cers 14.19 (Jan. 2022). Number: 19 Publisher: 454 Multidisciplinary Digital Publishing Institute, 455 p. 4704. ISSN: 2072-6694. DOI: 10.3390/ 456 cancers14194704. URL: https://www.mdpi. 457 com/2072-6694/14/19/4704 (visited on 458 02/22/2024). 459
- [4] X. Yu, S.-H. Wang, J. M. Górriz, X.-W. Jiang, 460
  D. S. Guttery, and Y.-D. Zhang. "PeMNet for 461
  Pectoral Muscle Segmentation". In: *Biology* 462
  11.1 (Jan. 14, 2022), p. 134. ISSN: 2079-7737. 463
  DOI: 10.3390/biology11010134. 464
- Y. Guo, W. Zhao, S. Li, Y. Zhang, and Y. 465
   Lu. "Automatic segmentation of the pectoral 466
   muscle based on boundary identification and 467
   shape prediction". In: *Physics in Medicine and* 468

 469
 Biology 65.4 (Feb. 19, 2020), p. 045016. ISSN:

 470
 1361-6560. DOI: 10.1088/1361-6560/ab652b.

- Z. Yang, I. Choi, J. Choi, J. Jung, M. Ryu, [6]471 and H. S. Yong. "Deep learning-based pec-472 toralis muscle volume segmentation method 473 from chest computed tomography image us-474 ing sagittal range detection and axial slice-475 based segmentation". In: PLOS ONE 18.9 476 (Sept. 5, 2023). Publisher: Public Library of 477 Science, e0290950. ISSN: 1932-6203. DOI: 10. 478 1371/journal.pone.0290950. URL: https: 479 //journals.plos.org/plosone/article? 480 id=10.1371/journal.pone.0290950 (visited 481 on 09/05/2024). 482
- W. Li, Y. Lu, K. Zheng, H. Liao, C. Lin, J. |7|483 Luo, C.-T. Cheng, J. Xiao, L. Lu, C.-F. Kuo, 484 and S. Miao. "Structured Landmark Detection 485 via Topology-Adapting Deep Graph Learning". 486 In: vol. 12354. 2020, pp. 266-283. DOI: 10. 487 1007/978-3-030-58545-7\_16. arXiv: 2004. 488 08190[cs]. URL: http://arxiv.org/abs/ 489 2004.08190 (visited on 07/12/2022). 490
- J. Kang, Q. Zhou, and H. Tong. "JuryGCN: 491 8 Quantifying Jackknife Uncertainty on Graph 492 Convolutional Networks". In: Proceedings of 493 the 28th ACM SIGKDD Conference on Knowl-494 edge Discovery and Data Mining. KDD '22. 495 New York, NY, USA: Association for Com-496 puting Machinery, Aug. 14, 2022, pp. 742-497 752. ISBN: 978-1-4503-9385-0. DOI: 10.1145/ 498 499 3534678.3539286. URL: https://dl.acm. org/doi/10.1145/3534678.3539286 (vis-500 ited on 10/21/2024). 501
- 502 [9] A. Kendall and Y. Gal. What Uncertainties Do
   503 We Need in Bayesian Deep Learning for Com 504 puter Vision? Oct. 5, 2017. DOI: 10.48550/
   505 arXiv.1703.04977. arXiv: 1703.04977 [cs].
   506 URL: http://arxiv.org/abs/1703.04977
   507 (visited on 09/12/2024).
- E. Hüllermeier and W. Waegeman. "Aleatoric [10]508 and epistemic uncertainty in machine learning: 509 an introduction to concepts and methods". 510 In: Machine Learning 110.3 (Mar. 1, 2021), 511 pp. 457–506. ISSN: 1573-0565. DOI: 10.1007/ 512 s10994-021-05946-3. URL: https://doi. 513 org/10.1007/s10994-021-05946-3 (visited 514 on 09/05/2024). 515
- 516 [11] S. Vashishth, P. Yadav, M. Bhandari, and P.
  517 Talukdar. Confidence-based Graph Convolu518 tional Networks for Semi-Supervised Learn519 ing. arXiv.org. Jan. 24, 2019. URL: https:
  520 //arxiv.org/abs/1901.08255v2 (visited
  521 on 09/06/2024).
- F. Wang, Y. Liu, K. Liu, Y. Wang, S. Medya, and P. S. Yu. Uncertainty in Graph Neural Networks: A Survey. Mar. 11, 2024. DOI:

10.48550/arXiv.2403.07185. arXiv: 2403. 525 07185[cs,stat]. URL: http://arxiv.org/ 526 abs/2403.07185 (visited on 09/06/2024). 527

- T.-M. Dutschmann, L. Kinzel, A. ter Laak, [13]528 and K. Baumann. "Large-scale evaluation of k-529 fold cross-validation ensembles for uncertainty 530 estimation". In: Journal of Cheminformatics 531 15.1 (Apr. 28, 2023), p. 49. ISSN: 1758-2946. 532 DOI: 10.1186/s13321-023-00709-9. URL: 533 https://doi.org/10.1186/s13321-023-534 00709-9 (visited on 09/05/2024). 535
- I. C. Moreira, I. Amaral, I. Domingues, A. 536
  Cardoso, M. J. Cardoso, and J. S. Cardoso. 537
  "INbreast: toward a full-field digital mammographic database". In: Academic Radiology 539
  19.2 (Feb. 2012), pp. 236–248. ISSN: 1878-4046. 540
  DOI: 10.1016/j.acra.2011.09.014. 541