Active Multi-Task Learning with Uncertainty Weighted Loss for Coronary Calcium Scoring

³ Bernhard Föllmer^{a)}, Federico Biavati, Christian Wald

- 4 Department of Radiology, Charité-Universitätsmedizin Berlin, corporate member of Freie
- 5 Universität Berlin and Humboldt-Universität zu Berlin, Berlin, Germany
- 6 Sebastian Stober
- 7 Artificial Intelligence Lab, Otto-von-Guericke-Universität, Magdeburg, Germany
- ⁸ Jackie Ma
- ⁹ Department of Artificial Intelligence, Fraunhofer Heinrich Hertz Institute, Berlin, Germany
- ¹⁰ Marc Dewey^{b)}
- ¹¹ Department of Radiology, Charité-Universitätsmedizin Berlin, corporate member of Freie Universität Parlin and Humboldt Universität zu Parlin, Cormany
- ¹² Universität Berlin and Humboldt-Universität zu Berlin, Berlin, Germany
- ¹³ Berlin Institute of Health and DZHK (German Centre for Cardiovascular Research), partner
- 14 site Berlin, Germany
- ¹⁵ Wojciech Samek^{b)}
- ¹⁶ Department of Artificial Intelligence, Fraunhofer Heinrich Hertz Institute, Berlin, Germany
- 17

19

Version typeset February 10, 2022 bernhard.foellmer@charite.de

Abstract

Purpose: Coronary artery calcium (CAC) scoring is an independent marker for the 20 risk of coronary heart disease events. Automatic methods for quantifying CAC could 21 reduce work load and assist radiologists in clinical decision making. However, large 22 annotated datasets must be acquired and labeled to achieve very good model per-23 formance, which is an expensive process and requires expert knowledge. By labeling 24 only most informative samples, active learning can reduce the number of training data 25 required. Multi-task learning techniques can improve model performance by joint 26 learning of multiple related tasks and extraction of shared informative features. 27

Methods: We propose an uncertainty weighted multi-task model for coronary cal-28 cium scoring in ECG-gated, non-contrast enhanced cardiac calcium scoring CT. The 29 model is trained to solve the two tasks of coronary artery region segmentation (weak 30 labels) and coronary artery calcification segmentation (strong labels) simultaneously 31 in an active learning scenario to improve model performance and reduce the num-32 ber of required training samples. We compare our model with a single-task U-Net 33 and a sequential-task model, as well as other state-of-the-art methods. The model is 34 evaluated based on 1275 individual patients of three different datasets (DISCHARGE, 35 CADMAN, or CaScore) and the relationship between performance and various influence-36 ing factors (image noise, metal artifacts, motion artifacts, image quality) is analyzed. 37 **Results:** The joint learning of multiclass coronary artery region segmentation and 38 binary coronary calcium segmentation improves calcium scoring performance. Since 39 shared complementary information can be learned from both tasks, the model reaches 40 optimal performance with only 12% of the training data and one-third of the labeling 41

- time in an active learning scenario. We identified image noise as one of the most important factors influencing model performance, along with anatomical anomalies and metal artifacts.
- 45 **Conclusions:** Our multi-task learning approach with uncertainty weighted loss im-
- ⁴⁶ proves calcium scoring performance by joint learning of shared features and reduces
- ⁴⁷ labeling costs when trained in an active learning scenario.
- 48

Keywords: Coronary Artery Calcium Scoring, Deep Learning, Neural Networks, Active
 Multi-Task Learning, Uncertainty Weighted Loss

51 Abbreviations

- ⁵² CAC Coronary artery calcification.
- ⁵³ CAR Coronary artery region.
- ⁵⁴ CCTA Contrast enhances CT.
- ⁵⁵ **CSCT** Calcium Scoring CT.
- 56 **CT** Computed tomography.
- 57 CVD Cardiovascular disease.
- 58 HU Hounsfield unit.
- ⁵⁹ **ICC** Interclass correlation coefficient.
- 60 LAD Left anterior descending artery.
- ⁶¹ LCX Left circumflex artery.
- 62 LM Left main artery.
- ⁶³ MTL Multi-task learning.
- ⁶⁴ **PPV** Positive predictive value.
- 65 **RCA** Right coronary artery.

66 I. Introduction

Cardiovascular disease (CVD) is the global leading cause of death.¹ Coronary calcium is 67 commonly associated with coronary atherosclerosis and its absence is associated with a very 68 low risk of adverse coronary events.² In clinical practice, semi-automatic software is used to 69 manually select coronary artery calcifications (CAC) in computed tomography (CT) image 70 slices from automatically labeled candidates, which is a tedious and time-consuming process 71 in large-scale studies.³ Typically, ECG-gated, non-contrast enhanced computed tomography, 72 known as calcium scoring computed tomography (CSCT), is used to identify CAC.⁴ The 73 Agatston score⁵ is the most common measure used to quantify CAC with the aim of defining 74 appropriate cardiac risk categories. In recent years, deep learning models such as convolu-75 tional neural networks have been used to automatically quantify CAC based on 2D slices ^{3,4} 76 or 2.5D/3D volumetric input data 6,7 . 77

⁷⁸ Methods have been developed for different examination types such as non-contrast ⁷⁹ enhanced ECG-gated calcium scoring CTs,⁸ contrast-enhanced coronary CT angiography (CCTA)⁹ or a combination of both.^{3,10,11} Since the segmentation of the cardiac tree is very challenging in CSCT, methods using non-contrast enhanced and contrast enhanced CT usually map spatial information about coronary arteries from the contrast enhanced CT to the non-contrast enhanced CT.^{10,12} Most methods perform segmentation of the calcified lesion to estimate the Agatston score and classify detected calcification based on the corresponding left anterior descending artery (LAD), left circumflex artery (LCX), and right coronary artery (RCA), but some also perform regression directly.¹³

Most of the current state-of-the-art methods only learn from sparse calcifications. Therefore, very large and heterogeneous datasets need to be acquired and labeled to train models that are robust and achieve satisfactory performance to be used in a clinical setting. Unfortunately, this is an expensive and time consuming process and requires expert knowledge. It calls for methods that can reduce labeling costs and improve performance by integrating the radiologist into the training process.

Active learning techniques are able to reach higher performance while using a smaller num-93 ber of annotated training samples by active sample selection and therefore reduce labeling 94 costs. In active learning, the learner (deep neural network) iteratively selects only the most 95 informative samples based on a selection criterion such as uncertainty sampling, query by 96 committee, expected error reduction, or expected model for labeling.¹⁴ The method inte-97 grates the radiologist into the training process and avoids labeling of uninformative samples. 98 Spatial information about coronary arteries and corresponding coronary calcifications is very 99 important to distinguish between coronary and extra-coronary calcifications. Calcifications 100 are usually very sparse, which makes it difficult to extract features with spatial information. 101 Extraction of spatial information about the coronary arteries in an auxiliary task can aid in 102 the localization of coronary calcifications. Multi-task learning (MTL) is a technique which 103 learns multiple related tasks together, to improve model performance by sharing complemen-104 tary information.¹⁵ In coronary calcium scoring, the spatial information of coronary arteries 105 is closely related to the calcium scoring task and therefore learning of coronary artery regions 106 serves as a good auxiliary task to support the original calcium scoring task. However, the 107 optimization of multiple loss functions for multi-task learning is a crucial factor and tuning 108 loss weighting by hand is difficult and computationally expensive. Many task balancing 109 approaches for dense predictions such as static weighting. GradNorm¹⁶, dynamic weight av-110 erage¹⁷, dynamic task prioritization¹⁸ or uncertainty weighted loss¹⁹ have been developed 111

and have shown that the best optimization method should be selected on a per case basis.²⁰ 112 The training of multi-task models in an active learning scenario can be challenging if the 113 dataset is very small. In this work, we exploit a multi-task model (MTL-model) with 114 uncertainty-weighted loss that outperforms a single-task U-Net and a sequential-model. The 115 model achieves very good performance on small training sets and can therefore be used in 116 active learning scenarios. The model performs as well as other state-of-the-art methods 117 and achieves similar results compared to our statically-weighted MTL-model with optimally 118 chosen weighting parameter. The contributions of this paper can be summarized as follows: 119

• We propose a novel learning paradigm for coronary calcium scoring by simultaneous learning of multiple related tasks to increase data efficiency and model performance by leveraging auxiliary information through shared informative features.

• We propose a multi-task encoder-decoder model for simultaneous coronary artery regions segmentation (multiclass) and coronary artery calcification segmentation (binary) to improve model performance compared to single-task models.

- We show that our model obtains optimal performance with substantially less training data (12%) and reduces annotation time to one-third in an active learning scenario compared to training on the full dataset.
- We demonstrate the importance of loss weighting for optimal model performance of our multi-task model and show how uncertainty weighted loss can facilitate active multi-task learning.
- We show that our model performs almost as well as the best state-of-the-art methods in terms of F1-score, intraclass correlation coefficient (ICC) and sensitivity of CAC volume, on a common benchmark dataset for coronary calcium scoring.

¹³⁵ II. Materials and methods

In this Section we present our multi-task model which performs simultaneous segmentation of coronary artery regions and segmentation of coronary artery calcifications. In Section II.A. we describe the used datasets and corresponding annotation strategies. We introduce the multiclass coronary artery region segmentation task II.B.1. and the binary lesion segmentation task II.B.2.. We propose the multi-loss optimization method using uncertainty weighted loss II.B.3. and give a detailed description of our implemented network architecture and training procedure II.B.4.. In Section II.B.5. we introduce a single-task U-Net and sequential model as
comparison models. In Section II.C. we introduce our active learning approach in which we
use only the most informative samples to decrease annotation costs and propose our hybrid
sampling strategy.

146 II.A. Datasets

For the evaluation of our multi-task learning approach, we compare the performance on three
different datasets. A detailed flowchart about the dataset selection process can be found in
the supplementary material.

The DISCHARGE-trial is a prospective multi-center randomized controlled trial to 150 examine for which patients with suspected coronary artery disease based on stable chest 151 pain, cardiac CT or cardiac catheterization is best suited.^{21,22} Our DISCHARGE dataset 152 consists of calcium scoring CTs (CSCT) from 1262 patients (708 male, 554 female) of the 153 trial. Image data were acquired from 26 clinical sites using 14 different scanner types. 154 Annotations for coronary artery calcification were acquired for all scans. Weak annotations 155 of coronary artery regions were only acquired for 215 randomly selected scans and randomly 156 divided into 140 CT scans (6721 slices) for training (65%) and 75 CT scans (3636 slices) for 157 validation (35%). All remaining 1047 CT scans (57452 slices) were used as test set. Only 158 one CSCT from each patient was selected for the dataset. The reconstruction of the CT 159 scans was performed using filtered back projection method (383) and iterative reconstruction 160 methods (879). To keep the data as close as possible to real-life clinical data, diagnostic CTs 161 including scans with metal artifacts (pace-maker, artificial valves, etc.), scans with severe 162 motion artifacts, high level of noise or anatomical abnormalities were not excluded. Note, 163 annotations of coronary artery regions (CAR) and coronary artery calcifications (CAC) are 164 available for training and validation set. For the test set, only CAC annotations are available. 165 The annotations were performed by two observers. Observer one was a trained physician 166 who annotated coronary calcifications, observer two was a trained medical imaging scientist 167 who annotated coronary artery regions. Available contrast-enhanced CT scans (CCTA) were 168 not included because the overall goal of the method is to predict coronary heart disease risk 169 without the need to inject a contrast agent. 170

171

The second dataset consists of CT scans from the publicly available orCaScore challenge

on (semi-)automatic coronary calcium scoring.²³ The framework provides 72 pairs of CSCT 172 and corresponding contrast-enhanced CT angiography (CCTA) from the same patient ac-173 quired at four academic hospitals. The data has been divided into a 32-scan training set 174 and a 40-scan test set. For the training set, a reference standard by two expert observers, a 175 radiologist with > 12 years of experience in CAC scoring and a research physician, are pro-176 vided. CT scans with anatomical abnormalities, intracoronary stents, and metal implants 177 as well as CTs showing severe motion artifacts or extremely high levels of noise determined 178 by visual inspection were excluded. Annotations of CAR were additionally acquired for the 179 training set. 180

The third dataset consists of CSCT from the single-center randomized controlled Coronary Artery Disease Management (CAD-Man) study.²⁴ The dataset consists of 156 CT scans and annotations were only acquired for CAC. The dataset serves as an additional test set. The reconstruction of the CT scans was performed using filtered back projection method. Reference standards are provided by one expert observer. The annotations of coronary calcifications and artery regions were performed by a trained medical imaging scientist.

¹⁸⁷ Differences between the three datasets regarding distribution of candidate lesions are ¹⁸⁸ shown in Table 1. A candidate lesion is defined as connected 3-D image voxels (6-¹⁸⁹ connectivity) with intensities greater than 130 HU.

	# geong		ICY		OTHED CAD	Candidates
	# scans		LUA	INCA		per scan
DISCHARGE Training	140	344	168	338	865k	6183
DISCHARGE Test	1047	2375	1042	1872	6254k	5978
DISCHARGE Validation	75	198	118	221	432k	5768
orCaScore Training	32	103	21	56	138k	3454
orCaScore Test	40	-	-	-	-	-
CADMAN Test	156	335	151	156	1400k	8980

Table 1: Number of candidate lesions (connected 3-D image voxels with intensities greater than 130 HU) distributed in the DISCHARGE training set, validation set and test set, orCaScore training set and CADMAN test set. Candidate lesions are distributed over calcified coronary artery lesions (LAD, LCX, RCA) and other structures such as bones or extra-coronary calcifications (OTHER_CAC). Since the orCaScore test set is not public, no information about distribution of candidate lesions is available.

¹⁹⁰ II.A.1. Annotation procedure

The coronary artery tree is divided into three sub-trees corresponding to the left anterior descending artery (LAD), left circumflex artery (LCX) and right coronary artery (RCA) including main branch (LM) and its side-branches. The left main artery (LM) is included in the sub-tree of the LAD. The training of the multi-task model requires annotations for the two task of coronary artery calcification (CAC) segmentation and coronary artery region (CAR) segmentation.

¹⁹⁷ II.A.2. Annotation of coronary artery calcifications (CAC)

The annotation of the CAC was performed by thresholding and highlighting all voxels with 198 a density above 130 HU (candidate lesions). The observer annotates all highlighted voxels of 199 calcified lesions and assigns a class according to the corresponding coronary artery ("LAD", 200 "LCX", "RCA"). Calcified lesions corresponding to multiple arteries (e.g. calcified lesions 201 in bifurcations) were divided by annotating the voxels according to the arteries. For model 202 evaluation, lesions where defined as connected voxels (6-connectivity) with a minimum lesion 203 volume of 1.5mm³ All remaining candidate lesions were annotated as "OTHER_CAC". The 204 annotations of coronary artery calcifications were performed using an in-house developed 205 semi-automatic segmentation module for 3D slicer²⁵. 206

²⁰⁷ II.A.3. Weak annotation of coronary artery regions (CAR)

The annotations of CARs are acquired using weak annotations (scribbles), since a precise 208 segmentation of the vessel tree is impossible in non-contrast CTs due to the missing contrast 209 between arteries and surrounding tissue. To overcome the problem of misleading labels, we 210 did not label regions between arteries and surrounding tissue which are difficult to distinguish 211 or a precise labeling of the boundary would be extremely time-consuming. To facilitate and 212 speed up the annotation process, we used an in-house semi-automatic segmentation mod-213 ule developed for 3D Slicer²⁵. At first, the annotator was using a scribble to annotate the 214 three main arteries in each slice. In the second step, an additional scribble (closed contour) 215 was used to surround the arteries and isolate the annotated artery scribble from the tissue. 216 In the third step, connected component analysis was performed to divide the annotations 217

into different components. The largest component (background) was joint with the closed
contour scribbles and labeled as OTHER_CAR. If no coronary artery was seen in the slice,
the annotator only placed a single scribble for OTHER_CAR in the image. Examples of the
performed annotations can be seen in Figure 2.

II.B. Multi-task segmentation network with uncertainty weighted loss

We propose a multi-task segmentation network following an encoder-decoder structure with 224 skip connections, inspired by the U-Net architecture.^{26,27} The multi-task network archi-225 tecture is illustrated in Figure 1 and performs multiclass coronary artery region (CAR) 226 segmentation and coronary artery calcifications (CAC) segmentation at the same time. Fea-227 tures extracted by the encoder are shared with the two decoders for the tasks of multiclass 228 coronary artery region segmentation (T_R) and binary segmentation of calcified lesions (T_L) . 229 Since the information of predicted coronary artery regions is a useful prior information for 230 segmentation of calcified lesions, feature maps extracted by the decoder for coronary artery 231 region segmentation are shared with the decoder for binary segmentation of calcified lesions. 232 To utilize this prior information about candidate lesions, we concatenate the image slice (512) 233 $px \ge 512 px$) with a candidate lesion mask to form the input tensor. The candidate lesion 234 mask was created by thresholding the image using a constant threshold of 130 HU. During 235 training, the losses of both task L_R (loss for task T_R) and L_L (loss for task T_L) are combined 236 using an uncertainty weighting loss,¹⁹ to jointly optimize the model parameters. 237

²³⁸ II.B.1. Coronary artery region segmentation task

The network aims to learn coronary artery regions from weakly labeled regions, as shown in Figure 2. Weak labels²⁸ are defined as segmentations, which are imprecise but less costly to obtain than pixel-level annotations. Since in non-contrast enhanced CT scans the spatial boundary between the vessels and the surrounding tissue cannot be determined precisely, pixels x of an image i are either annotated and belong to the annotated pixel set $\Omega_{R,i}$ with one of the CAR classes {LAD, LCX, RCA, OTHER_CAR}, or are not annotated. The pixel-wise softmax function ²⁶ and focal loss ²⁹ are used to deal with large class imbalance

II.B. Multi-task segmentation network with uncertainty weighted loss

247



Figure 1: Multi-task model for coronary artery calcification (CAC) scoring. The image and the CAC candidate lesion mask are concatenated to form the input tensor. The model consists of one encoder that shares feature maps with two decoders of the multiclass coronary artery region (CAR) segmentation task T_R and binary CAC segmentation task T_L . Predictions are combined by multiplying binary CAC segmentation with multiclass CAR segmentation to perform multiclass calcification segmentation

between background pixels (OTHER_CAR) and pixels of coronary artery regions.

$$\mathcal{L}_{R,i} = \sum_{x \in \Omega_{R,i}} \sum_{c_R=1}^{4} -w_{c_R} y_{c_R}(x) (1 - p_{c_R}(x))^{\gamma_R} \log(p_{c_R}(x))$$
(1)

The γ_R parameter smoothly adjusts the rate at which easily segmented pixels are downweighted and w_R balances the loss. $p_{c_R}(x)$ and $y_{c_R}(x)$ are the pixel-wise softmax output and the reference class of pixel $x \in \Omega_{R,i}$, respectively. The pixel set $\Omega_{R,i}$ contains all labeled pixels. Unlabeled pixels (gaps) $x \notin \Omega_{R,i}$, as shown in Figure 2, are ignored and not used for loss calculation. The parameter w_{c_R} is a weighting parameter that balances the importance of the classes and handles the data imbalance problem. The parameter c_R is the channel of the corresponding CAR class.

²⁵⁵ II.B.2. Binary lesion segmentation task

The lesion segmentation network performs a binary segmentation of candidate coronary artery lesions into the classes $\{CAC, OTHER_CAC\}$. Feature maps extracted by the decoder for the coronary artery region segmentation are shared with the decoder for the binary lesion segmentation. The binary focal loss $\mathcal{L}_{L,i}$ defined as (2) is calculated based on all voxels from candidate lesions grouped in the set $\Omega_{L,i}$.

$$\mathcal{L}_{L,i} = \sum_{x \in \Omega_{L,i}} \sum_{c_L=1}^{2} -w_{c_L} y_{c_L}(x) (1 - p_{c_L}(x))^{\gamma_L} \log(p_{c_L}(x))$$
(2)

The parameters
$$p_{c_L}(x)$$
, $y_{c_L}(x)$ as well as $x \in \Omega_{L,i}$ and w_{c_L} are defined analogouesly to

- the region segmentation in Subsection II.B.1.. The output of the binary CAC segmentation
- decoder is multiplied (channel-wise) with the output of the CAR segmentation decoder to

²⁶⁵ perform multiclass CAC segmentation.



Figure 2: Multi-task annotations of an image slice with coronary artery calcifications (CAC) in the left anterior descending artery (LAD), left circumflex artery (LCX) and right coronary artery (RCA) (A). Weak annotations of coronary artery regions (CAR) for the LAD - red, LCX - yellow, RCA - blue and OTHER_CAR - green (B). Strong annotations of coronary artery calcifications in the LAD - red, LCX - yellow, RCA - blue and other objects with density higher 130 HU (OTHER_CAC) - green (C).

²⁶⁶ II.B.3. Uncertainty based weighted loss

The performance of a multi-task network depends strongly on the weighting of the losses. The most commonly used loss weighting strategy for multi-task learning is static weighting, which computes a weighted sum of the losses using balancing parameters α_i . The static weighted loss of our multi-task model is the weighted sum of the losses for multiclass segmentation of coronary artery regions \mathcal{L}_R and the binary segmentation of coronary calcifications \mathcal{L}_L , as shown in Equation (3).

 $\mathcal{L}_{total}(\mathbf{W}) = \alpha \mathcal{L}_R(\mathbf{W}) + (1 - \alpha) \mathcal{L}_L(\mathbf{W})$ (3)

261

This method is simple but unfortunately computationally expensive to fine tune.³⁰ The determination of the optimal weighting parameter value α is even more challenging in active multi-task learning, since the model is initially trained with a very small number of annotated training data. Other methods are based on dynamic weight average (DWA) using task-specific feature-level attention¹⁷ or use gradient normalization¹⁶ to balance losses.

For our multi-task calcium scoring model, we use the uncertainty weighted loss method of 279 Cipolla et al.¹⁹. The uncertainty based weighting uses homoscedastic uncertainty to weight 280 loss functions of each task.¹⁹ We combine the outputs of the last layers (softmax output) 281 from the decoders based on the homoscedastic uncertainty. To model the uncertainty, we 282 introduce the positive scalar σ_R for coronary artery region segmentation task and σ_L for the 283 binary calcification segmentation task. The parameters can be interpreted as Boltzmann 284 distributions (also called Gibbs distribution) where the input is scaled by σ_R^2 and σ_L^2 respec-285 tively. The total loss \mathcal{L}_{total} in Equation (4) is an uncertainty weighted loss of \mathcal{L}_R and \mathcal{L}_L 286 where W represents the parameters of the multi-task network. A detailed deviation can be 287 found in the supplementary material. 288

289

290

$$\mathcal{L}_{total}(\mathbf{W}, \sigma_R, \sigma_L) = \frac{1}{\sigma_R^2} \mathcal{L}_R(\mathbf{W}) + \frac{1}{\sigma_L^2} \mathcal{L}_L(\mathbf{W}) + \log \sigma_R + \log \sigma_L$$
(4)

This loss is smoothly differentiable, and is well formed such that the task weights will not converge to zero. For practical reasons, we predict the log variance $\log \sigma^2$, which is more stable and avoids any division by zero.¹⁹

²⁹⁴ II.B.4. Multi-task network architecture and training procedure

To train the multi-task network, we oversample slices with calcifications to form balanced 295 mini batches (20% samples with calcifications, 80% without calcifications). The encoder con-296 sists of eight downsampling blocks, each block consists of two convolutional layers, dropout 297 layer, batch normalization ³¹ and ReLU activation function.³² The two decoders consist of 298 eight upsampling blocks, where each block consists of a bilinear upsampling layer and two 299 convolution layers, dropout layer, batch normalization layer and ReLU activation function, 300 respectively. The feature maps of the upsampling block of the coronary artery region (CAR) 301 segmentation decoder are shared with each upsampling block of the coronary calcification 302

segmentation decoder, but not vice versa, to follow the causal relation between coronary 303 artery regions and coronary artery calcifications. Skip connections between the encoder 304 and the two decoders are implemented as concatenations and used to share feature maps 305 from the respective downsampling block.²⁶ The model was trained with a batch size of 8, 306 the Adam³³ optimizer, an initial learning rate of 5e-04, learning rate decay of 0.95 after 307 every 5 epochs and L2 weight decay. During training, the dropout rate of the inner layer 308 between the encoder and decoder was set to 0.5, and the dropout rate for all other lay-309 ers was set to zero. During our experiments, we found that the convergence of the static 310 weighted loss MTL-model depends strongly on the initial learning rate but due to high 311 training times, we did not perform further detailed hyper parameter analysis. We use focal 312 loss for both training tasks. The focal loss parameters of task T_R were set to $\gamma_R = 2.0$, 313 $\alpha_{OTHER CAR} = 0.01, \alpha_{LAD} = 1.0, \alpha_{LCX} = 1.0, \alpha_{RCA} = 1.0$. The loss parameters of the coro-314 nary calcification segmentation task were set to $\gamma_L = 2.0$, $\alpha_{CAC} = 1.0$, α_{OTHER} $_{CAC} = 0.01$. 315 To train the network, we augmented the image slices by small translations to prevent over-316 fitting. To perform multiclass calcification segmentation, the predictions of the two tasks are 317 combined by multiplying the binarized calcification segmentation with the coronary artery 318 region segmentation. To avoid overtraining, we use early stopping based on the perfor-319 mance of the validation set. The training stopped after approximately 200k iterations (240 320 epochs), where one iteration corresponds to a batch of 8 slices. Training was performed 321 on an NVIDIA Tesla V100, 32GB and PyTorch framework. More details and a pretrained 322 model can be found at (https://github.com/Berni1557/MTAL-CACS). 323

³²⁴ II.B.5. Single-task model, sequential-task model and multi-task model

We compare our multi-task model with a single-task and a sequential-task model in Figure 3 325 to show that simultaneous training of related tasks can extract informative shared features 326 and improve model performance. The single-task model consists of a multiclass U-Net²⁶ with 327 the same downsampling and upsampling block architecture as in the multi-task network. 328 The last layer consists of four channels (OTHER_CAC, LAD, LCX, RCA) for multiclass 329 segmentation of coronary calcifications. The sequential model consists of two separated 330 models. The first model is trained for multiclass coronary artery region (CAR) segmentation. 331 After the training has finished, the predictions are used for the training of the coronary 332 calcification (CAC) segmentation network. Therefore, the CAR predictions are concatenated 333

with the image and CAC candidate mask and serve as input for the binary segmentation network for coronary calcifications. The goal of the sequential model is to follow the causal relation between CAR and CAC.



Figure 3: Single-task model, sequential-task model and multi-task model architecture comparison. The single-task model (lower, left) consists of a multiclass U-Net. The sequential model (upper) consists of a model for multiclass coronary artery region (CAR) segmentation whose predictions are used to train the coronary artery calcification (CAC) segmentation network. The multi-task model (lower, right) consists of one encoder and two decoders for the prediction of coronary artery region (CAR) segmentations and coronary calcification (CAC) segmentations which are combined for multiclass segmentation of CAC.

³³⁷ II.C. Active learning with uncertainty weighted multi-task model

Labeling of coronary calcifications in CT scans is a laborious and time-consuming task and requires significant expert knowledge.³⁴ Labeling for MTL methods tends to be more expensive since each task requires its own annotations. Active learning is able to reduce the costs

by iteratively labeling only most informative samples thus achieving optimal performance 341 with a smaller number of samples.³⁵ In multi-task learning with static weighting parameter 342 α , the best performing parameter has to be determined, which is a difficult and expensive 343 process¹⁹ and is often performed using grid search on the entire annotated dataset. In active 344 learning, the estimation of a static weighting parameter α is even more challenging to tune, 345 since the data distribution changes after each sampling round and thereby the optimal value 346 of parameter α changes as well. Moreover, the estimation on small datasets can be very 347 sensitive to the randomly drawn initial training samples. 348

We simulate active learning to investigate whether our uncertainty weighted loss model can overcome these problems. There are several approaches for active multi-task learning such as active learning via bandits,³⁶ active learning frameworks for adaptive filtering³⁷ or value of information based methods.³⁸

For our approach, we developed a sampling strategy based on uncertainty sampling and 353 random sampling, which we call the hybrid sampling strategy. First, we apply monte carlo 354 dropout (MCD)³⁹ during inference for all samples which are not in the training set, predict 355 segmentation maps and repeat this process $N_{MCD} = 10$ times. Dropout rate was set to 0.01 356 for all dropout layers. Based on the predictions, we estimate the MC sample variance⁴⁰ for 357 each pixel, corresponding to candidate calcifications (pixel with density values greater than 358 130 HU) and calculate the average variance for each sample. We sort all samples in descend-359 ing order and randomly sample from the top 20% with highest variance. Selected samples 360 and respective annotations are added to the training set. We use this simple strategy, since 361 it is not our goal to improve sampling strategies, but rather to investigate their general 362 applicability. This sampling strategy has low computational complexity and increases the 363 diversity of batch query samples.³⁵. We compare our hybrid sampling method with the 364 random sampling method, where we randomly select samples from the unlabeled dataset for 365 labeling and training. 366

³⁶⁷ III. Results

In this Section we first introduse used performance metrics compare the performance of our 368 proposed multi-task model with the singe-task U-Net and sequential-task model trained on 369 the full DISCHARGE training set III.B.. In Subsection III.C., we compare our multi-task 370 model with other state-of-the-art models. In Subsection III.D., we analyze our uncertainty 371 weighted multi-task model in an active learning scenario and show that the number of 372 required training samples and annotation time can be reduced compared to labeling the 373 full training set. Finally, we analyze the influence of image noise, metal artifacts, motion 374 artifacts and image quality on model performance in Subsection III.E. 375

³⁷⁶ III.A. Performance metrics

Mi

The performance of our models for coronary calcium scoring in Table 2 was evaluated on 377 volume and lesion level with binary and multiclass segmentation metrics⁴¹. For the eval-378 uation of the multiclass coronary artery region segmentation task T_R , we use the Micro 379 F1-score on volume level. The Micro F1-score is the harmonic mean of Micro precision and 380 Micro recall based on the the coronary arteries, excluding the OTHER_CAR class (5). For 381 the Micro precision and Micro recall, the number of true positives (TP_{sum}) is the number 382 of all correctly classified pixels of the coronary artery regions, excluding pixels of the class 383 OTHER_CAR. The number of false positives (FP_{sum}) is the number of pixels belonging 384 to the class OTHER CAR but being misclassified as one of the coronary arteries, plus all 385 pixels of coronary arteries that are incorrectly assigned to another artery. The number of 386 false negatives (FN_{sum}) is the number of pixels belonging to the coronary arteries but being 387 misclassified as OTHER CAR, plus all pixels of coronary arteries that are incorrectly as-388 signed to another artery. Therefore, misclassifications between arteries are counted as false 389 negatives and false positives. 390

391

$$cro F1-score = 2 * \frac{Micro-precision * Micro-recall}{Micro-precision + Micro-recall}$$
(5)

$$\text{Micro-recall} = \frac{TP_{sum}}{TP_{sum} + FN_{sum}} \tag{6}$$

392

393

$$\frac{\text{Iicro-precision} = \frac{TP_{sum}}{TP_{sum} + FP_{sum}}$$
(7)

394	The evaluation of the binary coronary calcification task T_L was evaluated based on
395	the positive predictive value (PPV), sensitivity and F1-score. The evaluation of the result-
396	ing multiclass calcification segmentation was evaluated based on the F1-score calculated
397	irrespective of the artery-specific label, to be comparable with other methods. For the com-
398	parison with other methods, we use to intraclass correlation (ICC), sensitivity and F1-score
399	in Table 4 and 3.
400	To evaluate our active learning method in Figure 5, we used the Micro-F1 score of the re-
401	sulting multi-class calcification segmentation.
402	The risk categorization performance in Table 5 was evaluated based on the linearly weighted

403 Cohen's kappa as a measure of agreement between the reference category and the catego-

⁴⁰⁴ rization based on the MTL-model.

⁴⁰⁵ III.B. Comparison of single-task, sequential-task and multi-task ⁴⁰⁶ model

We trained all three models described in Section II.B.5. on the full DISCHARGE training set 407 and evaluate the performance based on the DISCHARGE test set. In Table 2 we compare 408 results for coronary artery region (CAR) segmentation task T_R in terms of Micro F1-score 409 and binary coronary calcification (CAC) segmentation task T_L in terms of F1-score, positive 410 predictive value and sensitivity. The Micro F1-score is reported for the resulting multiclass 411 calcification segmentation. For T_R we report Micro F1-scores only for the validation set, 412 since annotations of the DISCHARGE test set were not available for CAR. For the static-413 MTL model we set weighting parameter to the optimal value $\alpha = 0.4$, determined based on 414 the maximum Micro F1-score for calcification segmentation using grid-search method. Un-415 certainty weighted loss MTL-model and static weighted MTL-model with optimal weighting 416 parameter value reached similar performance of F1-score=0.881 and F1-score=0.882, re-417 spectively. Both MTL-models (static weighted MTL and uncertainty weighted loss MTL) 418 outperform the single-task model (F1-score=0.804) and sequential model (F1-score=0.769). 419 The performance of the multi-task models is very good at the volume level, but lower at lesion 420 level, due to a false positive predictions caused by misclassification of noise. As expected, 421

	CAR	Bi	Binary calcification		
	segmentation	seg	segmentation		
	task T_R	tas	sk T_L		segmentation
	Micro F1	PPV	Sen.	F1	F1
	(Vol)	(Vol/Num)	(Vol/Num)	(Vol/Num)	(Vol)
Single-task		0.900	0.726	0.804	0.775
(U-Net)	-	(0.219)	(0.923)	(0.353)	0.775
Sequential-task	0.479	0.916	0.663	0.769	0.740
model	0.472	(0.231)	(0.871)	(0.365)	0.740
Static weighted		0.937	0.833	0.882	
MTL-model	0.459	(0.412)	(0.883)	(0.562)	0.850
$(\alpha = 0.4)$		(0.412)	(0.000)	(0.502)	
Uncertainty		0.024	0.949	0.991	
weighted loss	0.451	(0.924)	(0.890)	(0.561)	0.849
MTL-model		(0.413)	(0.880)	(0.562)	

Table 2: Performance comparison of the single-task model (U-Net), sequential-task model, static weighted MTL-model and uncertainty weighted loss MTL-model. Evaluation is based on Micro F1-score for coronary artery region segmentation task T_R (only available for DIS-CHARGE validation dataset) and F1-score, positive predictive value (PPV) and sensitivity (Sen.) for binary calcification segmentation task T_L and Micro F1-score for the combined multiclass calcification segmentation of the DISCHARGE test dataset on volume and lesion level.

the best performing model for coronary artery region segmentation task T_R is the sequential model (Micro F1-score=0.472), since the first of the two sequential networks performs only this task. Figure 4 shows an example for the predictions of CAR and CAC by the multi-task network. We compare the multi-task predictions for severe noise in Figure S3 (supplementary material) to show how the uncertainty weighted loss MTL-model outperformed the sequential model.

⁴²⁸ III.C. Performance comparison with other methods

To compare our model (uncertainty weighted loss MTL-model) against other methods, we evaluate the performance based on the orCaScore test set described in Section II.A.. The orCaScore dataset does not provide any reference annotations for the test set and is therefore well suited for model comparison. For a fair comparison we trained our model twice: once on the DISCHARGE training set and once on the orCaScore training set. Results on the orCaScore test set compared with other methods are shown in Table 3. The results of



Figure 4: Visualization of overlap between predicted coronary calcifications (CAC) and coronary artery regions (CAR). Coronary calcifications in the left anterior descending artery (LAD), left circumflex artery (LCX) and right coronary artery (RCA) (A). Predicted coronary artery regions of the LAD - red, LCX - yellow and RCA - blue (B) and 3D surface model of the segmented CAR (C).

our model evaluated on the DISCHARGE test set and CADMAN test set is compared with
methods evaluated on other non-public datasets in Table 4. Note that results are not directly
comparable due to unknown data distributions.

We can see that our model trained on the DISCHARGE training set (F1-score=0.958) 438 performs very good using only CSCTs. The best performing method of Gogin et al.⁷ (F1-439 score=0.975) is using an ensemble of 3D CNNs to perform calcium scoring. Other methods 440 are using CCTA to map segmentations from cardiac structures (heart, aorta, coronary ar-441 teries) in the CCTA to the CSCT 7 or use preprocessing by cylindrical cropping around 442 an initial automatic segmentation of the ascending aorta.³. In the work of D. Eng et al. 443 ⁴⁵, two deep learning models were used to automate CAC scoring using gated unenhanced 444 coronary CTs and non-gated unenhanced chest CTs, but reported performance metrics are 445

⁴⁴⁶ not comparable with those in Table 4.

⁴⁴⁷ Note that our model performs well on the full DISCHARGE test set (F1-score=0.881), ⁴⁴⁸ however due to the large variability in the dataset and inclusion of scans with motion and ⁴⁴⁹ metal artifacts, the performance is lower than on the orCaScore test set (F1-score=0.958). ⁴⁵⁰ Similar dataset dependent performance differences can be seen in Table 3 and 4 by Wolterink ⁴⁵¹ et al.³ and Zhang et al.⁸. In Section III.E. we analyze different influencing factors to find ⁴⁵² reasons for these surprising findings.

453 The per patient risk categories were predicted based on the estimated agatston score of

		D	TOO	a	
Mathada	Interaction	Dataset	ICC	Sen.	F1.
Methous	Interaction	# scans (train, test)	(Vol)	(Vol)	(Vol)
Observer 1 ³	Manual	CSCT (-, 40)	0.998	0.985	0.9860
Observer 2 ³	Manual	CSCT (-, 40)	0.984	0.998	0.975
Shahzad et al. ⁴²	Automatic	CSCT+CCTA (209, 40)	0.971	0.621	0.893
Yang et al. ¹⁰	Semi-Auto.	CSCT+CCTA (40, 40)	0.992	0.940	0.968
Kelm et al. ⁴³	Automatic	CSCT+CCTA (32, 40)	0.980	0.838	0.943
Kondo et al. ⁸	Semi-Auto.	CSCT+CCTA (32, 40)	0.621	0.513	0.623
Durlak et al. ⁴⁴	Automatic	CSCT (32, 40)	0.989	0.835	0.951
Wolterink et al. ⁴³	Automatic	CSCT (373, 40)	0.986	0.845	0.947
Zhang et al. ⁸	Automatic	CSCT (129, 40)	0.991	0.911	0.954
Gogin et al. ⁷	Automatic	CSCT (783, 40)	0.995	0.968	0.975
Proposed network	Automatic	CSCT (215, 40)	0.994	0.955	0.958
Proposed network (orCaScore train)	Automatic	CSCT (32, 40)	0.984	0.961	0.928

Table 3: Performance comparison between our model and other state-of-the-art methods for automated coronary calcium scoring in cardiac CT on the orCaScore test set. Comparison is based on interclass correlation coefficient (ICC), Sensitivity (Sen.) and F1-score for CAC volume. The first block shows the performance of the two observers on the orCaScore test set. The second block shows results of all methods using non-contrast enhanced CT (CSCT) and contrast-enhanced coronary CT angiography (CCTA) on the orCaScore test set. The third block shows results of all methods using only CSCT on the orCaScore test set.

Mathada	Dataset	ICC	Sen.	F1
Methods	# scans (train, test)	(Vol.)	(Vol.)	(Vol.)
Kurkure et al. ⁴⁶	CSCT (100, 105)	-	0.921	-
Išgum et al. ⁴⁷	CSCT (228, 76)	-	0.738	-
Brunner et al. ⁴⁸	CSCT (30, 30)	-	0.863	-
Shahzad et al. ⁴²	CSCT (209, 157)	-	0.839	-
Zhang et al. ⁸	CSCT $(129 \text{ with 5-fold CV})$	0.986	0.905	0.946
Wolterink et al. 43	CSCT (373, 530)	0.96	0.79	0.85
Vos et al. ¹³	CSCT (373, 530)	0.97	-	-
Zeleznik et al. ⁴⁹	CSCT (129, (441, 663, 4021))	0.89, 0.80, 0.792	-	-
Velzen et al. ⁵⁰	CSCT (373, 529)	0.970	-	-
Proposed network	CSCT (215, 1047)	0.055	0.941	0.001
(DISCHARGE train)	-DISCHARGE test	0.900	0.041	0.001
Proposed network	CSCT (215, 154)	0.847	0.041	0.822
(DISCHARGE train)	-CADMAN test	0.041	0.941	0.022

Table 4: Results of state-of-the-art methods for automated coronary calcium scoring in cardiac CT on non-public datasets. Results are compared in terms of interclass correlation coefficient (ICC), Sensitivity (Sen.) and F1-score. The method "Proposed network (DIS-CHARGE train)" refers to the proposed uncertainty weighted MTL-model trained on the DISCHARGE training set. The proposed network was evaluated on the DISCHARGE test set and CADMAN test set.

the CAC segmentations and compared with the risk categories based on the reference an-454 notations. The confusion matrices of risk category predictions and corresponding linearly 455 weighted Cohen's kappa (κ) for all three datasets are shown in Table 5. We use a linearly 456 weighted kappa because risk categories are on an ordinal rating scale and the deviations 457 are weighted differently depending on their size. It shows that κ is much higher for the 458 orCaScore dataset (κ =0.97) compared to the DISCHARGE (κ =0.80) or CADMAN dataset 459 $(\kappa=0.80)$. Misclassifications of the risk category occurs mainly between category I and II 460 because of false positive predictions. 461

a) DISCHARGE test set, $\kappa = 0.80$

	Au	Automated risk category					
Risk	Ι	II	III	IV	Total		
Ι	267	159	8	7	441		
II	3	284	21	5	313		
III	0	2	108	10	120		
IV	0	0	2	171	173		
Total	270	445	139	193	1047		
c) CADMAN test set, $\kappa = 0.80$							

	Aut	Automated risk category				
Risk	Ι	II	III	IV	Total	
Ι	39	16	4	0	59	
II	0	49	3	2	54	
III	0	1	18	2	21	
IV	0	0	1	21	22	
Total	39	66	26	25	156	

b) or CaScore training set, $\kappa = 0.97$

	Aι	Automated risk category					
Risk	Ι	II	III	IV	Total		
Ι	7	1	0	0	8		
II	0	8	0	0	8		
III	0	0	8	0	8		
IV	0	0	0	8	8		
Total	7	9	8	8	32		

Table 5: Confusion matrices show the agreement in CVD risk for the DISCHARGE test set (a), orCaScore training set (b) and CADMAN test set (c). Categorization is based on the total Agatston score with I: 0, II: [1,100), III: [100,300), IV: > 300.

462 III.D. Active learning evaluation with uncertainty weighted loss 463 MTL-model

We evaluate our uncertainty weighted loss MTL-model in an active learning scenario by conducting two experiments. In the first experiment, we analyzed the model performance trained in an active learning scenario for different loss weighting strategies (uncertainty weighted loss, static weighted loss) and sampling strategies (random sampling, hybrid sampling) described in Subsection II.C.. We initially trained the model with only 100 randomly

selected samples (slices) and double the number of samples in each sampling round. Instead 469 of retraining the model from scratch after each round, we continue training with the larger 470 dataset and a reduced initial learning rate to 1e-04 compared to 5e-04 for the initial training. 471 We useed early stopping based on the validation set to avoid overtraining in each sampling 472 round. The Micro F1-score of the multi-class calcification segmentation is used to compare 473 different models. As shown in Figure 5 with uncertainty weighted loss and hybrid sampling 474 method, the model required only three sampling rounds and 800 annotated slices (12% of)475 the training set) to achieve similar performance (Micro F1-score=0.846) as when trained on 476

the full training set (Micro F1-score=0.849).



Figure 5: Performance comparison between different loss weighting methods (static weighted loss and uncertainty weighted loss) as well as different sampling methods (random sampling and hybrid sampling) in an active learning scenario.

To compare the model performance based on the used loss weighting strategy and sampling method after three sampling rounds, we compare the model performance proportion (compared to uncertainty weighted model on the full dataset) in an active learning scenario in Table 6. It can be seen that the uncertainty weighted loss outperforms static weighting for random and hybrid sampling by 6.4% and 4.9%, respectively. This can be explained by the

- fact that the data distribution of the training set is changing in each sampling round and especially during the first sampling rounds. The uncertainty weighted loss method can com-
- ⁴⁸⁵ pensate for this distribution shift, but static weighted loss cannot. It can also be seen that
- 486 hybrid sampling outperforms random sampling for static weighted and uncertainty weighted
- loss by 5.4% and 4.0%, respectively. The hybrid sampling method selects only the most
- ⁴⁸⁸ informative image slices and can therefore reduce number of required samples.

	Random sampling	Hybrid sampling
Static-weighted loss	$89.28\% \ (0.758/0.849)$	$94.70\% \ (0.804/0.849)$
Uncertainty-weighted loss	$95.64\% \ (0.812/0.849)$	99.65% (0.846/0.849)

Table 6: Model performance proportion of the active learning model after three sampling rounds compared to the performance of the uncertainty weighted loss MTL-model trained on the full dataset.

The labeling of the additional coronary artery region annotations requires extra time, 489 even if the annotation process is an efficient semi-automatic process described in Sec-490 tion II.A.3.. An approximation of the required annotation time for 1) coronary calcifications, 491 2) coronary calcifications and coronary artery regions and 3) only informative slices of coro-492 nary calcifications and coronary artery region was investigated empirically and is shown in 493 Table 7. It shows that annotation of CAC and CAD with active learning reduces the an-494 notation cost to approximately one-third, compared to labeling of calcifications on the full 495 training set, even though labeling CAC and CAD is more time consuming. 496

	Annot. time	Number of	Annot. time	Improvement
	per slice [s]	labeled slices	training set [s]	ratio
CAC	4.0	6721	26884	1.0
CAC + CAR	12.0	6721	80652	3.0
CAC + CAR + AL	12.0	800	9600	0.36

Table 7: Approximated annotation time for annotation of coronary artery calcifications (CAC) compared to annotation of coronary calcifications and coronary artery regions (CAC+CAR) and annotation of coronary calcifications and coronary artery regions using active learning (CAC+CAR+AL).

- ⁴⁹⁷ In a second experiment, we analyzed the impact of the number of training samples on
- the estimated optimal weighting parameter α in Equation (3) using a grid search method.
- ⁴⁹⁹ Therefore, we trained the static weighted loss model multiple times with varying weighting
- parameter α from 0.1 to 0.9 and step size of 0.1 for a very small randomly selected dataset

(only 100 samples) and compared the results when trained on the full dataset. It turns out that the estimated optimal parameter of the full dataset $\alpha = 0.4$ does not match with the optimal weighting parameter of the small dataset $\alpha = 0.2$ because the small dataset does not represent the data distribution of the full dataset. If a non-optimal parameter value would have been selected after the first sampling round of the active learning method, optimal performance would not have been achieved. Alternatively, α could have been redetermined in each sampling round, but this would be computationally very expensive. A detailed analysis

⁵⁰⁸ can be found in Figure S2 in the supplementary material.

III.E. Influence of image noise, metal artifacts, motion artifacts and image quality on the model performance

The performance rises by 4.6% if CT scans with severe image noise, metal artifacts, motion 511 artifacts and image quality are excluded. We can see in Table 3 and Table 4 that the test 512 performance on the orCaScore test set (Micro F1-score = 0.961) is much higher compared to 513 the test performance on the DISCHARGE dataset (Micro F1-score = 0.881). To explain the 514 performance difference, we analyzed the influence of four factors 1) image noise, 2) metal 515 artifacts, 3) motion artifacts, 4) image quality. We estimated the noise using a method 516 similar to Christianson et al.⁵¹ First we segmented the CT image into the heart-related 517 tissue types (-200 to 140 HU), second, a noise image filter⁵² was applied to the segmented 518 region, third, a histogram was generated and the highest peak was selected as noise level.⁵¹ 519 The noise levels of all CT scans of the test dataset were normalized using z-score⁵³ and the 520 most noisy 20% were labeled as noisy CT scans. Metal artifacts and motion artifacts were 521 determined visually and scans were labeled according to their presence or absence. Image 522 quality was visually assessed by a high level of disturbance or anatomical abnormalities and 523 labeled as good or bad quality, accordingly. Note that none of the CT scans in the test 524 set were deemed as nondiagnostic (unsatisfactory for diagnosis) by a radiologist. Examples 525 for the four influencing factors are shown in Figure 6. It shows that the Micro F1-score 526 ranges from 0.881 including all scans in the test set to 0.927 if all noisy scans, scans with 527 metal or motion article and poor quality images are excluded. When only noisy images are 528 excluded, performance increases by 3.1%. Surprisingly, when we exclude images with severe 529 motion artifacts, performance drops only by 0.02%. This can by explained by the fact that 530 in motion artifacts, calcifications appear very large, resulting in a high number of "falsely" 531



Figure 6: Image examples with severe image noise (A), metal artifacts (B), motion artifacts (C), low image quality (abnormality provoked by hiatal hernias) (D).

⁵³² labeled true positives in the data set. Excluding samples with motion artifacts decreases the ⁵³³ number of true positives and thus the corresponding Micro F1-score. A detailed analysis ⁵³⁴ of the influencing factors and its influence on the model performance can be found in the ⁵³⁵ supplementary material.

536 IV. Discussion

In this paper, we have proposed an MTL-model with uncertainty weighted loss for coronary calcium scoring in ECG-gated, non-contrast enhanced cardiac CTs. The model can be trained in an active learning scenario and requires only 12% of the training data and approximately one-third of the annotation time to achieve the same performance as when trained with the

541 <mark>full dataset.</mark>

To the best of our knowledge, our model is the first that performs segmentation of coronary 542 artery regions (CAR) based on weak annotations and segmentation of coronary calcifications 543 (CAC) in an end-to-end framework. We compared our multi-task models with a single-544 task model (multiclass U-Net) and a sequential model. It can be seen in Table 2 that the 545 multi-task models outperform other models in terms of positive predictive value, sensitiv-546 ity, F1-score and Micro F1-score. The benefits of an MTL-model compared to a multiclass 547 U-Net and a sequential model are the shared information between coronary artery region 548 segmentation task T_R and calcification segmentation task T_L . In contrast to the multiclass 549 U-Net, the MTL-model is able to learn important spatial information from weakly labeled 550 samples and is able to transfer this knowledge for segmentation of coronary calcifications. 551 Explanation techniques such as layer wise relevance propagation 5^{4} could lead to a deeper 552 understanding about different prediction strategies but they are beyond the scope of this 553

554 paper.

To investigate the uncertainty weighted loss, we compared the performance with a MTLmodel trained using static weighted loss. When using optimal weighting parameter, the performance is similar, but it is important to note that the determination of the optimal weighting parameter value is a challenging and expensive process and even more difficult to estimate in an active learning scenario.

To reduce labeling costs, we investigated our uncertainty weighted multi-task network in an 560 active learning scenario and could show that our model reaches optimal performance with 561 substantial less training samples. The uncertainty weighted loss MTL-model is able to bal-562 ance losses when the data distribution is changing after each sampling round. We compared 563 different active learning scenarios and could show in Figure 5 that uncertainty weighted loss 564 outperforms static weighted loss in random sampling and hybrid sampling. The biggest dis-565 advantage of static weighting is the estimation of weighting parameter α , which is difficult 566 to obtain and sensitive to the size of the training set shown in Figure S2 (supplementary 567 material). In contrast to T. Gong et al.⁵⁵, we did not notice more instability issues when 568 our uncertainty weighted loss model was trained on small datasets. 569

We compared our uncertainty weighted MTL-model with other methods in Table 3 on the 570 orCaScore dataset and could show that our model performs very good in terms of F1-571 score, ICC and sensitivity using only CSCT. To compare the performance with respect to 572 the dataset, we tested our model on three different datasets. To our surprise, the model 573 trained on the DISCHARGE training set performed better on the orCaScore test set (Mi-574 cro F1-score=0.958) than on the DISCHARGE test set (Micro F1-score=0.881). The test 575 performance on the CADMAN dataset (Micro F1-score=0.822) was even lower than on the 576 DISCHARGE test set due to a higher number of false positive predictions. This can be 577 explained by a higher level of noise in the CADMAN dataset, since it contains only filtered 578 back projections. The influence of noise can also be reflected in the higher number of lesion 579 candidates per scan in Table 1 for the CADMAN and DISCHARGE dataset. 580

The predictions of cardiovascular disease (CVD) risk categories based on the segmentations in Table 5 show a very good agreement of $\kappa = 0.97$ for the orCaScore dataset but a lower agreement of $\kappa = 0.80$ for the DISCHARGE test set. Mislabeled noise leads to a high false positive rate between risk category I (total Agatston score is 0) and II (total Agatston score between 1 and 100) and similar findings have been made in R. Zeleznik et al..⁴⁹

We also trained our model on the orCaScore training set with additional annotations for 586 CAR and reached only slightly lower performance (Micro F1-score=0.928). To gain a better 587 understanding of the different influencing factors (exclusion criteria) related to model per-588 formance, we compared the performance after exclusion of scans due to 1) image noise, 2) 589 metal artifacts, 3) motion artifacts and 4) image quality). If all exclusion criteria were met, 590 the Micro F1-score increased from 0.886 to 0.931. We have also shown that image noise is 591 one of the most influencing factors on model performance beside metal artefacts and image 592 quality. Scans with motion artifact had no effect on performance, which can be explained 593 by visual expansion of the lesion area due to motion, mainly in the proximal RCA, leading 594 to overestimation of the lesion volume in both the labeling phase by the radiologist and the 595 prediction phase by the network. 596

⁵⁹⁷ IV.A. Limitations

We have seen that the convergence of the MTL-model trained with static weighted loss was more sensitive to changes of the learning rate compared to uncertainty weighted loss when trained on a small datasets. Nevertheless, training time requires several hours which makes tuning of the hyper-parameter challenging and limits the possibility to draw general conclusions.

- Our method is processing 2D axial CT slices. The usage of 3D-information might be beneficial as shown in Zhang et al. ⁸ and recently published methods based on 3D-CNN ensembles ⁷ achieved very good results. Since our active learning approach is based on the labeling of only the most informative slices the 3D-annotations would be sparse. Learning dense 3D segmentations from sparse annotations can be challenging in a multi-task network ⁵⁶ therefore we leave a 3D extension of our method for future work.
- Furthermore, reference standards for the DISCHARGE and CADMAN datasets were provided by only one experienced observer for coronary calcifications and coronary artery regions. Independent annotations from a second observer and clarification of discrepancies by consensus would improve the quality of the dataset but since the annotation process requires expert knowledge and is tedious and time-consuming, we leave this improvement of the dataset for a further research project.
- ⁶¹⁵ The analysis of influencing factors is limited to four factors (image noise, metal artifacts, mo-

tion artifacts, image quality), yet other factors such as reconstruction method, scanner type,
slice thickness or slice spacing are known to influence model performance but are beyond
the scope of this work and will be investigated in future work.

⁶¹⁹ IV.B. Further research directions

With respect to our results, we have to critically reflect the question which loss and per-620 formance metrics are best suited for risk prediction of coronary heart disease events. Our 621 model performs well on F1-scores, ICC and sensitivity of CAC volume but lacks precision 622 on CVD risk agreement. In further analysis we will investigate how a direct prediction of the 623 risk categories ¹³ can be integrated into our model to improve risk categorization. A major 624 focus will be on improving the prediction of patients with zero calcium score. We also plan to 625 extend our model from 2D input data to 3D to take advantage of 3D context information and 626 overcome our limitations. We evaluated the uncertainty weighted MTL-model in an active 627 learning scenario using our hybrid sampling method and believe that the model is also ap-628 plicable with other sampling strategies but leave further analysis as future work. Additional 629 future work may investigate how radiologist-in-the-loop frameworks might use explanations 630 to guide a more efficient active learning based labeling process for coronary calcium scoring. 631 A deeper understanding of the model behavior supported by explanations and quantification 632 of model uncertainties would enable the radiologist to understand predictions and assist in 633 medical decision making. 634

635 V. Conclusions

In this work we have proposed a multi-task model with uncertainty weighted loss for coronary 636 calcium scoring. The model improves calcium scoring performance by extracting shared 637 informative features from the two tasks of coronary artery region (CAR) segmentation and 638 coronary artery calcifications (CAC) segmentation. The model performance was evaluated 639 using a large multi-center dataset of the DISCHARGE trial (1047 CSCTs), a single-center 640 dataset of the CAD-Man study (156 CSCTs) and the multi-center or CaScore test set (40 641 CSCTs). When trained in an active learning scenario, the model achieves optimal perfor-642 mance with only 12% of the training samples, reduces annotation time to one-third and 643

enables the integration of the radiologist into the training loop. The good performance
and the reduction of required annotated image slices might enable the training of models
applicable in a clinical setting.

Acknowledgement 647

Acknowledgment should be provided to the DISCHARGE Trial Group (Napp et al.)²² 648 for the collection and provision of the data from the DISCHARGE project (603266-649 2, HEALTH-2012.2.4.-2) funded by the FP7 Program of the European Commission 650 (https://www.dischargetrial.eu): 651 Pál Maurovich-Horvat, M.D., Ph.D., M.P.H., Maria Bosserdt, Ph.D., Klaus F. Kofoed, M.D., 652 D.M.Sc., Nina Rieckmann, Ph.D., Theodora Benedek, M.D., Ph.D., Patrick Donnelly, M.D., 653 José Rodriguez-Palomares, M.D., Ph.D., Andrejs Erglis, M.D., Cvril Štěchovský, M.D., Gin-654

tarė Šakalytė, M.D., Nada Čemerlić Adić, M.D., Matthias Gutberlet, M.D., Ph.D., Jonathan 655

D. Dodd, M.D., Ignacio Diez, M.D., Gershan Davis, M.D., Elke Zimmermann, M.D., Cezary 656

Kępka, M.D., Radosav Vidakovic, M.D., Ph.D., Marco Francone, M.D., Ph.D., Małgorzata 657

Ilnicka-Suckiel, M.D., Ph.D., Fabian Plank, M.D., Ph.D., Juhani Knuuti, M.D., Rita Faria, 658

M.D., Stephen Schröder, M.D., Colin Berry, M.D., Luca Saba, M.D., Balazs Ruzsics, M.D., 659

Ph.D., Christine Kubiak, Ph.D., Iñaki Gutierrez-Ibarluzea, Ph.D., Kristian Schultz Hansen, 660

Ph.D., Jacqueline Müller-Nordhorn, M.D., M.P.H., Bela Merkely, M.D., Ph.D., Andreas 661

Dehlbæk Knudsen, M.D., Imre Benedek, M.D., Ph.D., Clare Orr, M.D., Filipa Xavier Va-662

lente, M.D., Ph.D., Ligita Zvaigzne, M.D., Vojtěch Suchánek, M.D., Antanas Jankauskas, 663

M.D., Filip Adić, M.D., Michael Woinke, M.D., Mark Hensey, M.B., B.Ch., B.A.O., Iñigo 664 Lecumberri, M.D., Erica Thwaite, M.D., Michael Laule, M.D., Mariusz Kruk, M.D., Aleksan-

dar N. Neskovic, M.D., Ph.D., Massimo Mancone, M.D., Donata Kuśmierz, M.D., Gudrun 666

Feuchtner, M.D., Mikko Pietilä, M.D., Ph.D., Vasco Gama Ribeiro M.D., Tanja Drosch, 667

M.D., Christian Delles, M.D., Gildo Matta, M.D., Michael Fisher, M.D., Ph.D., Bálint Szil-668

veszter, M.D., Ph.D., Linnea Larsen, M.D., Ph.D., Mihaela Ratiu, M.D., Ph.D., Stephanie 669

Kelly, Bruno Garcia del Blanco, M.D., Ph.D., Ainhoa Rubio, M.D., Zsófia D. Drobni, M.D., 670

Birgit Jurlander, M.D., Ph.D., Ioana Rodean, M.D., Susan Regan, Hug Cuéllar Calabria, 671

M.D., Ph.D., Melinda Boussoussou, M.D., Thomas Engstrøm, M.D., D.M.Sc., Roxana Ho-672

das, M.D., Adriane E. Napp, Ph.D., Robert Haase, M.D., Georg M. Schuetz, M.D., Sarah 673

Feger, M.D., Konrad Neumann, Ph.D., Henryk Dreger, M.D., Matthias Rief, M.D., Viktoria 674

Wieske, M.D., Melanie Estrella, Ph.D., Peter Martus, Ph.D., Marc Dewey, M.D. 675

676

665

Affiliations: 677

MTA-SE Cardiovascular Imaging Research Group (P.M.-H., B.S., Z.D.D.) at Heart and Vas-678 cular Center (B.M., M.B.) and Department of Radiology, Medical Imaging Center (P.M.-H.), 679 Semmelweis University, Budapest, Hungary; Department of Radiology (M.B., E.Z., A.E.N., 680 R. Haase, G.M.S., S.F., M. Rief, V.W., M.E., M.D.) and Cardiology (M.L., H.D.) and Insti-681 tute of Public Health (N.R., J.M.-N) and Institute of Biometry and Clinical Epidemiology 682 (K.N.), Charité - Universitätsmedizin Berlin, corporate member of Freie Universität Berlin 683 and Humboldt-Universität zu Berlin and DZHK (German Centre for Cardiovascular Re-684 search), partner site Berlin (H.D., M.D.) and Berlin Institute of Health (M.D.) and Berlin 685 University Alliance (M.D.), Berlin, Germany; Department of Cardiology (K.F.K., A.D.K., 686 T.E.) and Radiology (K.F.K., A.D.K.), Copenhagen University Hospital - Rigshospitalet & 687 Department of Clinical Medicine, Faculty of Health and Medical Sciences and Department 688 of Cardiology, Herlev-Gentofte Hospital (L.L.) and Department of Cardiology, Nordsjael-689 lands Hospital, (B.J.), University of Copenhagen, Copenhagen, Denmark; Department of 690 Internal Medicine, Clinic of Cardiology, (T.B., R. Hodas) and Department of Radiology and 691 Medical Imaging (M. Ratiu), University of Medicine and Pharmacy, Science and Technology 692 "G.E.Palade" and County Clinical Emergency Hospital Tirgu Mures (T.B.) and Center of 693 Advanced Research in Multimodality Cardiac Imaging, Cardio Med Medical Center (I.B., 694 I.R.), Tirgu Mures, Romania; Department of Cardiology, Southeastern Health and Social 695 Care Trust, (P.D., C.O., S.K., S.R.) Belfast, United Kingdom; Department of Cardiology 696 (J.R.-P, F.X.V., B.G.B.) and Radiology (H.C.C.), Hospital Universitari Vall d'Hebron, Insti-697 tut de Recerca (VHIR), Universitat Autònoma de Barcelona, Barcelona, Spain; Department 698 of Cardiology (A.E.) and Radiology (L. Z.), Paul Stradins Clinical University Hospital and 699 University of Latvia (A.E.), Riga, Latvia; Department of Cardiology (C.Š.) and Department 700 of Imaging Methods (V.S.), Motol University Hospital, Prague, Czech Republic; Department 701 of Cardiology (G.Š.) and Department of Radiology (A.J.), Medical Academy, Lithuanian 702 University of Health Sciences and Department of Cardiology, Hospital of Lithuanian Uni-703 versity of Health Sciences (G.S.), Kaunas, Lithuania; Faculty of Medicine, University of Novi 704 Sad, (N.Č.A., F.A.) and Department of Cardiology, Institute for Cardiovascular Diseases of 705 Vojvodina (N.Č.A., F.A.), Novi Sad, Sremska Kamenica, Serbia; Department of Radiology 706 (M.G.) and Department of Cardiology (M.W.), University of Leipzig Heart Centre, Leipzig, 707 Germany; Department of Radiology (J.D.D.) and Department of Cardiology (M.H.), St. 708 Vincent's University Hospital and School of Medicine (J.D.D.), University College Dublin, 709

Dublin, Ireland; Department of Cardiology (I.D., A.R.) and Department of Radiology (I.L.), 710 Basurto Hospital, Bilbao Spain; Department of Cardiology (G.D.) and Department of Ra-711 diology (E.T.), Aintree University Hospital, Liverpool and School of Medicine, University 712 of Central Lancashire (G.D.), Preston, United Kingdom; National Institute of Cardiology 713 (C. Kepka, M.K.), Warsaw, Poland; Department of Cardiology, Internal Medicine Clinic, 714 Clinical Hospital Center Zemun (R.V., A.N.N.) and Faculty of Medicine, University of Bel-715 grade (R.V., A.N.N.), Belgrade, Serbia; Department of Radiological, Pathological and On-716 cological Sciences (M. Francone) and Department of Clinical Internal, Anesthesiologic and 717 Cardiovascular Sciences (M.M.), Sapienza University of Rome, Rome, Italy; Department 718 of Cardiology (M.I.S.) and Department of Radiology (D.K.), Wojewodzki Szpital Specjal-719 istyczny We Wroclawiu, Wroclaw, Poland; Department of Internal Medicine III, Cardiology 720 (F.P.) and Department of Radiology (G.F.), Innsbruck Medical University, Innsbruck, Aus-721 tria; Turku PET Centre (J.K.) and Heart Center (M.P.), Turku University Hospital and 722 University of Turku and Administrative Centre, Health Care District of Southwestern Fin-723 land (M.P.), Turku, Finland; Department of Cardiology, Centro Hospitalar de Vila Nova 724 de Gaia/ Espinho (R.F., V.G.R.), Vila Nova de Gaia, Portugal; Department of Cardiology, 725 ALB FILS KLINIKEN GmbH (S.S., T.D.), Göppingen, Germany; Institute of Cardiovas-726 cular & Medical Sciences, University of Glasgow (C.B., C.D.), Glasgow and Golden Jubilee 727 National Hospital (C.B.), Clydebank, United Kingdom; Department of Radiology, Univer-728 sity of Cagliari (L.S.) and Department of Radiology, Azienda Ospedaliera Brotzu (G.M.), 729 Cagliari, CA, Italy; Department of Cardiology, Royal Liverpool University Hospital (B.R., 730 M.F.) and Institute for Cardiovascular Medicine and Science, Liverpool Heart and Chest 731 Hospital (B.R., M. Fisher), and Faculty of Health and Life Sciences, University of Liver-732 pool (M. Fisher), Liverpool, United Kingdom; ECRIN-ERIC (European Clinical Research 733 Infrastructure Network-European Research Infrastructure Consortium) (C. Kubiak), Paris, 734 France; Basque Foundation for Health Innovation and Research, Barakaldo / Bizkaia and 735 Basque Office for Health Technology Assessment (I.G.-I.), Vitoria - Gasteiz, Spain; Uni-736 versity of Copenhagen, Department of Public Health, Section for Health Services Research 737 (K.S.H.), Copenhagen, Denmark; Bavarian Cancer Registry, Bavarian Health and Food 738 Safety Authority (J.M.-N.), Munich, Germany; Department of Clinical Epidemiology and 739 Applied Biostatistics, Universitätsklinikum Tübingen, Tübingen (P.M.), Germany 740 Acknowledgment should also be provided for the CAD-MAN Trial Group (Dewey et al.)²⁴ 741

- ⁷⁴² for the collection and provision of the data from the CAD-MAN project (ClinicalTrials.gov
- ⁷⁴³ Identifier: NCT00844220) funded by a grant of the Heisenberg programme of the German
- ⁷⁴⁴ Research Foundation to Marc Dewey.
- The authors thank the organizers of the orCaScore Challenge for launching this internationaland open competition.
- 747 This work was funded by the German Research Foundation through the graduate pro-
- ⁷⁴⁸ gram BIOQIC (GRK2260, project-ID: 289347353), the priority program SPP-Radiomics
- ⁷⁴⁹ (SPP2177, project-ID: 402688427) and the DISCHARGE project (603266-2, HEALTH-
- ⁷⁵⁰ 2012.2.4.-2) funded by the FP7 Program of the European Commission.

751 Conflict of Interest

The author Marc Dewey declares relationships with the following companies: Prof. Dewey 752 has received grant support from the FP7 Program of the European Commission for the 753 randomized multicenter DISCHARGE trial (603266-2, HEALTH-2012.2.4.-2). He also re-754 ceived grant support from German Research Foundation (DFG) in the Heisenberg Program 755 (DE 1361/14-1), graduate program on quantitative biomedical imaging (BIOQIC, GRK 756 2260/1), for fractal analysis of myocardial perfusion (DE 1361/18-1), the Priority Pro-757 gramme Radiomics for the investigation of coronary plaque and coronary flow (DE 1361/19-1 758 [428222922] and 20-1 [428223139] in SPP 2177/1). He also received funding from the Berlin 759 University Alliance (GC SC PC 27) and from the Digital Health Accelerator of the Berlin 760 Institute of Health. Prof. Dewey is European Society of Radiology (ESR) Research Chair 761 (2019–2022) and the opinions expressed in this article are the author's own and do not rep-762 resent the view of ESR. Per the guiding principles of ESR, the work as Research Chair is 763 on a voluntary basis and only remuneration of travel expenses occurs. Prof. Dewey is also 764 the editor of Cardiac CT, published by Springer Nature, and offers hands-on courses on 765 CT imaging (www.ct-kurs.de). Institutional master research agreements exist with Siemens, 766 General Electric, Philips, and Canon. The terms of these arrangements are managed by 767 the legal department of Charité - Universitätsmedizin Berlin. Professor Dewey holds a 768 joint patent with Florian Michallek on dynamic perfusion analysis using fractal analysis 769 (PCT/EP2016/071551 and USPTO 2021 10,991,109). 770

771 Other authors declared no conflicts of interest.

772

⁷⁷³ ^{a)} Author to whom correspondence should be addressed:

- 774 Bernhard Föllmer
- 775 Charité Universitätsmedizin Berlin
- 776 Klinik für Radiologie
- 777 Campus Charité Mitte (CCM)
- ⁷⁷⁸ Charitéplatz 1, 10117 Berlin
- ⁷⁷⁹ E-Mail: bernhard.foellmer@charite.de
- 780

⁷⁸¹ ^{b)} These authors should be considered as joint senior authors.

783 782 784	\mathbf{R}_{1}	eferences WHO, Cardiovascular diseases, fact sheet 317. World Health Organization. (2021).
785 786	2	M. Dewey, Cardiac CT in Clinical Practice, in <i>Cardiac CT</i> , <i>Second Edition</i> , pages 33–42, Springer Berlin Heidelberg, 2014.
787 788 789 790 791	3	J. M. Wolterink, T. Leiner, B. D. De Vos, J. L. Coatrieux, B. M. Kelm, S. Kondo, R. A. Salgado, R. Shahzad, H. Shu, M. Snoeren, R. A. Takx, L. J. Van Vliet, T. Van Walsum, T. P. Willems, G. Yang, Y. Zheng, M. A. Viergever, and I. Išgum, An evaluation of automatic coronary artery calcium scoring methods with cardiac CT using the orCaScore framework, Medical Physics 43 , 2361–2373 (2016).
792 793 794	4	G. Santini, D. Della Latta, N. Martini, G. Valvano, A. Gori, A. Ripoli, C. L. Susini, L. Landini, and D. Chiappino, An automatic deep learning approach for coronary artery calcium segmentation, IFMBE Proceedings 65 , 374–377 (2017).
795 796 797	5	A. S. Agatston, F. W. R. Janowitz, F. J. Hildner, N. R. Zusmer, M. Viamonte, and R. Detrano, Quantification Coronary Artery Calcium Using Ultrafast Cumputed To- mography, Journal of the American College of Cardiology 15 , 827–832 (1990).
798 799 800	6	N. Lessmann, B. V. Ginneken, M. Zreik, P. A. D. Jong, B. D. D. Vos, M. A. Viergever, and I. Išgum, Automatic calcium scoring in low-dose chest CT using deep neural net- works with dilated convolutions, 37 , 615–625 (2018).
801 802 803 804 805	7	 N. Gogin, M. Viti, L. Nicodème, M. Ohana, H. Talbot, U. Gencer, M. Mekukosokeng, T. Caramella, Y. Diascorn, J. Y. Airaud, M. S. Guillot, Z. Bensalah, C. Dam Hieu, B. Abdallah, I. Bousaid, N. Lassau, and E. Mousseaux, Automatic coronary artery calcium scoring from unenhanced-ECG-gated CT using deep learning, Diagnostic and Interventional Imaging 102, 683–690 (2021).
806 807 808	8	W. Zhang, J. Zhang, X. Du, Y. Zhang, and S. Li, An end-to-end joint learning framework of artery-specific coronary calcium scoring in non-contrast cardiac CT, Computing 101 , 667–678 (2019).
809	9	J. M. Wolterink, T. Leiner, B. D. de Vos, R. W. van Hamersvelt, M. A. Viergever, and

I. Išgum, Automatic coronary artery calcium scoring in cardiac CT angiography using paired convolutional neural networks, Medical Image Analysis **34**, 123–136 (2016).

- ¹⁰ G. Yang, Y. Chen, X. Ning, Q. Sun, H. Shu, and J.-L. Coatrieux, Automatic coronary
 ⁸¹³ calcium scoring using noncontrast and contrast CT images., Medical physics 43, 2174
 ⁸¹⁴ (2016).
- ¹¹ S. G. M. V. Velzen, N. Hampe, and B. D. D. Vos, AI for Calcium Scoring, pages 1–22.
- ¹² R. Shahzad, L. van Vliet, W. Niessen, and T. V. Walsum, Automatic Classification of
 ⁸¹⁷ Calcification in the Coronary Vessel Tree, Orcascore.Isi.Uu.Nl.
- ¹³ B. D. de Vos, J. M. Wolterink, T. Leiner, P. A. de Jong, N. Lessmann, and I. Isgum, Direct Automatic Coronary Calcium Scoring in Cardiac and Chest CT., IEEE transactions
 on medical imaging 38, 2127–2138 (2019).
- ¹⁴ A. Smailagic, H. Y. Noh, P. Costa, D. Walawalkar, K. Khandelwal, M. Mirshekari,
 J. Fagert, A. Galdran, and S. Xu, MedAL : Accurate and Robust Deep Active Learning
 ⁸²³ for Medical Image Analysis, (2018).
- ¹⁵ S. Ruder, An Overview of Multi-Task Learning in Deep Neural Networks * arXiv : 1706
 . 05098v1 [cs . LG] 15 Jun 2017, (2017).
- ¹⁶ Z. Chen, V. Badrinarayanan, C. Y. Lee, and A. Rabinovich, GradNorm: Gradient
 ⁸²⁷ normalization for adaptive loss balancing in deep multitask networks, 35th International
 ⁸²⁸ Conference on Machine Learning, ICML 2018 2, 1240–1251 (2018).
- ¹⁷ S. Liu, E. Johns, and A. J. Davison, End-to-end multi-task learning with attention,
 ⁸³⁰ Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern
 ⁸³¹ Recognition 2019-June, 1871–1880 (2019).
- ¹⁸ M. Guo, A. Haque, D.-A. Huang, S. Yeung, and L. Fei-Fei, *Dynamic Task Prioritization* for Multitask Learning: 15th European Conference, Munich, Germany, September 8-14,
 ²⁰¹⁸, Proceedings, Part XVI, pages 282–299, 2018.
- ¹⁹ R. Cipolla, Y. Gal, and A. Kendall, Multi-task Learning Using Uncertainty to Weigh
 ⁸³⁶ Losses for Scene Geometry and Semantics, Proceedings of the IEEE Computer Society
 ⁸³⁷ Conference on Computer Vision and Pattern Recognition, 7482–7491 (2018).

838 ²⁰ 839 840	 S. Vandenhende, S. Georgoulis, W. Van Gansbeke, M. Proesmans, D. Dai, and L. Van Gool, Multi-Task Learning for Dense Prediction Tasks: A Survey, IEEE Transactions on Pattern Analysis and Machine Intelligence, 1–20 (2021).
841 21	M. Dewey, The Discharge trial, https://www.dischargetrial.eu/, 2021.
842 843 844 845	 A. E. Napp, R. Haase, M. Laule, G. M. Schuetz, M. Rief, H. Dreger, G. Feuchtner, G. Friedrich, Š. Miloslav, and M. Dewey, Computed tomography versus invasive coro- nary angiography : design and methods of the pragmatic randomised multicentre DIS- CHARGE trial, pages 2957–2968 (2017).
846 847	I. I. Jelmer M. Wolterink, Bob D. de Vos, Tim Leiner, Max A. Viergever, orCaScore, https://orcascore.grand-challenge.org/, 2021.
848 24 849 850 851 852	 M. Dewey, M. Rief, P. Martus, B. Kendziora, S. Feger, H. Dreger, S. Priem, F. Knebel, M. Böhm, P. Schlattmann, B. Hamm, E. Schönenberger, M. Laule, and E. Zimmermann, Evaluation of computed tomography in patients with atypical angina or chest pain clinically referred for invasive coronary angiography: Randomised controlled trial, BMJ (Online) 355 (2016).
853 25	F. A. Jolesz, Intraoperative Imaging and Image- Guided Therapy 123.
855 856	O. Ronneberger, P.Fischer, and T. Brox, U-Net: Convolutional Networks for Biomedical Image Segmentation, in <i>Medical Image Computing and Computer-Assisted Intervention</i> (<i>MICCAI</i>), volume 9351 of <i>LNCS</i> , pages 234–241, Springer, 2015.
857 27 858 859 860 861	R. Ke, A. Bugeau, N. Papadakis, P. Schuetz, and C. B. Schönlieb, Learning to Segment Microscopy Images with Lazy Labels, in <i>Lecture Notes in Computer Science (including</i> <i>subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)</i> , volume 12535 LNCS, pages 411–428, Springer Science and Business Media Deutschland GmbH, 2020.
862 28 863	G. Papandreou, Lc. Chen, K. Murphy, and A. L. Yuille, Weakly- and Semi-Supervised Learning of a Deep Convolutional Network for Semantic Image Segmentation.
864 29 865 866	T. Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollar, Focal Loss for Dense Object Detection, IEEE Transactions on Pattern Analysis and Machine Intelligence 42 , 318–327 (2020).

- ³⁰ S. Kongyoung, C. Macdonald, and I. Ounis, Multi-Task Learning using Dynamic Task
 Weighting for Conversational Question Answering, pages 17–26 (2020).
- ³¹ S. Ioffe and C. Szegedy, Batch normalization: Accelerating deep network training by
 reducing internal covariate shift, 32nd International Conference on Machine Learning,
 ICML 2015 1, 448–456 (2015).
- ³² M. J. Brown, L. A. Hutchinson, M. J. Rainbow, K. J. Deluzio, and A. R. De Asha,
 ⁸⁷³ Rectified Linear Units Improve Restricted Boltzmann Machines, Journal of Applied
 ⁸⁷⁴ Biomechanics 33, 384–387 (2017).
- ³³ D. P. Kingma and J. L. Ba, Adam: A method for stochastic optimization, 3rd International Conference on Learning Representations, ICLR 2015 - Conference Track Proceedings, 1–15 (2015).
- ³⁴ S. Budd, E. C. Robinson, and B. Kainz, A survey on active learning and human-inthe-loop deep learning for medical image analysis, Medical Image Analysis 71, 102062
 (2021).
- ³⁵ P. Ren, Y. Xiao, X. Chen, X. Wang, X. Chang, P.-Y. Huang, and Z. Li, A Survey of
 Deep Active Learning, Technical report, 2020.
- ³⁶ M. Fang and D. Tao, Active multi-task learning via bandits, SIAM International Conference on Data Mining 2015, SDM 2015, 505–513 (2015).
- ³⁷ A. Harpale and Y. Yang, Active learning for multi-task adaptive filtering, ICML 2010 ⁸⁸⁶ Proceedings, 27th International Conference on Machine Learning, 431–438 (2010).
- ³⁸ Y. Zhang, Multi-task active learning with output constraints, Proceedings of the Na tional Conference on Artificial Intelligence 1, 667–672 (2010).
- ³⁹ Y. Gal and Z. Ghahramani, Dropout as a Bayesian Approximation: Representing Model
 ⁸⁹⁰ Uncertainty in Deep Learning, in *Proceedings of the 33rd International Conference on International Conference on Machine Learning Volume 48*, ICML'16, pages 1050–1059,
 ⁸⁹² JMLR.org, 2016.
- ⁴⁰ T. Nair, D. Precup, D. L. Arnold, and T. Arbel, Exploring uncertainty measures in deep networks for multiple sclerosis lesion detection and segmentation, in *Lecture Notes*

- in Computer Science (including subseries Lecture Notes in Artificial Intelligence and
 Lecture Notes in Bioinformatics), volume 11070 LNCS, pages 655–663, Springer Verlag,
 2018.
- ⁴¹ M. Grandini, E. Bagli, and G. Visani, Metrics for Multi-Class Classification: an Overview, pages 1–17 (2020).
- ⁴² R. Shahzad, T. V. Walsum, M. Schaap, A. Rossi, S. Klein, A. C. Weustink, P. J. D.
 Feyter, L. J. V. Vliet, and W. J. Niessen, Vessel Specific Coronary Artery Calcium
 Scoring: An Automatic System, Academic Radiology 20, 1–9 (2009).
- J. M. Wolterink, T. Leiner, R. A. P. Takx, M. A. Viergever, and I. Išgum, Automatic
 Coronary Calcium Scoring in Non-Contrast-Enhanced ECG-Triggered Cardiac CT With
 Ambiguity Detection, Tmi 34, 1867–1878 (2015).
- ⁴⁴ F. Durlak, M. Wels, C. Schwemmer, M. Sühling, S. Steidl, and A. Maier, Growing a Ran⁹⁰⁷ dom Forest with Fuzzy Spatial Features for Fully Automatic Artery-Specific Coronary
 ⁹⁰⁸ Calcium Scoring, in *Machine Learning in Medical Imaging*, edited by Q. Wang, Y. Shi,
 ⁹⁰⁹ H.-I. Suk, and K. Suzuki, pages 27–35, Cham, 2017, Springer International Publishing.
- ⁴⁵ D. Eng et al., Automated coronary calcium scoring using deep learning with multicenter
 external validation, npj Digital Medicine 4 (2021).
- ⁴⁶ U. Kurkure, D. R. Chittajallu, G. Brunner, Y. H. Le, and I. A. Kakadiaris, A supervised classification-based method for coronary calcium detection in non-contrast CT,
 ⁹¹⁴ International Journal of Cardiovascular Imaging 26, 817–828 (2010).
- ⁴⁷ I. Išgum, A. Rutten, M. Prokop, and B. van Ginneken, Detection of coronary calcifications from computed tomography scans for automated risk assessment of coronary artery disease, Medical Physics 34, 1450–1461 (2007).
- ⁴⁸ G. Brunner, D. R. Chittajallu, U. Kurkure, and I. A. Kakadiaris, Toward the automatic
 ⁹¹⁹ detection of coronary artery calcification in non-contrast computed tomography data.,
 ⁹²⁰ The international journal of cardiovascular imaging 26, 829–838 (2010).
- ⁴⁹ R. Zeleznik et al., Deep convolutional neural networks to predict cardiovascular risk
 from computed tomography, Nature Communications (2021).

- ⁵⁰ S. G. M. V. Velzen, N. Lessmann, and B. K. Velthuis, Deep Learning for Automatic
 ⁹²⁴ Calcium Scoring in CT : Validation Using Multiple Cardiac CT and Chest CT Protocols,
 ⁹²⁵ Radiology (2020).
- ⁵¹ O. Christianson, J. Winslow, D. P. Frush, and E. Samei, Automated Technique to
 Measure Noise in Clinical CT Examinations, American Journal of Roentgenology 205,
 W93–W99 (2015).
- ⁵² SimpleITK, SimpleITK, https://simpleitk.org, 2021.
- ⁵³ D. Sylvan, Introduction to Mathematical Statistics, volume 23, 2013.
- ⁵⁴ S. Bach, A. Binder, G. Montavon, F. Klauschen, K.-R. Müller, and W. Samek, On
 ⁹³² Pixel-Wise Explanations for Non-Linear Classifier Decisions by Layer-Wise Relevance
 ⁹³³ Propagation, PLOS ONE **10**, e0130140 (2015).
- ⁵⁵ T. Gong, T. Lee, C. Stephenson, V. Renduchintala, S. Padhy, A. Ndirango, G. Keskin,
 ⁹³⁵ and O. H. Elibol, A Comparison of Loss Weighting Strategies for Multi task Learning
 ⁹³⁶ in Deep Neural Networks, IEEE Access 7, 141627–141632 (2019).
- J.-M. Bokhorst, H. Pinckaers, P. Van Zwam, I. Nagtegaal, J. Van Der Laak, F. Ciompi,
 and F. C. Nl, Learning from sparsely annotated data for semantic segmentation in
 histopathology images, Proceedings of Machine Learning Research 102, 84–91 (2019).

Active Multi-Task Learning with Uncertainty Weighted Loss for 940 **Coronary Calcium Scoring**

Supplemental Materials

Dataset selection process Ι. 943

Figure S1 shows the selection process for the DISCHARGE and CADMAN datasets. 944

For the DISCHARGE dataset, we considered all 3883 patients (Np = 3883) of the DIS-945 CHARGE trial as potential eligible patients. First, we excluded all patients who were not 946 part of the study CT cohort. Second, at the date of 2020-09-01, all patient without any non-947 contrast enhanced cardiac CTs (CSCT) of 3.00 mm were excluded. Third, for each patients 948 with multiple CSCT reconstructions, we randomly selected one of the reconstructions and 949 excluded the rest. Fourth, we excluded all CSCTs with a slice spacing $\neq 3.0$ mm. 950

For the CADMAN dataset, we considered all 340 patients of the CADMAN trial as potential 951

eligible patients. First, we excluded all patients who are not part of the study CT cohort. 952

Second, we excluded all patients without CSCT. Third, we excluded all CSCT scans with a 953 slice thickness $\neq 3.0$ mm or slice spacing $\neq 3.0$ mm.



(a) Flowchart of the dataset selection process for the DISCHARGE dataset

(b) Flowchart of the dataset selection process for the CADMAN dataset

Figure S1: Flowchart of the dataset selection process for the DISCHARGE dataset and CADMAN dataset.

941

942

954

⁹⁵⁵ II. Dependency between model performance and loss ⁹⁵⁶ weighting parameter

The estimation of the optimal weighting parameter is challanging when trained in a active learning scenario. In Figure S2 we see the dependency between model performance (Micro F1-score) and static weighting parameter value α , trained on a small, randomly selected training set (100 samples) and the full training set (6721 samples). The model was train with varying weighting parameter α from 0.1 to 0.9 and step size of 0.1. It shows, that the optimal parameter value $\alpha = 0.2$ of the small training does not match with the optimal parameter value $\alpha = 0.4$ of the full training set.



Figure S2: Dependency between model performance and loss weighting parameter value α for a very small training set with 100 samples and the full training with 6721 samples.

⁹⁶⁴ III. Comparison of model predictions between the se ⁹⁶⁵ quential model and the uncertainty weighted loss ⁹⁶⁶ MTL-model

To analyze different prediction strategies of the sequential model and the uncertainty weighted loss MTL-model, we compared the predition results of an image with severe image noise. The uncertainty weighted MTL-model performs both task simultaneously and extracts joint informative features which helps to avoid false positive predictions of calcifications in noisy image slices. Figure S3 shows that the sequential model has problems to distinguish noise from micro calcifications and predicts more false positive than the uncertainty weighted MTL-model.

Image

Uncertainty weighted loss MTL-model



Sequential model



Figure S3: Comparison of model predictions for the uncertainty weighted MTL-model and sequential model. The first row shows the image slice (A) with a high level of noise, predicted coronary artery regions (CAR) (B) and coronary calcifications (CAC) (C) of the uncertainty weighted MTL-model. The second row shows the predicted coronary artery regions (CAR) (D) and coronary calcifications (CAC) (E) of the sequential model. In contrast to the sequential mode, joint learning of the two tasks and extraction of shared features, supports the ability to predicts fewer false positive calcifications in noisy images.

⁹⁷⁴ IV. Dependency between model performance and im ⁹⁷⁵ age noise, metal artifacts, motion artifacts and im ⁹⁷⁶ age quality

The model performance of the test set is strongly dependent on the selection criteria of
the included CT scans. To analyze the influence of the four factors 1) image noise, 2)
metal artifacts, 3) motion artifacts, 4) image quality, we performed a subset analysis of the
DISCHARGE test set. Each row in Figure 1 corresponds to a subset of the DISCHARGE
test set. The columns Noisy scan, Metal artifact, Motion artifact and Image quality are
exclusion criteria according to which the scans are included (✓) or excluded (✗).

# Scans	Noisy	Metal	Motion	Image	F1-score Sen.	PPV	ĸ	
	scan	artifacts	artifacts	quality		Den.	11 V	n
1047	1	\checkmark	\checkmark	1	0.881	0.842	0.924	0.800
924	X	\checkmark	\checkmark	1	0.914	0.887	0.942	0.834
1030	1	X	\checkmark	1	0.886	0.851	0.923	0.804
910	X	X	\checkmark	1	0.920	0.899	0.941	0.839
1025	1	\checkmark	X	1	0.878	0.833	0.929	0.795
907	X	\checkmark	X	1	0.914	0.881	0.950	0.828
1009	1	X	X	\checkmark	0.884	0.843	0.929	0.798
893	X	X	X	1	0.922	0.895	0.950	0.833
1028	1	\checkmark	\checkmark	X	0.886	0.851	0.924	0.809
910	X	\checkmark	\checkmark	X	0.919	0.897	0.942	0.839
1012	1	X	\checkmark	X	0.891	0.861	0.924	0.813
897	X	X	\checkmark	X	0.925	0.909	0.941	0.845
1006	1	1	X	X	0.884	0.843	0.930	0.804
893	X	\checkmark	X	X	0.920	0.892	0.951	0.833
991	1	X	X	X	0.890	0.854	0.929	0.807
880	X	X	X	X	0.927	0.905	0.950	0.839

Table 1: Dependency between model performance and exclusion criteria (image noise, metal artifacts, motion artifacts, image quality) on volume level for the DISCHARGE test dataset (\checkmark - included, \varkappa - excluded).

982

⁹⁸³ V. Deviation of the uncertainty weighted total loss

We derive our total loss function based on the uncertainty weighted loss method of Cipolla et al.¹⁹. The outputs of the coronary calcification segmentation decoder and coronary artery region segmentation decoder are defined as $f_{c_L}^{\mathbf{W}}(\mathbf{x})$ and $f_{c_R}^{\mathbf{W}}(\mathbf{x})$, respectively. In the derivation of the total loss in Cipolla et al., a scaled version of the output $f^{\mathbf{W}}(\mathbf{x})$ is squashed through a softmax function.

$$p\left(\mathbf{y} \mid \mathbf{f}^{\mathbf{W}}(\mathbf{x})\right) = \operatorname{Softmax}\left(\frac{1}{\sigma^2}\mathbf{f}^{\mathbf{W}}(\mathbf{x})\right) = \frac{\exp(\frac{1}{\sigma^2}f_c^{\mathbf{W}}(\mathbf{x}))}{\sum_{c'}\exp(\frac{1}{\sigma^2}f_{c'}^{\mathbf{W}}(\mathbf{x}))}$$
(E1)

⁹⁹⁰ The log likelihood for class c can be writte as

$$\log p\left(\mathbf{y}=c \mid \mathbf{f}^{\mathbf{W}}(\mathbf{x}), \sigma\right) = \frac{1}{\sigma^2} f_c^{\mathbf{W}}(\mathbf{x}) - \log \sum_{c'} \exp\left(\frac{1}{\sigma^2} f_{c'}^{\mathbf{W}}(\mathbf{x})\right)$$
(E2)

with $f_c^{\mathbf{W}}(\mathbf{x})$ the *c*'th element of the vector $f^{\mathbf{W}}(\mathbf{x})$.

We define the likelihood to factorise over the two outputs with the scaling parameters σ_R and σ_L , hence

$$\mathcal{L}_{total}(\mathbf{W}, \sigma_R, \sigma_L) = -\log p\left(\mathbf{y}^{\mathbf{R}} = c_R, \mathbf{y}^{\mathbf{L}} = c_L \mid \mathbf{f}^{\mathbf{W}}(\mathbf{x})\right)$$

$$= \log(\operatorname{Softmax}(y_R = c_R \mid \mathbf{f}^{\mathbf{W}}, \sigma_R) \cdot \operatorname{Softmax}(y_L = c_L \mid \mathbf{f}^{\mathbf{W}}, \sigma_L))$$

$$= -\log \frac{\exp(\frac{1}{\sigma_R^2} f_{c_R}^{\mathbf{W}}(\mathbf{x}))}{\sum_{c'_R} \exp(\frac{1}{\sigma_R^2} f_{c'_R}^{\mathbf{W}}(\mathbf{x}))} - \log \frac{\exp(\frac{1}{\sigma_L^2} f_{c_L}^{\mathbf{W}}(\mathbf{x}))}{\sum_{c'_L} \exp(\frac{1}{\sigma_L^2} f_{c'_L}^{\mathbf{W}}(\mathbf{x}))}$$
(E3)

$$\frac{1}{\sigma_2} \sum_{c'} \exp\left(\frac{1}{\sigma_2^2} f_{c'}^{\mathbf{W}}(\mathbf{x})\right) \approx \left(\sum_{c'} \exp\left(f_{c'}^{\mathbf{W}}(\mathbf{x})\right)\right)^{\overline{\sigma_2^2}}$$
(E4)

996

995

989

991

⁹⁹⁷ With the simplification 19 in equation (E4) we get:

$$\begin{aligned} \mathcal{L}_{total}(\mathbf{W}, \sigma_R, \sigma_L) &= -\log \frac{\exp(\frac{1}{\sigma_R^2} f_{c_R}^{\mathbf{W}}(\mathbf{x}))}{\sigma_R(\sum_{c'_R} \exp(f_{c'_R}^{\mathbf{W}}(\mathbf{x})))^{\frac{1}{\sigma_R^2}}} - \log \frac{\exp(\frac{1}{\sigma_L^2} f_{c_L}^{\mathbf{W}}(\mathbf{x}))}{\sigma_L(\sum_{c'_L} \exp(f_{c'_L}^{\mathbf{W}}(\mathbf{x})))^{\frac{1}{\sigma_L^2}}} \\ &= -\frac{1}{\sigma_R^2} \log \frac{\exp(f_{c_R}^{\mathbf{W}}(\mathbf{x}))}{\sum_{c'_R} \exp(f_{c'_R}^{\mathbf{W}}(\mathbf{x}))} + \log(\sigma_R) \\ &- \frac{1}{\sigma_L^2} \log \frac{\exp(f_{c_L}^{\mathbf{W}}(\mathbf{x}))}{\sum_{c'_L} \exp(f_{c'_L}^{\mathbf{W}}(\mathbf{x}))} + \log(\sigma_L) \end{aligned}$$
(E5)

998

Analogous, we propose our total loss based on the focal losses defined in (1) and (2) claimed on the basis of our experiments.

$$\mathcal{L}_{total}(\mathbf{W}, \sigma_R, \sigma_L) = -\frac{1}{\sigma_R^2} \left(\log \frac{\exp(f_{c_R}^{\mathbf{W}}(\mathbf{x}))}{\sum_{c'_R} \exp(f_{c'_R}^{\mathbf{W}}(\mathbf{x}))} w_{c_R} \left(1 - \frac{\exp(f_{c_R}^{\mathbf{W}}(\mathbf{x}))}{\sum_{c'_R} \exp(f_{c'_R}^{\mathbf{W}}(\mathbf{x}))}\right)^{\gamma_R} \right) + \log(\sigma_R)$$

$$-\frac{1}{\sigma_L^2} \left(\log \frac{\exp(f_{c_L}^{\mathbf{W}}(\mathbf{x}))}{\sum_{c'_L} \exp(f_{c'_L}^{\mathbf{W}}(\mathbf{x}))} w_{c_L} \left(1 - \frac{\exp(f_{c_L}^{\mathbf{W}}(\mathbf{x}))}{\sum_{c'_L} \exp(f_{c'_L}^{\mathbf{W}}(\mathbf{x}))}\right)^{\gamma_L} \right) + \log(\sigma_L)$$

$$= \frac{1}{\sigma_R^2} \mathcal{L}_R(\mathbf{W}) + \frac{1}{\sigma_L^2} \mathcal{L}_L(\mathbf{W}) + \log\sigma_R + \log\sigma_L$$
(E6)

1001

1002 with the two losses:

$$\mathcal{L}_{R}(\mathbf{W}) = -\log(\operatorname{Softmax}(y_{R} = c_{R} \mid \mathbf{f}_{c_{R}}^{\mathbf{W}}))w_{c_{R}}(1 - \operatorname{Softmax}(y_{R} = c_{R} \mid \mathbf{f}_{c_{R}}^{\mathbf{W}}))^{\gamma_{R}}$$
(E7)

$$\mathcal{L}_{L}(\mathbf{W}) = -\log(\operatorname{Softmax}(y_{L} = c_{L} \mid \mathbf{f}_{c_{L}}^{\mathbf{W}}))w_{c_{L}}(1 - \operatorname{Softmax}(y_{L} = c_{L} \mid \mathbf{f}_{c_{L}}^{\mathbf{W}}))^{\gamma_{L}}$$
(E8)

The final weighted loss depends on the model parameters **W** and the two task specific scalars σ_R and σ_L .