# Cascaded learning in intravascular ultrasound: coronary stent delineation in manual pullbacks

**Tobias Wissel[1*], Katharina A. Riedl[3], Klaus Schaefers[2], Hannes Nickisch[1], Fabian J. Brunner[3], Nikolas D. Schnellbaecher[1], Stefan Blankenberg[3,4], Moritz Seiffert[3,4], Michael Grass[1]**

[1]Philips Research – Hamburg, Germany; [2]Philips Research – Eindhoven, The Netherlands
[3]Department of Cardiology, University Heart & Vascular Center Hamburg, Hamburg, Germany
[4]German Center for Cardiovascular Research (DZHK), Partner Site Hamburg/Lübeck/Kiel, Germany

**Abstract**

**Purpose:** Implanting stents to re-open stenotic lesions during percutaneous coronary interventions is considered a standard treatment for the acute or chronic coronary syndrome. Intravascular ultrasound (IVUS) can be used to guide and assess the technical success of these interventions. Automatically segmenting stent struts in IVUS sequences improves workflow efficiency but is non-trivial due to a challenging image appearance entailing manifold ambiguities with other structures. Manual, ungated IVUS pullbacks constitute a challenge in this context. We propose a fully data-driven strategy to first longitudinally detect and subsequently segment stent struts in IVUS frames.

**Approach:** A cascaded deep learning approach is presented. It first trains an encoder model to classify frames as "stent", "no stent", or "no use". A segmentation model then delineates stent struts on a pixel level only in frames with stent label. The first stage of the cascade acts as a gateway to reduce the risk for false positives in the second stage, the segmentation, trained on a smaller and difficult-to-annotate dataset. Training of the classification and segmentation model was based on 49,888 and 1,826 frames of 74 sequences from 35 patients, respectively.

**Results:** The longitudinal classification yielded Dice scores of 92.96%, 82,35%, and 94,03% for the classes "stent", "no stent" and "no use". The segmentation achieved a Dice score of 65.1% on the stent ground truth (intra-observer performance: 75,5%) and 43.5% on all frames (including frames without stent, with guidewires, calcium or without clinical use). The latter improved to 49.5% when gating the frames by the classification decision and further increased to 57.4% with a heuristic on the plausible stent strut area.

**Conclusions:** A data-driven strategy for segmenting stents in ungated, manual pullbacks was presented – the most common and practical scenario in the time-critical clinical workflow. We demonstrated a mitigated risk for ambiguities and false positive predictions.

**Keywords**: intravascular ultrasound (IVUS), coronary, stent, segmentation, detection.

*First Author, E-mail: tobias.wissel@philips.com

## 1 Introduction

Coronary artery disease remains one of the leading causes of death worldwide, accounting for more than 9 million deaths alone in 2016 according to the World Health Organization (WHO) [1]. The disease is caused by atherosclerosis – an accumulation of plaques in the intima of the arterial wall which decrease the effective vessel diameter and thus form a stenosis. Stenotic vessels impede the flow of oxygenated blood into the cardiac muscle causing chest pain (angina pectoris) and ultimately provoking myocardial infarction.

Percutaneous coronary interventions (PCIs) with balloon dilatation and implantation of coronary stents constitute the preferred strategy in most patients with acute myocardial infarction but may also be performed in patients suffering from chronic coronary syndromes to improve symptoms and outcome. Here, intravascular imaging can support several parts of the clinical workflow: Recent studies suggest that treatment planning for complex lesions can significantly benefit from stenting criteria based on intravascular ultrasound [2] and that IVUS-guided procedures improve long-term clinical outcome in patients with acute myocardial infarction [3]. During the PCI, the assessment of stent malapposition, stent underexpansion, stent strut fractures, post-dilation decisions, its placement with respect to other anatomical structures such as bifurcations or the guidance for complex stenting procedures at bifurcation lesions may be improved by intravascular imaging [4]. In particular, the former two are seen as risk factors for in-stent restenosis or thrombosis [5, 6]. In addition, complex PCIs treating bifurcation lesions require elaborated procedures such as the culotte or other techniques. Here, typically two guidewires are inserted into both branches, one is jailed during stent deployment in the first branch and later used as a guide for re-wiring the second branch through the wire cells of the first stent (e.g., by proximal or distal

57     cell technique) [7]. Only after successful rewiring, the second branch can also be stented, which is

58     greatly facilitated by intra-vascular imaging and knowing the location of the first stent and its cells.

59     ACC/AHA and ESC/EACTS guidelines mention intravascular ultrasound (IVUS) as a

60     complement to intravascular optical coherence tomography (IVOCT). Both imaging modalities

61     can potentially mitigate the limitations of X-ray angiography in interventional guidance, but also

62     have different strengths and shortcomings. While IVOCT exhibits a higher cross-sectional

63     resolution and typically better image contrast, it also has a limited radial field-of-view (FoV) which

64     – even though it captures the vascular lumen border – often reaches not far beyond, in particular

65     for larger vessels. Adequate IVOCT imaging also requires saline flushing of the lumen, which is

66     not the case for IVUS. Depending on the transducer frequency, the latter can achieve several

67     different, but compared to IVOCT inferior, cross-sectional resolutions. Despite the characteristic

68     speckle noise patterns, IVUS is well suited to also evaluate plaques and vessel wall compositions

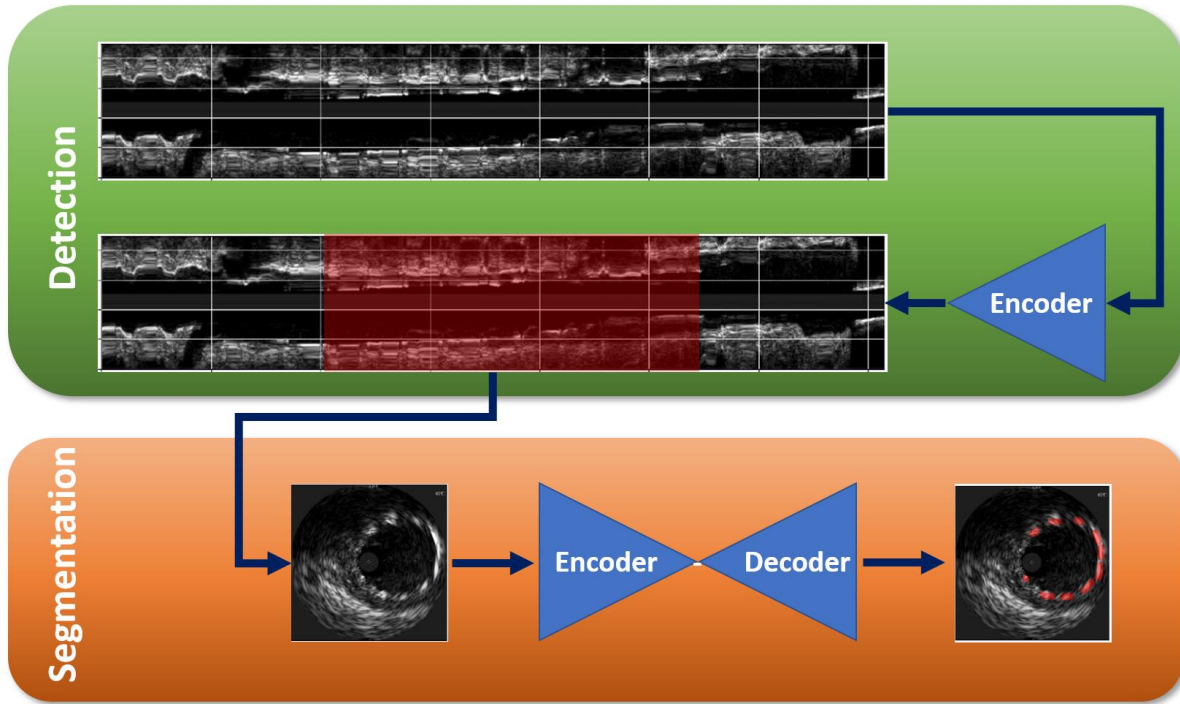69     as it acquires larger FoVs.

70     Advanced image interpretation often requires expertise and training [8, 9]. This establishes the

71     clinical need for making IVUS imaging more accessible and easier to use while minimizing, for

72     instance, procedure prolongation [4]. Both imaging techniques provide sequences of 2D

73     tomographic cross-sectional images when manually or automatically pulled along the artery. The

74     tube-like metal mesh of a stent therefore appears as a collection of spots typically located at the

75     inner circular border or sometimes within the intima. The number and shape of the visible spots

76     depends, for instance, on the inclination of the imaging plane. In IVUS, stent struts often even

77 appear merged into bright circular ring segments due to dense fibrous or micro-calcified tissue

78 close to the stent [26].

79 Motivated by the clinical needs above, research groups successfully detected stent struts in IVOCT

80 images [20, 21, 22]. Here algorithms such as the Tsantis et al. [22] can usually rely on the excellent

81 image quality, the higher image resolution as well as contrast inside the vascular lumen up to the

82 lumen boundary in IVOCT images. In their work they make use of characteristic intensity

83 distributions and wavelet-based matching, which would be rather challenging on IVUS images

84 with an inhomogeneous texture, speckle noise and many spot-like ambiguities.

85 Less groups targeted the more challenging task of extracting stents from IVUS images. Contour-

86 based predictions as proposed by Dijkstra [23] or Kitahara [24] and colleagues worked well, but

87 either contain semi-automatic steps, require well visible stent struts or neglect challenging frames

88 with acoustic shadowing, bifurcations or generally more than 90° signal attenuation in order to

89 work robustly. In clinical practice, frames of an IVUS pullback often contain artifacts caused, for

90 instance, by heart beat related pulsation of the surrounding anatomy, longitudinal swinging of the

91 transducer or simply rapid pulls, transducer-wall or transducer-device interactions during manual

92 pullbacks. Apart from ambiguities due to similar appearance, these artifacts and the low resolution

93 paired with speckle noise entail a high risk of false positives for fully automatic algorithms. This

94 has also been recognized by the research group around Balocco and Ciompi et al. who provided a

95 comprehensive approach to the problem of stent detection in automatic IVUS pullback data sets

96 [26].

97 They present a pipeline of three steps validated on automatic pullbacks which first performs an

98 image-based gating on the sequence [27], then detects stent struts and stent shape in cross-sectional

**Figure 1: Illustration of the cascaded concept. Frames of the manual pullback are first analyzed by an encoder network, which decides for one of three classes per frame: stent, no stent or no use. Only stent frames are then passed on to the encoder-decoder to segment the stent struts. Apart from the favorable training setup, this is also targeting a reduction of false positive predictions on frames that do not show a stent anyway.**

99   images and finally longitudinally localizes the positioned stent in pullback direction. They note

100  that the stent strut detection as proposed by Ciompi et al. [28, 29] is an essential part of the pipeline.

101  First, a 2-stage multi-class AdaBoost classifier generates pixel-wise label maps from handcrafted

102  appearance features. Based on heuristics including knowledge about the luminal area, the most

103  plausible elliptical stent shape is fitted to the strut class before a stent prototype filter confirms an

104  output set of likely struts locations.

105  In their most recent work, Balocco et al. define a likelihood function on the strut label mask, which

106  they successfully convert into a stent indicator variable along the vessel using the SAX algorithm

107  for solving the more challenging problem of stent localization [30]. As part of their conclusion,

108  they acknowledge that the approach could benefit from fully data-driven deep learning techniques.

109  They note that a thorough learning strategy is needed to deal with the otherwise substantial

110  requirement of carefully annotated data to solve this complex task with its high risk for false

111  positive predictions.

112  In this work, we present such a strategy. We propose a cascaded deep learning approach as

113  illustrated in Figure 1. The cascade reverses the pipeline from Balocco et al. [26] in that it first

114  solves the task of longitudinal detection with an encoder network, before an encoder-decoder

115  network segments the stent struts only in the detected stent ranges. This way, the encoder can be

116  trained on a huge amount of efficiently, frame-wise annotated data to solve a simpler task, while

117  the segmenter can be trained on a smaller amount of pixel-level annotated data to solve a more

118  challenging task. The latter type of data is often tedious to annotate, prone to label noise and hard

119  to get in larger quantities.

120  In contrast to previous work, the proposed concept has the ambition to work on single frames from

121  manual, ungated pullbacks as this is closest to time-critical clinical practice. We further recognize

122  that these practical scenarios impose a problem where – in contrast to other typical segmentation

123  tasks on clinical data – many images do not contain the actual target structure or instead contain

124  objects of similar appearance such as guidewires or calcium deposits. Due to the rather frequent

125  and sometimes rapid back and forth transducer motion, manual pullbacks also exhibit a particular

126  risk for artifacts and clinically irrelevant frames. Overall, these conditions require special attention

127  dedicated to false positive predictions during algorithm design. Therefore, the encoder will also be

128  trained against a "no use" class to serve as a better gateway for frames sent to the segmenter and

129  to reduce the risk for false positives.

130  Finally, our contributions can complement algorithmic features from earlier work and can be used

131  – together with existing solutions for lumen and vessel wall segmentation [31, 32, 33] – to assess

132  stent malapposition [26] or facilitate the rewiring during the treatment of bifurcation lesions. For

133  the two use cases, malapposition and rewiring – two different aspects of stent detection and

134  segmentation – are emphasized and will be addressed by dedicated metrics in this work: While the

135  first use case requires precise knowledge about the radial position of the stent mesh, the second

136  one rather asks for the angular occupation to identify mesh cells enabling guidewire passage.

137

## 2    Material and Methods

*2.1 Data and Annotations*

The IVUS data were acquired at the University Heart and Vascular Center (UHC) Hamburg Eppendorf using the Core or Core Mobile precision guided therapy system (Philips Healthcare, San Diego, USA). The pullbacks were conducted manually within the coronary vessels without ECG gating. The acquisitions from pre- or post-interventional stenting procedures used a 20 MHz phased array Eagle Eye Platinum probe (Philips Healthcare, San Diego, USA). Frames were given in Cartesian coordinates of size $500 \times 500$ pixels (isotropic resolution of 0.02mm, i.e. a field-of-view (FoV) of 10mm). They were downsampled to $224 \times 224$ pixels for the predictive task with a preserved FoV. All annotations were obtained in form of a consensus between two clinical experts routinely experienced with IVUS during stenting procedures. The study was approved by the corresponding ethics committee and institutional review board (IRB).

The clinical experts partially annotated 74 sequences from 35 patients using two types of annotations: (1) Intervals indicating the presence of a certain label in frames along the pullback yielding binary, frame-wise annotations per label, and (2) pixel-wise annotations done with a brush tool covering relevant areas in a 2D frame where a label is located.

In the first category, 20,020 frames (83 intervals) of stent, 17,532 frames (76 intervals) without stent and 12,336 frames (91 intervals) of no clinical use were annotated. The latter contained, for

156 instance, frames with artifacts due to rapid pulling, frames where the transducer was covered by

157 the catheter or ostial frames where the intima moved out of the FoV.

158 In the second category, label masks for 827 frames containing stent struts, for 619 calcium frames

159 and for 390 guidewire frames where annotated. In addition, 120 stent frames from another 19

160 pullbacks were annotated twice for assessing the intra-observer variability in the ground truth for
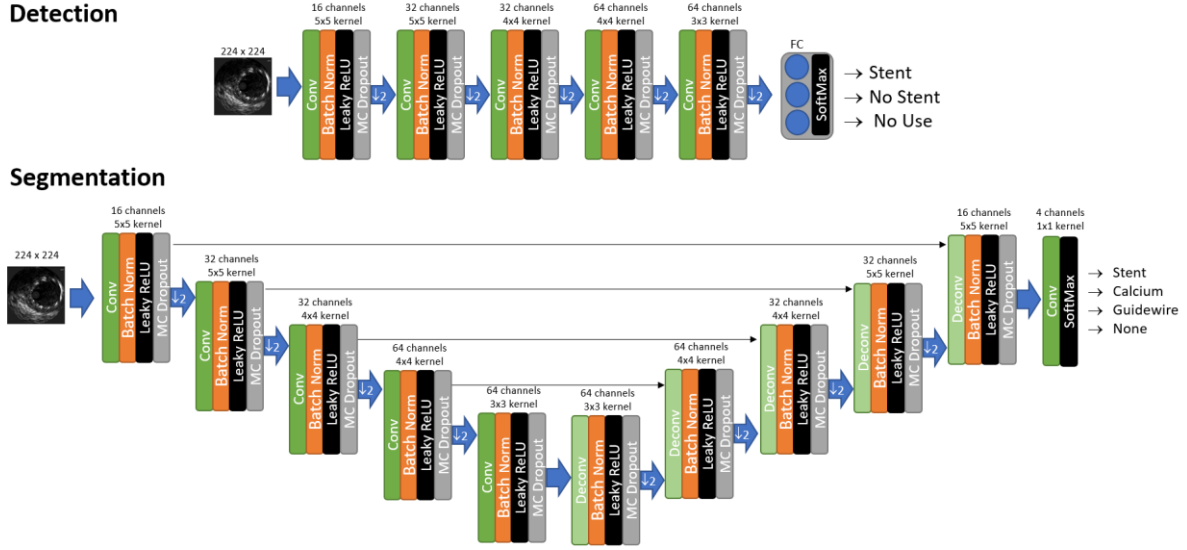
161 the main predictive task.

162 The intra-observer variability for the task of stent segmentation is quantified by a Dice coefficient

163 of 75.47%. This coefficient measures the agreement (in terms of true/false positive/negative pixels)

164 between the first and second annotation of the expert. Only frames with a first annotation (which

165 definitely contain a stent) were presented to the expert in the second round.

166 The data was split into five cross-validation (CV) folds for later evaluation. The partitioning was

167 carried out on patient level, i.e. all pullbacks and therefore all frames that came from the same

168 patients were always assigned to the same fold, such that a model tested on a particular fold would

169 be unbiased and would not have seen data from the test patients during training.

170 Furthermore, the data was partitioned into five folds such that the number of frames per fold were

171 roughly the same and labels per class were as balanced as possibly under these restrictions.

172 *2.2 Approach and Implementation*

173 The cascade makes use of two independently trained networks. During deployment, they

174 subsequently act upon the incoming pullback frames where the decision of the first part, the stent

175 detection, is passed forward to support the pixel-wise prediction of the second part, the stent

176 segmentation. This aims at achieving a twofold goal: (1) reducing the false positive rates on pixel

**Figure 2: Network architectures for stent detection (top, encoder network) and stent segmentation (bottom, encoder-decoder network with skip connections between both parts).**

177   level on frames that potentially do not contain a stent, and (2) allow for faster processing as the

178   segmentation network would only act on positive stent predictions preselected by the detector. To

179   implement this cascade concept, architectures belonging to the encoder-decoder family [12] as

180   well as the architecturally different and popular DeepLabV3 as an alternative option for the more

181   difficult segmentation task were used. The latter makes use of atrous spatial pyramid pooling

182   (ASPP) with dilated convolutions on top of a ResNet backbone [10].

183   The encoder-decoder family comprises several possibilities of variations, which differ by distinct

184   architectural features and give rise to well-known variants such as U-Net, DeconvNet or SegNet

185   for a segmentation task [11, 12, 13]. We therefore used the training set of the first cross-validation

186   fold to grid-search the benefit of such typical design features. These included: choice of the

187   activation function, normalization layer, different pooling/un-pooling variants (including max-

188   pooling, average pooling or strided (de)convolutions), usage of skip connections, Monte-Carlo

189   drop-out layers, padded/unpadded convolutions, or residual and squeeze&excite blocks in the

10

190    convolutional segments. We further optimized typical parameters such as the number of feature

191    maps per segment or the kernel size of the convolutional filters. The search was carried out for the

192    segmentation as well as detection task, whereas the latter only used the encoder.

193    The DeepLabV3 was used with a ResNet-50 backbone and explored in two versions: (1) pre-

194    trained weights on ImageNet, or (2) random weight initialization.

195    The subsequent paragraphs will describe our final choices from these two architecture families.

196    Decisions were made for the best performing design choice and in case of equal performance for

197    the simpler variants following the concept of Occam's razor by not multiplying choices beyond

198    necessity.

199    *2.2.1 Detector*

200    The detector as illustrated at the top of Figure 2 is a five-block encoder network, where each block

201    consists of a convolution, batch normalization, leaky ReLU and Monte-Carlo dropout layer. The

202    one-hot network output in terms of four mutually exclusive classes is realized by a fully connected

203    (fc) layer. The detector was trained for 80 epochs with a learning rate of 0.01 and by augmenting

204    the available number of frames with frame-level annotations by the same number of augmented

205    frames. Data augmentation used random rotation, scaling and axis flips as well as deformations

206    based on trigonometric displacement fields. Classes were equally weighted in a cross-entropy loss

207    term.

208    *2.2.2 Segmenter*

209    The segmentation network candidate from the encoder-decoder family consisted of 2×5 blocks for

210    the encoder and decoder respectively, which are connected via skip connections. The blocks are

211     equivalent to the detector network but did not make use of any weight sharing. The network outputs

212     probability maps for four mutually exclusive classes (stent, guidewire, calcium, none).

213     The encoder-decoder segmenter was trained for 120 epochs with a learning rate of 0.005 and by

214     augmenting the available number of frames with pixel-level annotations by twice the amount of

215     augmented frames. Reaching beyond point estimates for the respective predictions, Monte-Carlo

216     dropout was used (50 samples during deployment) to minimize the effect of individual trainings

217     on the performance and to increase repeatability over trainings [17]. Besides providing a measure

218     of predictive uncertainty, the dropout was also found to improve the segmentation quality when

219     using its mean predictions.

220     The DeepLabV3 candidate was explored in its original form, pre-trained on ImageNet and refined

221     on the IVUS data for 60 epochs as well as with randomly initialized weights and trained from

222     scratch for 120 epochs. The initial learning rate was taken from the original article.

223     All networks were trained using Adam [15], used a weight decay regularization ($\lambda = 0.0001$) and

224     batch normalization pooled across all pixels in a frame, and were initialized with the He method if

225     not stated otherwise [16]. The learning rate followed a schedule of hyperbolic decay such that half

226     the learning rate is reached after half the number of specified epochs. Training was done using the

227     generalized Dice loss (GDL) [14] where the weight of the stent class was doubled compared to the

228     other, auxiliary classes. Encoder-decoder implementations were done in CNTK [18] while the

229     DeepLabV3 variant also made use of PyTorch [19].

230     *2.3 Evaluation*

231    Typical evaluation metrics such as receiver-operator characteristics, Dice score, recall, precision,

232    false positive rates (FPR) and area-under-the curve (AUC) were used to compare performances of

233    the detector and segmenter.

234    As it was not feasible to single out individual stent struts in the majority IVUS frames (see

235    discussion of stent appearances in the introduction), the segmentation ground truth was given in

236    the form of general masks on pixel level. Although this could not give rise to a stent strut detection

237    rate as in [31], we defined two dedicated metrics which are tailored to the two major motivations:

238    (1) malapposition, and (2) rewiring.

239    *2.3.1 Angular Dice Score (DCA)*

240    Successful rewiring requires knowledge of angular occupation to identify cell passages through

241    the stent mesh. To quantify angular overlap, we do not compute the Dice score between the

242    predicted voxels p and the human annotation. Instead, we project both annotation and prediction

243    onto a ring of constant radius around the IVUS catheter tip and compute the Dice score between

244    the projections PR(p) and PR(a). The resulting metric is – like the plain Dice score -- symmetric,

245    between 0 and 1 but focuses on angular overlap while ignoring discrepancies along the radius.

246    *2.3.2 Symmetric Median Skeleton Distance (SMSD)*

247    Confirmation of stent apposition predominantly requires radial accuracy of the stent segmentation

248    to set it into relation to the intima contour of the vessel wall. To achieve this, we construct a metric

249    that is not a relative measure of performance but quantifies accuracy in pixels, i.e. multiples of the

250    image resolution (0.0446mm) and therefore de facto a metric length scale . We skeletonize the

251    prediction and the annotation masks and compute the 50% Hausdorff distance between the two

252 resulting skeletons. Formally, we compute the median Euclidean distance between all prediction

253 skeleton points and their respective closest point on the annotation skeleton. To guarantee

254 symmetry, the final SMSD is then the average between this median and the median from

255 annotation to skeleton

256 The two proposed metrics DCA and SMSD measure complementary aspects of the overlap

257 between prediction and annotation and hence provide a more detailed picture of the performance

258 of our pipeline than the Dice score alone.

259

## 3    Results

*3.1 Stent Detection*

The detector part of the cascade was trained on frame-level annotations to predict intervals which

contain a certain class. These classes were modeled to be mutually exclusive.

The resulting scores pooled from all test sets of the 5-fold cross-validation are shown in Figure 3.

All scores were computed based on the mean probabilities of the Monte-Carlo output distributions.

The normalized confusion matrix in Figure 3 shows a true positive rate above 90% for all classes

on the diagonal (probability threshold 0.5). Highest accuracies are achieved on the "no use" task,

i.e., when detecting frames with artifacts or frames of no clinical use. Lowest accuracy was

obtained for the detection of areas without stent. The same effect can also be observed from the

receiver-operator-characteristic (ROC) curves in Figure 4 (left) as well as from the Dice curves in

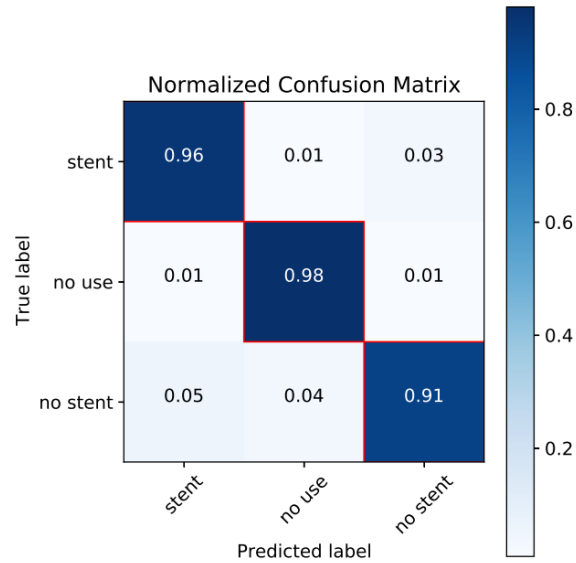Figure 4 (right). Both curves were generated by computing the scores based on a sweep across
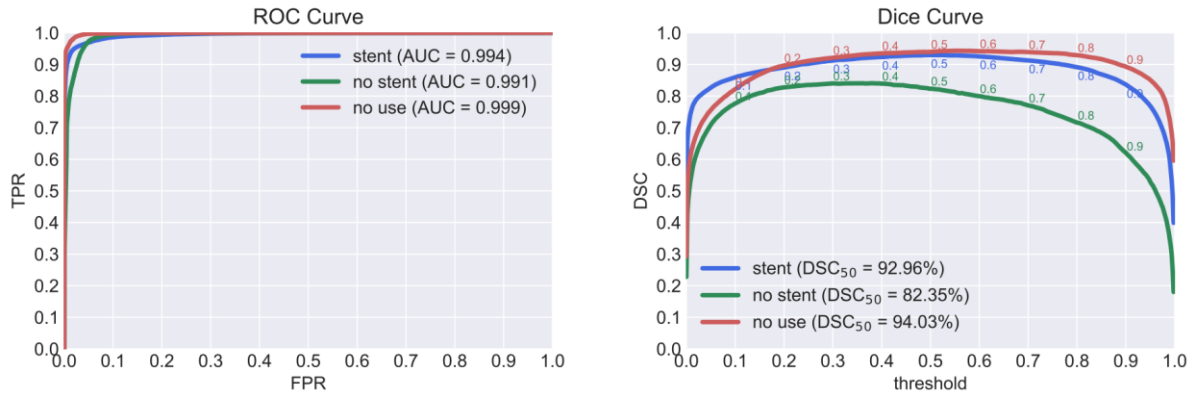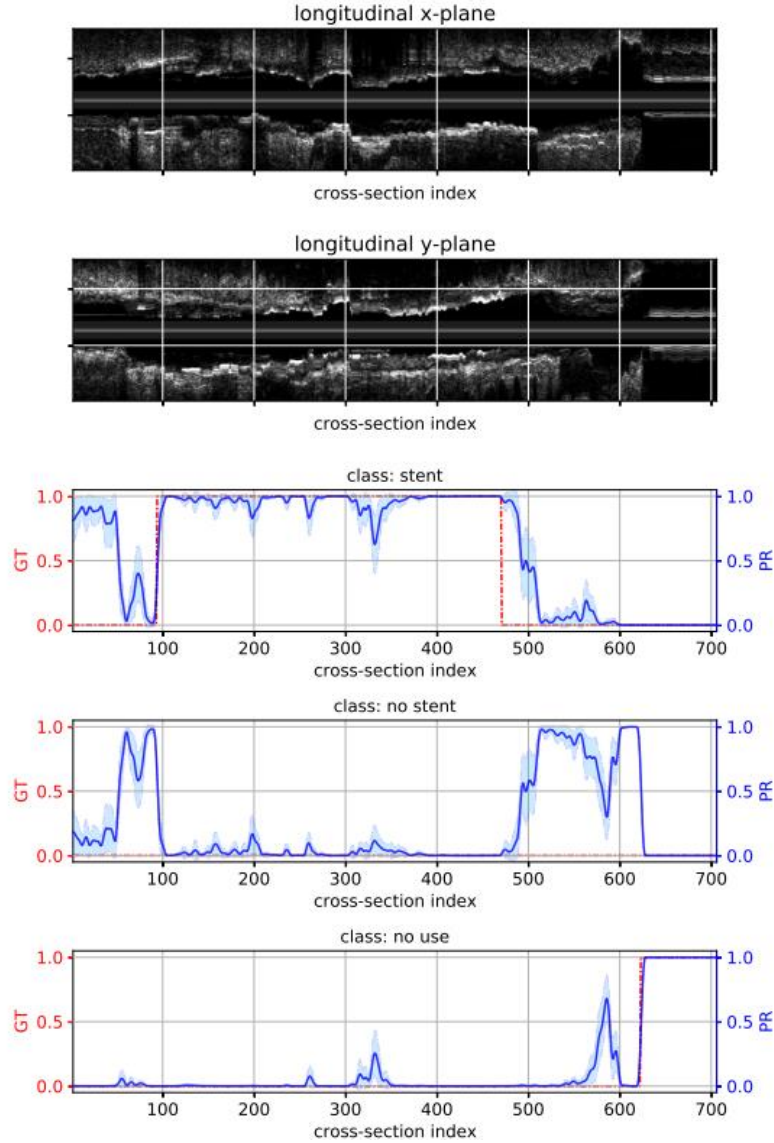


**Figure 3: Confusion matrix for the stent detection network after pooling the results on all five test folds.**

**Figure 4: Left: ROC curve for varying thresholds on the output probability maps. ROC-AUC values are listed for a threshold of t = 0.5. Right: Dice curves showing the dependence of the Dice score on the chosen threshold. A good compromise is achieved when thresholding all classes at 0.5.**

272     different output thresholds for the probability maps ($t \in [0, 1.0]$). The ROC curves yield high areas-

273     under-the-curve (AUCs).

274     Figure 5 shows the mean and standard deviation predictions along an example pullback sequence

275     for each of the first three classes longitudinally plotted.

**Figure 5: Example pullback along with encoder results. First two rows show orthogonal cuts through the pullback in longitudinal direction and the lower three graphs the ground truth (dash-dotted red) and predictions (solid blue line Monte-Carlo mean and blue shading Monte-Carlo estimate of the standard estimation) Regions where none of the three classes has a red ground truth value at 1 have not been annotated by the expert. Here, predictions cannot be compared with a target label. During the last frames the transducer was covered by the catheter, which is correctly recognized as "no use".**

From a cross-sectional perspective, Figure 6 presents example frames along with the corresponding detector decision. The number of true positives, false positives and false negatives displayed does not reflect their corresponding proportions in the overall predictions.
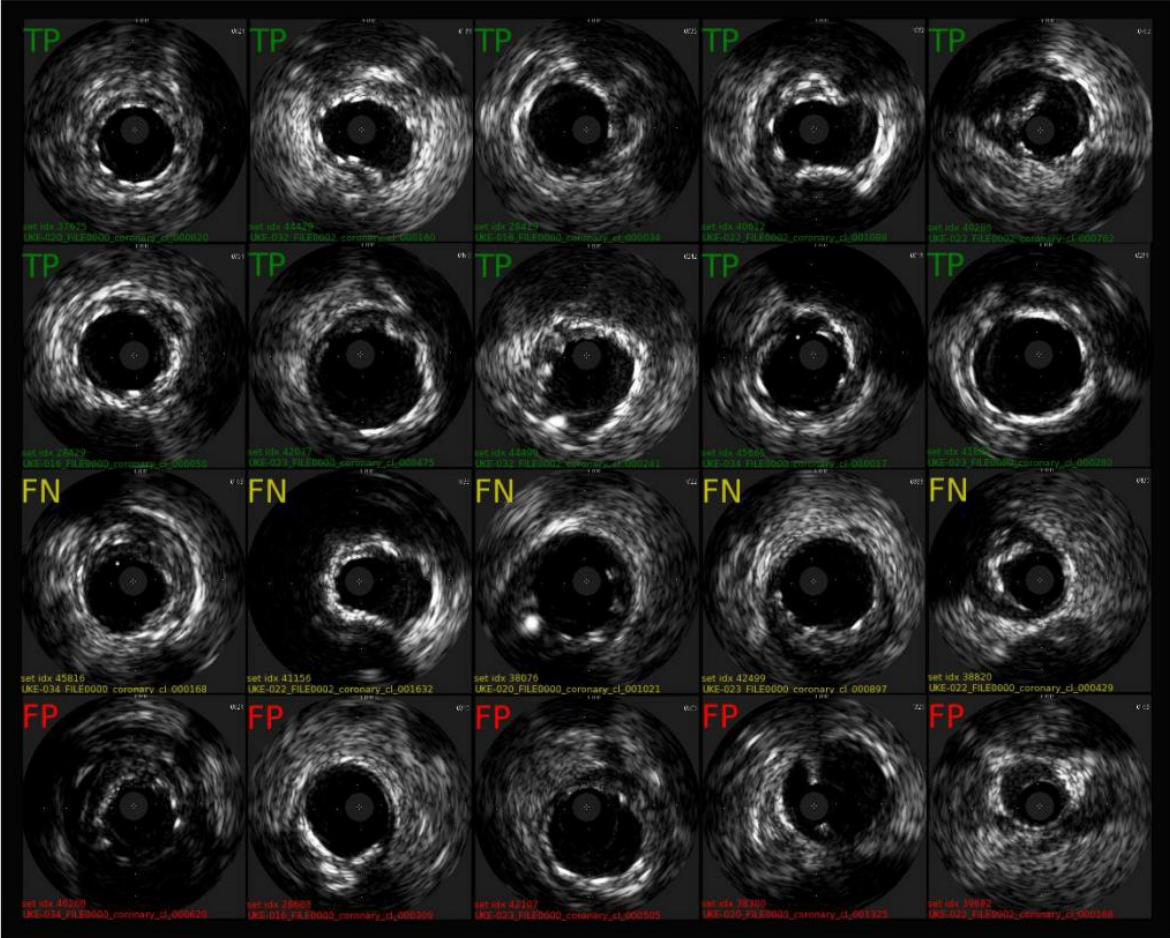
17

**Figure 6: Example frames labeled with their corresponding predictions from the stent output of the detector network: true positive decisions (top two rows, typical to challenging examples from left to right), false negative decisions (detector missed ground truth annotations, third row), and false positive detections (bottom row, wrong predictions without ground truth label).**

279    In the first two rows very typical stent examples can be seen on the left, where bright speckle spots

280    – indicating metal struts – are visible along the outer lumen contour. Further to the right, frames

281    are getting more challenging: stent struts can only be seen partially along the lumen border

282    (sometimes only two or three of them), calcium deposits are also present in the frame or smaller

283    calcifications are even attached to the stent making it more difficult to distinguish single struts and

284    merging them into arc-shaped structures.

285    Visual inspection of some of the false positive examples (first and third example in that row)

286    suggests that although ground truth was not labeling them as stent frames, they still seem to clearly

287  show stent struts. In that sense they are only tagged as false positives because they are evaluated

288  against remaining inconsistencies in the ground truth. The other false positives again fall into the

289  ambiguous borderline category described above, where a retrospective decision even for an expert

290  is challenging.

291  *3.2 Stent Segmentation*

292  The second stage of the cascade, the segmentation network, was implemented using an encoder-

293  decoder and a DeepLabV3 architecture trained on pixel-level annotations. The intra-observer

294  variability for the task of stent segmentation yielded a Dice coefficient of 75.47% indicating the

295  extent of agreement between the first and second annotation of the experts. For the second

296  annotation only frames from the first annotation (which definitely contain a stent) were presented

297  to the expert. The Dice metric does therefore not capture any errors the expert would make when

298  re-identifying these frames as stent frames or when ignoring frames without a stent (but maybe

299  calcifications or other ambiguities).

300

301  Given this reference context we evaluated the automatic segmentation as follows: The network

302  was exclusively trained on frames of the overall dataset, which are accompanied by pixel-level

303  annotations of the involved classes (stent, calcium, guidewire). All other frames were ignored as

304  previous scouting experiments did not provide any beneficial evidence for the involvement of other

Table 1:Comparison of metrics scoring the quality of stent segmentation on two different data supports: (1) all frames with available ground truth, (2) only frames containing stent pixel-level ground truth. Scores are Dice coefficient, false positive rate (FPR), angular Dice coefficient (DCA) and symmetric median skeleton distance (SMSD) in pixels, i.e. multiples of the image resolution 0.0446mm. Best scores are highlighted in bold.

| All | Dice [%] | FPR [%] | DCA [%] | SMSD [px] |
|---|---|---|---|---|
| **Encoder-Decoder** | 43.47 | **0.38** | 50.29 | 5.17 |
| **DeepLabV3, 120ep, no pretrain** | **46.65** | 0.71 | **54.54** | 5.12 |
| **DeepLabV3, 60ep, pretrain** | 45.42 | 0.9 | 53.83 | **5.06** |
| **Stent Loc** | | | | |
| **Encoder-Decoder** | **65.05** | **0.51** | **75.20** | **2.77** |
| **DeepLabV3, 120ep, no pretrain** | 62.50 | 1.09 | 72.95 | 3.47 |
| **DeepLabV3, 60ep, pretrain** | 61.96 | 1.37 | 73.69 | 3.23 |

305 additional labels (such as "no stent"). At each epoch a new balanced and shuffled set of frames

306 from this pre-selection was used for training.

307 In the following we tag this set of frames with the identifier "all" in contrast to the identifier "stent

308 loc", which refers to the set of frames on which for which a stent ground truth was available on

309 pixel level.

310 In a first step, we evaluated the performances of the encoder-decoder candidate as well as the

311 DeepLabV3 architecture on both sets. The results in Table 1 show similar scores with either

312 approach. While the encoder-decoder performed slightly better on the stent frames, DeepLabV3

313 had partly better scores when evaluated on all frames. Note, that the encoder-decoder variant was

314 insensitive to varying thresholds on its probability outputs and therefore used a simple argmax,

315 while the probability threshold on the stent class was optimized for DeepLabV3.

316 Overall, we consider differences in segmentation quality only as minor and hence proceed with

317 the encoder-decoder variant for further evaluations. With this decision we follow Occam's razor

318 by selecting the simpler of two equally performing approaches. Furthermore, we consider the

319 slightly better performance of the encoder-decoder network on the stent frames as more relevant

320 because the stent detector part of the later cascade will prevent the segmentation network from

321 being applied to all frames. Finally, even on all frames it showed the lowest FPR establishing

322 excellent conditions for the cascade.

323

324 In a second step, evaluation and testing of stent segmentation performance then had different foci

325 of interest:

326    1. Stent segmentation was evaluated on all frames carrying a stent pixel-level label

327        (equivalent to the intra-observer scenario, and similar to evaluations in other studies on

328        segmentation tasks).

329    2. Stent segmentation was evaluated on all frames of the above-mentioned pre-selection

330        (including all frames which did not contain a stent, but only one of the other labels).

331    3. Stent segmentation was evaluated based on all frames that passed a gating criterion (either

332        a simple heuristic, or a positive stent detector decision, or both).

333 The three scenarios above can be interpreted as follows. The first scenario would be the ideal case

334 where it was possible to perfectly identify which frames contain a stent and therefore avoiding any

335 false positives on frames which do not show a stent anyway. This constitutes an upper bound for

336 scores which can be achieved with the other two scenarios. The second scenario is the naive

337 solution: The segmentation network needs to identify frames without a stent on its own and avoid

338 highlighting any pixels in them as stent. Finally, the last scenario is the proposed one, where a

339  mechanism is used to sort out frames at which the segmentation network does not need to look,

340  because they contain no stent. This has the advantage of not only minimizing the risk for

341  unnecessary false positives, but it also reduces computation cost as the segmentation network is

342  only deployed if really required.

343  *3.2.1 Simple Gating Heuristic*

344  A simple heuristic can be constructed under two hypotheses. First, the majority of false positive

345  predictions on frames that do not show a stent are of rather small size and it is rare that the network

346  would claim large areas as stent although they are not. Second, frames which do show a stent

347  typically cover a certain minimal area as there should be a minimal number of struts visible in the
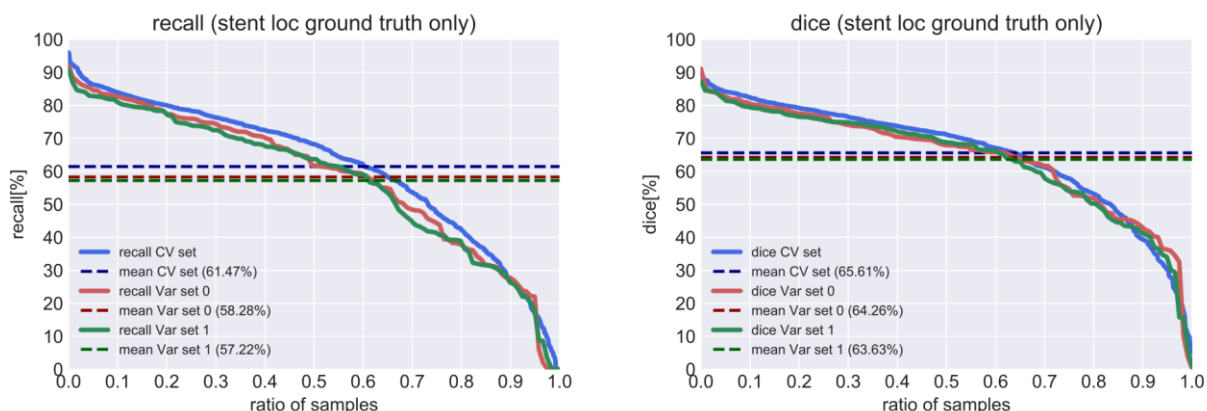
348  image, e.g. at least more than three [40].

349  Under these circumstances a threshold on size of the segmented area can be applied (in mm² or as

350  the image size is always the same in number of pixels). If the number of pixels segmented as stent

351  in a frame falls below the threshold, then the frame will be cleared from them.

352  A histogram on areas covered by the strut ground truth masks reveals that areas below 500 pixels

353  are very unlikely. In frames with 224×224 pixels, 125 pixels (0.249 mm²) roughly correspond to

354  the area of one stent strut. We therefore explore thresholds of either 400 or 500 pixels only.
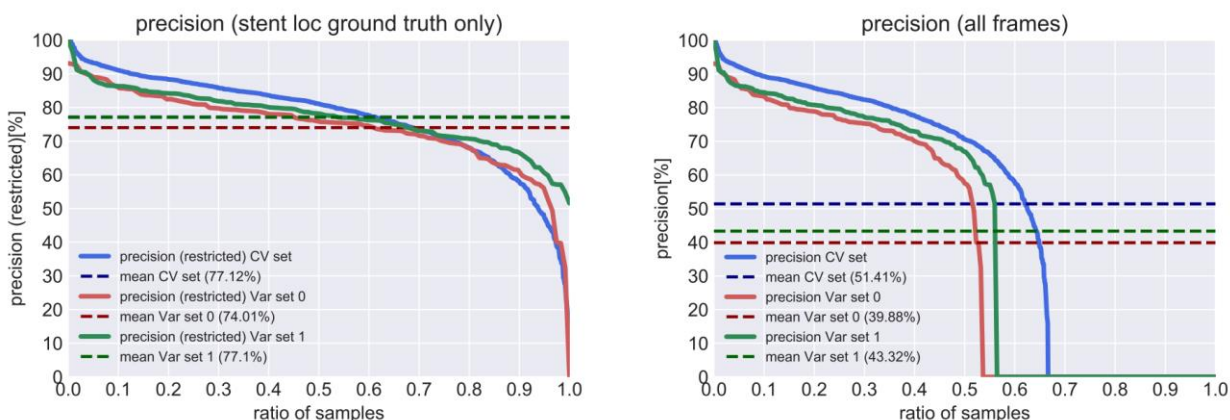
355  *3.2.2 Stent Detector Gating.*

356  In line with the initial proposal for the cascaded approach, we investigate the benefit of using the

357  detector network from Section 3.1 for gating the segmentation task. The segmentation network

358  will only act upon frames which have a positive decision for stent from the detector network.

359  The detector network was only applied on unseen frames: As the trainings for detector and

360  segmentation networks used the same cross-validation split on patient level, we always applied the

361  corresponding unbiased detector network to the segmentation test fold at hand.



**Figure 8: Sorted rank plots for recall and Dice scores computed per-frame. Dashed lines indicate average metric on the intra-observer variance set. Left: Recall ranking for all frames with annotated ground truth. Right: Dice score ranking for all frames with annotated ground truth.**



**Figure 7: Sorted rank plots for the precision score computed per-frame. Dashed lines indicate average metric on the intra-observer variance set. Left: Precision ranking for all samples with annotated ground truth. Right: Precision ranking for all frames containing automatic segmentations.**

362  Scores can be computed (1) either on the complete set of frames after the gating correction of the

363  detector took place, or (2) scores can be computed only on all frames with a positive detector

364  decision.

365  Both scores are expected to be similar, but not equal as the detector network still has a low risk for

366  false negatives outside the set of frames with positive detector decision.

367  *3.2.3 Overall Scores*



**Figure 9: Score matrices for the segmentation network. Scores are presented for different frame supports and post processing (pp) steps: no post-processing (no-pp), segmentation rejection based on a "< N pixel" threshold (Npx) and based on detector decisions (detect). Left: Dice scores for three different post-processing scenarios (vertical) on different frame supports (horizontal). The effect of the detector is evaluated on all frames (segmentations are corrected based on detector decision) and on frames with positive detector decision for stent only. Right: False positive rates (fpr) (listed as average number of pixels in a 224×224 frame) for four different post-processing scenarios (vertical) on different frame supports (horizontal). A false positive rate of 0.25% or 125 pixels (0.249 mm²) roughly corresponds to the area of one stent strut as pictured by the imaging modality.**

368  The first three graphs in Figure 8 and Figure 7 show the scores only computed on all frames that

369  have a pixel-level stent ground truth. Scores (recall, Dice, and precision) have been computed per

370  frame and ranked from the highest to the lowest value. Dashed horizontal lines indicate the average

371  scores on test sets from all folds joined as well as the average scores on the held-out intra-observer

372  set (both first as well as second annotation run of the expert).

373  Generally, the average scores for all sets are very similar and low scores are only found on a small

374  set of frames (steep fall on the right margin of the Dice and precision plots). The fourth plot shows

375  the precision evaluated on all frames. A precision can be computed on all frames on which the

376  number of predicted stent pixels is different from zero. It can be observed that there is a larger

24

377   number of frames with zero precision, i.e. frames which did not have a ground truth annotation but

378   were still assigned stent pixels by the segmentation network. Further investigations were carried

379   out to evaluate the characteristics of these predictions (e.g. whether they contain large or small

380   areas of false positives).

381   *3.2.4 Dice Scores and False Positive Rates*
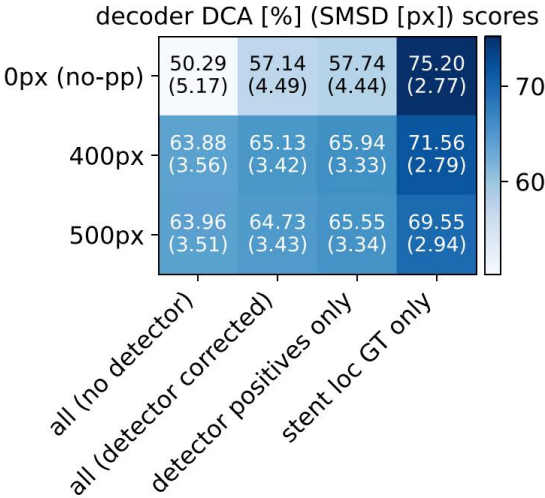
382   Figure 9 (left) presents Dice scores in a matrix for different frame supports (sets of frames on

383   which the score was computed) and different post-processing procedures. In the very right column,

384   the upper bound Dice score is shown. This bound is achieved when only evaluating on frames with

385   stent ground truth. After applying the thresholding heuristic at 500 px or 400 px a small decrease

386   can be observed. While the threshold can correct a moderate number of probable false positive

387   segmentations, it also bears the risk to erroneously reject actual true positives in rare cases.

388   On the very left column, we list the Dice scores when computing them on all frames. Scores are

389   generally smaller than for the previous case due to frames of zero precision (cf. Figure 7 right).

390   After applying the heuristic, more than 10% improvement in Dice can be observed. This indicates

391   that indeed most frames with zero precision are caused by only small false positively assigned

392   areas in the frame.

393   When using the detector-based gating, an increase in Dice score can be observed, too and when

394   applying both at the same time, an additional benefit of the detector given the already applied

395   heuristic can be seen. This is true for both scores, those computed on detector frames with positive

396   decision for stent as well as on all frames after detector gating has been applied. The latter is

397   slightly lower, which reflects the small number of false negatives from the detector network.

398   The Dice scores also indicate that a more conservative threshold of 400 px is more suitable than

399   the higher 500 px one. Overall, the segmentation achieves 86.2% of the human (intra-operator)

400 performance when only applied to frames with a stent ground truth (same scenario the human

401 expert faced) and 75.1% of the human performance on all frames when used within the cascade.



**Figure 10: Angular Dice (DCA) and symmetric median skeleton distance (SMSD, in brackets) scores of the segmentation network. The former is reported in percent and the latter in pixel, i.e. multiples of the image resolution 0.0446mm. Equivalent to the Dice scores in Figure 9 results are shown without post-processing (no-pp) and after applying the 400px or 500px heuristic. Again, we evaluate on all frames, stent-only frames as well as on the two cascading scenarios**

402 Figure 9 (right) provides more details on specific false positive rates. The matrix illustrates average

403 false positive rates for different post-processing scenarios (no post-processing, the thresholding

404 heuristic, detector gating or both) and frame supports (arranged according to the expert label they

405 carry and irrespective of whether they have been part of the cross validation sets or not). Again,

406 the combination of detector and heuristic post-processing yields best results.

407 *3.2.5 Angular Dice and Median Skeleton Distance*

408 Equivalent to Figure 9, Figure 10 reports the results for the angular Dice (DCA) and the symmetric

409 median skeleton distance. The overall performance of the segmentation network in terms of these

410 scores confirms the trend observations made on the Dice score earlier. The angular Dice exceeds

411 the full Dice by approximately 8-10% indicating a good angular coverage of the stent ground truth.

412    The median skeleton distance yielded values in the lower single digit regime which also confirms

413    a segmentation which is well aligned with the ground truth in radial direction. Both scores are

414    constructed to be insensitive with respect to inaccuracies in the thickness of the stent segmentation.

415    *3.2.6 Example Segmentations*

416    Example frames with their corresponding segmentations can be seen in Figure 11. Both images

417    illustrate random examples drawn equally spaced from the score rankings discussed earlier (Dice

418    score in the top and precision in the bottom image). Green areas indicate true positive

419    segmentations (network was correct), red indicates false positive areas (network predicted stent on

420    pixels where ground truth said the opposite), and yellow indicates false negatives (stent areas

421    missed by the column-wise (from left-to-right and top-to-bottom). Scores for the examples

422    decrease from left to right due to the way they are sampled from the ranking plots. Therefore, the

423    first frames give an impression on very good segmentations, while frames on the right indicate

424    some of the poorer frames where the stent is visible in terms of small bright struts along the outer

425    lumen contour.

**Figure 11: Example frames drawn in steps of constant frame proportions from the score rankings. Therefore, the plots provide a representative selection of frames covering the full range of scores achieved per frame (scores decrease from left to right). Green: true positive pixels, Yellow: false negative pixels, Red: false positive pixels. Top: Frames drawn at equal spacings from the Dice rankings (Figure 8, right). Bottom: Frames drawn at equal spacings from the precision rankings (Figure 7, right).**

## 4 Discussion

*4.1 Stent Detection*

With a true positive rate beyond 95% stent detection performance was found to be excellent and hardly justifies the usage of more complex and heavy-weight architectures. Increasing complexity by elevating the number of feature maps, has even been found to have a negative effect on the performance.

Possible reasons for the slightly weaker results on the "no stent" class are twofold. On the one hand, it is the class with the least training examples available. This could introduce a bias with respect to the other two classes even though samples are balanced out prospectively by the minibatch sampling. However, the pool of input image variation from which the balanced set can be sampled in each epoch is still higher for the two larger classes.

Moreover, frames without a stent were maybe the most difficult to annotate. While for artifact or stent regions the expert would always choose the indicative example intervals, it is more difficult to define the typical "no stent" interval. Such parts of the pullback can still contain a lot of variation such as guidewires, calcium, bifurcations etc. Furthermore, intervals between clear stent intervals, for instance, could still accidentally contain a (for the human eye) poorly visible stent at the pullback transition into the stent.

The former also appears to be one of the main reasons for false negative stent predictions: Most of the examples in row three of Figure 6 exhibit thickened bright outlines at the lumen border which are easily confused between calcification, stent or calcified stent – even for a human expert. Only rarely a few single struts stand out clearly. Finally, the bottom row of the same Figure illustrates label noise originating from stent frames that were occasionally missed by the expert during the

449    time-consuming annotation process. While most of these frames only contain a few single struts,

450    they still impact the performance measures of the detector in terms of false positive predictions.

451    *4.2 Stent Segmentation*

452    Despite the very promising segmentation quality as evaluated by visual inspection (see for instance

453    Figure 11), the absolute Dice values may seem only moderate. This is a well-known phenomenon

454    which is frequently observed with the segmentation of smaller structures as opposed to segmenting

455    larger organs from images. One example is given by lymph node segmentations from CT images

456    where Dice scores of around 50% are typically achieved [40]. Like stent struts, these structures

457    only occupy smaller areas with a low number of pixels. The label noise is typically higher as they

458    are more difficult to recognize or to precisely annotate by the clinical expert. Apart from

459    overlooking some of them, the annotation brush as used in this work easily misses or erroneously

460    annotates a few pixels. Furthermore, the transition between stent strut and surrounding areas is not

461    always clear due to the speckle noise and point spread function of the image acquisition which

462    operate on a similar scale. Finally, the amount of data and the extent of the annotation burden

463    naturally contribute to this effect and the resulting human inconsistencies. The segmentation

464    network will find a systematic way of labelling stent struts, but inaccuracies of a few pixels will

465    significantly impact the Dice score. The Dice achieved in the intra-observer study confirms this

466    effect without involvement of the network.

467    Dice scores of the encoder-decoder candidate were found to be like those of the DeepLabV3

468    candidate. As a possible reason we remark that due to the number of convolutional segments and

469    larger kernel sizes the effective field size of the optimized encoder-decoder network covers almost

470    the full field-of-view. We found this to be an important aspect as only the global view onto the

471    image will provide information about systematic patterns in the composition of a typical IVUS

472     frame and the typical location of a stent within it. The very low number of false positives found

473     further away from the lumen border supports this hypothesis. At the same time the large effective

474     field size prevents the network from getting lost in the interpretations on finer scales comprising

475     ambiguous speckles. Finally, the ASPP blocks with their dilated convolutions in DeepLabV3 may

476     play a similar role for the other candidate. It however does not suffice for a unique advantage as

477     the encoder-decoder candidate can compensate for this with its effective field size. This is also

478     support by the scores achieved. Finally, the lower complexity encoder-decoder network also serves

479     as a regularization against an over-interpretation of the image (see the lower FPR). This has

480     particular importance for predictive tasks on IVUS data on which the risk of misinterpreting

481     speckle as the target seems to be higher than missing a uniquely identifying feature in the

482     appearance of stent struts.

483

484     In relative comparisons, the Dice scores are highest when only considering frames that definitely

485     show stent struts. Any segmentation decision on the other frames of the set can only degrade this

486     score and depart from this upper bound or in the best case keep it unchanged. The latter can only

487     be achieved if the segmentation network will not decide for stent in any pixel of these frames. To

488     avoid this decision in first place, the cascade proposes to shift it to either the introduced heuristic,

489     the detector network or both. Indeed, best performances are achieved when combining the 400px

490     heuristic with the detector network. The positive impact of the former confirms that the decrease

491     of the Dice in the larger frame sets mainly originates from a small number of falsely labelled pixels

492     such as bright ambiguously looking speckles.

493     Yet, the false positive rate is highest in frames containing a stent label. This can be expected as the

494     typical speckle appearance of IVUS frames makes it difficult to decide where a strut starts and

495    where it ends. We have also seen in Figure 6 that for parts of the bright outer lumen contour it is

496    sometimes difficult to judge whether it still is a (e.g. micro-calcified) stent. This also leads to the

497    moderate Dice score between the two expert annotations.

498    In general, the low average rates for the FPR again confirm that the zero precision frames in Figure

499    7 are mainly caused by small erroneously segmented areas. The rates are naturally lowest on "no

500    stent" frames, and higher on calcified frames, frames of no clinical use (e.g. with artifacts) or

501    frames showing guidewires.

502    Once calcifications enter the frame or start to mix with the stent struts, the task gets more difficult

503    for both expert and network. In the bottom image of Figure 11, the very right columns present zero

504    precision frames. Again, we observe only small areas – mainly single, sometimes stretched, bright

505    spots which are mistaken as stent by the network. Some of them are again debatable from a ground

506    truth perspective (such as first row, second from the left).

507

508    Finally, the DCA and SMSD confirm the trends discussed above. Moreover, their motivation is

509    twofold: First, they are tailored to the two main clinical use cases identified earlier: the detection

510    of stent malapposition which requires accurately locating the stent contour in radial direction, and

511    re-wiring branches in bifurcation lesions which require knowing which angular segments are

512    occupied by the stent mesh and which parts could allow guidewire passage. The former is targeted

513    by SMSD, which evaluates the median distance between the ground truth and segmentation

514    skeletons of the stent masks. We compute the median as opposed to larger percentiles or the

515    maximum (as done in Hausdorff distances) to focus on the radial errors as much as possible. We

516    expect higher percentiles to be more influenced by relative angular shifts between ground truth

517  and segmentation. This is however already captured by the angular Dice (DCA), which is more

518  dedicated to the re-wiring use case.

519  As a second motivation, both scores aim at decreasing the impact originating from the "thickness"

520  of the annotated or segmented stents. In addition to the influential factors discussed above, also

521  the brush size favored by the clinical expert may have introduced a bias into the annotations, which

522  should not enter the evaluation. The existence of this bias seems to be confirmed by an angular

523  Dice which is up to 10% higher than the original Dice score while observing a small SMSD at the

524  same time.

525  *4.3 Cascading Concept for a Data-Driven Approach to Stent Segmentation*

526  In the cascaded concept, the encoder successfully served as a gateway for triggering the

527  segmentation network and can thus mitigate the high risk for false positives. This demonstrates

528  the value and superiority of the cascade approach with respect to a conventional pure segmentation.

529  The segmentation network faces the challenge of solving the harder task of pixel-wise stent strut

530  localization while learning from a very limited set of difficult-to-annotate frames. A risk for false

531  positive segmentations then naturally stems from manifold ambiguities in the IVUS image which

532  are difficult to resolve without a more global view on the data – even for the human expert. The

533  disproportionally higher number of frames without a stent compared to those that in fact contain a

534  stent contributes to this. On frames with a stent ground truth the segmentation network achieved

535  86.2% of the human expert performance. This constitutes the fairest comparison because the expert

536  faced only frames during intra-observer variability analysis of which it was known that they show

537    stent struts. Yet, the detector-segmenter cascade applied to on all frames still reached 75.1% of

538    this previous human performance.

539    This aspect gains even higher relevance when considering that all predictions operate on a per-

540    frame basis as all pullbacks were ungated and conducted manually unlike the automatic pullbacks

541    from previous work [26, 29]. These are carried out at constant directed speed and allow the

542    algorithm to exploit longitudinal context. This was not the case here, i.e. the transducer could, for

543    instance, remain at a certain location for some time, go back and forth, be pulled at varying speeds,

544    or even get caught in calcified twists and be subject to rapid jumps. In addition, anatomical and

545    device-related ambiguities with similarly bright contrast have also been identified earlier: Liu at

546    al. [34], for instance, report the visibility of the pericardial border close to the vessel, guidewires

547    or varying contrast due to non-orthogonal reflections of the acoustic wave as misleading factors

548    for calcium assessment. These and others certainly also play a role in automatic stent detection as

549    struts typically have a similar localization in the vessel, exhibit similar contrast and texture features

550    and may finally be entangled with other structures such as spotty calcifications or dense fibers

551    [26].

552    **5   Conclusions and Future Directions**

553    In our work, we successfully demonstrated a data-driven deep learning strategy for segmenting

554    stent struts in IVUS frames. We used a cascaded approach which reverses the order of in-frame

555    strut segmentation and longitudinal stent detection as it was proposed in previous work on the

556    same topic which still made use of a handcrafted processing chain [26, 40]. Using a cascade over

557    a conventional frame-wise segmentation along the pullback yielded a successful strategy for

558    opening up the clinical problem to modern data-driven learning. The results confirm our rationale:

559 Training an encoder network for the simpler task of frame classification on a large set of easy-to-

560 annotate frames yields promising results which can be leveraged to facilitate strut segmentation of

561 a subsequent encoder-decoder segmentation network. Being able to address the problem of stent

562 segmentation by modern learning techniques in such a way will substantially facilitate future

563 development cycles reacting, for instance, on covariate shift stemming from different IVUS

564 transducer types, transducer frequencies, varying post-processing steps in the image acquisition

565 chain or the application to different vasculature. Furthermore, our work sets the stage for easily

566 integrating further predictions in terms of multi-task learning using the same strategy but pushing

567 the exploitation of synergies between related tasks.

568 We also demonstrated that additional regularization and post-processing can further increase

569 consistency in the segmentation output. Here, we applied a simple heuristic based on the total area

570 covered by the strut prediction in a frame. Although it was not our focus here, future work can

571 make use of more elaborate heuristics such as number and size of connected components in the

572 thresholded output probability mask, or the position of the segmented strut areas with respect to

573 other identified anatomical structures such as the outer lumen contour (internal elastic lamina) or

574 plaques in intima as done in other studies [26, 29].

575 In this context also other valuable findings from previous studies can complement our work.

576 Inspired by recent successes of combining radiomics and data-driven approaches, the handcrafted

577 features identified by Ciompi et al. [29] can be integrated and possibly further increase robustness

578 of our approach. In addition, the SAX algorithm proposed by Balocco et al. [30] can be

579 investigated as a replacement for the simple thresholding strategy which we used for the encoder

580 part. As the encoder, however, also classifies frames of no use and without a stent in a mutually

581 exclusive manner, the feasibility of a multivariate variant needs to be considered carefully. Our

582 more intuitive likelihood function compared to [30] at the encoder output (defined in [0, 1]) may

583 already sufficiently solve the problem of mapping output probabilities to discrete class decisions.

584 Our Dice curves confirm minimal dependence on the choice of the output threshold.

585 Finally, future work aims at improving the robustness of our approach e.g. exploiting prior

586 knowledge about the data domain [36] and common error sources [35] or synthetic data generation

587 approaches [37]. Furthermore, while our approach already solves other problems such as the

588 identification of clinically useful frames, e.g. for interventional or retrospective navigation, further

589 predictive models can be combined with the cascade. A natural goal is an extension by existing

590 lumen/vessel wall segmentation models [35] such that stent malapposition can be automatically

591 assessed or a combination with bifurcation and side branch detection to facilitate rewiring by

592 intelligent guidance in complex PCIs.

593

## Disclosures

Dr. Seiffert reports non-financial support from Boston Scientific as well as grants and personal fees from Philips outside the submitted work.

## Acknowledgments

## References

1. Alexandra N. Nowbar, Mauro Gitto, James P. Howard, Darrel P. Francis, Rasha Al-Lamee, Mortality From Ischemic Heart Disease, Analysis of Data From the World Health Organization and Coronary Artery Disease Risk Factors From NCD Risk Factor Collaboration, Circulation: Cardiovascular Quality and Outcomes. 2019; 12:e005375, 2019

2. Zhang J, Gao X, Kan J, Ge Z, Han L, Lu S, Tian N, Lin S, Lu Q, Wu X, Li Q, Liu Z, Chen Y, Qian X, Wang J, Chai D, Chen C, Li X, Gogas BD, Pan T, Shan S, Ye F, Chen SL. Intravascular Ultrasound Versus Angiography-Guided Drug-Eluting Stent Implantation: The ULTIMATE Trial. J Am Coll Cardiol. 2018 Dec 18;72(24):3126-3137. doi: 10.1016/j.jacc.2018.09.013. Epub 2018 Sep 24. PMID: 30261237.

3. Ik Jun Choi, Sungmin Lim, Eun Ho Choo, Byung-Hee Hwang, Chan Joon Kim, Mahn-Won Park, Jong-Min Lee, Chul Soo Park, Hee Yeol Kim, Ki-Dong Yoo, Doo Soo Jeon, Ho Joong Youn, Wook-Sung Chung, Min Chul Kim, Myung Ho Jeong, Youngkeun Ahn, & Kiyuk Chang (2021). Impact of Intravascular Ultrasound on Long-Term Clinical Outcomes in Patients With Acute Myocardial Infarction. *JACC: Cardiovascular Interventions, 14(22), 2431-2443.*

4. Gao XF, Ge Z, Kong XQ, Kan J, Han L, Lu S, Tian NL, Lin S, Lu QH, Wang XY, Li QH, Liu ZZ, Chen Y, Qian XS, Wang J, Chai DY, Chen CH, Pan T, Ye F, Zhang JJ, Chen SL; ULTIMATE

617      Investigators. 3-Year Outcomes of the ULTIMATE Trial Comparing Intravascular Ultrasound Versus

618      Angiography-Guided Drug-Eluting Stent Implantation. JACC Cardiovasc Interv. 2021 Feb

619      8;14(3):247-257. doi: 10.1016/j.jcin.2020.10.001. Epub 2020 Oct 29. PMID: 33541535.

620   5.   Trabattoni D, Bartorelli AL. IVUS in bifurcation stenting: what have we learned? EuroIntervention.

621      2010; 6(suppl J): J88 – J93. https://doi.org/10.4244/EIJV6SUPJA14.

622   6.   Yoon HJ, Hur SH. Optimization of stent deployment by intravascular ultrasound. Korean J Intern

623      Med. 2012; 27 (1):30 – 38. https://doi.org/10.3904/kjim.2012.27.1.30.

624   7.   Fadi J. Sawaya, Thierry Lefèvre, Bernard Chevalier, Phillipe Garot, Thomas Hovasse, Marie-Claude

625      Morice, Tanveer Rab, & Yves Louvard (2016). Contemporary Approach to Coronary Bifurcation

626      Lesion Treatment. *JACC: Cardiovascular Interventions, 9(18), 1861-1878.*

627   8.   Neumann F-J, Sousa-Uva M, Ahlsson A, Alfonso F, Banning AP, Benedetto U, et al., 2018

628      ESC/EACTS Guidelines on myocardial revascularization. European heart journal. 2019;40(2):87-165

629   9.   ACC/AHA/SCAI Writing Committee. American college of cardiology/american heart association

630      task force on practice guidelines. Journal of the American College of Cardiology, 47(1):e1–e121,

631      2006.

632   10.  Liang-Chieh Chen and George Papandreou and Florian Schroff and Hartwig Adam, 2017.

633      "Rethinking Atrous Convolution for Semantic Image Segmentation". CoRR, abs/1706.05587.

634   11.  Ronneberger, Olaf; Fischer, Philipp; Brox, Thomas, 2015. "U-Net: Convolutional Networks for

635      Biomedical Image Segmentation". arXiv:1505.04597

636   12.  Noh, H.; Hong, S. & Han, B. Learning Deconvolution Network for Semantic Segmentation

637      Proceedings of the 2015 IEEE International Conference on Computer Vision (ICCV), IEEE

638      Computer Society, 2015, 1520-1528

639   13.  Vijay Badrinarayanan and Alex Kendall and Roberto Cipolla, 2015. SegNet: A Deep Convolutional

640      Encoder-Decoder Architecture for Image Segmentation. CoRR, abs/1511.00561.

641   14.  Carole H. Sudre, Wenqi Li, Tom Vercauteren, Sebastien Ourselin, and M. Jorge Cardoso.

642      Generalised Dice overlap as a deep learning loss function for highly unbalanced segmentations. In

643     Deep Learning in Medical Image Analysis and Multimodal Learning for Clinical Decision Support,

644     pages 240–248, 2017.A. Harris et al., "Free-space optical wavelength diversity scheme for fog

645     mitigation in a ground-to-unmanned-aerial-vehicle communications link," *Opt. Eng.* **45**(8), 086001

646     (2006) [doi:10.1117/1.2338565].

647 15. Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In International

648     Conference on Learning Representations (ICLR), 2015.

649 16. Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing

650     human-level performance on imagenet classification. 2015.

651 17. Yarin Gal, Riashat Islam, & Zoubin Ghahramani (2016). Deep Bayesian Active Learning with Image

652     Data. In *Bayesian Deep Learning workshop, NIPS*.

653 18. Microsoft Corporation. Microsoft cognitive toolkit, release 2.3, 2017.

654 19. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., … Chintala, S., 2019. PyTorch:

655     An Imperative Style, High-Performance Deep Learning Library. In Advances in Neural Information

656     Processing Systems 32 (pp. 8024–8035). Curran Associates

657 20. Zhao Wang, Jenkins, M. W., Linderman, G. C., Bezerra, H. G., Fujino, Y., Costa, M. A., Wilson, D.

658     L., & Rollins, A. M. (2015). 3-D Stent Detection in Intravascular OCT Using a Bayesian Network

659     and Graph Search. IEEE transactions on medical imaging, 34(7), 1549–1561

660 21. Hong Lu, Madhusudhana Gargesha, Zhao Wang, Daniel Chamie, Guilherme F. Attizzani, Tomoaki

661     Kanaya, Soumya Ray, Marco A. Costa, Andrew M. Rollins, Hiram G. Bezerra, and David L. Wilson,

662     "Automatic stent detection in intravascular OCT images using bagged decision trees," Biomed. Opt.

663     Express 3, 2809-2824 (2012)

664 22. Tsantis S, Kagadis GC, Katsanos K, Karnabatidis D, Bourantas G, Nikiforidis GC. Automatic vessel

665     lumen segmentation and stent strut detection in intravascular optical coherence tomography. Med

666     Phys. 2012 Jan;39(1):503-13. doi: 10.1118/1.3673067. PMID: 22225321.

667    23. Jouke Dijkstra, Gerhard Koning, Joan C Tuinenburg, Pranobe V Oemrawsingh, & Johan H.C Reiber

668        (2003). Automatic stent border detection in intravascular ultrasound images. International Congress

669        Series, 1256, 1111-1116.

670    24. Kitahara H, Kobayashi Y, Yock PG, Fitzgerald PJ, Honda Y. Deep learning-based intravascular

671        ultrasound segmentation for the assessment of coronary artery disease. Int J Cardiol. 2021 Jun

672        15;333:55-59. doi: 10.1016/j.ijcard.2021.03.020. Epub 2021 Mar 16. PMID: 33741429.

673    25. Iuga, AI., Carolus, H., Höink, A.J. *et al.* Automated detection and segmentation of thoracic lymph

674        nodes from CT using 3D foveal fully convolutional neural networks. *BMC Med Imaging* **21,** 69

675        (2021). https://doi.org/10.1186/s12880-021-00599-z

676    26. Simone Balocco, Francesco Ciompi, Juan Rigla, Xavier Carrillo, Josepa Mauri, Petia Radeva,

677        Chapter 10 - Computer-Aided Detection of Intracoronary Stent Location and Extension in

678        Intravascular Ultrasound Sequences, Editor(s): Simone Balocco, Intravascular Ultrasound, Elsevier,

679        2020, Pages 159-183, ISBN 9780128188330, https://doi.org/10.1016/B978-0-12-818833-0.00010-2.

680    27. Gatta C, Balocco S, Ciompi F, Hemetsberger R, Rodriguez Leor O, Radeva P. Real-time gating of

681        IVUS sequences based on motion blur analysis: method and quantitative validation. Med Image

682        Comput Comput Assist Interv. 2010;13(Pt 2):59-67. doi: 10.1007/978-3-642-15745-5_8. PMID:

683        20879299.

684    28. Francesco Ciompi, Rui Hua, Simone Balocco, Marina Albert, iOriol Pujol, Carles Caus, Josepa

685        Mauri, Petia Radeva (2013) Learning to Detect Stent Struts in Intravascular Ultrasound. In: Sanches

686        J.M., Micó L., Cardoso J.S. (eds) Pattern Recognition and Image Analysis. IbPRIA 2013. Lecture

687        Notes in Computer Science, vol 7887. Springer, Berlin, Heidelberg

688    29. Ciompi, F., Balocco, S., Rigla, J., Carrillo, X., Mauri, J. and Radeva, P. (2016), Computer-aided

689        detection of intracoronary stent in intravascular ultrasound sequences. Med. Phys., 43: 5616-5625.

690        https://doi.org/10.1118/1.4962927

691   30. Balocco, S., Ciompi, F., Rigla, J., Carrillo, X., Mauri, J. and Radeva, P. (2019), Assessment of

692       intracoronary stent location and extension in intravascular ultrasound sequences. Med. Phys., 46: 484-

693       493. https://doi.org/10.1002/mp.13273

694   31. Balocco, S.; Gatta, C.; Ciompi, F.; Wahle, A.; Radeva, P.; Carlier, S.; Unal, G.; Sanidas, E.; Mauri,

695       J.; Carillo, X.; Kovarnik, T.; Wang, C.-W.; Chen, H.-C.; Exarchos, T. P.; Fotiadis, D. I.; Destrempes,

696       F.; Cloutier, G.; Pujol, O.; Alberti, M.; Mendizabal-Ruiz, E. G.; Rivera, M.; Aksoy, T.; Downe, R. W.

697       & Kakadiaris, I. A., Standardized evaluation methodology and reference database for evaluating

698       IVUS image segmentation, Computerized Medical Imaging and Graphics , 2014, 38, 70 - 90

699   32. Katouzian, A, Angelini, ED, Carlier, SG, Suri, JS, Navab, N, Laine, AF (2012). A state-of-the-art

700       review on segmentation algorithms in intravascular ultrasound (IVUS) images. IEEE Trans Inf

701       Technol Biomed, 16, 5:823-34.

702   33. Bargsten L, Schlaefer A. SpeckleGAN: a generative adversarial network with an adaptive speckle

703       layer to augment limited training data for ultrasound image processing. Int J Comput Assist Radiol

704       Surg. 2020 Sep;15(9):1427-1436. doi: 10.1007/s11548-020-02203-1. Epub 2020 Jun 18. PMID:

705       32556953; PMCID: PMC7419454

706   34. Shengnan Liu, Tara Neleman, Eline M.J. Hartman, Jurgen M.R. Ligthart, Karen T. Witberg, Antonius

707       F.W. van der Steen, Jolanda J. Wentzel, Joost Daemen, & Gijs van Soest (2020). Automated

708       Quantitative Assessment of Coronary Calcification Using Intravascular Ultrasound. Ultrasound in

709       Medicine & Biology, 46(10), 2801-2809.

710   35. Bargsten, L., Riedl, K. A., Wissel, T., Brunner, F. J., Schaefers, K., Sprenger, J., Grass, M., Seiffert,

711       M., Blankenberg, S., & Schlaefer, A. (2021). Tailored methods for segmentation of intravascular

712       ultrasound images via convolutional neural networks. In N. V. Ruiter & B. C. Byram (Eds.), Medical

713       Imaging 2021: Ultrasonic Imaging and Tomography. SPIE. https://doi.org/10.1117/12.2580720

714   36. Bargsten, L., Riedl, K. A., Wissel, T., Brunner, F. J., Schaefers, K., Grass, M., Blankenberg, S.,

715       Seiffert, M & Schlaefer, A. (2021). Attention via Scattering Transforms for Segmentation of Small

716       Intravascular Ultrasound Data Sets. In Proceedings of Machine Learning Research 1–14, 2021.

717     37. Bargsten, L., Schlaefer, A. SpeckleGAN: a generative adversarial network with an adaptive speckle

718         layer to augment limited training data for ultrasound image processing. Int J CARS 15, 1427–1436

719         (2020). https://doi.org/10.1007/s11548-020-02203-1

720

721 **Tobias Wissel** is a research scientist at Philips Research Hamburg. He received his PhD in robotics

722 and cognitive systems from the University of Lübeck. His research interests include signal and

723 medical image processing as well as machine learning for diagnostics and image-guided therapy

724 with a focus on cardiovascular diseases.

725

726 **Katharina A. Riedl** is a physician at the Department of Cardiology at the University Heart &

727 Vascular Center Hamburg. She received her MD in ultra-high field MRI from the University of

728 Würzburg. Her clinical and research interest include intravascular imaging, cardiovascular

729 magnetic resonance imaging and coronary computed tomography with focus on coronary artery

730 disease.

731

732 **Klaus Schaefers** has studied information management and worked in several research and

733 development positions in organizations such as IBM and Fraunhofer. Since 2016 Klaus is working

734 with Philips Research as a software architect, where he focusses his work on big data and AI

735 related projects.

736

737 **Hannes Nickisch** is a senior research scientist at Philips Research Hamburg. He obtained a PhD

738 from Technical University Berlin and Max Planck Institute in probabilistic modeling. His research

739 interests include medical image analysis and biophysical simulations – in particular for

740 cardiovascular applications – as well as probabilistic machine learning methods.

741

742 **Fabian J Brunner** is a cardiologist at the Department of Cardiology at the University Heart &

743 Vascular Center Hamburg, Germany. His clinical and research interest include the prevention of

744 coronary artery disease as well as its interventional treatment.

745

746 **Nikolas D. Schnellbächer** is a research scientist at Philips Research Hamburg. He obtained a PhD

747 in Physics from Heidelberg University.

748 His research interests cover statistical physics, signal processing and machine learning and their

749 application for image processing and image reconstruction for medical imaging.

750

751 **Stefan Blankenberg** is a full Professor of Medicine. He is Director of the Clinic for Cardiology

752 and Medical Director at the University Heart & Vascular Center Hamburg.

753

754 **Moritz Seiffert** is head senior physician at the Department of Cardiology at the University Heart

755 & Vascular Center Hamburg. His clinical focus and research interests include complex coronary

756 interventions and percutaneous treatment of valvular heart disease.

757

758 **Michael Grass** is a Principal Scientist at Philips Research Hamburg and a lecturer on Medical

759 Imaging Systems at the Hamburg University of Technology, Germany. He holds a PhD in physics

760 and a Dr. habil. and Privatdozent in engineering for his work on Medical Imaging Systems.

761 Research interests cover diagnostic and interventional medical imaging, tomographic image

762 reconstruction, and artificial intelligence in medical imaging.

763

764

## List of Captions

**Figure 1: Illustration of the cascaded concept. Frames of the manual pullback are first analyzed by an encoder network, which decides for one of three classes per frame: stent, no stent or no use. Only stent frames are then passed on to the encoder-decoder to segment the stent struts. Apart from the favorable training setup, this is also targeting a reduction of false positive predictions on frames that do not show a stent anyway.**

**Figure 2: Network architectures for stent detection (top, encoder network) and stent segmentation (bottom, encoder-decoder network with skip connections between both parts).**

**Figure 3: Confusion matrix for the stent detection network after pooling the results on all five test folds.**

**Figure 4: Left: ROC curve for varying thresholds on the output probability maps. ROC-AUC values are listed for a threshold of t = 0.5. Right: Dice curves showing the dependence of the Dice score on the chosen threshold. A good compromise is achieved when thresholding all classes at 0.5.**

**Figure 5: Example pullback along with encoder results. First two rows show orthogonal cuts through the pullback in longitudinal direction and the lower three graphs the ground truth (dash-dotted red) and predictions (solid blue line Monte-Carlo mean and blue shading Monte-Carlo estimate of the standard estimation) Regions where none of the three classes has a red ground truth value at 1 have not been annotated by the expert. Here, predictions cannot be compared with a target label. During the last frames the transducer was covered by the catheter, which is correctly recognized as "no use".**

**Figure 6: Example frames labeled with their corresponding predictions from the stent output of the detector network: true positive decisions (top two rows, typical to challenging examples from left to right), false negative decisions (detector missed ground truth annotations, third row), and false positive detections (bottom row, wrong predictions without ground truth label).**

**Figure 7: Sorted rank plots for the precision score computed per-frame. Dashed lines indicate average metric on the intra-observer variance set. Left: Precision ranking for all samples with annotated ground truth. Right: Precision ranking for all frames containing automatic segmentations.**

**Figure 8: Sorted rank plots for recall and Dice scores computed per-frame. Dashed lines indicate average metric on the intra-observer variance set. Left: Recall ranking for all frames with annotated ground truth. Right: Dice score ranking for all frames with annotated ground truth.**

**Figure 9: Score matrices for the segmentation network. Scores are presented for different frame supports and post processing (pp) steps: no post-processing (no-pp), segmentation rejection based on a "< N pixel" threshold (Npx) and based on detector decisions (detect). Left: Dice scores for three different post-processing scenarios (vertical) on different frame supports (horizontal). The effect of the detector is evaluated on all frames (segmentations are corrected based on detector decision) and on frames with positive detector decision for stent only. Right: False positive rates (fpr) (listed as average number of pixels in a 224×224 frame) for four different post-processing scenarios (vertical) on different frame supports (horizontal). A false positive rate of 0.25% or 125 pixels (0.249 mm²) roughly corresponds to the area of one stent strut as pictured by the imaging modality.**

**Figure 10: Angular dice (DCA) and symmetric median skeleton distance (SMSD, in brackets) scores of the segmentation network. The former is reported in percent and the latter in pixel, i.e. multiples of the image resolution 0.0446mm. Equivalent to the dice scores in Figure 9 results are shown without post-processing (no-pp) and after applying the 400px or 500px heuristic. Again, we evaluate on all frames, stent-only frames as well as on the two cascading scenarios**

**Figure 11: Example frames drawn in steps of constant frame proportions from the score rankings. Therefore, the plots provide a representative selection of frames covering the full range of scores achieved per frame (scores decrease from left to right). Green: true positive pixels, Yellow: false negative pixels, Red: false positive pixels. Top: Frames drawn at equal spacings from the Dice rankings (Figure 7, right). Bottom: Frames drawn at equal spacings from the precision rankings (Figure 8, right).**