
Encoder-Decoder Architectures for Clinically Relevant Coronary Artery Segmentation

João Lourenço Silva¹ Miguel Nobre Menezes² Tiago Rodrigues² Beatriz Silva²

Fausto J. Pinto²

Arlindo L. Oliveira¹

¹INESC-ID / Instituto Superior Técnico, University of Lisbon
{joao.lourenco.silva, arlindo.oliveira}@tecnico.pt

²Cardiology Department, CAML, CCUL, Lisbon School of Medicine, University of Lisbon
{mmenezes.gm, tiagoerodrigues, beatrizsilvae}@gmail.com
faustopinto@medicina.ulisboa.pt

Abstract

Coronary X-ray angiography is a crucial clinical procedure for the diagnosis and treatment of coronary artery disease, which accounts for roughly 16% of global deaths every year. However, the images acquired in these procedures have low resolution and poor contrast, making lesion detection and assessment challenging. Accurate coronary artery segmentation not only helps mitigate these problems, but also allows the extraction of relevant anatomical features for further analysis by quantitative methods. Although automated segmentation of coronary arteries has been proposed before, previous approaches have used non-optimal segmentation criteria, leading to less useful results. Most methods either segment only the major vessel, discarding important information from the remaining ones, or segment the whole coronary tree based mostly on contrast information, producing a noisy output that includes vessels that are not relevant for diagnosis. We adopt a better-suited clinical criterion and segment vessels according to their clinical relevance. Additionally, we simultaneously perform catheter segmentation, which may be useful for diagnosis due to the scale factor provided by the catheter's known diameter, and is a task that has not yet been performed with good results. To derive the optimal approach, we conducted an extensive comparative study of encoder-decoder architectures trained on a combination of focal loss and a variant of generalized dice loss. Based on the EfficientNet and the UNet++ architectures, we propose a line of efficient and high-performance segmentation models using a new decoder architecture, the EfficientUNet++, whose best-performing version achieved average dice scores of 0.8904 and 0.7526 for the artery and catheter classes, respectively, and an average generalized dice score of 0.9234. Source code will be available at: <https://github.com/jlcsilva/EfficientUNetPlusPlus>.

1 Introduction

Coronary arteries are the blood vessels that carry blood to the heart. Coronary Artery Disease (CAD), also known as Coronary Heart Disease (CHD) or Ischemic Heart Disease (IHD), is the narrowing or blockage of these arteries, caused by the build-up of atherosclerotic plaques inside them, which can lead to limited blood flow and consequent damage to the heart muscle. CAD is the cause of roughly 16% of global deaths every year [1].

X-ray coronary angiography (CAG) is one of the main procedures for CAD diagnosis and treatment. Patients submitted to CAG are catheterized and have their arteries filled with a radio-opaque contrast agent that makes them visible in X-ray images. Traditionally, physicians used these images to assess the presence and severity of stenosis (i.e., artery narrowing) through visual inspection. However, this method is highly subjective and potentially unreliable, which led to the development of Quantitative Coronary Angiography (QCA), a diagnostic support tool. Using semi-automatic edge-detection algorithms, QCA reports the vessel diameter at user-specified locations and the point of stenosis. Nevertheless, the low contrast and resolution of CAG images, the uneven contrast agent distribution and the presence of artifacts such as pacemakers, the spine and the catheter itself make this task very challenging. Thus, QCA still requires manual correction of the vessel boundaries before calculating the stenosis percentage, limiting its use in clinical practice. Indeed, in everyday practice, the severity of stenosis is still assessed visually in most cases, rather than with QCA software.

Recently, coronary artery segmentation performance in CAG images has been significantly improved by deep learning methods. Most of them either segment only the major vessel or try to segment the whole coronary tree based primarily on contrast differences. The procedures, however, may not be clinically optimal. The former discards potentially damaged vessels whose lesions may not be negligible, and the latter includes secondary vessels that may not be relevant for either diagnostic or therapeutic purposes, potentially distracting physicians from the important ones. We circumvent these shortcomings by adopting a better-suited clinical criterion, with the aid of expert physicians, in which a vessel is only segmented if it is important to properly assess other vessels' lesions or contains a non-negligible lesion itself. More specifically, we only segment arteries no smaller than 2 mm at their origin, as thinner vessels are generally deemed inadequate for revascularization.

Additionally, with more complex lesion assessment and anatomical feature extraction in mind, we simultaneously segment the catheter, whose known diameter provides a scale factor that may come to play an important role in current and future diagnostic methods. To the best of our knowledge, simultaneous catheter and coronary artery segmentation in CAG images has only been performed in one previous work [2], reporting dice score coefficients (DSC) of 0.54 and 0.69 for the artery and catheter classes, respectively. We obtained DSCs of 0.8904 and 0.7526 for these same classes.

Due to the adopted segmentation criteria and the inclusion of a catheter class, the problem we aim to solve is more complex than those addressed in previous work. Instead of simply learning to distinguish between vessels - or main vessels, easily identifiable by their larger width - and background, our model must also learn to segment catheters and determine the vessels' clinical relevance. To address these challenges, we adopt a loss function that combines a variant of Generalized Dice Loss (GDL) [3, 4] and Focal Loss (FL) [5]. The former provides global segmentation quality information and is designed to handle class imbalance. The latter uses pixel-wise information to force models to focus on hard, misclassified examples, leading to improved segmentation of less common classes and challenging parts of images, such as class boundaries.

To determine the best architecture for this segmentation task, we conducted an extensive comparative study of existing encoders and decoders, which provided insights into the best architectural patterns for this and, presumably, other medical image segmentation problems. Based on these findings, we propose a new efficient and high-performance segmentation architecture, the EfficientUNet++, based on the EfficientNet family of models [6] and the UNet++ [7] decoder architecture. Using this architecture, we achieved DSC of 0.8904 and 0.7526 for the artery and catheter classes, respectively, and a generalized dice score of 0.9234.

Overall, the main contributions of this paper are as follows:

1. We propose an approach to perform simultaneous coronary artery and catheter segmentation in CAG images, using a new and better-suited clinical criterion, in which vessels are only labeled as such if they are deemed relevant for diagnostic and therapeutic purposes;
2. We perform an extensive quantitative and qualitative comparison of the performance of existing encoders and decoders, which may provide valuable insights for other medical image segmentation tasks;
3. Based on the findings of our study, we propose a line of efficient and high-performance segmentation models based on EfficientNet backbones and the new EfficientUNet++ decoder architecture, enabling practitioners to choose a trade-off between model size and model performance, according to the available hardware and the clinical needs.

2 Related Work

Major vessel segmentation. Previous work has shown that major vessel segmentation can be improved by replacing the U-Net’s encoder with popular image classification backbones, either pre-trained on ImageNet [4, 8] or trained from scratch on a relatively small dataset composed of 3200 CAG images [9]. Additionally, it has also been shown that the use of a modified generalized dice loss function that uses weights to offset class imbalance and features a tunable penalty for false positives and false negatives could further improve segmentation performance [4]. In the sequence of these findings, we train our models using a combination of the proposed loss function and focal loss and compare their segmentation quality when using different state-of-the-art encoders.

Other authors have proposed a U-Net-based nested encoder-decoder architecture, named T-Net [10]. To simplify the optimization process, the authors replaced the U-Net’s blocks with residual ones. In addition, to enable feature reuse, they arranged the pooling and up-sampling operations to make all the feature maps extracted by the encoder available to every layer of equal or greater depth of the decoder, in a DenseNet-like fashion. These modifications enhanced information flow through the network and enabled it to outperform a standard U-Net. The use of dense connections is also present in the U-Net++ [7], which we explore in our work.

Whole vessel body segmentation. One of the main challenges of the coronary artery segmentation task is the discrimination between vessels and artifacts. Given that arteries are only visible in the presence of contrast, previous work has used the images acquired before contrast injection as a second-channel input to help a U-Net discern between the vessels and the background [11]. However, for this approach to be effective, it must be coupled with an image alignment algorithm to compensate for the motion caused by heartbeat and respiration. Furthermore, it requires the entire angiographic sequence to be acquired with minimal table motion, which can be hard to achieve, as standard clinical practice involves moving the patient table to follow the flow of dye within the vessels. For this reason, we do not use this technique in our work.

In line with what we propose in this paper, some authors have also attempted to use different segmentation architectures to achieve better performance than what is possible with the commonly used U-Net. Specifically, they proposed using a pre-trained PSPNet [12] and a U-Net++ combined with a feature pyramid network to improve multi-scale feature detection [13]. Other proposals include a fully convolutional encoder-decoder network specifically designed for vessel segmentation, featuring Gaussian convolutions and trained with deep supervision [14].

Other authors have proposed to use a multi-layer perceptron to produce segmentation masks based on the multi-scale features extracted by multi-scale Gaussian matched and Gabor filters [15, 16], with one of them achieving state-of-the-art performance on the coronary tree segmentation task [15]. However, this method does not suit our needs, as the task we aim to solve requires the use of higher-level features that make it possible to discern between catheter and vessels and determine the clinical relevance of each vessel. More specifically, we only segment arteries no smaller than 2 mm at their origin, as thinner vessels are generally deemed inadequate for revascularization.

A number of proposals [17, 18, 19] have explored the use of temporal information to complement the 2D data of the frame to be segmented. Some have adopted simple approaches such as replacing the input blocks of 2D segmentation models with 3D blocks capable of processing inputs composed of multiple frames [18] or using a temporal filter to compute the weighted average of successive output frames to mitigate the presence of artifacts and inter-frame flicker. More complex alternatives have also been studied, such as the use of a U-Net-based architecture with a 3D encoder and a 2D decoder [19]. Regardless of the strategy adopted, the use of temporal information led to quantitative improvements in segmentation performance. Nevertheless, it has not been experimentally confirmed that the obtained results are clinically more relevant. Due to the uneven distribution of contrast during the angiographic sequence and the motion caused by respiration and heartbeat, it is usual for some parts of the vessels and lesions to be more visible in some frames than others. Thus, temporal context may help mitigate false positives caused by the presence of artifacts and false negatives induced by hard to perceive vessels. However, there is also a non-negligible risk that it may lead some lesions to go unnoticed if the models attribute more importance to frames in which lesions are not perceivable than those in which they are evident. Due to this concern, we consider that the use of temporal context should be subject to in-depth analysis and leave its study for future work.

Orthogonal work has featured a coarse-to-fine approach in which the segmentation mask produced by a Fully Convolutional Network (FCN) is sparsely enhanced by the use of a region-based segmentation model [20]. Other authors have used a dual-path sliding-window CNN to combine local information and global context from two different sized patches to produce a vessel probability map, which they then refined using a similar network and an edge map computed from the angiogram [21]. Even though both methods perform well on the coronary tree segmentation task, their reliance on local information makes them unsuitable for our task, which requires global context to determine the clinical relevance of each vessel.

To address the scarcity of labeled data and alleviate the burden on annotators, previous work has introduced a framework for weakly supervised learning that allows generic segmentation models to learn from pseudo labels generated by automatic vessel enhancement [22]. The model is trained using a self-paced scheme, in which it learns from pixels in descending difficulty order. In each iteration, according to the combined uncertainty of the model and the vessel enhancement algorithm for each part of the image, local online manual annotation refinement is performed. This procedure allows the pseudo labels to be improved and the model to learn well without as much effort from the annotators as in fully supervised methods. We do not use weakly nor semi-supervised learning in this work, but it could be a future line of research, as it addresses label scarcity, one of the main issues in medical image segmentation tasks.

Catheter and coronary tree segmentation. To the best of our knowledge, simultaneous catheter and coronary tree segmentation has only been addressed in a single previous work, in which a U-Net-based Siamese architecture was trained with automatically generated annotations [2]. Using noisy low-level binary segmentation and optical flow, the authors generated multi-class annotations that were successively improved in a multistage segmentation approach. While this task is similar to ours, in the sense that it performs both catheter and artery segmentation, it does not address the clinical relevance issue and aims to segment every vessel of the coronary tree.

Other segmentation criteria. To the best of our knowledge, an intermediate segmentation criterion, in which secondary arteries with diameter inferior to 1 mm at their origin are not segmented, has been proposed only once [23]. In our work, we use a similar criterion but only segment vessels no smaller than 2 mm at their origin, as thinner vessels are generally deemed inadequate for revascularization. To perform this task, the authors couple an Angiographic Processing Network (APN), trained to learn the best possible preprocessing filter to improve segmentation performance, with the U-Net and DeepLabV3+ [24] architectures. The results show that the APN module helps both models to achieve better results. While the APN and preprocessing, in general, can improve segmentation performance in multiple tasks, including the one we are addressing, in this paper we focus on the segmentation network’s architecture and leave the study of preprocessing methods to future work.

3 Architecture

Given the profusion of high-performing segmentation models, we decided to conduct an extensive comparative study of existing encoders and decoders to determine the best architectures for this task and what makes them perform better than the others. In this section, we describe the metrics used for model comparison and the experiments we conducted. As a result of our study, we propose a new efficient and high-performing architecture, the EfficientUNet++, based on the EfficientNet image classification models and the UNet++ decoder.

3.1 Evaluation Metrics

We evaluate the segmentation quality of each class using DSC, precision and recall. The overall segmentation quality of each image is evaluated using DSC and generalized dice score (GDS).

Let C be the number of classes, N be the number of pixels, G be the one-hot encoded ground-truth, with $g_{ci} \in \{0, 1\}$ denoting whether pixel i belongs to class c or not, and P be the predicted probabilistic map, with $p_{ci} \in [0, 1]$ representing the probability of pixel i belonging to class c . Then,

$$GDS = 2 \frac{\sum_{c=1}^C w_c \sum_{i=1}^N g_{ci} p_{ci}}{\sum_{c=1}^C w_c \sum_{i=1}^N g_{ci} + p_{ci}}, \tag{1}$$

where w_c is the weight assigned to class c . When all weights are set to 1, GDS is equivalent to DSC. By setting all weights to $w_c = 1/(\sum_{i=1}^N g_{ci})^2$, each class’s contribution to the score is corrected by the inverse of its area, reducing the correlation between region size and dice score [25]. Consequently, GDS attributes the same importance to all classes independently of their frequency, making it a fairer metric for multi-class segmentation performance than DSC.

3.2 Loss Function

The problem we aim to solve can be interpreted as the combination of two sub-tasks: a macro-level and a micro-level one. The former consists in identifying the vessels and catheters, distinguishing them from each other and from artifacts, and determining which arteries are clinically relevant. The latter concerns the precise delineation of class contours, which is crucial for a reliable diagnosis based on the produced segmentation masks.

To address these tasks, we propose the use of a loss function composed of a variant of generalized dice loss [3, 4], pGDL, which provides information on global segmentation quality, and focal loss [5], FL, which provides a pixel-wise evaluation focused on the harder pixels, which are usually the ones belonging to less common classes, near class boundaries and in artifact regions. The loss function used can be defined as:

$$Loss = pGDL + \lambda FL, \quad (2)$$

where λ is a hyperparameter that controls the weight given to each term. For simplicity, we attribute the same importance to both terms and use $\lambda = 1$ in all our experiments.

Generalized Dice Loss. As its name suggests, GDL is an extension of the dice loss (DL) function. By attributing weights to the segmentation classes, it provides invariance to different label set properties. Using the notation defined in Section 3.1, GDL takes the form:

$$GDL = 1 - GDS = 1 - 2 \frac{\sum_{c=1}^C w_c \sum_{i=1}^N g_{ci} p_{ci}}{\sum_{c=1}^C w_c \sum_{i=1}^N g_{ci} + p_{ci}}, \quad (3)$$

where GDS is the function defined in Eq. 1. When all class weights, w_c , are set to 1, GDL is equivalent to DL. We adopt an even more generalized version of the DL function, pGDL, which adds a penalty for false positives and false negatives to GDL [4]. It is defined by:

$$pGDL = \frac{GDL}{1 + k(1 - GDL)}, \quad (4)$$

where k is a hyperparameter that controls the weight of this additional penalty. When k is set to 0, pGDL is equivalent to GDL. In this work, we set k to 0.75, as we empirically verified that this value worked well for all models and led to better performance than $k = 0$, for most of them. Concerning the class weights, we follow Sudre et al. [3] and set them to $w_c = 1/(\sum_{i=1}^N g_{ci})^2$. This weighting corrects the contribution of each class by the inverse of its area [25], reducing the correlation between region size and dice score and consequently mitigating the negative effects of class imbalance.

Focal Loss. Focal loss is based on the cross-entropy (CE) loss function and was originally designed for one-stage object detection scenarios with extreme imbalance between foreground and background classes [5]. For notational convenience and based on the notation above, we define:

$$p'_{ci} = \begin{cases} p_{ci} & \text{if } g_{ci} = 1 \\ 1 - p_{ci} & \text{if } g_{ci} = 0 \end{cases} \quad (5)$$

and rewrite the CE loss function as $CE(p'_{ci}) = -\log(p'_{ci})$. Thus, FL takes the form:

$$FL(p'_{ci}) = -\alpha'(1 - p'_{ci})^\gamma \log(p'_{ci}), \quad (6)$$

where $\alpha \in [0, 1]$ is a weighting factor that balances positive and negative examples, with α' defined analogously to p'_{ci} , and $(1 - p'_{ci})^\gamma$ is a modulating factor, controlled by $\gamma \geq 0$, that down-weights easy examples and forces the model to focus on and learn from hard ones. When $\alpha' = 1$ and $\gamma = 0$, FL is equivalent to CE loss. As γ increases, so does the importance given to hard examples compared to easy ones. For simplicity, we followed the FL function’s authors [5] and used $\gamma = 2$ and $\alpha = 0.25$, as these values allowed to obtain good results.

3.3 Encoder Comparison

Previous work [4, 8, 9] has shown that the U-Net’s performance can be enhanced by replacing its encoder with more sophisticated image classification architectures, both when using transfer learning from a large dataset, like ImageNet [4, 8], and when training from scratch on a relatively small dataset composed of 3200 XRA images [9], suggesting that the U-Net’s original backbone is too rudimentary and not an as good feature extractor as those designed for image classification.

On this assumption, we trained multiple models to perform our segmentation task, using image classification architectures pre-trained on ImageNet as encoders. To avoid overfitting the encoders to our small dataset and to evaluate the quality of the visual representations learned from ImageNet, their weights were frozen during decoder training, with the additional benefit of shortening the networks’ training time. Furthermore, to investigate the existence of synergies between certain encoder-decoder pairs, we trained each backbone with multiple decoders. In particular, we used the U-Net [26], commonly used for medical image segmentation, the UNet++ [7], which has been shown by its authors to outperform the U-Net in multiple medical image segmentation tasks, and the DeepLabV3+ [24], a state-of-the-art semantic segmentation architecture.

Figures 1a, 1b and 1c show the measured segmentation performances as a function of the total number of FLOPS, when using encoders from the RegNetY [27], ResNeXt [28] and ResNet [29] families. Graphs of performance as a function of the number of parameters are included in the appendix. Notably, for every decoder, the best performance is achieved using an EfficientNet backbone. Furthermore, for the same performance, EfficientNet encoders are always more efficient than other backbones, parameter-wise but especially computation-wise. Due to their compound scaling, EfficientNets are generally thinner than other encoders at each scale, i.e., use fewer channels to represent information. Thus, as decoder computation scales with feature map dimension, EfficientNets allow building much more computationally efficient systems than wider encoders. This is particularly evident when using complex decoders, such as the UNet++, that heavily process extracted features.

We observe that, in general, better image classification architectures allow higher segmentation performance, leading us to the intuitive and widely accepted premise that image classification performance is correlated with feature extraction capabilities. However, we notice that, for some decoders and encoder families, segmentation performance starts degrading after a certain encoder complexity is surpassed. As the encoder weights are not updated, and the same decoders converge to better solutions when using simpler encoders of the same family, we conclude that the source of degradation is decoder overfitting. Degradation is especially noticeable when using the UNet++ decoder. We hypothesize that, combined with their heavy processing of extracted features, their high number of parameters leads them to overfit the training set. The same happens for the U-Net, but not as noticeably, as it has fewer parameters and performs less processing of extracted features. Decoder overfitting could probably be mitigated or even solved using regularization techniques, a larger dataset than the one we use, composed of 270 CAG, or both. Interestingly, EfficientNet encoders seem to have a regularizing effect on decoders, suggesting they generalize and transfer better to new tasks than other models, which can be related with their thinner feature maps.

Besides all these advantages, the low number of channels used by EfficientNets to represent extracted features improves memory efficiency, allowing relatively larger training batches, which can be very important for researchers with limited hardware resources.

3.4 Decoder Comparison

Given the EfficientNet models’ superior performance and efficiency, we used them as encoders for decoder comparison experiments, with each decoder being trained using all EfficientNet backbones. Figure 1d shows the segmentation performance as a function of the number of FLOPS for the DeepLabV3+ [24], FPN [30], LinkNet [31], MANet [32], PAN [33], PSPNet [12], ResUNet [34],

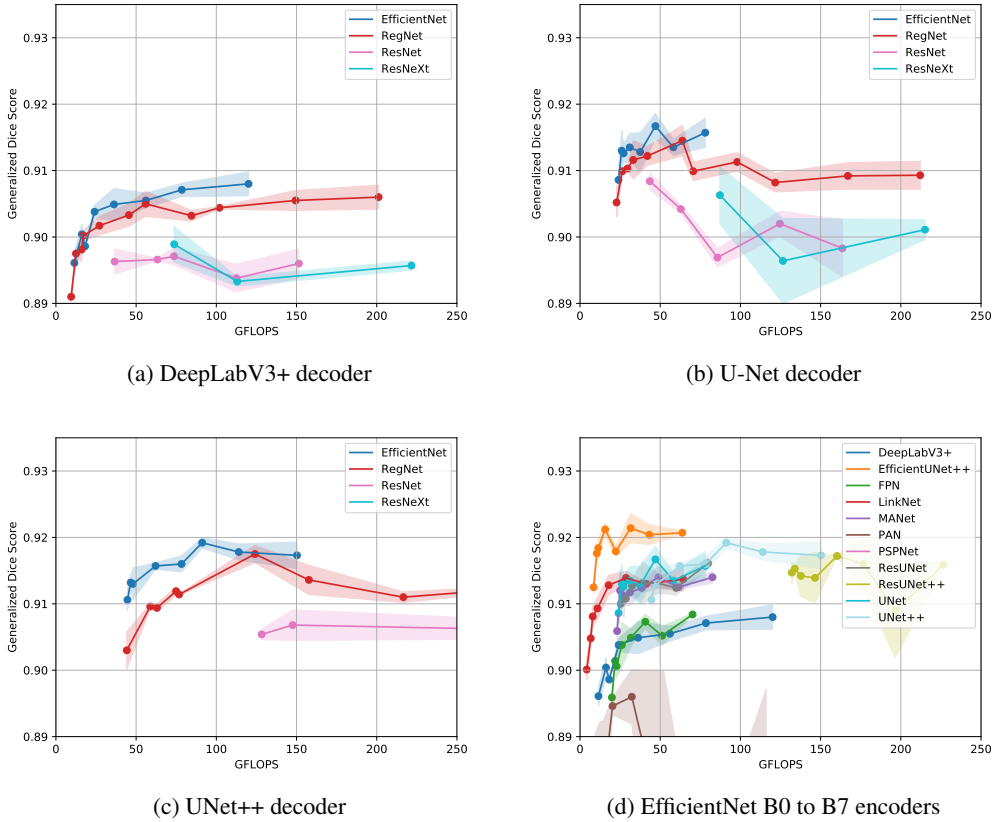


Figure 1: Segmentation performance, measured by generalized dice score (GDS), as a function of the number of FLOPS. Figures (a), (b) and (c) show the performance of different encoders combined with the (a) DeepLabV3+, (b) U-Net and (c) UNet++ decoders. Figure (d) shows the performance of different decoders combined with the EfficientNet B0 to B7 encoders. Each polygonal line corresponds to an encoder family. The dots represent the following models, in ascending order of FLOPS: EfficientNet - B0, B1, B2, B3, B4, B5, B6, B7; RegNet - Y2, Y4, Y6, Y8, Y16, Y32, Y40, Y64, Y80, Y120, Y160; ResNet - 18, 34, 50, 101, 152; ResNeXt - 50_32x4d, 101_32x4d, 101_32x8d. Above 250 GFLOPS, performance keeps degrading and is omitted. Models with GDS below 0.89 are also omitted.

ResUNet++ [35], U-Net [26] and UNet++ [7] decoder architectures. Graphs of performance as a function of the number of parameters are included in the appendix.

The results of our experiments confirm the importance of the skip connections between encoder and decoder featured in U-Net-based models. The UNet++ and ResUNet++ achieve the best performances among all models, and the LinkNet, MANet and ResUNet achieve good results, similar to the U-Net's. However, the PAN decoder also uses skip connections but has poor performance, suggesting that the attention mechanisms it uses in its blocks and bridge module are prejudicial for this task.

The role played by attention mechanisms is not very clear. While they seem to be what harms the PAN's performance, in the ResUNet++ they appear to be beneficial, and in the MANet they do not seem to have any effect, as it performs very similarly to the U-Net in which it is based. The importance of residual connections is also unclear, as they reduce the ResUNet's performance compared to the U-Net but work well in the ResUNet++. Given that attention mechanisms and residual connections alone do not seem to improve performance in the MANet and ResUNet, respectively, we hypothesize that it is the combination of both that allows the ResUNet++ to perform so well. However, confirming this theory would require an in-depth comparison of the attention mechanisms used in the ResUNet++, MANet and PAN, and their synergies with residual connections, which we leave for future work.

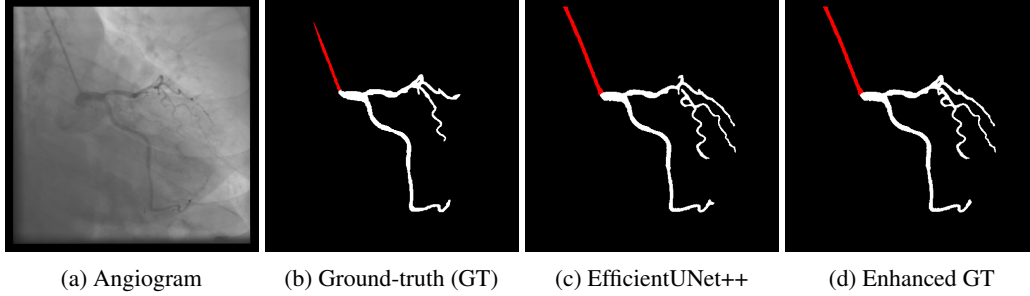


Figure 2: In some cases, the segmentation model produces better segmentation masks than humans. Here, the automatic segmentation mask, (c), was used by physicians to enhance the original ground-truth, (b), and produce the mask shown in (d).

Interestingly, the UNet++, which performs similarly to the ResUNet++ but more efficiently, both parameter and computation-wise, does not use residual connections nor attention. Instead, it uses densely connected nested decoder sub-networks, which promote feature reuse and allow it to extract more information from the encoder’s feature maps at each scale.

Architectures based on pyramid pooling and feature pyramids are the worst-performing ones. Due to the lack of skip connections to provide accurate localization information, these models produce coarser segmentation maps than those of the U-Net-based models. While this allows them to perform well on generic segmentation tasks, it harms their performance when applied to medical images, which require fine segmentation. This is the case of the DeepLabV3+, which uses atrous convolutions and an atrous spatial pyramid pooling (ASPP) module in the encoder, the FPN, whose segmentation is based on feature pyramids, and the PSPNet, which obtains local information by applying pyramid pooling to the output of an encoder with dilated convolutions.

3.5 EfficientUNet++ Architecture

When coupled with EfficientNet backbones, the UNet++ achieves high segmentation performance with reasonable parameter and computational efficiency. However, while the number of parameters is not a major concern, the computation required for inference can be prohibitive of widespread clinical usage, as it requires expensive hardware to be run promptly for entire angiographic sequences, usually comprised of about a hundred frames.

To address this, we propose a new architecture, the EfficientUNet++, that reduces computational complexity by replacing the UNet++’s blocks with residual bottlenecks with depthwise convolutions. Furthermore, to enhance performance, we process the bottleneck feature maps with concurrent spatial and channel squeeze and excitation (scSE) blocks [36], which combine the channel attention of squeeze and excitation (SE) blocks [37] with spatial attention.

As shown in Figure 1d, when combined with EfficientNet encoders, our decoder architecture establishes a line of efficient and high-performance segmentation models, which obtained a maximum average generalized dice score of 0.9239 when using the EfficientNetB5 backbone. For the artery and catheter classes, this model obtained average dice scores of 0.8904 and 0.7526, respectively.

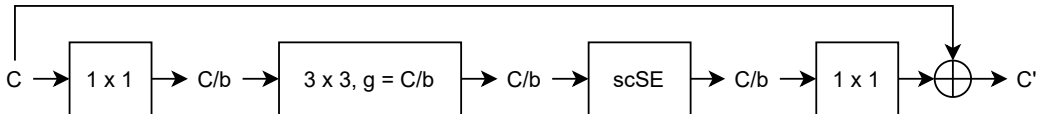


Figure 3: EfficientUNet++’s convolutional block. Each convolution is followed by batch normalization [38] and Hardswish [39], except for the last 1×1 convolution, which is not activated. C and C' are the numbers of input and output channels. Feature map height and width are not altered. We set the bottleneck ratio, b , to 1, and the number of convolution groups, g , to be equal to the number of input channels, making the 3×3 convolution depthwise. The scSE block uses a squeeze ratio of 1.

4 Implementation Details

Training details. We used encoders pre-trained on ImageNet and froze their weights during decoder training. Kaiming [40] and Xavier initialization [41] were used to initialize decoder weights in hidden and output layers, respectively. The models were trained on 4 Tesla V100S GPUs for 150 epochs each, using Adam with $\beta_1 = 0.9$, $\beta_2 = 0.999$, a mini-batch size of 8, no weight decay, and an initial learning rate of 0.001. The learning rate was divided by 10 at the 50th and 100th epochs.

In most experiments, we used public PyTorch implementations [42] under an MIT license. To keep comparisons fair, we extended the repository with ResUNet, ResUNet++ and EfficientUNet++ implementations. To obtain average scores and standard deviations, each model was trained and tested three times.

Dataset. Our dataset is exclusively composed of retrospective data, whose use was approved by the ethics committee. It comprises 270 anonymized 512×512 pixels CAG images annotated by three expert cardiologists. The images were acquired from multiple viewing angles of the left (LCA) and right coronary arteries (RCA) of 47 patients. The dataset was split into a training and a test set, composed of 237 and 33 images, respectively. The images of each set were carefully chosen to keep them representative of the original one, having approximately identical distributions regarding the observed arteries, viewing angles and number of images annotated by each physician. In total, the training set and test set contain 165/72 and 23/10 images of the LCA/RCA, respectively.

Data augmentation. Our image augmentation policy consists of the sequential application of the following transformations: 1) rotation of the image at a random angle between $\pm 20^\circ$, as this is the approximate imaging angle range; 2) horizontal and vertical shifts at random rates within $\pm 10\%$; 3) zoom at a random rate between -10% and 10% ; 4) brightness change at a random ratio within $\pm 40\%$, to account for the variability of the angiography images' brightness across acquisition devices.

Data augmentation is performed online. Each time an image is sampled in a batch, three augmentations are created, increasing the number of images seen per epoch by a factor of 4. Due to the near-infinite number of possible distinct augmentations, it is very unlikely for the model to see the same augmentation more than once. Thus, the number of different images seen by the model rises by a factor equal to the number of training epochs compared to offline data augmentation.

5 Discussion and Future Work

In this work, we propose a new and better-suited clinical criterion for simultaneous catheter and artery segmentation in CAG images. Whereas most previous approaches either segment only the major vessel or the whole coronary tree, based mostly on contrast information, we segment arteries according to their clinical relevance, as assessed by expert physicians. To derive the best approach for this task, we compared multiple encoder and decoder architectures, which we trained on a combination of focal loss and a variant of generalized dice loss, which aggregates information on global and pixel-wise segmentation quality, handles class imbalance, and forces models to focus on hard pixels in class boundary and artifact regions.

Among the existing architectures, we found the EfficientNet and the UNet++ to be the best encoder and decoder, respectively. Due to their compound scaling, EfficientNets are efficient not only in terms of parameters and computation, but also in the way they represent features, generally using fewer channels at each scale than other models, with the latter's benefit being threefold: 1) it requires less computation from decoders; 2) it seems to have a regularizing effect; 3) it slightly reduces memory use during training. Based on the EfficientNet and the UNet++, we propose a new decoder architecture, the EfficientUNet++. By replacing the UNet++'s blocks residual bottlenecks with depthwise convolutions and using scSE spatial and channel attention blocks, our architecture outperforms all other we tested and is much more computationally efficient than the UNet++.

This paper gives origin to multiple directions of future work. Regarding model architecture, it would be interesting to study the impacts of different attention mechanisms in segmentation performance. The findings of such research could then be used to build models even better than the EfficientUNet++. Another promising line of research is the application of semi-supervised learning to this and other medical image segmentation tasks to take advantage of the great amount of existing unlabeled data.

Broader Impact

This paper presents a new and better-suited clinical criterion for segmentation of coronary X-ray angiography images and proves that current state-of-the-art encoder and decoder architectures can achieve high levels of performance in this task. With these findings, we hope to open the way for a new line of research on clinically relevant coronary artery segmentation, which has not yet been explored and whose results may be of great usefulness for clinical practice.

Additionally, to derive the best approach for this task, we conducted an extensive comparative study of existing encoder and decoder architectures, whose findings will hopefully be useful for other medical image segmentation tasks. As a result of this study, we propose a new segmentation architecture, the EfficientUNet++, that combines the EfficientNet encoders with a more efficient and better-performing version of the UNet++ decoder to build a line of efficient and high-performance segmentation models. By providing all the code and implementation details, we hope other researchers and engineers will further improve it and use it in their segmentation tasks.

Having been designed to be used as a support diagnostic tool, any possible failure of this system will be safeguarded by existing diagnostic methods and will not put patients' lives at risk.

Our method does not leverage biases in the data and, to the best of our knowledge, no one will be put at disadvantage by the results of this research.

References

- [1] K. E. Rudd, S. C. Johnson, K. M. Agesa, K. A. Shackelford, D. Tsoi, D. R. Kievlan, D. V. Colombara, K. S. Ikuta, N. Kissoon, S. Finfer *et al.*, "Global, regional, and national sepsis incidence and mortality, 1990–2017: analysis for the global burden of disease study," *The Lancet*, vol. 395, no. 10219, pp. 200–211, 2020.
- [2] A. Vlontzos and K. Mikolajczyk, "Deep segmentation and registration in x-ray angiography video," *arXiv preprint arXiv:1805.06406*, 2018.
- [3] C. H. Sudre, W. Li, T. Vercauteren, S. Ourselin, and M. J. Cardoso, "Generalised dice overlap as a deep learning loss function for highly unbalanced segmentations," in *Deep learning in medical image analysis and multimodal learning for clinical decision support*. Springer, 2017, pp. 240–248.
- [4] S. Yang, J. Kweon, and Y.-H. Kim, "Major vessel segmentation on x-ray coronary angiography using deep networks with a novel penalty loss function," in *International Conference on Medical Imaging with Deep Learning—Extended Abstract Track*, 2019.
- [5] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2980–2988.
- [6] M. Tan and Q. Le, "Efficientnet: Rethinking model scaling for convolutional neural networks," in *International Conference on Machine Learning*. PMLR, 2019, pp. 6105–6114.
- [7] Z. Zhou, M. M. R. Siddiquee, N. Tajbakhsh, and J. Liang, "Unet++: A nested u-net architecture for medical image segmentation," in *Deep learning in medical image analysis and multimodal learning for clinical decision support*. Springer, 2018, pp. 3–11.
- [8] S. Yang, J. Kweon, J.-H. Roh, J.-H. Lee, H. Kang, L.-J. Park, D. J. Kim, H. Yang, J. Hur, D.-Y. Kang *et al.*, "Deep learning segmentation of major vessels in x-ray coronary angiography," *Scientific reports*, vol. 9, no. 1, pp. 1–11, 2019.
- [9] Z. Xian, X. Wang, S. Yan, D. Yang, J. Chen, and C. Peng, "Main coronary vessel segmentation using deep learning in smart medical," *Mathematical Problems in Engineering*, vol. 2020, 2020.
- [10] T. J. Jun, J. Kweon, Y.-H. Kim, and D. Kim, "T-net: Nested encoder–decoder architecture for the main vessel segmentation in coronary angiography," *Neural Networks*, vol. 128, pp. 216–233, 2020.
- [11] J. Fan, J. Yang, Y. Wang, S. Yang, D. Ai, Y. Huang, H. Song, A. Hao, and Y. Wang, "Multichannel fully convolutional network for coronary artery segmentation in x-ray angiograms," *Ieee Access*, vol. 6, pp. 44 635–44 643, 2018.
- [12] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia, "Pyramid scene parsing network," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 2881–2890.

- [13] C. Zhao, H. Tang, J. Tang, C. Zhang, Z. He, Y.-P. Wang, H.-W. Deng, R. Bober, and W. Zhou, "Semantic segmentation to extract coronary arteries in fluoroscopy angiograms," *medRxiv*, 2020.
- [14] P. M. Samuel and T. Veeramalai, "Vssc net: vessel specific skip chain convolutional network for blood vessel segmentation," *Computer Methods and Programs in Biomedicine*, vol. 198, p. 105769, 2021.
- [15] K. Santhi and A. R. M. Reddy, "An automated framework for coronary analysis from coronary cine angiograms using machine learning and image analysis techniques," *INFORMATION TECHNOLOGY IN INDUSTRY*, vol. 9, no. 1, pp. 1406–1412, 2021.
- [16] F. Cervantes-Sanchez, I. Cruz-Aceves, A. Hernandez-Aguirre, M. A. Hernandez-Gonzalez, and S. E. Solorio-Meza, "Automatic segmentation of coronary arteries in x-ray angiograms using multiscale analysis and artificial neural networks," *Applied Sciences*, vol. 9, no. 24, p. 5507, 2019.
- [17] D. Liang, J. Qiu, L. Wang, X. Yin, J. Xing, Z. Yang, J. Dong, and Z. Ma, "Coronary angiography video segmentation method for assisting cardiovascular disease interventional treatment," *BMC medical imaging*, vol. 20, no. 1, pp. 1–8, 2020.
- [18] L. Wang, D. Liang, X. Yin, J. Qiu, Z. Yang, J. Xing, J. Dong, and Z. Ma, "Coronary artery segmentation in angiographic videos utilizing spatial-temporal information," *BMC Medical Imaging*, vol. 20, no. 1, pp. 1–10, 2020.
- [19] D. Hao, S. Ding, L. Qiu, Y. Lv, B. Fei, Y. Zhu, and B. Qin, "Sequential vessel segmentation via deep channel attention network," *Neural Networks*, vol. 128, pp. 172–187, 2020.
- [20] L. N. L. Thuy, T. D. Trinh, L. H. Anh, J. Y. Kim, H. T. Hieu *et al.*, "Coronary vessel segmentation by coarse-to-fine strategy using u-nets," *BioMed Research International*, vol. 2021, 2021.
- [21] E. Nasr-Esfahani, N. Karimi, M. H. Jafari, S. M. R. Soroushmehr, S. Samavi, B. Nallamothu, and K. Najarian, "Segmentation of vessels in angiograms using convolutional neural networks," *Biomedical Signal Processing and Control*, vol. 40, pp. 240–251, 2018.
- [22] J. Zhang, G. Wang, H. Xie, S. Zhang, N. Huang, S. Zhang, and L. Gu, "Weakly supervised vessel segmentation in x-ray angiograms by self-paced learning from noisy labels with suggestive annotation," *Neurocomputing*, vol. 417, pp. 114–127, 2020.
- [23] K. Iyer, C. P. Najarian, A. A. Fattah, C. J. Arthurs, S. R. Soroushmehr, V. Subban, M. A. Sankardas, R. R. Nadakuditi, B. K. Nallamothu, and C. A. Figueroa, "Angionet: A convolutional neural network for vessel segmentation in x-ray angiography," *medRxiv*, 2021.
- [24] L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam, "Encoder-decoder with atrous separable convolution for semantic image segmentation," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 801–818.
- [25] W. R. Crum, O. Camara, and D. L. Hill, "Generalized overlap measures for evaluation and validation in medical image analysis," *IEEE transactions on medical imaging*, vol. 25, no. 11, pp. 1451–1461, 2006.
- [26] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *International Conference on Medical image computing and computer-assisted intervention*. Springer, 2015, pp. 234–241.
- [27] I. Radosavovic, R. P. Kosaraju, R. Girshick, K. He, and P. Dollár, "Designing network design spaces," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 10 428–10 436.
- [28] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He, "Aggregated residual transformations for deep neural networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1492–1500.
- [29] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [30] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie, "Feature pyramid networks for object detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 2117–2125.

- [31] A. Chaurasia and E. Culurciello, “Linknet: Exploiting encoder representations for efficient semantic segmentation,” in *2017 IEEE Visual Communications and Image Processing (VCIP)*. IEEE, 2017, pp. 1–4.
- [32] R. Li, S. Zheng, C. Duan, C. Zhang, J. Su, and P. Atkinson, “Multi-attention-network for semantic segmentation of fine resolution remote sensing images,” *arXiv preprint arXiv:2009.02130*, 2020.
- [33] H. Li, P. Xiong, J. An, and L. Wang, “Pyramid attention network for semantic segmentation,” *arXiv preprint arXiv:1805.10180*, 2018.
- [34] Z. Zhang, Q. Liu, and Y. Wang, “Road extraction by deep residual u-net,” *IEEE Geoscience and Remote Sensing Letters*, vol. 15, no. 5, pp. 749–753, 2018.
- [35] D. Jha, P. H. Smedsrud, M. A. Riegler, D. Johansen, T. De Lange, P. Halvorsen, and H. D. Johansen, “Resunet++: An advanced architecture for medical image segmentation,” in *2019 IEEE International Symposium on Multimedia (ISM)*. IEEE, 2019, pp. 225–2255.
- [36] A. G. Roy, N. Navab, and C. Wachinger, “Concurrent spatial and channel ‘squeeze & excitation’ in fully convolutional networks,” in *International conference on medical image computing and computer-assisted intervention*. Springer, 2018, pp. 421–429.
- [37] J. Hu, L. Shen, and G. Sun, “Squeeze-and-excitation networks,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 7132–7141.
- [38] S. Ioffe and C. Szegedy, “Batch normalization: Accelerating deep network training by reducing internal covariate shift,” in *International conference on machine learning*. PMLR, 2015, pp. 448–456.
- [39] A. Howard, M. Sandler, G. Chu, L.-C. Chen, B. Chen, M. Tan, W. Wang, Y. Zhu, R. Pang, V. Vasudevan *et al.*, “Searching for mobilenetv3,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 1314–1324.
- [40] K. He, X. Zhang, S. Ren, and J. Sun, “Delving deep into rectifiers: Surpassing human-level performance on imagenet classification,” in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 1026–1034.
- [41] X. Glorot and Y. Bengio, “Understanding the difficulty of training deep feedforward neural networks,” in *Proceedings of the thirteenth international conference on artificial intelligence and statistics*. JMLR Workshop and Conference Proceedings, 2010, pp. 249–256.
- [42] P. Yakubovskiy, “Segmentation models pytorch,” https://github.com/qubvel/segmentation_models.pytorch, 2020.

A Performance vs. Computation Trade-Off

In clinical practice, depending on the available hardware resources and clinical needs, it may be necessary to make a trade-off between performance and computational efficiency. To help practitioners make that choice, we present the Pareto frontier of all tested models in Figure 4, with the number of FLOPS as a function of performance, measured by generalized dice score. Using these axes, each model on the Pareto frontier is the most efficient at each performance level, and best-performing at each computation regime. Thus, all other models need not be considered when making a trade-off between performance and computational efficiency.

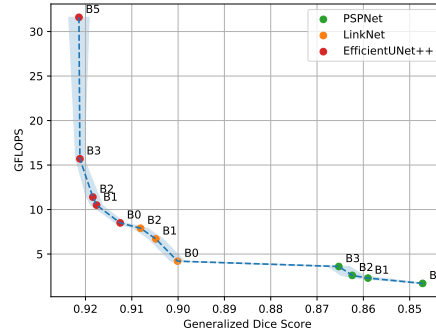


Figure 4: FLOPS as a function of performance, measured by generalized dice score. The dashed polygonal line corresponds to the Pareto frontier. Each dot represents a model: the text labels indicate the EfficientNet backbone, from B0 to B7, and the colors denote the decoder architecture.

In Figure 4, the lower a point is, the less computation it requires for inference, and the more to the left it is, the better its performance. Therefore, the gentler the slope between a model A and a better-performing model B, the more significant is the relative merit of A compared to B. Taking that and the low performance of PSPNet decoders into account, we suggest using the slightly more computationally demanding LinkNet-based architectures when there is a limit on computational complexity.

B Attention Blocks Ablation

Figures 5a and 5b show the performance of the EfficientUNet++ decoder architecture without attention mechanisms and when using SE, sSE and scSE attention blocks, which perform channel attention, spatial attention, and concurrent channel and spatial attention. From the graph, we conclude that, for high computation regimes, combining channel and spatial attention improves performance at the cost of only a slight increase in computation and parameters. At low computation regimes, all models have similar performance and computational efficiency. Thus, scSE blocks are an overall better choice, and we use them regardless of the computation regime.

Notably, even without attention mechanisms, the best EfficientUNet++ and UNet++ models perform similarly, showing that replacing the UNet++’s blocks with residual bottlenecks with depthwise convolutions increases efficiency without performance loss.

We hypothesize that, as EfficientUNet++ models using scSE blocks generally outperform those using SE blocks, architecture variants using other attention mechanisms could exceed the performance of both, possibly with less computation and fewer parameters. Therefore, we argue that the study of attention in decoder architectures may be a promising line of future research.

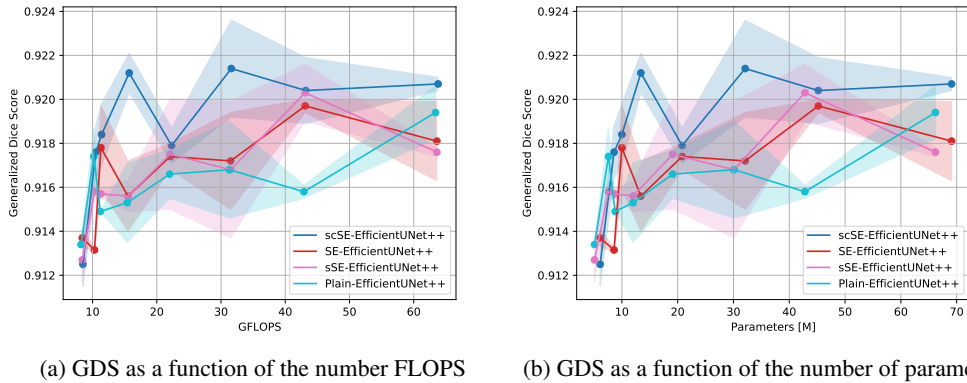


Figure 5: EfficientUNet++ performance without attention mechanisms (in cyan) and when using sSE spatial attention blocks (in pink), SE channel attention blocks (in red), and scSE concurrent channel and spatial attention blocks (in blue). The dots represent the following EfficientNet models, in ascending order of FLOPS and parameters: B0, B1, B2, B3, B4, B5, B6, B7.

C Qualitative Decoder Comparison

Figures 6 and 7 show the segmentation masks obtained by each decoder when coupled with the encoder it performs best with, for a left and a right coronary artery, respectively. Visual inspection of these images supports the results of quantitative comparison. The U-Net-based models, except for the PAN, generally perform well. The FPN and DeepLabV3+ produce coarser yet very reasonable results, and the PSPNet performs very poorly, having trouble distinguishing catheters from arteries.

The four best-performing architectures, i.e., the EfficientUNet++, ResUNet++, UNet++ and U-Net, output very good segmentation masks with few differences between them. Given the qualitative similarity between these models, the main advantage of the EfficientUNet++ is its higher computational efficiency. Also, we hypothesize that due to the attention mechanisms, it has more room for improvement than other architectures if trained on a larger dataset.

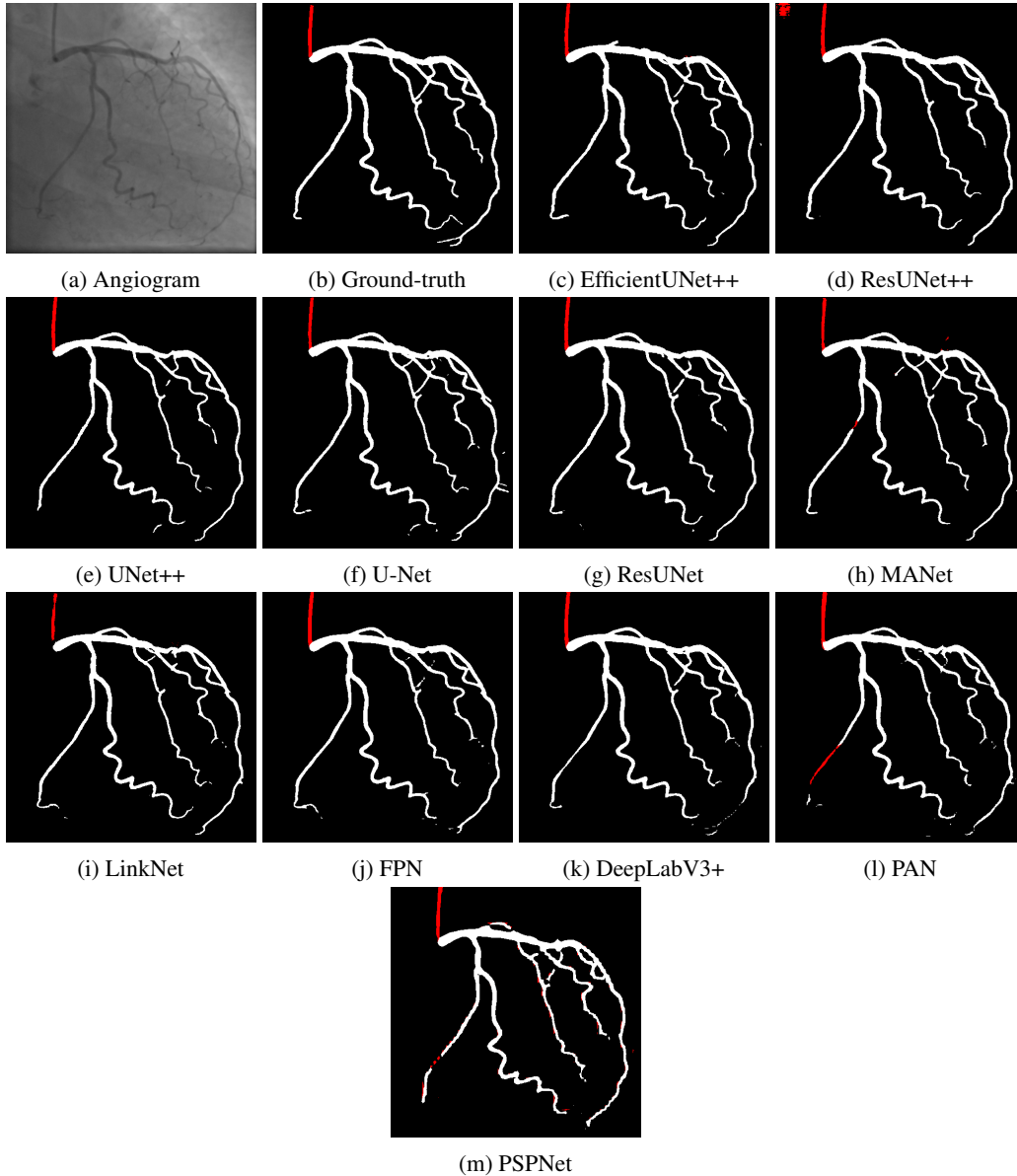


Figure 6: Segmentation of a left coronary artery. Each decoder was coupled with the encoder it performs best with. The masks are sorted in descending order of the models' average performances.

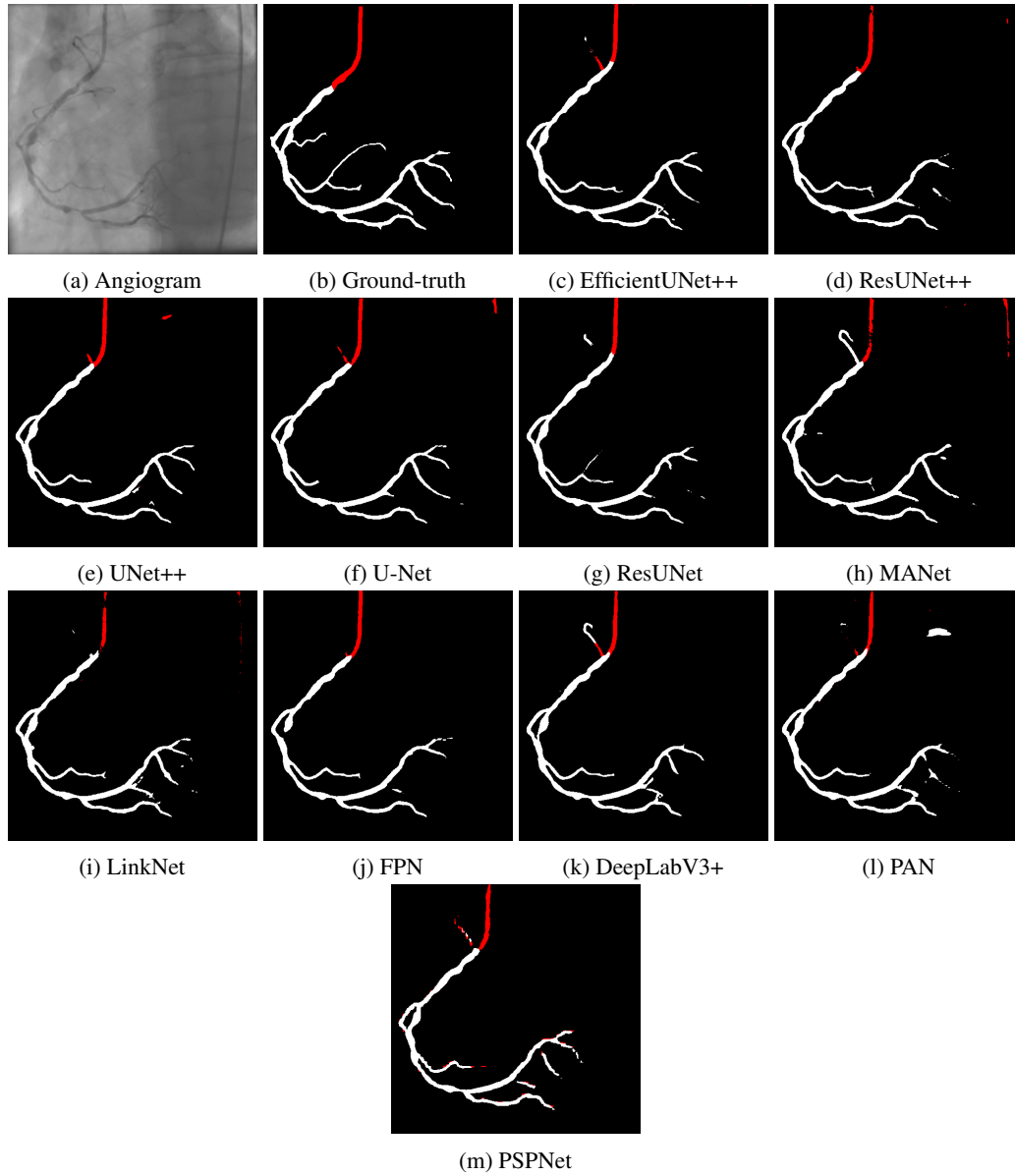
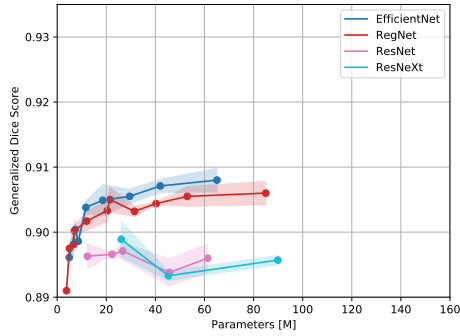
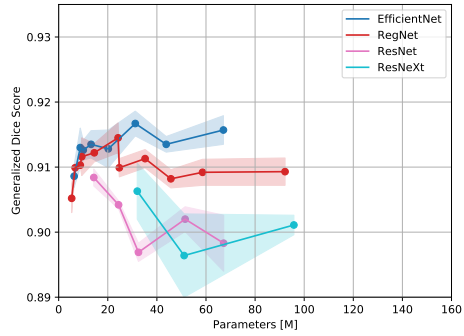


Figure 7: Segmentation of a right coronary artery. Each decoder was coupled with the encoder it performs best with. The masks are sorted in descending order of the models' average performances.

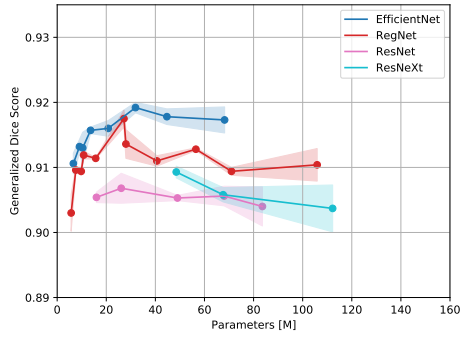
D Encoder and Decoder Performance Comparison Graphs



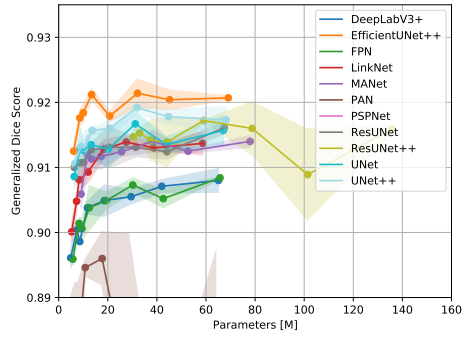
(a) DeepLabV3+ decoder



(b) U-Net decoder



(c) UNet++ decoder



(d) EfficientNet B0 to B7 encoders

Figure 8: Segmentation performance, measured by generalized dice score (GDS), as a function of the number of parameters. Figures (a), (b) and (c) show the performance of different encoders combined with the (a) DeepLabV3+, (b) U-Net and (c) UNet++ decoders. Figure (d) shows the performance of different decoders combined with the EfficientNet B0 to B7 encoders. Each polygonal line corresponds to an encoder family. The dots represent the following models, in ascending order of parameters: EfficientNet - B0, B1, B2, B3, B4, B5, B6, B7; RegNet - Y2, Y4, Y6, Y8, Y16, Y32, Y40, Y64, Y80, Y120, Y160; ResNet - 18, 34, 50, 101, 152; ResNeXt - 50_32x4d, 101_32x4d, 101_32x8d. Models with GDS below 0.89 are also omitted.