# Automating Carotid Intima-Media Thickness Video Interpretation with Convolutional Neural Networks*

Jae Y. Shin*, Nima Tajbakhsh*, R. Todd Hurst, Christopher B. Kendall, and Jianming Liang‡

## Abstract

*Cardiovascular disease (CVD) is the leading cause of mortality yet largely preventable, but the key to prevention is to identify at-risk individuals before adverse events. For predicting individual CVD risk, carotid intima-media thickness (CIMT), a noninvasive ultrasound method, has proven to be valuable, offering several advantages over CT coronary artery calcium score. However, each CIMT examination includes several ultrasound videos, and interpreting each of these CIMT videos involves three operations: (1) select three end-diastolic ultrasound frames (EUF) in the video, (2) localize a region of interest (ROI) in each selected frame, and (3) trace the lumen-intima interface and the media-adventitia interface in each ROI to measure CIMT. These operations are tedious, laborious, and time consuming, a serious limitation that hinders the widespread utilization of CIMT in clinical practice. To overcome this limitation, this paper presents a new system to automate CIMT video interpretation. Our extensive experiments demonstrate that the suggested system significantly outperforms the state-of-the-art methods. The superior performance is attributable to our unified framework based on convolutional neural networks (CNNs) coupled with our informative image representation and effective post-processing of the CNN outputs, which are uniquely designed for each of the above three operations.*
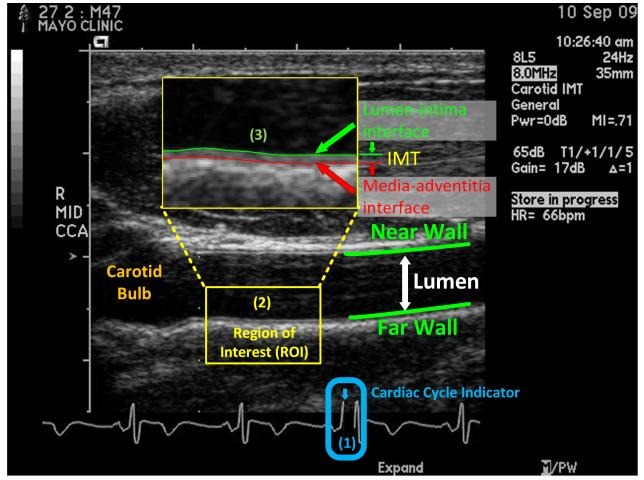
Figure 1: Longitudinal view of the carotid artery in an ultrasound B-scan image. CIMT is defined as the distance between the lumen-intima interface and the media-adventitia interface, measured approximately $1\,cm$ distal from the carotid bulb on the far wall of the common carotid artery at the end of the diastole; therefore, interpreting a CIMT video involves three operations: (1) select three end-diastolic ultrasound frames (EUFs) in each video (the cardiac cycle indicator, a black line, shows to where in the cardiac cycle the current frame corresponds); (2) localize a region of interest (ROI) approximately $1\,cm$ distal from the carotid bulb in the selected EUF; (3) measure the CIMT within the localized ROI. This paper aims to automate these three operations simultaneously through a unified framework based on convolutional neural networks.

## 1. Introduction

Given the clinical significance of carotid intima-media thickness (CIMT) as an early and reliable indicator of cardiovascular risk, several methods have been developed for CIMT image interpretation. The CIMT is defined as the distance between the lumen-intima and media-adventitia interfaces at the far wall of the carotid artery (Figure 1). Therefore, to measure CIMT, the lumen-intima and the media-adventitia interfaces must be identified. As a re-

sult, the earlier approaches are focused on analyzing the intensity profile and distribution, computing the gradient [24, 27, 5], or combining various edge properties through dynamic programming [13, 3, 25]. Recent approaches [16, 4, 23, 28, 7, 2] are mostly based on active contours (aka, snakes) or their variations [9]. Some of these approaches require user interaction, while other approaches aim for complete automation through integrating with various image processing algorithms, such as Hough transform [21] and dynamic programming [25]. Most recently, Menchn-Lara et al. employed a committee of standard multilayer perceptrons in [18] and a single standard multilayer perceptron with an auto-encoder in [19] for CIMT image interpretation, but both methods did not outperform the snake-based methods from the same research group [2, 1]. For a more complete survey of methods for automatic CIMT measurements, please refer to the review studies conducted by Molinari et al.[22] and Loizou et al. [15].

However, nearly all the aforementioned methods are focused on only the third operation: CIMT measurement, ignoring the two preceding operations, i.e., frame selection and ROI localization. To our knowledge, the only system that simultaneously automates the three operations is the work [26], an extension of [29], which automatically selects the EUF frame, localizes the ROI in each selected EUF frame, and provides the CIMT measurement in the selected ROI. However, as with other works, this method is based on hand-crafted algorithms, which often lack the desired robustness for routine clinical use, a weakness that we aim to overcome in this paper.

A key contribution of this paper is a new system that accelerates CIMT video interpretation by automating all the three operations in a novel unified framework based on convolutional neural networks (CNNs). We will show that with proper pre-processing and post-processing, our proposed CNN-based approach can significantly outperform the existing methods in all aspects of CIMT image interpretation including frame selection, ROI localization, and CIMT measurements, making the following specific contributions:

- A unified framework based on CNNs that automates the entire CIMT interpretation process. This is in contrast to the prior works where only the very last step of the CIMT interpretation process was automated. The performance of the suggested system significantly outperforms the hand-crafted approach [26], which, to our knowledge, is the only system in the literature that aimed to automate all the above three tasks.
- A novel frame selection method based on the ECG signals at the bottom of ultrasound frames. The suggested method utilizes effective pre-processing of patches and post processing of CNN outputs, enabling a significant increase in the performance of a baseline CNN.
- A new method that localizes the ROI for CIMT in-

terpretation. The suggested method combines the discriminative power of a CNN with a contextual constrain to accurately localize the ROIs in the selected frames. We demonstrate that the suggested contextually-constrained CNN outperforms the performance of a baseline CNN.
- A framework that combines CNNs with active contour models for accurate boundary segmentation. Specifically, given a localized ROI, the CNN initializes two open snakes, which further deform to acquire the shapes of intima-media boundaries. We show that the segmentation accuracy of the suggested method is far higher than the state-of-the-art methods.
- Extensive evaluation of each stage of the suggested CIMT interpretation system. Specifically, we perform leave-one-patient-out cross-validation[1] using only the training CIMT videos to tune the parameters of the suggested system, and then thoroughly evaluate the performance of our system using a large number of independent test CIMT videos.

## 2. CIMT Protocol

The CIMT exams utilized in this paper were performed with B-Mode ultrasound using an 8-14MHz linear array transducer utilizing fundamental frequency only (Acuson Sequoia[TM], Mountain View, CA, USA) [6]. The carotid screening protocol begins with scanning bilateral carotid arteries in a transverse manner from the proximal aspect to the proximal internal and external carotid arteries. The probe is then turned to obtain the longitudinal view of the distal common carotid artery. The sonographer optimizes the 2D images of the lumen-intima and media-adventitia interfaces at the level of the common carotid artery by adjusting overall gain, time gain, compensation and focus position. Once the parameters are optimized, the sonographer captures two CIMT videos focused on the common carotid artery from two optimal angles of incidence. The same procedure is repeated for the other side of neck, resulting in a total of 4 CIMT videos for each subject.

## 3. Method

Our goal is to automate the three operations in CIMT video interpretation, i.e, given a CIMT video, our method will automatically identify three EUFs (Section 3.1), localize an ROI in each EUF (Section 3.2), and segment the lumen-intima and media-adventitia interfaces within each ROI (Section 3.3).

### 3.1. Frame Selection

We select the EUFs based on the ECG signal embedded at the bottom part of a CIMT video. The cardiac cycle indi-

---

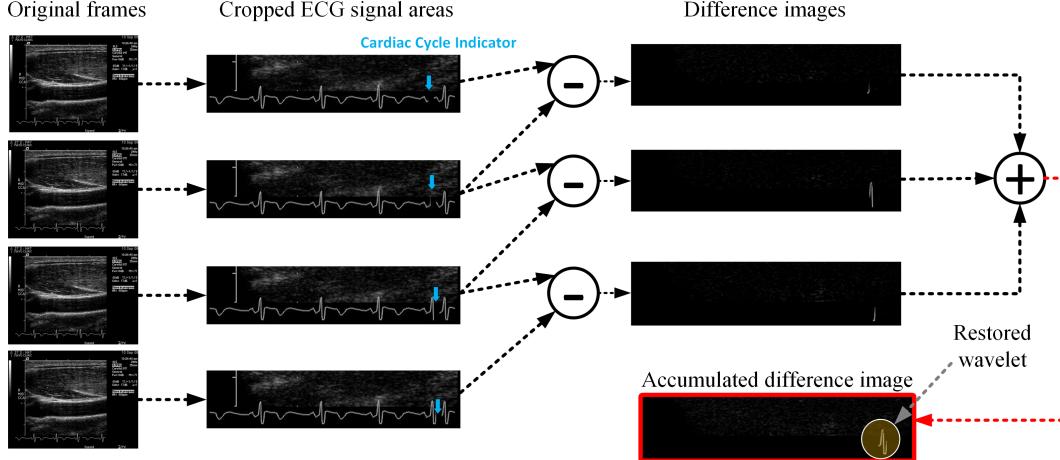[1]We leave all the videos from one patient out for validation.

Figure 2: An accumulated difference image is generated by adding up three neighboring difference images.

cator is represented by a moving black line in each frame. Since the ECG signal is overlaid on the ultrasound image, there is quite bit of noise around the indicator. The challenge is to reconstruct the original ECG signal from noisy frames and to detect the R peaks from the ECG signal, as the R-peaks correspond to the EUFs. To do so, we introduce accumulated difference images that carry sufficient information for CNN to learn and distinguish R-peaks from non-R-peaks.

**Training Phase:** Let $I^t$ denote an image subregion selected from the lower part of an ultrasound frame so that it contains the ECG signal. We first construct a set of difference images $d^t$ by subtracting every consecutive pairs of images, $d^t = |I^t - I^{t+1}|$, and then form accumulated difference images by adding up every three neighboring difference images, $D^t = \sum_{i=0}^{2} d^{t-i}$. Accumulated difference image $D^t$ can capture the cardiac cycle indicator at frame $t$. Figure 2 illustrates how an accumulated difference image is generated.

Next, we determine the location of the restored wavelet in each accumulated difference image. For this purpose, we find the weighted centroid $c = [c_x, c_y]$ of each accumulated difference image $D^t$ as follows:

$$c = \frac{1}{Z_t} \sum_{p \in D^t} D^t(p_x, p_y) \times p$$

where $p = [p_x, p_y]$ is a pixel in the accumulated difference image and $Z_t = \sum_{p \in D^t} D^t(p_x, p_y)$ is a normalization factor that ensures the weighted centroid stays within the image boundary. Once centroids are identified, we extract patches of size $32 \times 32$ around the centroid locations. Specifically, we extract patches with up to 2 pixel translations from each centroid. However, we do not scale the patches in data augmentation, because doing so would inject label noise in the training set. For instance, a small restored wavelet may take the appearance of an R-peak after expanding or an R-peak
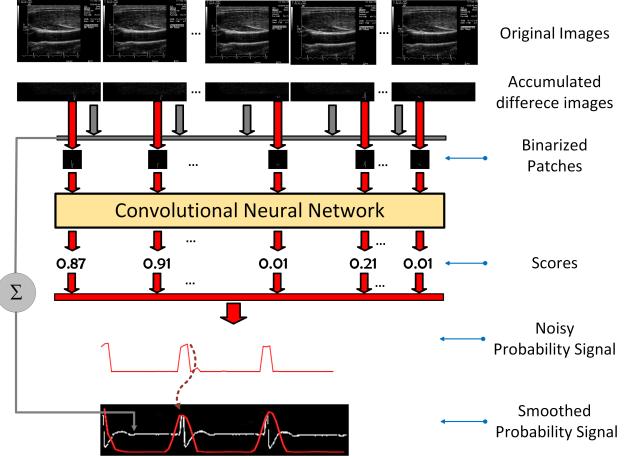


Figure 3: The test stage of our automatic frame selection scheme.

wavelet may look like a non-R-peak wavelet after shrinking. Nor do we perform rotation-based patch augmentation, because we do not expect the restored wavelets to appear with rotation in the test image patches. Once collected, patches are binarized using Otsu's method. In Section 4, we discuss the choice of binarization method through an extensive set of experiments. Each binary patch is then labeled as positive if it corresponds to an EUF (i.e., an R-peak); otherwise negative. Basically, given a patch, we first determine the accumulated difference image from which the patch is extracted. We then trace back to the underlying difference images and check whether they are related to the EUF or not. Once the patches are labeled, we form a stratified set with 96,000 patches to train a 2-way CNN for frame selection.

**Testing Phase:** Figure 3 shows our frame selection system given a test video. We first compute an accumulated differ-

ence image for each frame in the video. We then extract image patches from the weighted centroids of the accumulated difference images. The probability of each frame being the EUF is measured as the average probabilities assigned by the CNN to the corresponding patches. By concatenating the resulting probabilities for all the frames in the video, we obtain a probability signal whose local maxima indicate the locations of the EUFs. However, the generated probability signals often exhibit abrupt changes, which can cause too many local maxima along the signal. We therefore first smooth the probability signal using a Gaussian function, and then find the EUFs by locating the local maxima of the smoothed signals. In Figure 3, for illustration purposes, we have also shown the reconstructed ECG signal, which is computed as the average of the accumulated difference images, $\frac{1}{N}\sum_{t=1}^{N} D^t$ with $N$ being the number of frames in the video. As seen, the probability of being the EUF reaches its maximum around the R peaks of the QRS complexes (as desired) and then smoothly decays as it distances from the R peaks. By mapping the locations of the local maxima to the frame numbers, we can identify the EUFs in the test video.

### 3.2. ROI Localization

Accurate localization of the ROI is challenging, because, as seen in Figure 1, there are no significant differences that can be observed in image appearance among the ROIs on the far wall of the carotid artery. To overcome this challenge, we utilize the location of the carotid bulb as a contextual constraint. We choose this constraint for two reasons: 1) the carotid bulb appears as a distinct dark area in the ultrasound frame and thus can be uniquely identified; 2) according to the consensus statement of American society of Electrocardiography for cardiovascular risk assessment, the ROI should be placed approximately $1\,\mathrm{cm}$ from the carotid bulb on the far wall of the common carotid artery. While the former motivates the use of the carotid bulb location as a constraint from a technical point of view, the latter justifies this constraint from a clinical standpoint.

**Training Phase:** We incorporate this constraint in the suggested system by training a 3-way CNN that simultaneously localizes both ROI and carotid bulb, and then refines the estimated location of the ROI given the location of the carotid bulb. Figure 9 in the supplementary material illustrates how the image patches are extracted from a training frame. We perform data augmentation by extracting the training patches within a circle around the locations of the carotid bulbs and the ROIs. The negative patches are extracted from a grid of points sufficiently far from the locations of the carotid bulbs and the ROIs. Note that the above translation-based data augmentation is sufficient for this application, because our database provides a relatively large number of training EUFs, from which a large set of training patches can be collected. Once the patches are collected, we form

a stratified training set with approximately 410,000 patches to train a 3-way CNN for constrained ROI localization.

**Testing Phase:** Referring to Figure 4, during the test stage, the trained CNN is applied to all the pixels in the EUF, generating two confidence maps with the same size as the EUF. The first confidence map shows the probability of a pixel being the carotid bulb and the second confidence map shows the probability of a pixel being the ROI. One way to localize the ROI is to find the center of the largest connected component within the ROI confidence map without considering the detected location of the carotid bulb. However, this naive approach may fail to accurately localize the ROI. For instance, a long-tale connected component along the far wall of the carotid artery may cause substantial ROI localization error. To compound the problem, the largest connected component of the ROI confidence map may appear far from the actual location of the ROI, resulting in a complete detection failure. To overcome these limitations, we constraint the ROI location $l_{roi}$ by the location of the carotid bulb $l_{cb}$. For this purpose, we first determine the location of the carotid bulb as the centroid of the largest connected component within the first confidence map, and then localize the ROI using the following formula

$$l_{roi} = \frac{\sum_{p \in C^*} M(p) \cdot p \cdot I(p)}{\sum_{p \in C^*} M(p) \cdot I(p)} \tag{1}$$

where $l_{roi}$ denotes the ROI location, $l_{cb}$ denotes the center of the carotid bulb, $M$ denotes the confidence map of being the ROI, $C^*$ is the largest connected component in $M$ that is the nearest to the carotid bulb, and I(p) is an indicator function for pixel $p = [p_x, p_y]$ that is defined as

$$I(p) = \begin{cases} 1, & \text{if } \|p - l_{cb}\| < 1\,\mathrm{cm} \tag{2} \\ 0, & \text{otherwise} \tag{3} \end{cases}$$

The indicator function I(p) is binary function that simply includes pixel when the value is 1 as in Eq. 2, otherwise excludes pixel when the value is 0 as in Eq. 3.

### 3.3. Intima-Media Thickness Measurement

Measuring intima-media thickness require a continuous and one-pixel precise boundary for lumen-intima and media-adventitia. Lumen-intima is relatively easier to detect because of strong gradient change at the border, however, detecting media-adventitia interface is quite challenging due to its subtle image gradients and noise around its border. We approach this problem as a 3-way classification task: 1) lumen-intima interface, 2) media-adventitia interface, and 3) background.

**Training Phase:** To train 3-way CNN, we collected sparse background patches and then pixel-by-pixel image patches around lumen-intima interface and media-adventitia interface with additional patches $\pm 3$ pixels from the ground
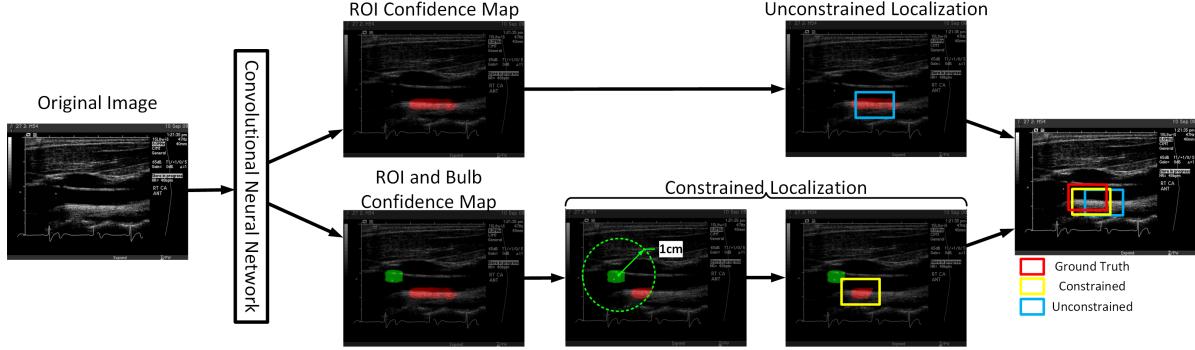
Figure 4: The test stage of our ROI localization method. In the unconstrained scenario, we only use the ROI confidence map, which results in relatively large localization error. In the constrained mode, given the estimated location of the carotid bulb, we localize the ROI more accurately.

truth. Using ±3 pixels for additional patches around intima-media boundary was necessary to balance number of patches with background patches and produced better results. Figure 10 in the supplementary material illustrates how the training patches are collected from an ROI.

**Testing Phase:** Figure 5 illustrates the testing process. The 3-way trained CNN is applied in a sliding-window fashion for a given test ROI and generates two confidence maps (Figure 5(b)) with the same size as the ROI. Since confidence map is thicker than a pixel, we choose the maximum response column-by-column and generate a new binary image as shown in Figure 5(c). Finally, we use two active contour models (a.k.a, snakes) [12] for segmenting lumen-intima and media-adventitia interfaces. Figure 5(d) shows two final converged snakes and we take measurements as the average vertical distance between the two snakes.

## 4. Experiments

We use a database of 92 CIMT videos captured from 23 subjects with 2 CIMT videos from the left and 2 CIMT videos from the right carotid artery of each subject. The ground truth for each video contains the EUF number, the locations of ROI, and the segmentation of lumen-intima and media-advantitia interfaces. For consistency, we use the same training set and the same test set (no overlap with training) for all three tasks. Our training set contains 48 CIMT videos of 12 subjects with a total of 4,456 frames and our test set contains 44 CIMT videos of 11 subjects with a total of 3,565 frames. For each task, we perform leave-one-patient-out cross-validation based on the *training subjects* to tune the parameters, and then evaluate the performance of the tuned system using the test subjects.

**Architecture:** As shown in Table 1, we employ a CNN architecture with 2 convolutional layers, 2 subsampling layers, and 2 fully connected layers (see Section 5 for our justifications). We also append a softmax layer to the last

fully connected layer so as to generate probabilistic confidence values for each class. Our CNN architecture has input patches of size 32x32, and we resize the collected patches to 32x32 prior to the training process. For the CNNs used in our experiments, we employ a learning rate of $\alpha = 0.001$, a momentum of $\mu = 0.9$, and a constant scheduling rate of $\gamma = 0.95$.

**Pre- and post-processing for frame selection:** We have experimentally found out that binarized image patches improve the quality of convergence and accuracy of frame selection. Furthermore, we have observed that the standard deviation of the Gaussian function used for smoothing the probability signals, can also substantially influence frame selection accuracy. Therefore, we have conducted leave-one-patient-out cross-validation based on the training subjects to find the best binarization method and the optimal standard deviation of the Gaussian function. For binarization, we have considered a fixed set of thresholds and adaptive thresholding using Otsu's method. For smoothing, we have considered a Gaussian function with different standard deviation ($\sigma_g$) as well as the scenario where no smoothing is applied. For each configuration of parameters, we have done a free-response ROC (FROC) analysis. We consider a selected frame a true positive, if it is found within one frame from the expert-annotated EUF; otherwise, a false positive.

Our leave-one-patient-out cross-validation study, summarized in Figure 11 in the supplementary material, indicates that the use of a Gaussian function with $\sigma_g = 1.5$ for smoothing the probability signals and adaptive thresholding using Otsu's method achieve the highest performance. Figure 6 shows the FROC curve of our system for the test subjects using the above parameters. For comparison, we have also shown the operating point of the hand-crafted approach [26], which is significantly outperformed by the suggested system.

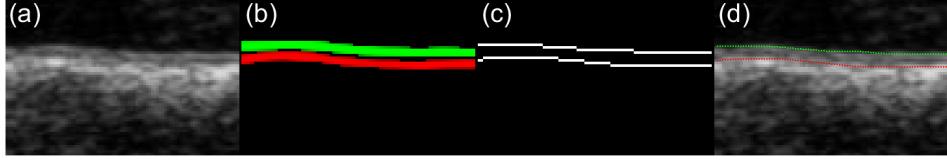**Constrained ROI Localization:** We conduct a leave-one-

Figure 5: The test stage of lumen-intima and media-adventitia interface detection. (a) a test ROI. (b) The trained CNN generates a confidence map where the green and red colors indicate the likelihood of lumen-intima interface and media-adventitia interface, respectively. (c) The thick probability band around each interface is thinned by selecting the largest probability for each interface in each column. (d) The step-like boundaries are refined through two open snakes.

Table 1: The CNN architecture used in our experiments. Note that $C$ is the number of classes, which is 2 for frame selection and 3 for both ROI localization and intima-media thickness measurements.

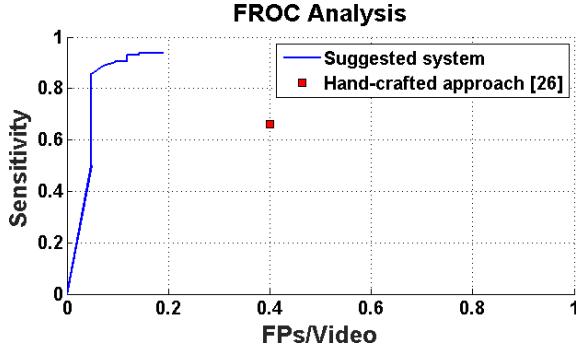| layer | type | input | kernel | stride | pad | output |
|---|---|---|---|---|---|---|
| 0 | input | 32x32 | N/A | N/A | N/A | 32x32 |
| 1 | convolution | 32x32 | 5x5 | 1 | 0 | 64x28x28 |
| 2 | max pooling | 64x28x28 | 3x3 | 2 | 0 | 64x14x14 |
| 1 | convolution | 64x14x14 | 5x5 | 1 | 0 | 64x10x10 |
| 2 | max pooling | 64x10x10 | 3x3 | 2 | 0 | 64x5x5 |
| 2 | fully connected | 64x5x5 | 5x5 | 1 | 0 | 250x1 |
| 2 | fully connected | 250x1 | 1x1 | 1 | 0 | Cx1 |



Figure 6: FROC curve of our frame selection system for the test subjects using the tuned parameters. For comparison, we have also shown the operating point of the prior hand-crafted approach [26], which is significantly outperformed by the suggested system.

patient-out cross-validation study based on the training subjects to find the optimal size of the training patches. Our cross-validation analysis, summarized in Figure 12 in the supplementary material, indicates that the use of $1.8 \times 1.8$ cm patches achieves the most stable performance, yielding low ROI localization error with only a few outliers. Fig-

ure 7 shows the ROI localization error of our system for the test subjects using the optimal size of training patches. To demonstrate the effectiveness of our constrained ROI localization method, we have also included the performance of the unconstrained counterpart. In the constrained mode, we use Eq. 1 for ROI localization whereas in the unconstrained mode we localize the ROI as the center of the largest connected component in the corresponding confidence map without considering the location of the carotid bulb. Our method achieves an average localization error of $0.19\,\mathrm{mm}$ and $0.35\,\mathrm{mm}$ in the constrained and unconstrained modes, respectively. The decrease in localization error is statistically significant ($p < 0.01$). Also as seen in Figure 7, our method in the unconstrained mode has resulted in 3 complete localization failures (outliers), which have been corrected in the constrained mode. Furthermore, compared with the hand-crafted approach [26], our system in the constrained mode shows a decrease of $0.1\,\mathrm{mm}$ in ROI localization error, which is statistically significant ($p < .00001$).

**Intima-Media Thickness Measurement:** We determined the optimal image patch size by leave-one-patient-out cross-validation using various image patch sizes and found that $360 \times 360\ \mu m$ achieved slightly lower localization error and fewer outliers (see Figs. 12-13 in the supplementary ma-
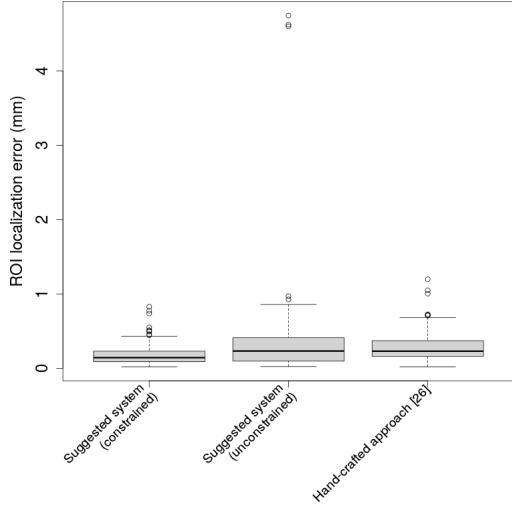
Figure 7: ROI localization error for the test subjects. Our method in the constrained mode outperforms both the unconstrained counterpart and the prior hand-crafted approach [26].
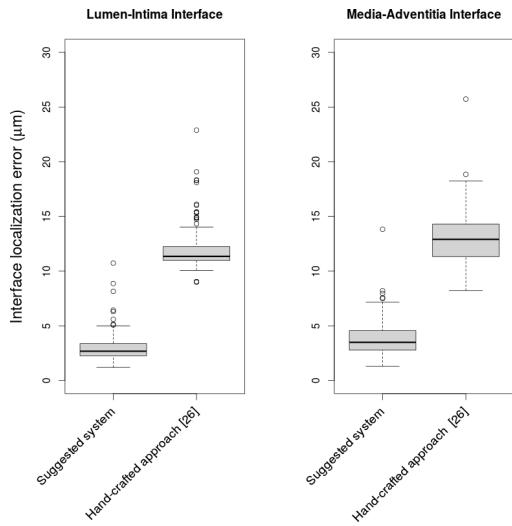


Figure 8: Localization error of the lumen-intima and media-adventitia interfaces for the suggested system and the prior hand-crafted approach [26]. The results are obtained for the test subjects.

terial). Figure 8 shows the interface localization error of our system on the test subjects, where we break down the overall localization error for lumen-intima and that of the media-adventitia interface as well as the hand-crafted approach [26] for each interface. We further analyzed agreement between our system and the expert with the Bland-Altman plot (see Figure 14 in the supplementary material).

## 5. Discussions

In Section 4, we investigated how the choice of patch binarization and degree of Gaussian smoothing affect the accuracy of frame selection. Here, we would like to discuss our findings and provide insights about our choices. We choose to binarize the patches, because it reduces appearance variability and suppress the low-magnitude noise content in the patches. Without patch binariztion, one can expect a large amount of variability in the appearance of wavelets can deteriorate the performance of the subsequent CNN (see Figure 11 in the supplementary material). The choice of binarization threshold is another important factor. The use of a high threshold results in the partial appearance of the wavelets in the resulting binary patches, reducing the discriminatory appearance features of the patches. A low threshold, on the other hand, can intensify noise content in the images, which decreases the quality of training samples and consequently a drop in classification performance. According to our analyses, it is difficult to find a fixed threshold that can both suppress the noise content and keep the shapes of the restored wavelets intact in all the collected patches. Otsu's method seems to overcome this limitation by adaptively selecting a binarization threshold according to the intensity distribution of each individual patch. For patches with intensity values between 0 and 1, the adaptive thresholds have a mean of 0.15 and standard deviation of 0.05. The wide range of adaptive thresholds explains why a constant threshold may not perform as desirably.

Gaussian smoothing of the probability signals is also essential for accurate frame selection. This is because the probability signals prior to smoothing exhibit high frequency fluctuations, which may complicate the localization of the local maxima in the signals. The first cause of such high frequency changes is patch misplacement in the accumulated difference images. Recall that we extract the patches around the weighted centroids of the accumulated difference images. However, a large amount of noise content in the difference images may cause the weighted centroid to deviate from the center of the restored wavelet. In this case, the extracted patch may partially or completely miss the restored wavelet. This can manifest itself as a sudden change in the CNN output and as a result in the corresponding probability signal. The second cause of high frequency changes is the inherited high variance of CNNs. Use of ensemble of CNNs and data augmentation can alleviate this problem at a significant computation cost. Alternatively, we choose to mitigate these undesirable fluctuations using Gaussian smoothing for computational efficiency.

As described in Section 3.2, we constrain our ROI localization method by the location of the carotid bulb. This is because the bulb area appears as a relatively distinct dark area in the ultrasound frame. The distinct appearance of the carotid bulb is also confirmed by our experiments, where

we obtain the average bulb localization error of $0.26\,\mathrm{mm}$ for the test subjects with only one failure case, which is more favorable than the average unconstrained ROI localization error of $0.38\,\mathrm{mm}$ with 3 failure cases. Therefore, the localization of the bulb area can be done more reliably than the localization of the ROI, which motivates the use of the bulb location as a guide for more accurate ROI localization. We integrate this constraint into our localization system through a post-processing mechanism (see Eq. 1). Alternatively, we could train a regression CNN where each pixel in the image directly votes for the location of the ROI. However, this approach may be hindered by lack of stable anatomical structures in noisy ultrasound images. We will explore a regression CNN for ROI localization as future work.

In Section 4, we showed a high level of agreement between our system and the expert for the assessment of intima-media thickness. The suggested system achieves a mean absolute error of $2.8\,\mu\mathrm{m}$ with a standard deviation of $2.1\,\mu\mathrm{m}$ for intima-media thickness measurements. However, this level of measurement error cannot hurt the interpretation of the vascular age, because there exists a minimum difference of $400\,\mu\mathrm{m}$ between the average intima-media thickness of healthy and high-risk population ($600\,\mu\mathrm{m}$ for healthy and $\geq 1000\,\mu\mathrm{m}$ for high-risk population) [8]. To further put the performance of our system into perspective, in Table 2, we have compared the accuracy of intima-media thickness measurements produced by our system with those of the other automatic methods recently suggested in the literature. As seen, our method yields a lower level of mean absolute error and smaller standard deviation.

We used a LeNet-like CNN architecture in our study, but it does not limit the suggested framework to this architecture. In fact, we have experimented with deeper CNN architectures such as AlexNet [10] in both training and fine-tuning modes; however, we did not observe any significant performance gain. This was probably because the higher level semantic features detected by the deeper networks are not very relevant to the tasks in our CIMT applications. Meanwhile, the concomitant computational cost of deep architectures may hinder the applicability of our system, because it lowers the speed—a key usability factor of our system. We also do not envision that a shallower architecture can offer the performance required for clinical practice. This is because a network shallower than the LeNet has only one convolutional layer and thus limited to learning primitive edge like features. Detecting the carotid bulb and the ROI, and segmenting intima-media boundaries are relatively challenging tasks, requiring more than primitive edge-like features. Similarly, for frame selection, classifying the restored wavelets into R-peak and non-R-peak categories is similar to digit recognition, for which LeNet is a common choice of architecture. Therefore, LeNet-like CNN architecture seems to represent an optimal balance be-

Table 2: CIMT error for our system and the other state-of-the-art methods.

| Author | Year | Thickness error ($\mu\mathrm{m}$) |
|---|---|---|
| Current work | — | $2.8 \pm 2.1$ |
| Bastida-Jumillacite [1] | 2015 | $13.8 \pm 31.9$ |
| Ilea [7] | 2013 | $80 \pm 40$ |
| Loizou [14] | 2013 | $30 \pm 30$ |
| Molinari [20] | 2011 | $43 \pm 93$ |

tween efficiency and accuracy for CIMT video analysis.

We should note that throughout this paper, all performance evaluations were performed without involving any user interactions. However, our goal is not to exclude the user (sonographer) from the loop rather to relieve him from the three tedious, laborious, and time consuming operations by automating them while still offering the user a highly, user-friendly interface to bring his indispensable expertise onto CIMT interpretation through refining the automatic results easily at the end of each of the automated operations. For instance, our system is expected to automatically locate a EUF within one frame, which is clinically acceptable, but in case the automatic selected EUF is not the exact one as desired, the user can simply press an arrow key to move one frame forward or backward. From our experience, the automatically localized ROI is acceptable even if there is a small distance from the ground truth location, but the user still can easily drag the ROI and move it around as desired. Finally, in refining the automatically identified lumen-intima and media-adventitia interfaces, the original snake formulation comes with spring forces for user interaction [9], but given the small distance between the lumen-intima and media-adventitia interfaces, we have found that "movable" hard constraints as proposed in [12] are far more effective than the spring forces in measuring CIMT.

## 6. Conclusion

In this paper, we presented a unified framework to fully automate and accelerate CIMT video interpretation. Specifically, we suggested a computer-aided CIMT measurement system with three components: (1) automatic frame selection in CIMT videos, (2) automatic ROI localization within the selected frames, (3) automatic intima-media boundary segmentation within the localized ROIs. We based each of the above components on a CNN with a LeNet-like architecture and then boosted the performance of the employed CNNs with effective pre- and post-processing techniques. For frame selection, we demonstrated that how patch binarization as a pre-processing step and smoothing the probability signals as a post-processing step improve the results generated by the CNN. For ROI localization, we experimentally proved that the location of the carotid bulb, as a constraint in a post-processing setting, significantly improves

ROI localization accuracy. For intima-media boundary segmentation, we employed open snakes as a post processing step to further improve the segmentation accuracy. We compared the results produced by the suggested system with those of the major prior works, demonstrating more accurate frame selection, ROI localization, and CIMT measurements. This superior performance is attributed to the effective use of CNNs coupled with pre- and post- processing steps, uniquely designed for the three CIMT tasks.

# References

[1] M. Bastida-Jumilla, R. Menchón-Lara, J. Morales-Sánchez, R. Verdú-Monedero, J. Larrey-Ruiz, and J. Sancho-Gómez. Frequency-domain active contours solution to evaluate intima–media thickness of the common carotid artery. *Biomedical Signal Processing and Control*, 16:68–79, 2015.

[2] M. C. Bastida-Jumilla, R. M. Menchón-Lara, J. Morales-Sánchez, R. Verdú-Monedero, J. Larrey-Ruiz, and J. L. Sancho-Gómez. Segmentation of the common carotid artery walls based on a frequency implementation of active contours. *Journal of digital imaging*, 26(1):129–139, 2013.

[3] D.-C. Cheng and X. Jiang. Detections of arterial wall in sonographic artery images using dual dynamic programming. *Information Technology in Biomedicine, IEEE Transactions on*, 12(6):792–799, 2008.

[4] S. Delsanto, F. Molinari, P. Giustetto, W. Liboni, S. Badalamenti, and J. S. Suri. Characterization of a completely user-independent algorithm for carotid artery segmentation in 2-d ultrasound images. *Instrumentation and Measurement, IEEE Transactions on*, 56(4):1265–1274, 2007.

[5] F. Faita, V. Gemignani, E. Bianchini, C. Giannarelli, L. Ghiadoni, and M. Demi. Real-time measurement system for evaluation of the carotid intima-media thickness with a robust edge operator. *Journal of Ultrasound in Medicine*, 27(9):1353–1361, 2008.

[6] R. T. Hurst, R. F. Burke, E. Wissner, A. Roberts, C. B. Kendall, S. J. Lester, V. Somers, M. E. Goldman, Q. Wu, and B. Khandheria. Incidence of subclinical atherosclerosis as a marker of cardiovascular risk in retired professional football players. *The American journal of cardiology*, 105(8):1107–1111, 2010.

[7] D. E. Ilea, C. Duffy, L. Kavanagh, A. Stanton, and P. F. Whelan. Fully automated segmentation and tracking of the intima media thickness in ultrasound video sequences of the common carotid artery. *Ultrasonics, Ferroelectrics, and Frequency Control, IEEE Transactions on*, 60(1), 2013.

[8] D. Jacoby, I. Mohler, EmileR., and D. Rader. Noninvasive atherosclerosis imaging for predicting cardiovascular events and assessing therapeutic interventions. *Current Atherosclerosis Reports*, 6(1):20–26, 2004.

[9] M. Kass, A. Witkin, and D. Terzopoulos. Snakes: Active contour models. *International journal of computer vision*, 1(4):321–331, 1988.

[10] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.

[11] B. B. Le Cun, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel. Handwritten digit recognition with a back-propagation network. In *Advances in neural information processing systems*. Citeseer, 1990.

[12] J. Liang, T. McInerney, and D. Terzopoulos. United snakes. *Medical image analysis*, 10(2):215–233, 2006.

[13] Q. Liang, I. Wendelhag, J. Wikstrand, and T. Gustavsson. A multiscale dynamic programming procedure for boundary detection in ultrasonic artery images. *Medical Imaging, IEEE Transactions on*, 19(2):127–142, 2000.

[14] C. Loizou, T. Kasparis, C. Spyrou, and M. Pantziaris. Integrated system for the complete segmentation of the common carotid artery bifurcation in ultrasound images. In H. Papadopoulos, A. Andreou, L. Iliadis, and I. Maglogiannis, editors, *Artificial Intelligence Applications and Innovations*, volume 412 of *IFIP Advances in Information and Communication Technology*, pages 292–301. Springer Berlin Heidelberg, 2013.

[15] C. P. Loizou. A review of ultrasound common carotid artery image and video segmentation techniques. *Medical & biological engineering & computing*, 52(12):1073–1093, 2014.

[16] C. P. Loizou, C. S. Pattichis, M. Pantziaris, and A. Nicolaides. An integrated system for the segmentation of atherosclerotic carotid plaque. *Information Technology in Biomedicine, IEEE Transactions on*, 11(6):661–667, 2007.

[17] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. *arXiv preprint arXiv:1411.4038*, 2014.

[18] R.-M. Menchón-Lara, M.-C. Bastida-Jumilla, A. González-López, and J. L. Sancho-Gómez. Automatic evaluation of carotid intima-media thickness in ultrasounds using machine learning. In *Natural and Artificial Computation in Engineering and Medical Applications*, pages 241–249. Springer, 2013.

[19] R.-M. Menchón-Lara and J.-L. Sancho-Gómez. Fully automatic segmentation of ultrasound common carotid artery images based on machine learning. *Neurocomputing*, 151:161–167, 2015.

[20] F. Molinari, K. Meiburger, G. Zeng, A. Nicolaides, and J. Suri. Caudles-ef: Carotid automated ultrasound double line extraction system using edge flow. *Journal of Digital Imaging*, 24(6):1059–1077, 2011.

[21] F. Molinari, K. M. Meiburger, L. Saba, G. Zeng, U. R. Acharya, M. Ledda, A. Nicolaides, and J. S. Suri. Fully automated dual-snake formulation for carotid intima-media thickness measurement a new approach. *Journal of Ultrasound in Medicine*, 31(7):1123–1136, 2012.

[22] F. Molinari, G. Zeng, and J. S. Suri. A state of the art review on intima–media thickness (imt) measurement and wall segmentation techniques for carotid ultrasound. *Computer methods and programs in biomedicine*, 100(3):201–221, 2010.

[23] S. Petroudi, C. Loizou, M. Pantziaris, and C. Pattichis. Segmentation of the common carotid intima-media complex in ultrasound images using active contours. *Biomedical Engineering, IEEE Transactions on*, 59(11):3060–3069, 2012.

[24] P. Pignoli and T. Longo. Evaluation of atherosclerosis with b-mode ultrasound imaging. *The Journal of nuclear medicine and allied sciences*, 32(3):166–173, 1987.

[25] A. C. Rossi, P. J. Brands, and A. P. Hoeks. Automatic localization of intimal and adventitial carotid artery layers with noninvasive ultrasound: a novel algorithm providing scan quality control. *Ultrasound in medicine & biology*, 36(3):467–479, 2010.

[26] H. Sharma, R. G. Golla, Y. Zhang, C. B. Kendall, R. T. Hurst, N. Tajbakhsh, and J. Liang. Ecg-based frame selection and curvature-based roi detection for measuring carotid intima-media thickness. In *SPIE Medical Imaging*, pages 904016–904016. International Society for Optics and Photonics, 2014.

[27] P.-J. Touboul, P. Prati, P.-Y. Scarabin, V. Adrai, E. Thibout, and P. Ducimetière. Use of monitoring software to improve the measurement of carotid wall thickness by b-mode imaging. *Journal of hypertension*, 10:S37–S42, 1992.

[28] X. Xu, Y. Zhou, X. Cheng, E. Song, and G. Li. Ultrasound intima–media segmentation using hough transform and dual snake model. *Computerized Medical Imaging and Graphics*, 36(3):248–258, 2012.

[29] X. Zhu, C. B. Kendall, R. T. Hurst, and J. Liang. A user friendly system for ultrasound carotid intima-media thickness image interpretation. In *SPIE Medical Imaging*, pages 79681G–79681G. International Society for Optics and Photonics, 2011.

# Supplementary material

## Convolutional Neural Networks

As with multi-layer prceptrons, convolutional neural networks are trained using the back-propagation algorithm. If $D$ denotes a set of training images, $W$ denotes a matrix containing the weights of the convolutional layers, and $f_W(D^{(i)})$ denotes the loss for the $i^{th}$ training image, the loss over the entire training set is then computed as

$$\mathcal{L}(W) = \frac{1}{|D|} \sum_i^{|D|} f_W(X^{(i)}) \qquad (4)$$

Gradient descent is commonly used for minimizing the above loss function with respect to the unknown weights $W$. However, the modern massively parallelized implementations of CNNs are limited by the amount of memory on GPUs; therefore, one cannot evaluate the loss function based on the entire training set $D$ at once. Instead, the loss function is approximated with the loss over the mini-batches of training images of size $N << |D|$. A common choice of the mini-batch size is 128, which is a reasonable trade-off between loss noise suppression and memory management. Given the size of mini-batches, one can approximate the loss function as $\mathcal{L}(W) \approx \frac{1}{N} \sum_{i=1}^N f_W(X^{(i)})$, and iteratively update the weights of the network with the following equations:

$$\gamma_t = \gamma^{\left\lfloor \frac{tN}{|D|} \right\rfloor}$$
$$V_{t+1} = \mu V_t - \gamma_t \alpha \Delta L(W_t)$$
$$W_{t+1} = W_t + V_{t+1} \qquad (5)$$

where $\alpha$ is the learning rate, $\mu$ is the momentum that indicates the contribution of the previous weight update in the current iteration, and $\gamma$ is the scheduling rate that decreases learning rate $\alpha$ at the end of each epoch.
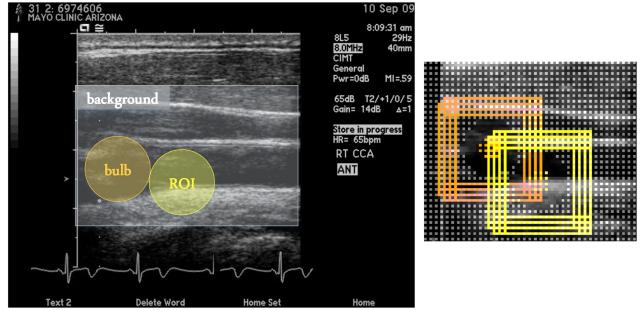
# Figures



Figure 9: For constrained ROI localization, we use a 3-way CNN whose training image patches are extracted from a grid of points on the background and around the ROI and the carotid bulb locations.
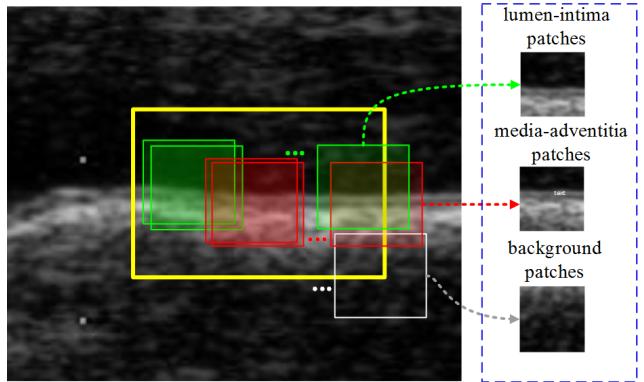


Figure 10: For lumen-intima and media-adventitia interface segmentation, we use a 3-way CNN whose training image patches are extracted from the background and around the lumen-intima and media-adventitia interfaces.
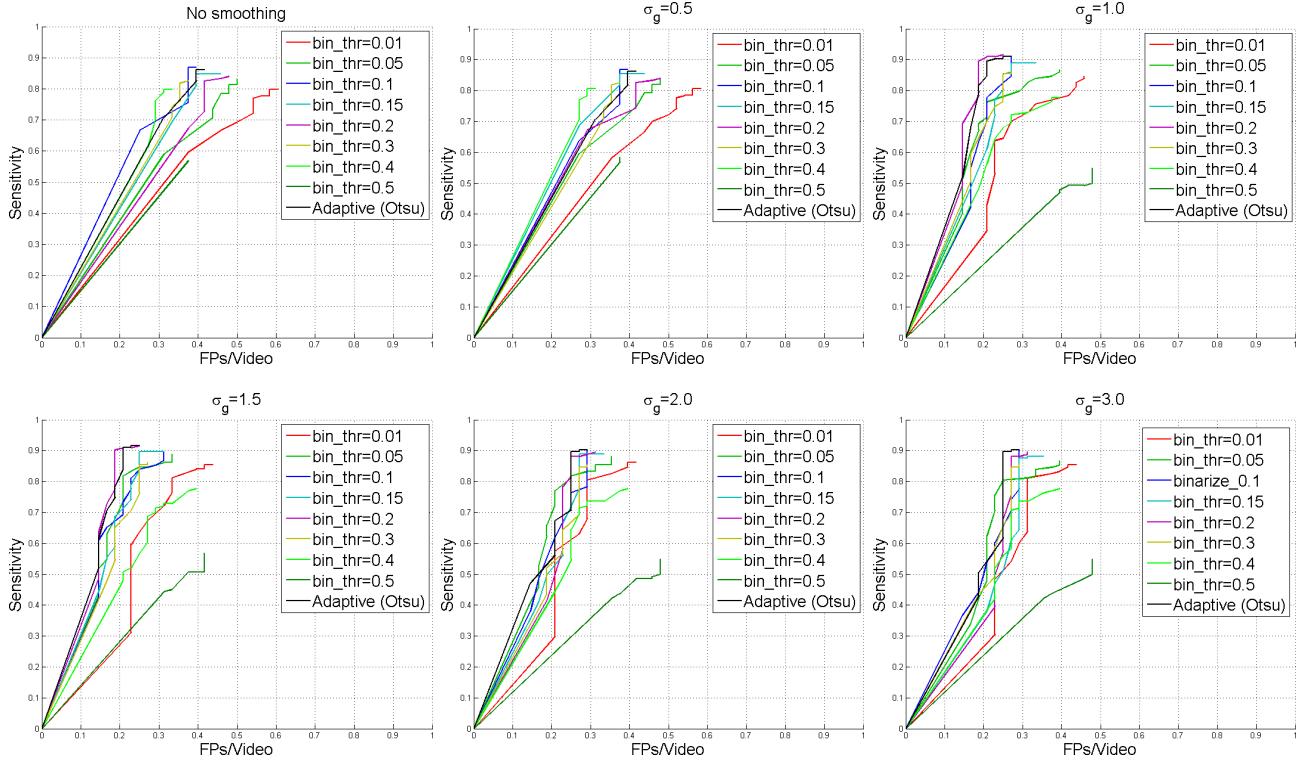
Figure 11: FROC curves of our system for automatic frame selection. Each plot shows FROC curves for different binarization thresholds and different levels of Gaussian smoothing. The results are obtained using leave-one-patient-out cross-validation based on the training subjects. As seen, no smoothing or a small degree of Gaussian smoothing leads to relatively low frame selection accuracy. This is because a trivial level of smoothing may not properly handle the fluctuations in the probability signals, causing a large number of false positives around an EUF. On the other hand, a large degree of smoothing may decrease the sensitivity of frame selection as the locations of the local maxima may be found more than one frame away from the expert-annotated EUFs. We therefore use a Gaussian function with $\sigma_g = 1.5$ for smoothing the probability signals. Our results also indicate that the adaptive thresholding method and a fixed threshold of 0.2 achieve the highest frame selection accuracy. However, we choose to use adaptive thresholding because it decreases the parameters of our system by one and that it performs more consistently at different levels of Gaussian smoothing.
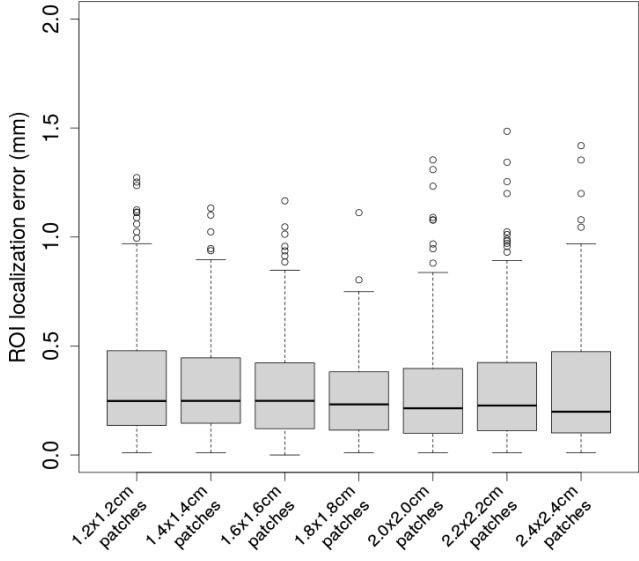
Figure 12: ROI localization error of our system for different sizes of patches. The results are obtained using leave-one-patient-out cross-validation based on the training subjects. In our analyses, we measure the localization error as the Euclidean distance between the estimated ROI location and the one provided by the expert. As can be seen, the use of $1.8 \times 1.8$ cm patches achieves the most stable performance, yielding low ROI localization error with only a few outliers.
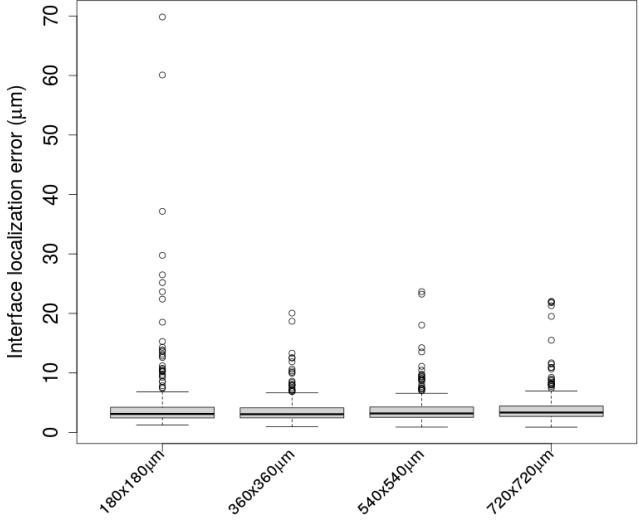


Figure 13: Combined interface localization error for different sizes of patches. The results are produced through a leave-one-patient-out cross-validation study based on the training subjects. Each box plots show the combined localization error of lumen-intima and media-adventitia interfaces for a different size of patches. In our analyses, we determine the localization error as the average of absolute vertical distances between our detected boundaries and the expert-annotated boundaries for the interfaces. As can be seen, while our system shows a high degree of robustness against different sizes of input patches, the use of patches of size $360 \times 360$ μm achieves slightly lower localization error and fewer outliers. Furthermore, this choice of patches yields higher computational efficiency compared to the larger counterpart patches.

## Agreement Analysis

We further analyze agreement between our system and the expert for the assessment of intima-media thickness. To this end, we use the Bland-Altman plot, which is a well-established technique to measure agreement between different observers. We have shown the Bland-Altman plot for the test subjects in Figure 14, where each circle represents a pair of thickness measurements, one from our method and one from the expert. As seen, the majority of circles fall within 2 standard deviations from the mean error, which suggests a large agreement between the automatically computed thickness measurements and those of the expert. Furthermore, Pearson product-moment correlation coefficient for the average and difference measurements is -0.097, indicating that the agreement between our method and the expert does not depend on intima-media thickness.
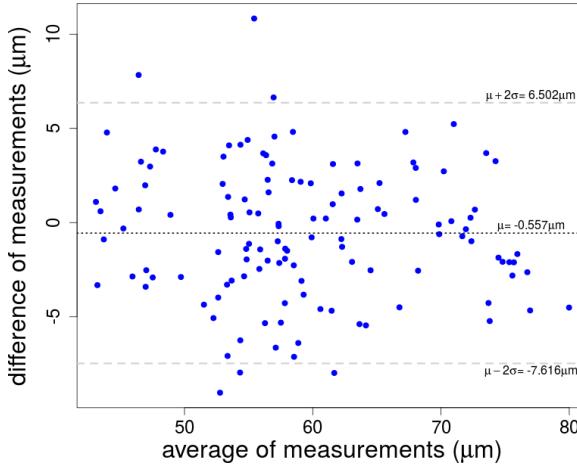


Figure 14: The Bland-Altman plot shows high agreement between our system and the expert for the assessment of intima-media thickness. Each circle in this plot represents a pair of thickness measurements from our method and the expert for a test ROI. In this plot, we have a total of 126 circles corresponding to 44 test videos.