
Enabling better pregnancy monitoring: The case of point-of-care diagnosis in fetal echocardiography

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 A major challenge in pre-natal healthcare delivery is the lack of devices and clin-
2 icians in several areas of the developing world. While the advent of portable
3 ultrasound machines and more recently, handheld probes, have brought down the
4 capital costs, the shortage of trained manpower is a serious impediment towards
5 ensuring the mitigation of maternal and infant mortality. Diagnosis of pre-natal
6 ultrasound towards several key pre-natal health indicators can be modelled as an im-
7 age analysis problem amenable to present day state-of-the art deep learning based
8 image and video understanding pipelines. However, deep learning based analysis
9 typically involves memory intensive models and the requirement of significant
10 computational resources, which is a challenging prospect in point-of-care health-
11 care applications. With the advent of portable ultrasound systems, it is increasingly
12 possible to expand the reach of automated prenatal health diagnosis. To accomplish
13 that, there is a need for lightweight architectures that can perform image analysis
14 tasks without a large memory or computational footprint. We propose a lightweight
15 convolutional architecture for assessment of ultrasound videos, suitable for those
16 acquired using mobile probes or converted from a DICOM standard from portable
17 machines. As exemplar of approach, we validated our pipeline for fetal heart
18 assessment (a first step towards identification of congenital heart defects) inclusive
19 of viewing plane identification and visibility prediction in fetal echocardiography.
20 This was attempted by models using optimised kernel windows and the construction
21 of image representations using salient features from multiple scales with relative
22 feature importance gauged at each of these scales using weighted attention maps
23 for different stages of the convolutional operations. Such a representation is found
24 to improve model performances at significant economization of model size, and
25 has been validated on real-world clinical videos.

26 1 Introduction

27 A key aspect of the UN Sustainable Development Goals relates to improving reproductive, maternal,
28 newborn and child health. A primary angle to the improvement of maternal and pre-natal health is the
29 adequate monitoring and assessment of fetal growth and abnormalities, so as to devise prognostic and
30 diagnostic measures in the event of possible adverse outcomes such as birth anomalies and congenital
31 diseases, the management and cure for some of which require advanced pre-planning even prior to
32 birth due to the technological and capital requirements involved in the management and redressal of
33 several such birth anomalies. Fetal ultrasound is the primary technique for prenatal health monitoring
34 and diagnosis, with several other modalities being restricted. Particularly, Congenital Heart Diseases
35 (CHDs) are responsible for driving infant mortality with rates being 8 in 1000 live births [4], and is
36 therefore a good case study for assessing the efficacy of automated point-of-care systems for pre-natal
37 healthcare delivery. Despite the universally acknowledged applications for ultrasound, systems of

38 image acquisition continue to be expensive and trained manpower is in short supply. Thus, usage of
39 automated image analysis systems built using machine learning algorithms is a potential avenue for
40 improving fetal health monitoring. In recent years, as deep learning based approaches became popular
41 for image processing applications, the size, computational requirements and complexity of models
42 along with data requirements remained a bottleneck towards deployment for point-of-care applications
43 for medical image analysis, despite rapid development in hardware for acquiring ultrasound scans with
44 the help of portable probes and mobile devices. An important clinical step in fetal heart ultrasound
45 characterisation, and essential for prognosis relevant to detection and management of CHDs, is the
46 visibility inference (whether or not the heart is visible in the frame) and the standard viewing plane
47 (4-chamber, 3-Vessel or Left Ventricular Outflow Tract/LVOT) identification. While deep learning
48 methods rely on end-to-end classifications by the feature learning and aggregation capabilities of
49 convolutional networks, we propose to leverage the presence of specific objects and anatomical
50 features defining a viewing plane at multiple scales through a measure of relevance imposed by
51 progressive attention modules [2]. This self-contained measure of importance of features in input
52 allows the models to train only on the features most relevant to the classes under study at the expense
53 of background, thereby reducing the size of the parameter space for such characterisation. This
54 idea leads us to explore the possibility of using attention layers to improve predictive accuracies of
55 lightweight architectures developed for mobile vision application [2,3].

56 2 Methodology

57 We attempt to improve the state of mobile ultrasound interpretation by constructing memory efficient
58 mobile deep learning architectures and augmenting the capacity and classification accuracies of the
59 lightweight models so developed by incorporation of an element of hierarchical prioritization of
60 information in the feature space through the use of stage-wise attention maps in the convolution
61 architecture. The idea is that while the usage of customised convolutional layers that use sets of 1x1
62 and 3x3 filters, with the former serving to impose separation in the depth level of the feature maps,
63 can reduce model size and computational cost, this comes at a reduction in the number of parameters
64 (not necessarily redundant). Such a reduced parameterization without controlling for parameter
65 importance to network decisions adversely affect performance for the given task. This performance
66 loss is reduced in the presented approach by the use of weighted attention mechanisms, where the
67 input images are partitioned into zones that are subsequently weighed to evaluate their contribution
68 towards the final classwise conditional likelihood for the whole image. Such attention based weighing
69 allows improvement in classification without reliance on extraneous model parameters. The role of
70 attention mechanisms in visual understanding of CNNs have been an area of active research. We
71 attempt to identify spatial cues that are most salient in informing the decisions by the convolutional
72 network on the given input. With the parameter budget being constrained for model efficiency,
73 we draw inspirations from the human mind's ability to extract relevant information from a scene
74 towards forming representative knowledge. This is replicated by having a weighted parameterisation
75 of obtained attention maps so as to magnify the impact of the most relevant features in the input
76 space and subdue the background towards final classification probability distributions obtained at the
77 softmax probability layer (Fig. 1). This is effectively a trainable mobile attention module, and can be
78 used at multiple locations in the architecture. The base architecture is inspired by the aggregation
79 of squeeze and excite modules introduced in [1] by substituting larger kernels with 1x1 kernels in
80 multiple layers and using 1x1 plus 3x3 kernels in alternate layers with the proportion of 3x3 filters
81 gradually increased to account for the complexity of neighborhood fine information in higher levels.
82 Each layer, represented as a set $s \in \{1, \dots, S\}$, is developed by a set of 1x1 and 3x3 filters that generate
83 the corresponding feature maps for every member of s as $F^s = \{f_1^s, f_2^s, \dots, f_n^s\}$. This specific
84 manner of representing feature maps is due to our interpretation that every member of the feature
85 map set, f_i^s , encodes the activations of spatial location i in layer s (each spatial location i is a square
86 region of a 100 x100 grid overlaid on a 2D feature map, so $1 \leq i \leq n$, $n=100$). With different feature
87 map dimensions across layers, the vector F^s has a variable length dimensions for constituent region
88 based encodings. This is resolved using a linear mapping for each of the three sets of F^s obtained to
89 map them to the dimension of that obtained using the final fully-connected layer F'^s , followed by a
90 dot product evaluation of each member of F^s with F'^s . This rationalisation with respect to the final
91 fully connected layer has an additional effect of capturing the overall global representation of the
92 input image as well. To obtain weighted attention over multiple layers, a softmax operation is applied
93 over the region-wise dot products with the final encoding. The attention weights so obtained are

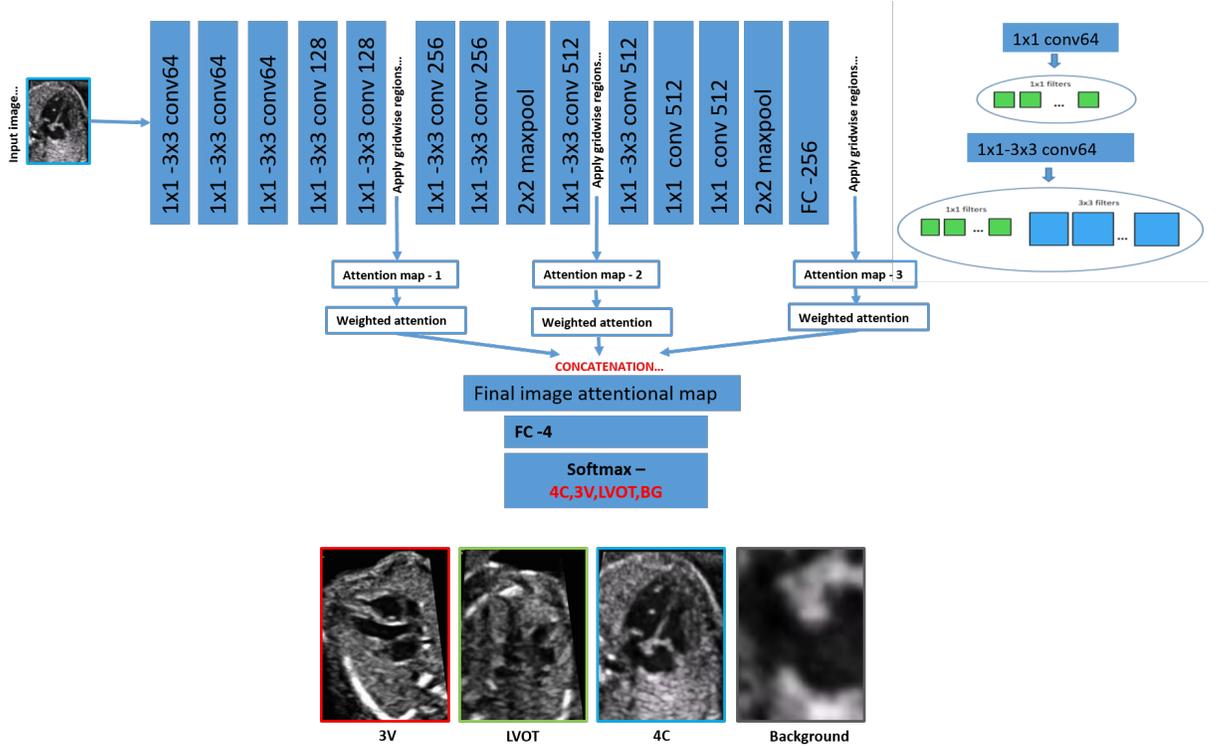


Figure 1: The overall architecture with attention maps at different stages. This is a representation of the configuration SN-att-2. For SN-att-1, the only attention map is after FC-512. There is a dot product with itself before being converted to the attention weight vector in that case, and this is followed by global concatenation towards creating the attention based representation.

Table 1: Performance of our attention driven models

Classification Accuracy (percentage)					
Method	4C	3V	LVOT	Non-standard/BG	Overall
Baseline [1]	85.42	70.14	65.71	80.13	75.35
SN-att-1	86.38	78.20	66.12	84.32	78.76
SN-att-2	88.60	78.95	69.34	83.52	80.10

94 assigned to each grid region defined for the attention map, and thus are a measure of the contribution
 95 of such a region to the overall loss function.

96 $a_i^s = \exp(\langle f_i^s, F'^s \rangle) / \sum_j^n \exp(\langle f_j^s, F'^s \rangle)$, where \langle, \rangle represents the dot product operation, $1 \leq i \leq n$,
 97 $n=100$ here

98 The attention weights $\{a_1^s, a_2^s, \dots, a_n^s\}$ so obtained are used to construct the global weighted
 99 attention per feature map $m_a^s = \sum_i^n a_i^s f_i^s$ and concatenated to obtain a final layer $M = [m_a^1, m_a^2$
 100 $a, \dots, m_a^S]$ called the Final Image Attention Map in Fig.1, ($S=3$ in our case). This is followed by a
 101 fully-connected layer FC-4 in Fig.1 with C nodes (C is the number of classes considered, $C=4$ here)
 102 for operations of softmax classification loss functions to obtain a classwise probability map. Thus, in
 103 effect, the concatenation of the weighted attention maps is used as a substitute for the fully-connected
 104 layer driven global image representation for image classification. This operation ensures that different
 105 feature regions at multiple scales of processing in convolutional stages are weighted directly and used
 106 to inform the final softmax cross-entropy classifier, instead of just using the fully-connected layer
 107 obtained by sequential convolutional operations.

108 3 Results

109 We start with a limited number of 91 cardiac screening videos from 12 subjects with gestational ages
110 ranging from 20 and 35 weeks during routine clinical scans. The duration of each video is between 2
111 and 10 seconds and a frame rate of 25 to 76 frames per second (39556 frames in total). It contained
112 one or more of the three views of the fetal heart and some background frames. Videos from 10
113 patients are used for training, and the remaining 2 for test experiments. For training, we split available
114 videos into frames and apply data augmentation by an updown and a top-bottom flipping. Individual
115 frames of size 430 x 510 are cropped into 224 x 224 centred about the heart centre, which was known
116 in the ground truth annotations. The models are trained using a batch size of 25 with a learning rate
117 of 0.001. A training:test split of 80:20 is used. The base architecture has no pooling layers till the
118 fifth 1x1-3x3 layer module to make feature maps available at a higher resolution to make regional
119 attention proposals optimally informative. We derive our attention maps from layer modules 5, 9
120 and 14. These maps are obtained as a set of encodings from a grid of regions obtained by dividing
121 the two-dimensional feature map in a 100 x 100 patch set. The encodings have associated weights
122 parameterized by a weight matrix. In the absence of established baselines in prior work for mobile
123 based classification in fetal echocardiography datasets, we compare the results obtained for with a
124 standard SqueezeNet architecture [1] adapted for handling our ultrasound image data, and our base
125 model with attention (SN-att-1 and SN-att-2 with SN-att-1 aggregating the attention layer from the
126 final fully connected layer and assuming different sections as representative of grid-regions in the prior
127 feature maps) for a classification of visibility and viewing planes in fetal echocardiography images.
128 The attention-based approach yields a notable performance improvement despite a negligible addition
129 to the model size (1.90 MB in baseline without attention vs 2.24 MB in SN-att-2), with the overall
130 baseline SqueezeNet accuracy of 75.35 exceeded by both versions of our attention based architectures
131 (78.76 and 80.10). The original SqueezeNet model adapted for this architecture is a heavier model
132 as well. Additionally, the inclusion of weighted attention improves performance in case of difficult
133 classes like 3V (78.20 and 78.95 vs 70.14 in baseline) and LVOT (66.12 and 69.34 vs 65.71). This is
134 because the weighted attention model allows enhanced reliance on finegrained discriminative features
135 and relatively ignores less-important features in the classification stages. The strategy to include
136 attention layers from different sections of the network as different sections learn different attributes of
137 the image is proven to enable better aggregation of salient features through the improvement by from
138 SN-att-1 (78.76) to SN-att-2 (80.10). To conclude, the ability of attentive classification to focus on
139 relevant features and diminish the role of the background effectively is reflected in the improved top-1
140 accuracies listed. Such an improvement without a large model complexity addition is of importance
141 in low-compute environments as in mobiles and EDGE devices in the clinical ultrasound space.
142 It is worth considering comparisons with quantization models, direct classification baselines from
143 deeper architectures and attention grids with variable resolutions. As of now, this work has been
144 attempted with competitive accuracies on actual clinical echocardiography videos, after conversion
145 from the DICOM standard to standard avi formats, which are then processed in our pipelines (the
146 video preparation and the learning/inference stages are therefore separate here). This conversion, is
147 integrated into our method. As future extension to the validations presented, it would be worthwhile
148 to port the pre-trained models directly onto mobile or other devices along with integration to support
149 input video streams derived using connected probes, similar to the demonstrations attempted for
150 ultrasound to mobile video conversions using handheld probes by industry players like Butterfly
151 Network Inc, Clarius and so on. That way, the processing and diagnosis step can be integrated with
152 the real-time acquisition and the whole pipeline can be used end-to-end, with a possible forward
153 integration to cloud services for later quality and diagnosis checks by qualified physicians located
154 physically away from the patient locations.

155 References

- 156 [1] Iandola, Forrest N., et al. "Squeezenet: Alexnet-level accuracy with 50x fewer parameters and < 0.5 mb
157 model size." arXiv preprint arXiv:1602.07360 (2016).
- 158 [2] Seo, Paul Hongsuck, et al. "Hierarchical attention networks." CoRR, abs/1606.02393 (2016).
- 159 [3] Jetley, Saumya, et al. "Learn to pay attention." arXiv preprint arXiv:1804.02391 (2018).
- 160 [4] N. Archer and N. Manning. Fetal Cardiology. Oxford University Press, 2009