

Research Article

# Temporal Feature Weaving for Neonatal Echocardiographic Viewpoint Video Classification

Satchel French<sup>1</sup>, Faith Zhu<sup>2</sup>, Amish Jain<sup>2</sup>, Naimul Khan<sup>1</sup>

1. Electrical, Computer, and Biomedical Engineering, Toronto Metropolitan University, Toronto, Canada; 2. Pediatrics, Mount Sinai Hospital, Toronto, Canada

**Automated viewpoint classification in echocardiograms can help under-resourced clinics and hospitals in providing faster diagnosis and screening when expert technicians may not be available. We propose a novel approach towards echocardiographic viewpoint classification. We show that treating viewpoint classification as video classification rather than image classification yields advantage. We propose a CNN-GRU architecture with a novel temporal feature weaving method, which leverages both spatial and temporal information to yield a 4.33% increase in accuracy over baseline image classification while using only four consecutive frames. The proposed approach incurs minimal computational overhead. Additionally, we publish the Neonatal Echocardiogram Dataset (NED), a professionally-annotated dataset providing sixteen viewpoints and associated echocardiography videos to encourage future work and development in this field. Code available at: <https://github.com/satchelfrench/NED>**

## 1. Introduction

This work aims to solve the problem of neonatal echocardiogram viewpoint classification by utilizing modern machine learning algorithms. This could broaden access to non-specialized technicians, provide improved training experiences and overall aid in the acquisition of higher quality echocardiograms. In a scan the anatomy is observed from a variety of perspectives, aptly referred to as “viewpoints”. Some viewpoints are easy to identify for a novice, but several are only separable by minor details. Additionally, viewpoints can appear different between patients, pathology’s and scan environment. Prior approaches to viewpoint classification have treated it as a simple image

classification problem<sup>[1]</sup>, utilizing off-the-shelf CNN-based networks such as VGG16. Those techniques have been effective in establishing a solid baseline, but they suffer from confusion on harder to distinguish viewpoints and have greater sensitivity to input variation, leaving considerable room for improvement. Additionally, though there are solutions explored for adult echocardiograms or alternate anatomies in neonates, there is a lack of literature supporting the use case for neonatal echocardiograms. Considering the anatomy of a heart is dynamic, a particular viewpoint may appear different during an alternate position of the cardiac cycle. Prior single-image classification approaches have neglected this salient information by not considering the movement of structures within the viewpoint. We present a multi-frame approach using a gated recurrent unit (GRU) network, <sup>[2]</sup> enabling the network to factor temporal information into its classification output. We demonstrate that this fundamental consideration of temporal-spatial relationship is a key factor in producing state of the art accuracy in viewpoint classification. key contributions are:

- The first of its kind, professionally labelled, open-source neonatal echocardiogram dataset featuring sixteen viewpoints, in both video and image formats.
- A modified ResNet-GRU architecture with our Temporal Feature Weaving (TFW) method which produces a 4.53% improvement in F1-Score compared to baseline image classification.

Our final proposed method with temporal feature weaving comes at no additional cost in model size or compute. The final model features approximately 30 million parameters, which is capable of being executed in real-time on a modern smartphone, furthering the accessibility of neonatal echocardiography.

## 2. Proposed Approach

### 2.1. Dataset

Our work introduces the first publicly available Neonatal Echocardiogram Dataset (NED). The dataset contains a total of 16 classes across 40 patient cases, totaling 1049 videos approximately 1 second in length. The scans are acquired organically from real patient scans, and as such contain a natural class imbalance. A subset of 12 common classes has been chosen as the primary focus of this work given their clinical relevance, scan homogeneity and sample size. These two subsets are referred to as the NED-16 and NED-12 respectively. The class balance of both datasets can be seen in Figure 1, where it should be highlighted that the three classes with the least samples, were removed to improve balance.

A fourth class, representing the Ductal Cut viewpoint is removed as it presents with a high variation between scans and is less commonly captured. The remaining subset of 12 classes and views are selected to be consistent with a generally standard capture procedure<sup>[3]</sup>.

Each scan in the dataset has been labelled by a physician with expertise and familiarity in the field and individual frames are organized respectively by patient, viewpoint, and video clip. Our dataset is accompanied by open-source PyTorch data-loaders that support loading frames as consecutive or evenly spaced sequences. Additionally, it includes sequence-wide augmentation transforms for rotation, scale, shift, flips and contrast adjustment.<sup>1</sup>

Neonatal Echocardiogram Dataset (NED) offers the following 16 classes for viewpoints: Apical-4-Chamber (*APICAL\_4C\_LVRV*), Apical-5-Chamber (*APICAL\_5C*), Apical-3-Chamber Right Ventricle (*APICAL\_3C\_RV*), Apical-3-Chamber Left Ventricle (*APICAL\_3C\_LV*), Apical-2-Chamber (*APICAL\_2C\_LV*), Parasternal Long Axis (PLAX) Left Ventricle (*PLAX\_LV*), PLAX Right Ventricle (RV) In-flow Focus (*PLAX\_RV\_IN*), PLAX Right Ventricle (RV) Out-flow focus (*PLAX\_RV\_OUT*), Parasternal Short Axis (PSAX) Apex (*PSAX\_APEX*), PSAX Papillary Muscle (*PSAX\_PAPS*), PSAX Mitral Valve (*PSAX\_MV*), PSAX Aortic Valve (*PSAX\_AV*), Pulmonary Artery Branch (*BRANCH\_PA*, NED-16 only), Ductal Cut (*DUCTAL\_CUT*, NED-16 only), Aortic Arch (*ARCH*, NED-16 only), and Subcostal Inferior Vena Cava (*SUBCOSTAL\_IVC*, NED-16 only).

Patient demographics available in Table 1 highlights the group statistic of dataset participants. Notable statistics include a moderate skew toward male patients at 61.8%, and median gestational age at time of scan of 32.15 weeks.

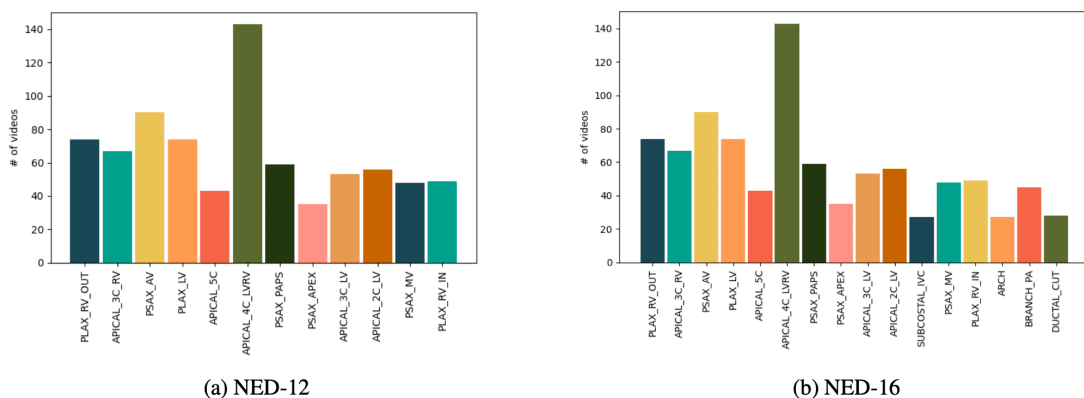


Figure 1. Class Distribution of NED

Metric	Value (std)
Patient Count	34
Gestational age at birth (weeks)	25.87 (1.84)
Gestational age at scan (weeks)	32.15 (5.2)
Age at scan (days)	46 (40)
Male sex (%)	21 (61.8)
Weight at scan (kg)	1.67 (1.01)
Systolic blood pressure (mmHg)	72 (13)
Diastolic blood pressure (mmHg)	42 (10))

**Table 1.** Patient Demographics

## 2.2. Model Architecture

The foundational ResNet-50 by He et al.<sup>[4]</sup> is used for both the image classifier and backbone of GRU based models. Given efficacy of VGG style networks for feature extraction<sup>[5][11]</sup>, ResNet is expected to provide same or better performance with fewer layers given its identity mapping via skip connections. In the context of GRU networks, ResNet enables the selection of both high-level and low-level features across time, providing richer feature combinations.

To compare image classifier performance, mean-voting of class probabilities across several frames was employed. It produces marginal improvement overall, and performs better on spaced sequences of frames, compared to a consecutive sequence. This is expected behaviour as given the model's conviction in its predictions, improvement would be seen only in cases where classes have a similar confidence, and would be better differentiated at alternate points during the cardiac cycle where anatomy is more dissimilar or imaging is less noisy.

### 2.2.1. Self-Attention Methods

To better understand the efficacy of this approach, a benchmark using an attention mechanism<sup>[6]</sup> is employed. This is a sensible approach given the success of attention models for sequence learning in

relevant tasks<sup>[7][8]</sup>. Several architectures were tested including the use of GRU with attention, but the best results were ultimately produced by splitting the 2048 length CNN feature outputs of four frames into a 512x16 sequence. This was fed into four attention heads, and the attended sequence was fed into an fully connected layer of 512x2048, before going to the 2048x12 class output layer.

### 2.3. Temporal Feature Weaving

To capture more meaningful temporal-spatial relationships we introduce the concept of temporal weaving. Temporal weaving divides the flattened feature vectors of each frame  $X_n$  into  $K$  segments (Equation 1), and reconstructs a new feature vector,  $W_k$ , by concatenating the respective chunks,  $C$ , between frames as seen in (Equation 2). Visually this is captured in the architecture Figure 2. In this paper we employ a simplistic approach to temporal weaving where  $K = N$ , meaning the weaved feature vector  $W$ , represents the same chunk evenly across time.

$$X_n = [C_{n1}, C_{n2}, C_{n3}, \dots, C_{n(K-1)}, C_{nK}] \quad (1)$$

$$W_k = [C_{1k}, C_{2k}, C_{3k}, \dots, C_{(N-1)k}, C_{Nk}] \quad (2)$$

Our approach produces  $N$  weaved feature vectors, which are given as a sequence input to the GRU. Predictions are made at each time-step  $t$ , with final class probabilities being taken from the final time-step.

Conceptually, this encourages the model to give attention to the same patch of the image across time, highlighting temporal information. Weaving the features this way provides a temporal signature in each sequence input, which can easily be obtained by contrasting chunks from different frames. This is thought to improve upon frame-sequence input as the ResNet backbone produces solely spatial features which are effectively compressed into smaller subspaces at each time-step in the GRU's hidden layer, making it unlikely that the minor differences across time persist as salient features. By weaving the features, the network is forced to make predictions on a single spatial-temporal feature patch, which when viewed as independent events, are implicitly chained together by the persistence of information in an RNN hidden layer. The end-of-The weaved features may be thought of as a special case where more relevant information persists to the end-of-time prediction.

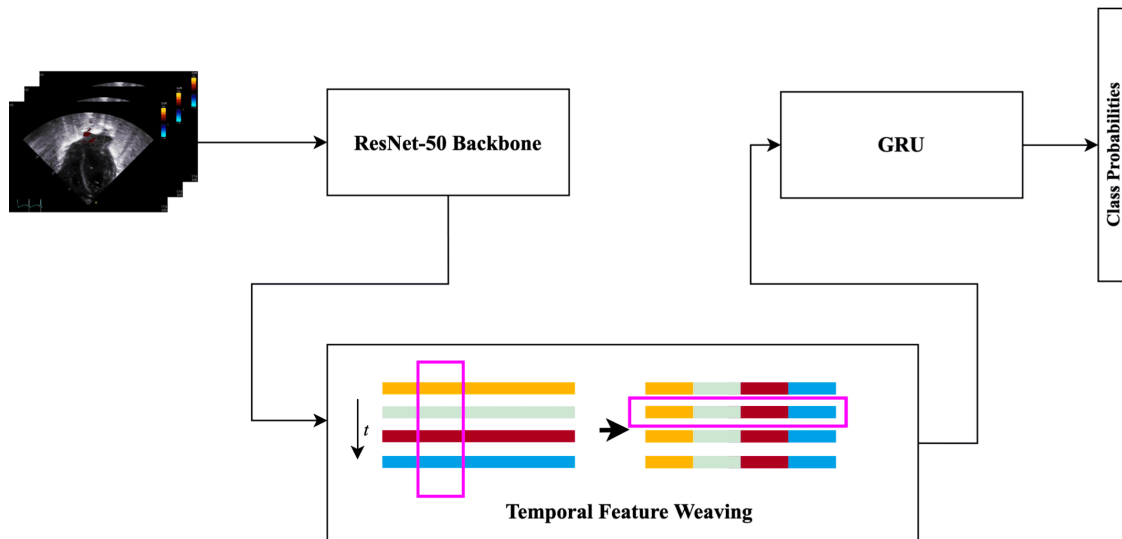


Figure 2. ResNet-50-GRU-Temporal-Weave Architecture

### 3. Experiments

#### 3.1. Training Parameters and Baseline Models

Data augmentation proved to be essential in producing quality results without over fitting the dataset. However, augmentation for echocardiograms requires more careful consideration to ensure that the transforms produced clinically viable samples. Our research indicated that a random rotation of [0, 25] degrees and auto contrasting yielded the best results. These findings are sensible from a clinical perspective and are consistent with literature<sup>[9]</sup>.

As a baseline image classifier architecture, ResNet-50<sup>[4]</sup> is trained from scratch on NED-12 (pre-training pm ImageNet provided no benefit). The model was trained for 50 epochs using cyclical learning rate scheduling with a maximum learning rate of ( $lr=0.08$ ) and batch size of ( $N=16$ ). It uses a modified gradient descent algorithm, with a momentum of ( $m=0.75$ ), L2 regularization of ( $d=5e-4$ ), and dropout of ( $p=0.6$ ). Finally, samples are generated by loading a single frame from random patients scans, resized to  $232 \times 232$ , cropped to  $224 \times 224$ , randomly augmented, and sampled with no-replacement for the epoch.

For the video classification model there was an obvious need to provide interleaved frames. While the model performs well with several sequential frames, by processing spaced frames it encodes temporal

information that may be unique to a viewpoint. Additionally, viewpoints may appear dramatically different throughout the cardiac cycle, with some parts of the anatomy moving slowly early in the scan, and then quickly moving during another segment.

When training the GRU based model, initial results failed to make meaningful improvements as training continued, likely due to the vanishing gradient problem recurrent neural networks suffer from. Although it's worth noting superior training performance was observed with GRU networks compared to LSTM, the reduction in loss was minimal. To better address these concerns, weights from the trained resnet-50 classifier were used, with outputs from the second-last fully connected layer being flattened into a feature vector of 2048 length. These features are generated for each frame in the sequence and weaved temporally before being fed into the model as a sequence of feature vectors. The GRU based model was also trained using cyclic learning rate scheduling with a maximum learning rate of ( $lr=0.003$ ) and a batch size of ( $N=16$ ) for 60 epochs. It also used Adam for gradient descent, with a momentum of ( $m=0.75$ ), L2 regularization of ( $d=1e-3$ ) and dropout ( $p=0.5$ ).

Models were initially evaluated on a random 90/10 train test split, and cross validated with ( $K=5$ ) folds. Folds are split patient wise to mitigate data leakage from any consistencies in a patients scan. Metrics are calculated during each fold and finally aggregated into accuracy, precision and recall scores.

### *3.2. Results*

Table 2 summarizes the results of experiments on the NED-12, including the proposed ResNet-GRU-TFW network, as well as comparisons to attention-based and mean-voting approaches. Our proposed solution of ResNet-GRU with Temporal Feature Weaving (TFW), produces an accuracy of 93.8%, and an F1-Score of 93.7% when evaluated on spaced frames outperforming all other models.

Model	Accuracy	Precision	Recall	F1
ResNet-50-Image-Classifer	89.55	88.99	89.55	89.20
Mean-Voting-4-Consecutive	89.73	90.66	89.73	90.01
Mean-Voting-4-Spaced	89.99	90.95	89.99	90.27
ResNet-GRU-4-Consecutive	88.46	88.84	88.46	88.44
ResNet-Attention-512-Consecutive	92.13	92.14	91.36	91.19
(Proposed) ResNet-GRU-TFW-4-Spaced	92.55	92.72	92.55	92.51
(Proposed) ResNet-GRU-TFW-4-Consecutive	93.88	93.65	93.88	93.73

Table 2. Results on NED-12

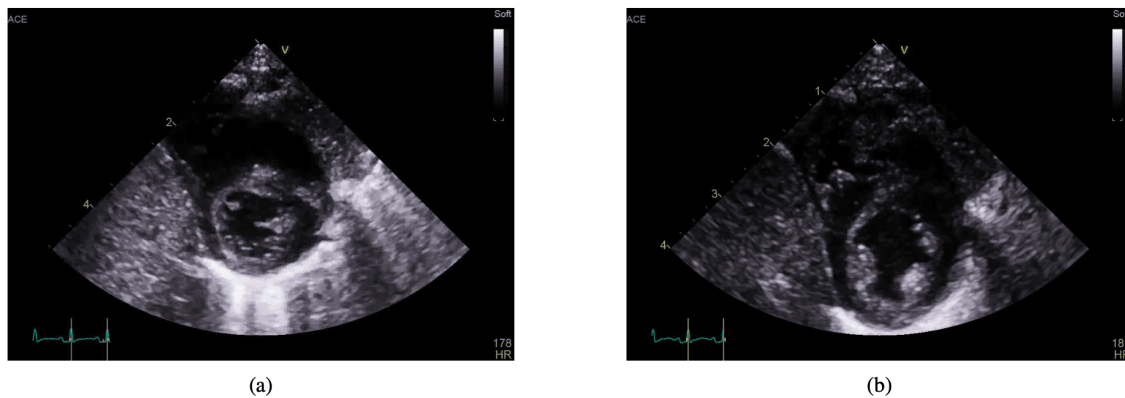


Figure 3. (a) Mitral-valve (*PSAX\_MV*) (b) Papillary muscle (*PSAX\_PAPS*)

## 4. Discussion

As seen in Table 2, there is a general progression of improvement in accuracy from single frame to multi-frame sequential and spaced configurations. Spaced-frame voting offers a better perspective throughout the cardiac cycle, where structures may appear meaningfully different. The ResNet-GRU architecture without temporal feature weaving performs worse across all metrics, showing a 0.76% decrease in F1 Score compared to the baseline. This reduced accuracy may stem from the model's



inability to memorize subtle differences between sequential frames, focusing instead on higher-level features. Recurrent neural networks compress input features by selecting which sequence elements to remember, potentially causing a lossy compression where certain discriminants are lost. The GRU may no longer access high-definition spatial features from the ResNet architecture, relying instead on high-level semantic features that underperform the baseline.

Substantial improvement is observed with the CNN-Attention approach, and weaving attended features made no difference as expected. This suggests the attention mechanism effectively identifies relevant patches for comparison across time, achieving the intended result. However, it still underperforms the TFW method by 2.4% in F1-Score. Additionally, multi-head attention adds training complexity and an extra 1 million parameters, whereas TFW introduces no such overhead.

Introducing TFW improves the F1 score by 3.31% and 4.53% for spaced and consecutive models, respectively. These gains over mean voting approaches support the hypothesis that considering the temporal signature of individual viewpoints enhances classification. Notably, the sequential case outperforms the spaced case with temporal feature weaving, contrary to the mean-voting scenario. This indicates an alternate prediction mechanism, where the sequential case better captures structural movements, suggesting that a unique temporal signature contributes to the model's predictions.

This effect is most evident when differentiating between static views of the parasternal short axis (PSAX) mitral valve and papillary muscle viewpoints. As shown in Figure 3, the viewpoints appear remarkably similar, differing mainly in a cropped focus on the papillary muscle and a clearer view of the mitral valve. The primary discriminator is the presence of the mitral valve, or "Fish Lips." However, since the valve is dynamic, it is not consistently identifiable by this feature alone from a single frame.

## 5. Conclusion

We propose a novel temporal feature weaving (TFW) approach to classify neonatal ultrasound videos into different viewpoints. Our work shows that tracking features across time adds value to viewpoint classification, which has not been explored in literature thoroughly. We also provide a public neonatal viewpoint classification dataset. Future work will focus on addressing limitations such as dataset size, using multiple labelers, and exploring other TFW approaches.

## Statements and Declarations

### *Compliance with ethical standards*

This study was performed in line with the principles of the Declaration of Helsinki. Approval was granted by the Research Ethics Board of the Mount Sinai Hospital (19-0065-E) and Toronto Metropolitan University (2019-389).

### *Acknowledgements*

This work was supported by the Natural Sciences and Engineering Research Council of Canada (Alliance Grant 546302-19 and Discovery Grant 2020-05471).

## Footnotes

<sup>1</sup> Github link to dataset and code to be provided upon acceptance.

## References

1. <sup>a, b</sup>Madani A, Arnaout R, Mofrad M, et al. Fast and accurate view classification of echocardiograms using deep learning. *Nature npj Digital Med.* 2018; 1(6). doi:[10.1038/s41746-017-0013-1](https://doi.org/10.1038/s41746-017-0013-1).
2. <sup>Δ</sup>Chung J, Gulcehre C, Cho K, Bengio Y (2014). "Empirical evaluation of gated recurrent neural networks on sequence modeling." *arXiv.* [arXiv:1412.3555 \[cs.NE\]](https://arxiv.org/abs/1412.3555).
3. <sup>Δ</sup>Groves AM, Singh Y, Dempsey E, Molnar Z, Austin T, El-Khuffash A, de Boode WP (2018). "Introduction to neonatologist-performed echocardiography." *Pediatric research.* 84 (Suppl 1): 1–12.
4. <sup>a, b</sup>He K, Zhang X, Ren S, Sun J (2015). "Deep residual learning for image recognition." *arXiv.* Available from: <https://arxiv.org/abs/1512.03385>.
5. <sup>Δ</sup>Simonyan K, Zisserman A (2015). "Very Deep Convolutional Networks for Large-Scale Image Recognition." *arXiv.* Available from: <https://arxiv.org/abs/1409.1556>.
6. <sup>Δ</sup>Vaswani A, Shazeer N, Parmar N, Uszkoreit J, Jones L, Gomez AN, Kaiser L, Polosukhin I (2023). "Attention Is All You Need." *arXiv.* Available from: <https://arxiv.org/abs/1706.03762>.
7. <sup>Δ</sup>Arnab A, Deghani M, Heigold G, Sun C, Lučić M, Schmid C (2021). "ViViT: A Video Vision Transformer." *arXiv.* [arXiv:2103.15691 \[cs.CV\]](https://arxiv.org/abs/2103.15691).

8. <sup>^</sup>Feng Z, Sivak JA, Krishnamurthy AK. "Two-stream attention spatio-temporal network for classification of echocardiography videos." In: 2021 IEEE 18th International Symposium on Biomedical Imaging (ISBI). 2021. p. 1461-1465. doi:[10.1109/ISBI48211.2021.9433773](https://doi.org/10.1109/ISBI48211.2021.9433773).
9. <sup>^</sup>Tupper A, Gagné C. "Analyzing Data Augmentation for Medical Images: A Case Study in Ultrasound Images." 2024. Available from: <https://arxiv.org/abs/2403.09828>.

## **Declarations**

**Funding:** This work was supported by the Natural Sciences and Engineering Research Council of Canada (Alliance Grant 546302-19 and Discovery Grant 2020-05471).

**Potential competing interests:** No potential competing interests to declare.