



Published in final edited form as:

Med Phys. 2022 June ; 49(6): 4071–4081. doi:10.1002/mp.15654.

Generative learning approach for radiation dose reduction in X-ray guided cardiac interventions

Fariba Azizmohammadi¹, Iñaki Navarro Castellanos², Joaquim Miró², Paul Segars³, Ehsan Samei³, Luc Duong¹

¹Interventional Imaging Lab, Department of Software and IT Engineering, École de technologie supérieure, Montreal, Canada

²Department of Pediatrics, CHU Sainte-Justine, Montreal, Canada

³Department of Radiology, Carl E. Ravin Advanced Imaging Laboratories, Duke University Medical Center, Durham, North Carolina, USA

Abstract

Background: Navigation guidance in cardiac interventions is provided by X-ray angiography. Cumulative radiation exposure is a serious concern for pediatric cardiac interventions.

Purpose: A generative learning-based approach is proposed to predict X-ray angiography frames to reduce the radiation exposure for pediatric cardiac interventions while preserving the image quality.

Methods: Frame predictions are based on a model-free motion estimation approach using a long short-term memory architecture and a content predictor using a convolutional neural network structure. The presented model thus estimates contrast-enhanced vascular structures such as the coronary arteries and their motion in X-ray sequences in an end-to-end system. This work was validated with 56 simulated and 52 patients' X-ray angiography sequences.

Results: Using the predicted images can reduce the number of pulses by up to three new frames without affecting the image quality. The average required acquisition can drop by 30% per second for a 15 fps acquisition. The average structural similarity index measurement was 97% for the simulated dataset and 82% for the patients' dataset.

Conclusions: Frame prediction using a learning-based method is promising for minimizing radiation dose exposure. The required pulse rate is reduced while preserving the frame rate and the image quality. With proper integration in X-ray angiography systems, this method can pave the way for improved dose management.

Keywords

cardiac interventions; radiation dose reduction; X-ray angiography

Correspondence Fariba Azizmohammadi, Interventional Imaging Lab, Department of Software and IT Engineering, École de technologie supérieure, Montreal H3C 1K3, Canada. fariba.azizmohammadi.1@ens.etsmtl.ca.

CONFLICT OF INTEREST

The authors have no conflicts to disclose.

1 | INTRODUCTION

Congenital heart disease (CHD) affects 1% of the population and is the most common type of birth malformation worldwide.¹ Patients with CHD are exposed to substantial amounts of ionizing radiation from diagnostic and treatment procedures.² In recent years, the number of complex, long-duration pediatric cardiac interventions has risen significantly. Consequently, the risks associated with radiation exposure among patients have also increased, which is why solutions must be found to reduce the radiation dose to as low as reasonably achievable while maintaining the required image quality.³ Minimizing radiation exposure in pediatric cardiology is paramount in interventional cardiology. Patients are subjected to either deterministic outcomes, such as skin necrosis, which is most commonly related to tissue rebounds, or stochastic effects, such as an increased risk of radiation-induced cancer and brain tumors.⁴ Moreover, complex CHDs must be catheterized repeatedly, thereby increasing the risk of radiation-induced cancer not only for patients but also for medical staff.⁵ Radiation exposure is, therefore, a major concern for pediatric populations, and determining the optimal dose for each patient is a highly relevant research topic in pediatric cardiology.

1.1 | Radiation dose reduction in X-ray angiography

Currently, X-ray angiography is widely accepted for minimally invasive interventions and provides adequate spatial and temporal image resolution. Fluoroscopy and fluorography are the two main fluoroscopically guided intervention modes in X-ray imaging. In fluoroscopy mode, the X-ray images are generated instantaneously and continuously to observe moving objects by capturing the motion. The images in this mode are not recorded and used to navigate the medical devices to specific locations within the patient in real-time. Fluorography mode requires a higher radiation exposure to generate and record high-resolution images for interpretation after the termination of the exposure.⁶ The required radiation dose for each acquisition mode is a function of the required image quality, the patient's size, and the time required to perform the procedure. Fluoroscopy time comprises the total time spent using fluoroscopy for image acquisition and is considered as one of the effective parameters for the final patient dosage.⁷

Previously, conventional analog X-ray equipment was used to deliver X-ray energy in a continuous dose. Recently, some strategies are applied to mitigate the radiation dose to the patients such as using the lowest possible fluoroscopic dose rate during live fluoroscopy, use of low frame rates (if possible), and use of multiple short fluoroscopic exposures instead of keeping the fluoroscope on continuously and minimizing the beam-on time for the fluoroscopy imaging.^{6,8}

Modern X-ray systems are equipped to deliver energy in pulses that can be adjusted to 7.5, 10, 15, and 30 frames per second (fps). In pulsed fluoroscopic imaging, the X-ray beam is switched on and off for every fluoroscopic frame, and thus the pulse width, or time duration of each frame, is lower than the time required in continuous fluoroscopy imaging. This allows for reducing the fluoroscopy time by replacing the continuous exposure with a pulsed beam delivery. However, images are temporally averaged and moving objects look unsharp and flicking. A sequence of pulsed images, including moving objects, appears more

continuous and less flickering at high pulse rates or frequencies based on the critical flicker frequency. At low frame rates, gap filling by replicating each acquired frame multiple times is applied to avoid flicker and minimize blurriness of moving targets. The term frame rate describes the number of frames that are generated per second, while the term pulse rate refers to the output of the fluoroscope, specifically the number of bursts of radiation that are emitted per second.⁹

Reducing the pulse rate during complex invasive cardiovascular procedures results in a considerable reduction of the total energy and the patient dose required for X-ray imaging.¹⁰ The average required dose rate scales as the square root of the frame rate, with an equal noise perception for the operator's eyes in pulsed fluoroscopy imaging.^{9,11} Hence, if the frame rate is reduced from 15 to 7.5 fps, the required dose rate is reduced by 30%, while doubling the frame rate from 15 to 30 fps increases the required dose rate by about 40%.^{9,11} One common approach for reducing the fluoroscopy time in X-ray fluoroscopy systems, involves the last image hold technique.¹²

1.2 | Relationship between motion estimation and the dose reduction for cardiac interventions

To keep the radiation dose as low as possible during the diagnostic and interventional procedures, motion compensation and prediction techniques are required to reduce potential misinterpretations caused by motion while preserving the image quality. Cardio-respiratory motion prediction has always been preferred in cardiac applications as it facilitates more accurate navigation procedures.

Deep learning architectures such as recurrent neural network (RNN) models are popular in cardiac imaging and in predicting the cardiorespiratory motion in diagnostic and interventional imaging processes.^{13–15} In these approaches, motion features (temporal and spatial) are extracted from image frames and memorized by the RNN model to predict upcoming images. However, predicting and generating realistic images and motion in an end-to-end system continues to present issues using existing models. Generative adversarial networks (GANs) are the tools used for learning deep representations. They can be used for both supervised and semisupervised learning by implicitly modeling high-dimensional data distribution. The main structure of GANs is based on training a pair of networks competing against each other. These two networks are generators and discriminators. The generator is like an art forger and produces realistic synthetic samples like images using a distribution. The discriminator acts as an art expert to distinguish the real sample from the synthetic generated one. These two networks are trained at the same time, allowing them to improve in their respective abilities until the discriminator is unable to tell the real and synthetic samples apart.¹⁶ Recently, GANs have been used as an advent method for video frame prediction. Prediction quality has been improved considerably using GANs, and the combination with RNNs has made it possible to predict multiple frames as well.¹⁷

1.3 | Proposed contribution

The contribution of this study is to predict dynamic X-ray angiography sequences using a generative model. A video frame prediction model is introduced to predict new X-ray

angiography frames. We introduced a new loss function to predict the temporal and spatial information of the arteries in angiography sequences. To minimize the vesselness structure differences between the predicted and ground truth images, a multiscale Hessian-based loss term is added to the loss function presented by Mathieu et al.¹⁸ Then, a predictive RNN-based motion model is trained to estimate the motion and content of single and/or multiple future frame(s) based on previously acquired frames in an end-to-end system.

This work is organized as follows: Section 2.1 describes the data used, Section 2.2 presents the X-ray frame prediction, while Section 2.3 presents the model architecture. The results and discussion are presented in Sections 3 and 4, respectively.

2 | MATERIALS AND METHODS

2.1 | Data description

We developed and validated our method using both simulated and patient X-ray angiography datasets from Sainte-Justine Hospital. Simulated X-ray sequences generated from realistic XCAT computational phantoms with cardiorespiratory motion¹⁹ were first investigated. The simulated motion included the beating heart and respiratory motions. The simulated dataset includes 56 different patients (32 male and 24 female) and 112 sequences (two sequences per patient, showing either the left coronary artery or the right coronary artery). All the generated sequences had a length of 75 frames and were acquired at 15 fps. The patient X-ray angiography database comprises 52 different patients with contrasted coronary arteries. This study was reviewed and approved by the Institutional Review Board of Sainte-Justine Hospital. Each patient presents a different number of sequences, with varying lengths. There is a total of 340 sequences, respectively, with a minimum and maximum length of 15 and 70 frames. All the data were acquired at 15 fps.

2.2 | X-ray angiography frame predictions

In this section, the effects of frame predictions on dose reduction are assessed in terms of the required dose rate and the total fluoroscopy time. The quantitative results of this assessment illustrate that reducing the total fluoroscopy time can have a considerable impact on cumulative radiation exposure reduction.

2.2.1 | Assessment of the impact of pulse rate reduction on the total radiation dose reduction—In our approach, we assumed that for any specific frame rate (7, 15, 30, 60 fps) the number of pulses required can be reduced during an X-ray imaging process such that the predicted frames can replace the real X-ray frames. Depending on the X-ray manufacturers, the dose for a given exposure duration is directly related to the pulse rate.^{20,21} or it can scale as the square root of the frame rate for uniform noise perceived by the operator's eyes.^{9,11} In this work, we considered the square root model.

According to this approach, for the same frame rate, a smaller pulse rate (i.e., dose rate) is required since T frames are predicted (Figure 1a). Considering K as the number of previously generated and visited frames and T as the number of predicted frames at each prediction mode, for every $K + T$ frames, T frames are predicted. Thus, the number of pulses

required at every second can be reduced by $FR \times (\frac{T}{K+T})$. Hence, the required dose rate (RDR) scales proportionally as

$$RDR \propto \sqrt{FR \times \frac{K}{K+T}}, \quad (1)$$

where the FR is the selected frame rate for the intervention or acquisition (7, 15, 30, 60 fps).

Given the parameter K , which is the number of previously generated and visited frames contributing to the prediction of the new frame per second, the X-ray exposure can pause at each predicting mode and resume in acquisition mode. Assuming t_T as the required time for T frames prediction, t_w as the required time window for $K+T$ acquisitions, and FT as the entire required fluoroscopy time (in seconds), the \widehat{FT} is the reduced fluoroscopy time:

$$\widehat{FT} = FT - \left(\left\lfloor \frac{FT}{t_w} \right\rfloor \times t_T \right) + t_r, \quad (2)$$

In any time window (t_w), the exposure time is reduced by the amount of time that is required to acquire T frames (t_T). The t_r is the remaining time in the X-ray angiography sequence ($t_r = t_{total} \bmod t_w$, $t_r \in W$) (Figure 1b).

Figure 1b is an example showing the difference between conventional continuous fluoroscopy, pulsed fluoroscopy, and our method, in terms of fluoroscopy time. For the pulsed fluoroscopy with frame prediction, the $\widehat{FT} = \Sigma(t_w - t_T) + t_r = \Sigma f t_i$ while $t_r \in W$. In pulsed fluoroscopy, less energy is exposed as compared to continuous fluoroscopy. In our approach, the X-ray device is supposed to pause at each prediction mode and resume in each acquisition mode. Thus, the total amount of fluoroscopy required in an X-ray imaging process is reduced.

2.2.2 | Cardiorespiratory motion and content estimation in X-ray sequences

—The prediction of upcoming frames of a video sequence requires two components, namely, the visual content and pixel displacement through time or motion. Thus, the proposed network learns the internal representation of image evolution through the sequence based on its content and motion. The model in this work consists of two different encoders: one for the visual content and a second one for the motion of the image sequence. These two key components need to be decomposed among the images and predicted separately. The motion features are extracted by an RNN-based encoder with long short-term memory (LSTM) and convolutional neural network (CNN), while the visible content features are only extracted from the last visited image with a CNN-based model. Deep learning methods have been applied successfully for video frame prediction in the literature.^{22–24}

2.3 | Model architecture

A generative model is built on an encoder–decoder framework. To extract the motion and content features of the images in sequences, a CNN model is used, in combination with an LSTM network. The LSTM cells are used to memorize the periodic aspect of the complex cardiorespiratory motion in the angiography sequences. According to our previous

work,¹³ the LSTM structure is robust enough to deal with different motion patterns in the cardiorespiratory motion signals during prediction. Therefore, an LSTM–CNN combination is used for a general motion estimator. The motion and content are predicted independently, using two encoders. Thus, the spatial and temporal dynamic features of the X-ray images are extracted and encoded separately. The model architecture also includes a concatenating section that combines the outputs of these encoders, as well as a multiscale residual that is used to avoid information loss before pooling in the network. The last part is the decoder, which reconstructs the predicted images. Figure 2 shows the complete structure of the model.

2.3.1 | Motion encoder—A convolutional LSTM (ConvLSTM) extracts the dynamic features in X-ray sequences. While the pixel-level features are extracted by a CNN, the sequential information is provided by the LSTM cells in the motion encoder. The motion encoder captures the local motions from one frame to the next in X-ray sequences. The cardiorespiratory movements of the objects (arteries, devices, catheters, wires, stents, etc.) are predicted directly (without using a surrogate object) and independently in the sequences.

The original presented motion encoder in Villegas et al.²² takes the element-wise image subtraction between (x_t and x_{t+1}) as an input. Since there are background movements in angiography images, the subtraction of original frames includes a lot of artifacts. In our approach, we filtered the input images by vesselness filter first and then subtracted the filtered input images to overcome the artifact caused by the background movement. Thus, the motion encoder tracks only the contrasted arteries' movement to encode the temporal dynamics of transformed images through the sequence (d_t). The output of the motion encoder is a function of filtered time frames subtraction ($x_{v(t+1)} - x_{v(t)}$), memory cell c_t and d_t .

2.3.2 | Content encoder—The content encoder extracts the essential spatial features from the visible contents, such as contrasted moving objects (arteries) and the background (ribs, bones, and devices) in the images. It takes the last observed frame x_t as input and encodes the spatial information in the image (CE_t) using a CNN network. The last observed frame has the most recent and important information that is required for the prediction of the future frame(s).

2.3.3 | Final prediction using the content and motion encoders' outputs—A multiscale encoder residual is used to compute the residual Res_t at each scale or layer just before the pooling layers of both motion and content encoders. The outputs of both encoders are concatenated and combined with the residual outputs (d_t , CE_t , Res_t) to perform pixel-level predictions in the decoder. These predictions can represent one or more frames in the future. The output of the model²² is as follows:

$$ME = [d_t, c_t] = f^{\text{motion}}(x_{v(t)} - x_{v(t-1)}, d_{t-1}, c_{t-1}) \quad (3)$$

$$CE = f^{\text{content}}(x_t) \quad (4)$$

$$Res_t^h = f^{\text{residual}}([CE^h, d_{t-1}^h]) \quad (5)$$

$$Output_t = f^{\text{combination}}([d_t, CE]) \quad (6)$$

$$\hat{x}_{t+1} = f^{\text{decoder}}(Output_t, Res_t), \quad (7)$$

where ME and CE are the motion and content encoder outputs, respectively. Res^h is the residual link at layer h being used to avoid information loss after pooling for each layer, and $Output_t$ represents the combination layer that concatenates the outputs of both motion and content encoders. The new frame is generated as the output of the decoder going through a $\tanh(\cdot)$ activation function.

2.3.4 | Loss function—A combination of terms (image space and generator loss terms) is minimized in this approach. We adjusted this loss function to predict the cardiac angiography sequences, given that the targets to track and predict are contrasted arteries. The total loss function is calculated as below considering the α and β as constant weights:

$$L_{\text{Total}} = \alpha L_{\text{IM}} + \beta L_{\text{GAN}}, \quad (8)$$

where L_{IM} represents the image space loss as a combination of terms that match the average pixel intensities with L_P , gradient difference to sharpen the predictions, and the new added subloss called vesselness²⁵ difference $L_{V_{SS}}$.

$$L_{\text{IM}} = \alpha L_{gdl} + \beta L_P + \gamma L_{V_{SS}}. \quad (9)$$

We penalized the difference between the second derivative of the Gaussian filter applied on the predicted and ground truth images with six different scales (vesselness σ range: 0.5–3 with step size: 0.5). The output of the vesselness filter on the images is the vesselness response image. The second derivatives encode the shape information, and the eigenvector corresponding to the smallest eigenvalue is the direction of the blood vessel locally. Hence, the $L_{V_{SS}}$ is applied to minimize the local differences between the predicted and ground truth images, which refer to the shape of the arteries.

The gradient difference term L_{gdl} ^{18,22} is applied to sharpen the generated images. This term directly assesses the gradient discrepancy of the ground truth and the predictions. The gradient difference between the ground truth image Y and the prediction \hat{Y} is given by

$$L_{gdl}(\hat{Y}, Y) = \sum_{i,j} (||Y_{i,j} - Y_{i-1,j}| - |\hat{Y}_{i,j} - \hat{Y}_{i,j-1}||^\lambda + ||Y_{i,j-1} - Y_{i,j}| - |\hat{Y}_{i,j} - \hat{Y}_{i,j}^\lambda|), \quad (10)$$

where λ is an integer greater or equal to 1 (here the $\lambda = 1$) and $|\cdot|$ is the absolute value function.¹⁸ The new vesselness difference term $L_{V_{SS}}$ matches the vesselness responses of

the predicted and ground truth images. The vesselness difference between the ground truth image Y and the prediction \hat{Y} is given by

$$L_{V_{ss}}(\hat{Y}, Y) = \sum_{i,j} |I_Y - I_{\hat{Y}}|. \quad (11)$$

To generate images correctly and avoid having the images being blurred by time, the generator loss in adversarial training L_{GAN} is added to solve the blurriness problem and induces realism in the image sequences, in addition to sharpening the images.¹⁸

$$L_{GAN} = -\log D([x_{1:t}, G(x_{1:t})]), \quad (12)$$

while $D(\cdot)$ represents the discriminator in adversarial training and $x_{1:t}$ is the input images concatenation. The adversarial discriminator loss (L_d) is defined by

$$L_d = -\log D([x_{1:t}, x_{t+1:t+T}]) - \log(1 - D([x_{1:t}, G(x_{1:t})])) \quad (13)$$

the concatenation of future ground truth images, and all of the predictions are represented as $x_{t+1:t+T}$ and $G(x_{1:t}) = \hat{x}_{t+1:t+T}$, respectively.^{18,22}

3 | RESULTS AND VALIDATIONS

The parameters for the X-ray angiography sequences were optimized for both the simulated and patient datasets. The number of iterations was evaluated between 1000 to 1500 for the simulated dataset and between 2000 and 2500 for the patient datasets. We divided the dataset into two parts: 80% of the dataset for training and 20% for testing. The model was evaluated on each dataset separately. Each sequence was divided into time slots or time windows of minimum $(K + T)$ frames. A single frame was repeatedly predicted at a time, and the prediction was included through the time slot while the previous predicted frame(s) contributed to new predictions. The number of previously generated and visited frames ($K = 7, 10$) contributing to predict the future frame(s) for the motion encoder was set based on capturing a complete heart cycle in time (0.8–1 s) and on the length of the shortest sequences in our dataset. All the parameters and hyperparameters were selected based on different experiments. The hyperparameters α , β , and γ were set to 1, 0.02, and 0.01, respectively, based on the experiments.

The quality of the predicted images was reduced by increasing the number of predictions. The visual quality of the predicted images using our method (the vesselness-based MCnet) was assessed as compared to the original MCnet in terms of certain similarity measurement metrics such as peak signal-to-noise-ratio (PSNR) and structural similarity index measurement (SSIM) (Tables 1–2). In our experiments, we predicted up to three frames with over 60% SSIM for both the original and vesselness-based MCnets (Tables 1–2).

According to the experiments, the quality of the predicted images is reduced by increasing the number of predictions. With the simulated data, the first three frames were well-predicted with 24–29 PSNR and between 87% to 97% SSIM (Table 1). For the patient

dataset, the best results refer to $K = 10$ in which the first three predicted frames reach between 68% and 82% SSIM (Table 2). Our experiments show that the parameter K must be equal to or greater than the number of frames required to cover a cardiorespiratory cycle. Moreover, the values for the parameter K in our experiments depend on the length of the shortest sequences in our patient dataset such that the $K + T$ must be equal to or less than the length of the shortest sequence in our dataset (13 frames). Based on the overall experiments with a patient and simulated datasets (Tables 1 and 2), the first three predicted frames have over 60% SSIM and the vesselness structure is clearly visible. Thus, at each second during the X-ray imaging process, the patients can be exposed to three fewer pulses while keeping the same frame rate (15 fps). The required frame acquisition (i.e., pulses) for a 15-fps sequence can drop by 23–30% (for $K = 10$ and, $K = 7$ respectively), and according to (1), the average required dose rate for 15 fps imaging on every second can be reduced by 0.63–0.47, as compared to real acquisition. Figure 4a,b show the samples of prediction with $K = 7$ and 10, respectively, and Figure 5 shows the overlay of the manually segmented ground truth arteries (in green) and the predictions.

To evaluate the motion prediction, we applied optical flow to estimate the motion between consecutive predicted frames as well as the ground truth frames. Optical flow is one common approach to detect the motion of moving objects in an image sequence, and it is defined as the distribution of visible velocities of moving objects in an image. Figure 3 shows the estimated movements between the four consecutive frames with optical flow. In the first row, the motion arrows are extracted from the ground truth sequence, and in the second and third rows the motion arrows are extracted from the predicted images using the vesselness MCnet and the original MCnet, respectively. The optical flow fields between each moving frame and the previous (source) frame are overlaid by moving frames ($F = 7, 8, 9$). The motion vectors in the frames predicted using the vesselness MCnet have mostly the same directions and same intensities in the region of interest (arteries) as the ground truth in all frames, while the intensities and directions of the detected motion vectors are different in the predicted frames using original MCnet.

From the test dataset, we randomly selected 30% of the sequences to evaluate the predicted content of the generated images with $K = 7$ (visited frames) and $T = 3$ (predicted frames). Coronary arteries were segmented in three consecutive frames of each selected sequence in both groups (ground truth and predictions) by a trained operator. From the resultant masks, we computed the Dice coefficients and Euclidean distances between the ground truth and the predicted images. Euclidean distance was calculated between the extracted centerlines of the segmented masks. Additionally, we reported results of a conventional gap-filling method (baseline) in the selected dataset (Table 3). The gap-filling method copied multiple times the last visited frame instead of being predicted. The Euclidean distance and Dice coefficients of the predicted images in our method and the ground truth were computed and compared with the Euclidean distance and Dice coefficients of the ground truth and the copied frames used as the gap filling.

The average computed Dice coefficients (between the ground truth and predicted images) over three predictions using our method was 0.78 ± 0.07 , while this value for the conventional gap filling was 0.63 ± 0.05 . Table 3 shows the comparison of our approach

and gap filling in terms of the computed Euclidean distances between the centerlines of the ground truth and the predictions. Based on these evaluations, for three consecutive frames, the results of the frame prediction with our approach outperform the baseline method (gap filling).

4 | DISCUSSION

This work presents a novel radiation dose management approach for pediatric interventional cardiology using a generative learning-based video frame prediction approach. This study can also facilitate the navigation of X-ray-guided interventions given the intrinsic motion compensation strategy it has in the frame predictions.

In our approach, a predictive model was introduced rather than an interpolation approach since interpolation methods require both future and former information. In frame prediction using this model, the idea is to extract the cyclic cardiorespiratory motion features from the previous frames and combine them with the visual content of the last visited frame.

The correlations between spatial and temporal features extracted from the previous frames allow self-supervision of the prediction of single or multiple frame(s) in an end-to-end system. This model can be transferable to adult patients by performing training on clinical data from adults. Additionally, the presented model can be fully adaptive to different patients with distinct respiratory and cardiac motion patterns. Compared to other video frame applications, X-ray sequences have less inherent uncertainty and variety when it comes to estimating upcoming frames since their grayscale images include limited objects for tracking, and the cardiorespiratory motion is periodic. However, the main challenge with X-ray sequence prediction in comparison to natural video prediction lies in the moving background, which makes motion prediction more complex in the former. In this work, we applied a new loss function and changed the input of the motion encoder using a vesselness filter to overcome the artifacts caused by the moving background.

Obtaining a minimum required image quality in X-ray angiography is highly challenging since different types of interventions may require different image qualities. Our results show the potential of our method for reducing the fluoroscopy time for pediatric cardiac interventions. In this work, we only focused on the pulse rate and fluoroscopy time reduction since our dataset was retrospective. Other dose indicators such as cumulative air kerma should be considered along with fluoroscopy time in our future work.

Significant efforts have been invested in improving the new generation of X-ray devices, given the importance of radiation dose reduction not only for pediatric patients with high potential risks of cancer but also for adult patients, cardiologists, and medical staff.^{26–28} This study can thus pave the way for the next generation of X-ray imaging devices, as it allows to optimize the induced radiation dose for patients and staff.

Future work will consider incorporating the heart cycle information using the ECG signal for more accurate motion estimation. Other model-based or hybrid approaches can be investigated to improve the accuracy of motion prediction. Additionally, video

superresolution methods can be included in the content predictor to improve the image quality of predictions.

5 | CONCLUSION

This work presents a novel radiation dose management approach for pediatric interventional cardiology using a learning-based video frame prediction. Such a prediction can reduce the amount of accumulated radiation dose for patients and staff by exposing them to fewer pulses while preserving the frame rate and the image quality.

ACKNOWLEDGMENTS

The authors would like to thank Sainte-Justine Hospital Cath lab technicians and Duke university (CVIT) for their time and valuable advice. This work was supported in part by the NSERC Discovery grant and by the National Institutes of Health biomedical resource grant P41-EB028744. The Quadro RTX6000 used for this research was donated by the NVIDIA Corporation.

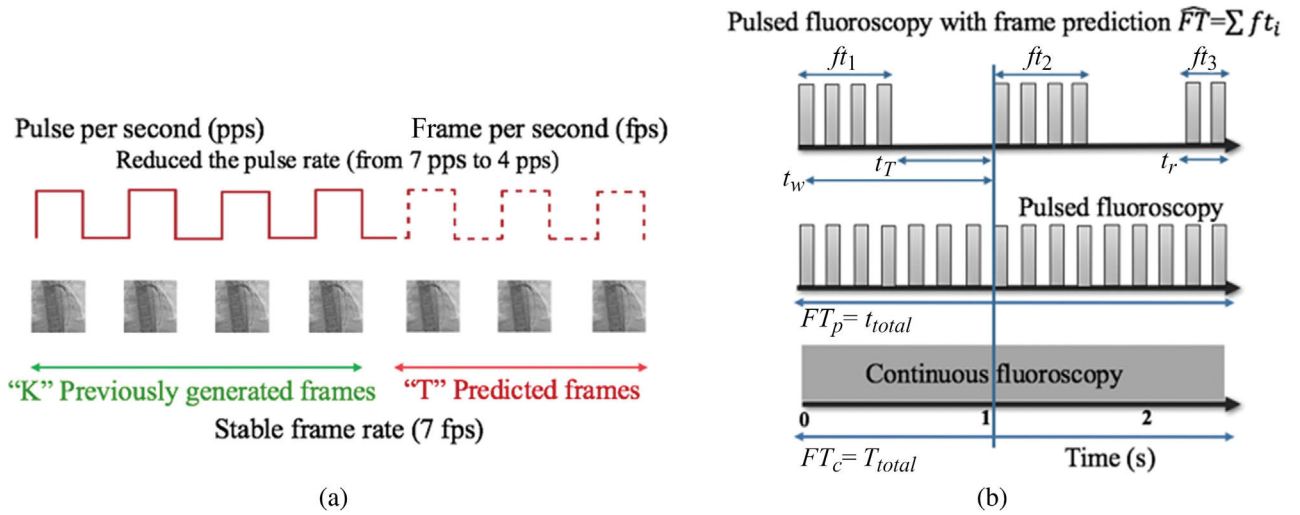
Funding information

Natural Sciences and Engineering Research Council of Canada (NSERC), Grant/Award Number: RGPIN-2021-03078

REFERENCES

1. Liu Y, Chen S, Zühlke L, et al. Global birth prevalence of congenital heart defects 1970–2017: updated systematic review and meta-analysis of 260 studies. *Int J Epidemiol*. 2019;48:455–463. [PubMed: 30783674]
2. Haas NA, Happel CM, Mauti M, et al. Substantial radiation reduction in pediatric and adult congenital heart disease interventions with a novel X-ray imaging technology. *Int J Cardiol Heart Vasc*. 2015;6:101–109. [PubMed: 28785634]
3. Gislason-Lee AJ, Kumcu A, Kengyelics SM, et al. How much image noise can be added in cardiac X-ray imaging without loss in perceived image quality? *J Electron Imaging*. 2015;24:051006.
4. Hunold P, Vogt FM, Schmermund A, et al. Radiation exposure during cardiac CT: effective doses at multi-detector row CT and electron-beam CT. *Radiology*. 2003;226:145–152. [PubMed: 12511683]
5. Chida K, Ohno T, Kakizaki S, et al. Radiation dose to the pediatric cardiac catheterization and intervention patient. *Am J Roentgenol*. 2010;195:1175–1179. [PubMed: 20966324]
6. Dauer LT. Radiation Dose Management for Fluoroscopically-Guided Interventional Procedures. National Council on Radiation Protection and Measurements. Radiation dose management for fluoroscopically-guided interventional medical procedures. Bethesda, MD: National Council on Radiation Protection and Measurements, 2010: Report 168. 2011.
7. Yamagata K, Aldhoon B, Kautzner J. Reduction of fluoroscopy time and radiation dosage during catheter ablation for atrial fibrillation. *Arrhythm Electrophysiol Rev*. 2016;5:144. [PubMed: 27617094]
8. Hirshfeld JW, et al. ACCF/AHA/HRS/SCAI clinical competence statement on physician knowledge to optimize patient safety and image quality in fluoroscopically guided invasive cardiovascular procedures: a report of the American College of Cardiology Foundation/American Heart Association/American College of Physicians Task Force on Clinical Competence and Training. *J Am Coll Cardiol*. 2004;44:2259–2282. [PubMed: 15582335]
9. Balter S Fluoroscopic frame rates: not only dose. *Am J Roentgenol*. 2014;203:W234–W236. [PubMed: 25148178]
10. Pyne CT, Gadey G, Jeon C, Piemonte T, Waxman S, Resnic F. Effect of reduction of the pulse rates of fluoroscopy and CINE-acquisition on X-ray dose and angiographic image quality during invasive cardiovascular procedures. *Circ Cardiovasc Interv*. 2014;7:441–446. [PubMed: 25006174]

11. Aufrichtig R, Xue P, Thomas CW, Gilmore GC, Wilson DL. Perceptual comparison of pulsed and continuous fluoroscopy. *Med Phys.* 1994;21:245–256. [PubMed: 8177157]
12. Wilson DL, Xue P, Aufrichtig R. Perception of fluoroscopy last-image hold. *Med Phys.* 1994;21:1875–1883. [PubMed: 7700194]
13. Azizmohammadi F, Martin R, Miro J, Duong L. Model-free cardiorespiratory motion prediction from X-ray angiography sequence with LSTM network. In: 2019 41st Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC). IEEE; 2019:7014–7018.
14. Fang H, Li H, Song S, et al. Motion-flow-guided recurrent network for respiratory signal estimation of X-ray angiographic image sequences. *Phys Med Biol.* 2020;65(24):245020. [PubMed: 32590382]
15. Lyu Q, Shan H, Xie Y, Li D, Wang G. Cine cardiac MRI motion artifact reduction using a recurrent neural network. arXiv preprint arXiv:2006.12700. 2020.
16. Creswell A, White T, Dumoulin V, Arulkumaran K, Sengupta B, Bharath AA. Generative adversarial networks: an overview. *IEEE Signal Process Mag.* 2018;35:53–65.
17. Hu Z, Wang JT. Generative adversarial networks for video prediction with action control. In: International Joint Conference on Artificial Intelligence. Springer; 2019:87–105.
18. Mathieu M, Couprie C, LeCun Y. Deep multi-scale video prediction beyond mean square error. arXiv:1511.05440. 2015.
19. Segars WP, Sturgeon G, Mendonca S, Grimes J, Tsui BM. 4D XCAT phantom for multimodality imaging research. *Med Phys.* 2010;37:4902–4915. [PubMed: 20964209]
20. Kobayashi T, HirshfeldJr JW. Radiation exposure in cardiac catheterization: operator behavior matters. 2017.
21. Mahesh M Fluoroscopy: patient radiation exposure issues. *Radiographics.* 2001;21:1033–1045. [PubMed: 11452079]
22. Villegas R, Yang J, Hong S, Lin X, Lee H. Decomposing motion and content for natural video sequence prediction. arXiv:1706.08033. 2017.
23. Hsieh J-T, Liu B, Huang D-A, Fei-Fei LF, Niebles JC. Learning to decompose and disentangle representations for video prediction. in *Adv Neural Inf Process Syst.* 2018:517–526.
24. Tulyakov S, Liu MY, Yang X, Kautz J. Mocogan: decomposing motion and content for video generation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. IEEE. 2018:1526–1535.
25. Frangi AF, Niessen WJ, Vincken KL, Viergever MA. Multiscale vessel enhancement filtering. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. Springer; 1998:130–137.
26. Gislason-Lee AJ, Keeble C, Malkin CJ, et al. Impact of latest generation cardiac interventional X-ray equipment on patient image quality and radiation dose for trans-catheter aortic valve implantations. *Br J Radiol.* 2016;89:20160269. [PubMed: 27610932]
27. Gislason-Lee AJ, Keeble C, Egleston D, Bexon J, Kengyelics SM, Davies AG. Comprehensive assessment of patient image quality and radiation dose in latest generation cardiac x-ray equipment for percutaneous coronary interventions. *J Med Imaging.* 2017;4:025501.
28. McNeice AH, Brooks M, Hanratty CG, Stevenson M, Spratt JC, Walsh SJ. A retrospective study of radiation dose measurements comparing different cath lab X-ray systems in a sample population of patients undergoing percutaneous coronary intervention for chronic total occlusions. *Catheter Cardiovasc Interv.* 2018;92:E254–E261. [PubMed: 29411518]

**FIGURE 1.**

(a) The sequence at 7 fps frame rate is acquired partially with exposed pulses and partially with predictions such that the pulse rate gets reduced while the frame rate remains constant ($K = 4$ and $T = 3$). (b) An example of three different fluoroscopy techniques. Less fluoroscopy time is required for pulsed discrete fluoroscopy by pausing the radiation beam after K acquired images for a prediction time t_T in each time window t_w compared to other methods ($\widehat{FT} < FT_p < FT_c$)

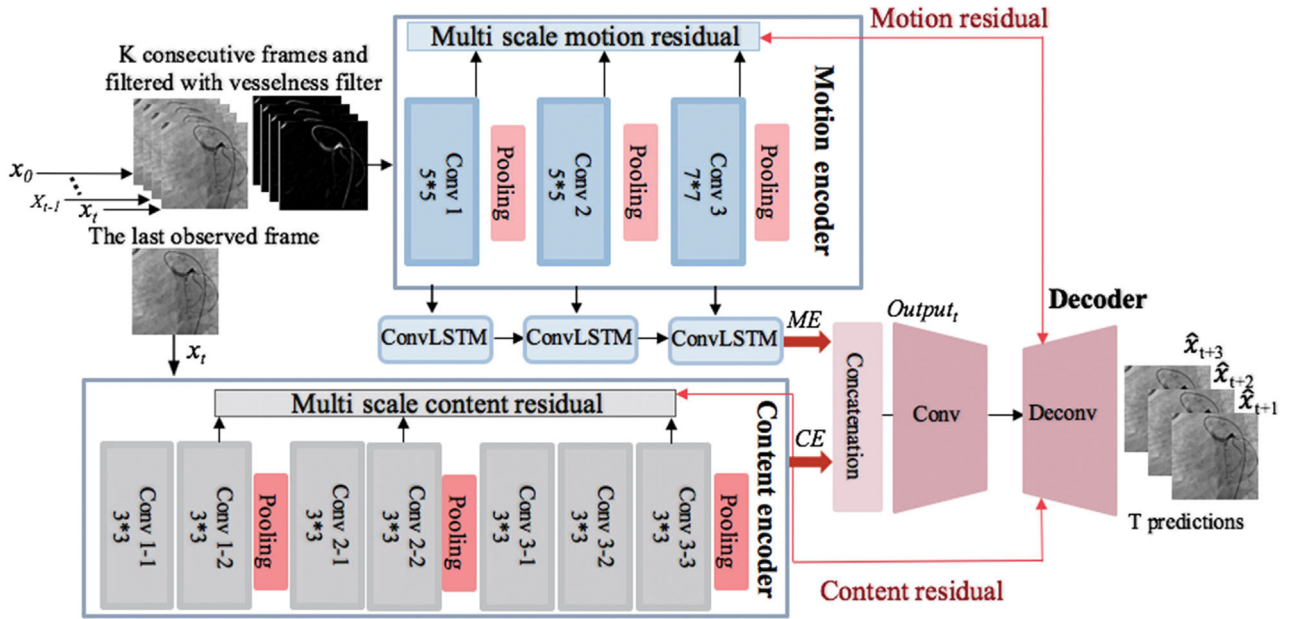


FIGURE 2.

The motion-content model structure. Two encoders extract the motion and content features separately (*ME* and *CE*). The input for the motion encoder is a subsequence of previously acquired and visited frames filtered by the vessellness filter. The input for the content encoder is the last visited frame. The outputs of these two encoders are concatenated to be decoded as a subsequence of predictions. The motion and content residuals are added to avoid information loss

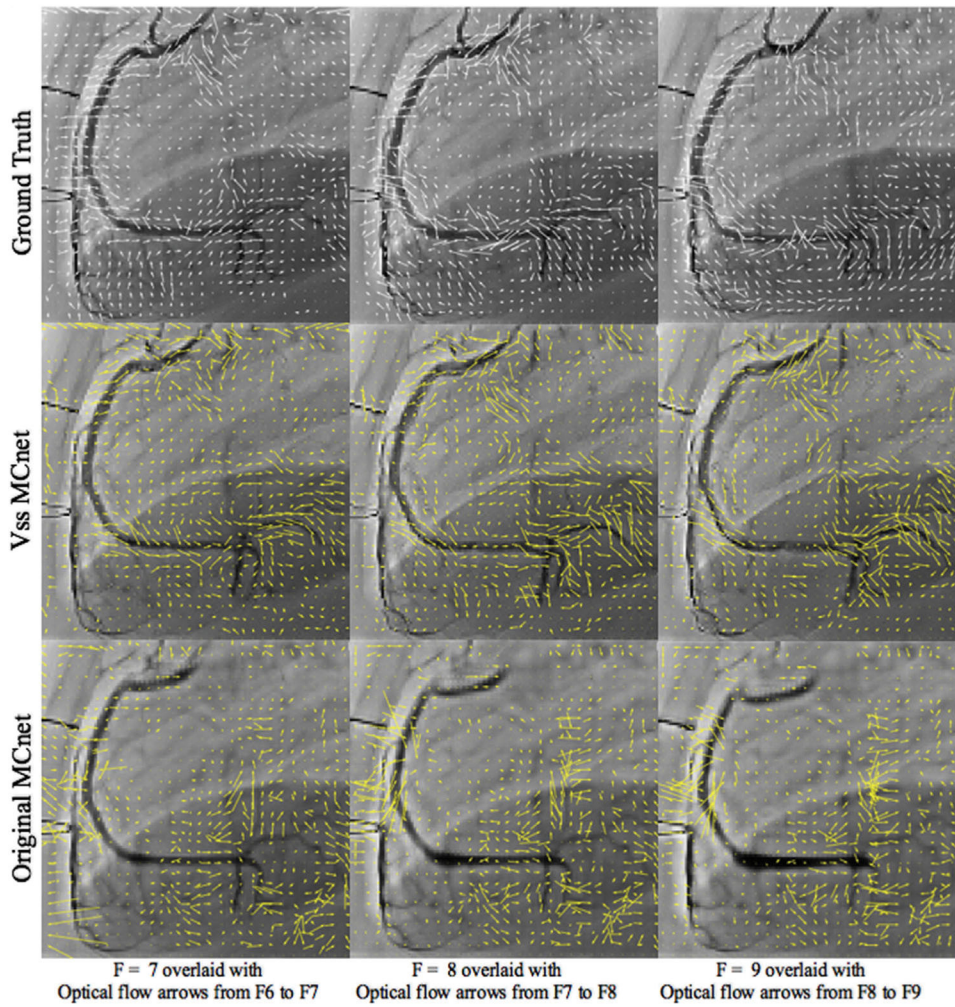


FIGURE 3.

Optical flow estimated motion fields of the ground truth sequence (top-white arrows) and the generated frames with vesselness-based MCnet on the second row and original MCnet on the third row (yellow arrows). The optical flow fields are overlaid to the predicted frames F7–F9

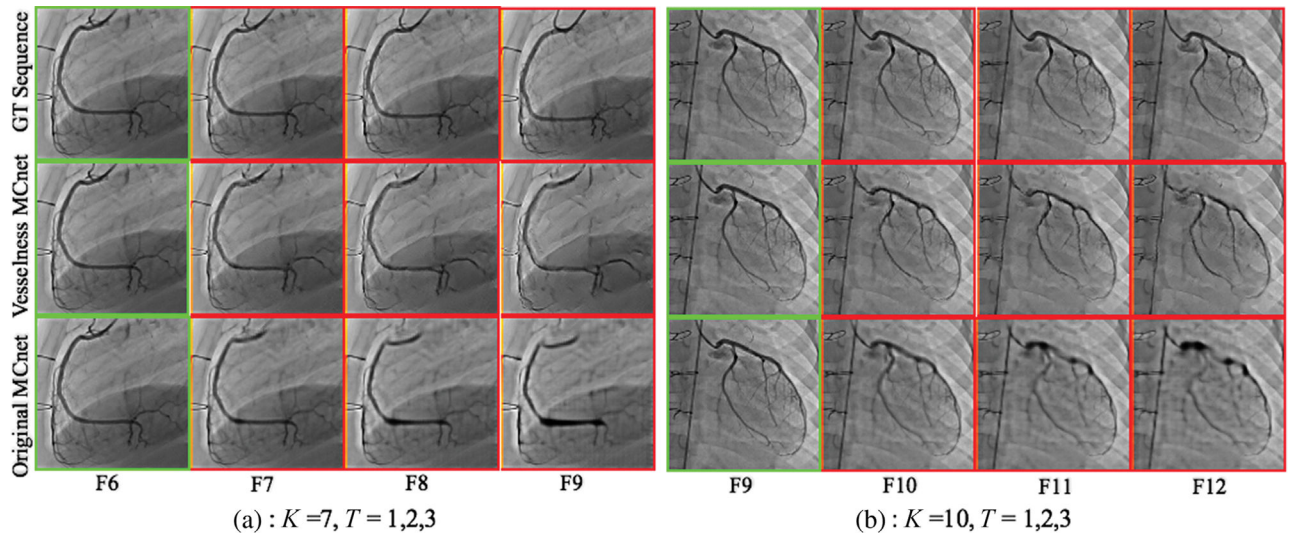


FIGURE 4.

The first row shows the ground truth sequence. The second and third rows show the results of vesselness-based MCnet and original MCnet, respectively. The predicted images are identified with a red outline and the last visited frame with a green outline

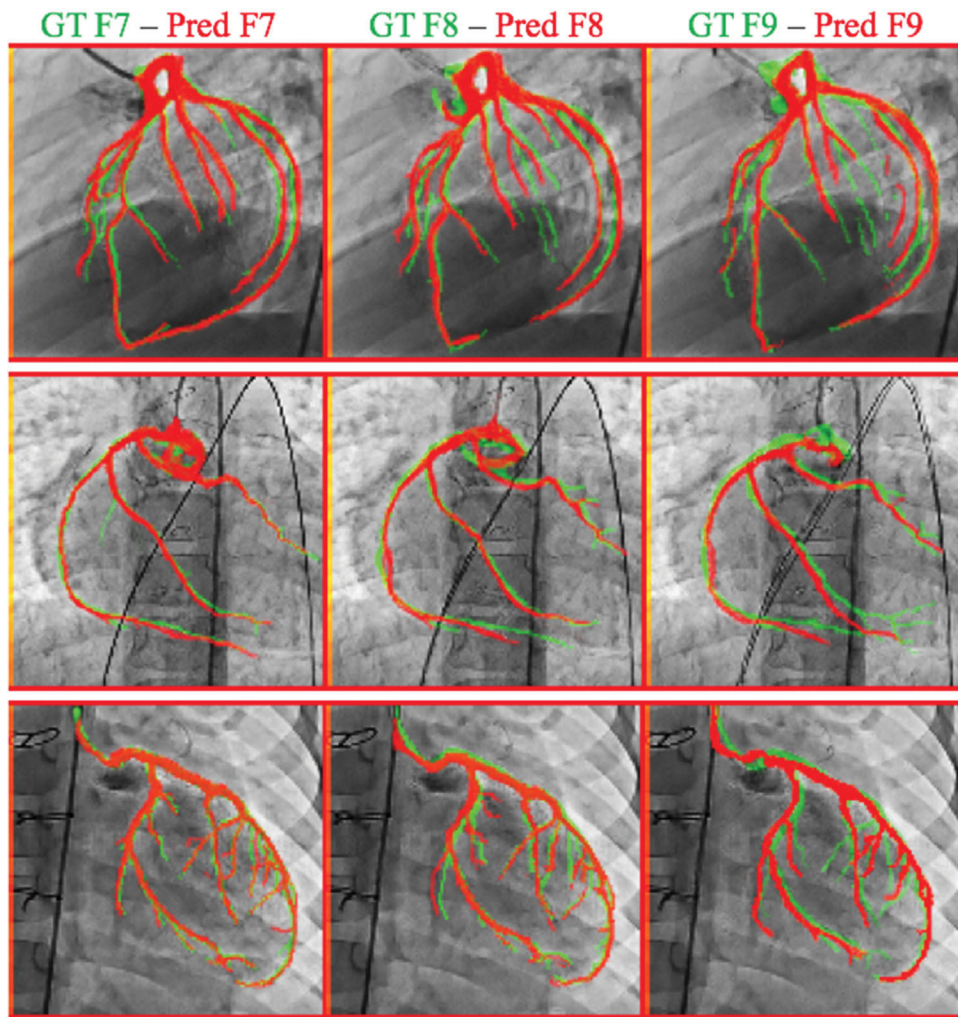


FIGURE 5.
An overlay of the manual segmentation masks for the ground truth in green and predicted sequences in red

TABLE 1

Average similarity measurements of the predicted images over the testing data on three predicted images for the simulated dataset

Frame	PSNR	Original MCnet PSNR	SSIM	Original MCnet SSIM
Simulated data $K = 7, T = 1,2,3$				
Frame 1	28.28	27.98	0.94	0.89
Frame 2	25.47	24.85	0.92	0.85
Frame 3	23.90	23.01	0.88	0.82
Simulated data $K = 10, T = 1,2,3$				
Frame 1	29.13	28.82	0.97	0.86
Frame 2	27.65	25.10	0.93	0.83
Frame 3	24.14	23.12	0.87	0.81

Abbreviations: PSNR, peak signal-to-noise-ratio; SSIM, structural similarity index measurement.

TABLE 2

Average similarity measurements of the predicted images over the testing data on three predicted images for the patient dataset

Frame	PSNR	Original MCnet PSNR	SSIM	Original MCnet SSIM
Patient data $K = 7, T = 1,2,3$				
Frame 1	27.10	26.75	0.79	0.80
Frame 2	24.42	23.59	0.68	0.70
Frame 3	23.10	21.54	0.61	0.61
Patient data $K = 10, T = 1,2,3$				
Frame 1	27.97	26.80	0.82	0.78
Frame 2	25.65	24.62	0.74	0.69
Frame 3	24.14	23.32	0.68	0.63

Abbreviations: MCnet, Motion Content network; PSNR, peak signal-to-noise-ratio; SSIM, structural similarity index measurement.

Euclidean distance between the centerlines of arteries in the predicted frames and ground truth for the frame prediction and gap filling

TABLE 3

Euclidean distance (mm)						
# Frame	Frame prediction			Gap filling		
	Mean	Maximum	SD	Mean	Maximum	SD
Frame 1	0.28 mm	0.76 mm	(±) 0.19 mm	0.33 mm	0.79 mm	(±) 0.22 mm
Frame 2	0.30 mm	0.78 mm	(±) 0.20 mm	0.39 mm	0.85 mm	(±) 0.31 mm
Frame 3	0.32 mm	0.84 mm	(±) 0.21 mm	0.51 mm	0.93 mm	(±) 0.35 mm