



### Frontiers in AI and Its Applications in Medical Physics

# GROWING SIGNIFICANCE OF DEEP LEARNING IN MODERN MEDICAL IMAGING

### **Rongping Zeng**

Division of Imaging, Diagnostics and Software Reliability (DIDSR)

Office of Science and Engineering Laboratories Center for Devices and Radiological Health U.S. Food and Drug Administration

July 2021

## Disclaimer



The mention of commercial products, their sources, or their use in connection with materials reported herein is not to be construed as either an actual or implied endorsement of such products by the Department of Health and Human Services.

## Outline

FDA

- Deep learning (DL) and medical imaging
- DL in image reconstruction/processing
  - Applications
  - Performance evaluation
- A few DIDSR research studies
  - MRI reconstruction
  - Synthetic digital mammography
  - Low-dose CT image denoising
- Summary

### **Growth of DL in medical imaging**

FDA



## **DL in modern medical imaging**

FDA





## **DL** in image reconstruction/processing

FDA

### \* Table 1 | Recent deep learning tomographic reconstruction methods for various biomedical imaging modalities based on their technical contributions

Modality	Architecture	Applications and citations					
X-ray CT	Image domain	Low-dose <sup>15-17,48,70,97</sup> , sparse view <sup>29,113</sup> , limited angle <sup>50</sup> , metal artefact <sup>52</sup> , dual-energy <sup>143</sup>					
	Unrolling	Sparse view <sup>29,58,59,61,62</sup> , low-dose <sup>60</sup>					
	Sensor domain	Metal artefact <sup>85,86</sup> , sinogram analysis <sup>87</sup>					
	Model-based/plug-and-play	Low-dose <sup>60,68,70,82,144</sup> , sparse view <sup>58,59</sup> , interior <sup>83</sup> , cone-beam artefact <sup>84</sup>					
	Domain transform	Low-dose and sparse view <sup>77-79</sup> , interior tomography <sup>83</sup>					
	GAN, unsupervised	Low-dose <sup>48,97/10</sup> , metal artefact <sup>130</sup>	*Wang, G., Ye, J.C. & De Man, B. <b>Deep learning for</b>				
MRI	Image domain	Accelerated MRI <sup>24,30,49</sup>	tomographic image reconstruction. Nat Mach				
	Unrolling	Accelerated MRI <sup>23,25,53-56</sup>	Intell <b>2</b> , 737–748 (2020)				
	Sensor domain	Accelerated MRI <sup>88,90,91</sup> , artefact removal <sup>89</sup>					
	Model-based/plug-and-play	Accelerated MRI <sup>26,63,72</sup> , dynamic MRI <sup>65,75</sup>	Potential to achieve comparable image				
	Domain transform	Accelerated MRI <sup>76</sup>					
	GAN, unsupervised	Accelerated MRI <sup>49,109</sup> , contrast synthesis <sup>106,107</sup>	performance as model-based iterative image				
PET/SPECT	Image domain	Low-dose PET31,32,145,146	reconstruction methods				
	Model-based/plug-and-play	Low-dose PET <sup>31,67-69</sup>	Faster computation speed				
	Domain transform	Low-dose PET <sup>80,81</sup>	New challenges on image quality assessment				
	GAN, unsupervised	Low-dose PET <sup>111</sup> , attenuation correction <sup>104,105</sup>					
Ultrasound	Image domain	Artefact removal <sup>45</sup> , photoacoustic artefact removal <sup>46,47</sup>					
	Sensor domain	Radiofrequency interpolation96					
	Model-based/plug-and-play	Photo-acoustic imaging <sup>64</sup>					
	Domain transform	Beamformer <sup>92,94,95</sup>					
Optical	Image domain	Super resolution <sup>36-28,103</sup> , ghost imaging <sup>42</sup> , scattering medium imaging <sup>40,41</sup> , mobile phone microscopy <sup>43</sup>					
	GAN, unsupervised	Deconvolution microscopy <sup>100,102,103</sup>					

## Image quality (IQ)

### Performance testingDiagnostic accuracy

### • Global IQ metrics:

- MSE, PSNR, SSIM etc.
- Standard system IQ metrics
  - MTF, NPS, CNR
  - Non-uniformity, CT number accuracy, etc.
- Task-based IQ metrics
  - Lesion detectability, estimation of lesion shape and lesion size, etc.

Less predictive

Predictive for linear systems

Most clinically relevant

## **Generalizability Performance**

- DL performance highly depends on the training data
- Real-world data is much broader
  - Imaging system/acquisition protocol
  - Patient population: sex, body size, disease, etc.
  - Data site
- Understanding of the generalizability performance improves regulatory review
  - More accurate definition of the intended use of a device
  - Appropriate design of the testing data variety

#### Key Considerations for Authors, Reviewers, and Readers of AI/ML Manuscripts in Radiology

#### Key Considerations

Are all three image sets (training, validation, and test sets) defined? (non overlap, independent) Is an *external* test set used for final statistical reporting? Have multivendor images been used to evaluate the AI algorithm? Are the sizes of the training, validation, and test sets justified? Was the AI algorithm trained using a standard of reference that is widely accepted in our field? Was preparation of images for the AI algorithm adequately described? Were the results of the AI algorithm compared with radiology experts and/or pathology? Was the manner in which the AI algorithm makes decisions demonstrated? Is the AI algorithm publicly available? Note.—AI = artificial intelligence, ML = machine learning.

Assessing Radiology Research on Artificial Intelligence, Bluemke et. al., Radiology, Dec., 2019

## A few DIDSR research studies

FDA

- MRI reconstruction
  - Effect of training data on MRI AUTOMAP reconstruction
- Synthetic digital mammography
  - Tasked-based evaluation of a DL-based synthetic digital mammography algorithm
- Low-dose CT image denoising
  - 1. Framework for DL model optimization
  - 2. Generalizability performance Evaluation



ImageNet: 50,000 natural

images

Metrics:  $\cap$ 

٠

- Global IQ (e.g., MAE)
- Task IQ (e.g., Cortical thickness)



#### Synthetic Digital Mammography (DM)

#### Virtual clinical trial for task-based evaluation of a deep learning synthetic mammography algorithm

Badal A, Cha K, Divel S, Graff CG, Zeng R, Badano A, SPIE medical imaging conference 2020

**Purpose**: Evaluate a synthetic DM method through a DL-based denoising of the digital breast tomosynthesis (DBT) 0° projection view.

- Network structures: CNN3, DnCNN, REDCNN
- IQ metrics
  - Global IQ: PSNR;
  - Task IQ: Lesion detectability (spiculated mass/calcification cluster)



- Using the VICTRE\* pipeline to create 120 breast phantom data
  - Training: 100
    - 25-projection DBT and high-dose DM (x25 dose)
    - Not contain lesions Mass
    - Testing: 3000
      - Contain lesions



\* <u>https://github.com/DIDSR/VICTRE</u>

FDA

#### Synthetic Digital Mammography (DM) Results







- Diagnostic performance of post-processed images does not always correlate with PSNR. •
- Lesion detectability might decrease even if PSNR in image increases. •

### Low-dose CT (LDCT) image denoising

• The DL network is trained to process a low-dose CT image slice (e.g., reconstructed by FBP) with the intend to output a cleaner image.

FDA



 Pairs of image patches extracted from full and quarter-dose image slices

#### LDCT - Study 1

#### Deep neural networks-based denoising models for CT imaging and their efficacy

Kc P, Zeng R, Farhangi M, Myers KJ, SPIE Medical Imaging 2021

**Purpose**: Holistic framework for optimizing DL denoising models using both global IQ and bench performance IQ metrics.

– CNN3, REDCNN, DnCNN, Du-DnCNN, GAN



• CNN3: 3-layer CNN, from Chen et el., BMOE 2017



REDCNN: 10-layer residual encoder-decoder CNN, from Chen

et al., IEEE-TMI, 2017



FDA

**DnCNN**: 17-layer CNN, from Zhang et al., IEEE\_TIP, 2017



 Du-DnCNN: Dilated U-shaped DnCNN, form Cheng et al., "Image-denoising-with-deep-cnns", 2020 <u>https://github.com/lychengr3x/Image-Denoising-with-Deep-CNNs.</u>



• **GAN**: RESNET as generator, DnCNN as discriminator, based on Goodfellow et al., "Generative adversarial nets", 2014

#### LDCT - Study 1 Results

• HU accuracy and PSNR

120		-35	
<mark>3</mark> 40			
990	•	•	



HU disk	-35	120	340	990	PSNR	
FBP	-34.91	119.97	338.94	987.88		
CNN3, Norm	-70.46	74.80	293.09	925.85	35.99	
CNN3, No Norm	-35.16	119.45	339.38	991.85	35.61	
REDCNN, Norm	-59.47	96.87	336.36	959.66	35.26	
REDCNN, No Norm	-33.83	120.02	339.89	992.39	35.76	
DnCNN, Norm	-41.01	115.67	334.25	982.51	35.13	
DnCNN, No Norm	-55.21	99.87	318.53	967.52	34.99	
Du-DnCNN, Norm	-49.36	102.99	316.01	960.31	34.82	

- Data normalization affects the HU accuracy performance differently for different network structures.
- Consider using bench performance to complement the global metrics when optimizing the DL denoising models.

#### LDCT - Study 2

#### Generalizability performance of a deep learning-based CT image denoising method

FD/

Zeng R, Lin Y, Li Q, Lu J, Skopec M, Fessler JA, Myers KJ, The 6<sup>th</sup> X-ray CT conference 2020

**Purpose**: Investigate the generalizability performance of a CT image denoising network (REDCNN) in data acquired with different scan parameters

- Dose, reconstruction filter, slice thickness
- Identify underlying data property shift that may affect the generalizability

		Training data groups	"	Trained models DL <i>kernel-thickness-dose"</i>		Cross testing	
Dose level effect Reconstruction Kernel effect	•	Smooth kernel / 3 mm thickness / 25% dose level Smooth kernel / 3 mm thickness / Mixed dose levels Sharp kernel / 3 mm thickness / mixed dose levels smooth kernel / 3 mm thickness / mixed dose levels		DLsmooth-3mm- <mark>25</mark> % DLsmooth-3mm-mix% DL <mark>sharp</mark> -3mm-mix% DL <mark>smooth</mark> -3mm-mix%	• • •	MSE Contrast-dependent MTF NPS	
Thickness effect	•	Smooth kernel / <b>1 mm thickness</b> / mixed dose levels Smooth kernel / <b>3 mm thickness</b> / mixed dose level		DLsmooth- <mark>1mm</mark> -mix% DLsmooth- <mark>3mm</mark> -mix%			



#### LDCT - Study 2 **Results: reconstruction kernel effect**

FDA



**OSEL** Accelerating patient access to innovative, safe, and effective medical devices through best-in-the-world regulatory science

Test:

Test:

#### LDCT - Study 2 Results: generalizability



- Performance of REDCNN, a slice-wise CT image denoising network REDCNN
  - Generalized poorly between the sharp and smooth reconstruction kernel;
  - Not different between the 1mm and 3mm slice thickness data;
  - More robust when trained with mixed-dose data.
- The NPS property in a CT image may be used to as an underlying data property to predict whether the denoising performance is generalizable to the data or not.



### **Summary**



- DL applies to every stage in a modern medical imaging chain
  - From image acquisition, image reconstruction, to image interpretation
- DL-based image reconstruction/processing
  - Consider multiple types of IQ metrics in algorithm development and performance evaluation
  - Investigate the generalizability performance
    - Help with more efficient training and testing data design
    - Ensure a DL method is used within its generalizable range

### Acknowledgements

- My FDA colleagues
  - Dr. Andreu Badal
  - Dr. Berkman Sahiner
  - Dr. Jiang Lu
  - Dr. Kyle J Myers
  - Dr. Kyoko Fujimoto
  - Dr. Prabhat Kc
- Host
  - Dr. Lei Xing



FDA

Rongping.zeng@fda.hhs.gov FDA/CDRH/OSEL/DIDSR