# Deep Learning Methods for Medical Image Processing

Fatemeh Zabihollahy, Ph.D.
Postdoctoral Research Fellow
University of California, Los Angeles

# Overview of the Presentation

- Introduction on image analysis methods
- Deep Learning (DL)
- Convolutional Neural Network (CNN)
- U-Net
- GAN
- Common challenges in applying DL-based method for medical image analysis
- Regularization
- Interpretability of DL models

# Image Analysis Problems

The science of analyzing medical problems based on different imaging modalities and digital image analysis techniques.

- Classification

- Object Localization

- Object Detection

- Segmentation

- Registration ( i.e., comparing different modalities/patients)

- Reconstruction (in CT refers to a mathematical process that generates tomographic images from X-ray projection data acquired at many different angles around the patient.)
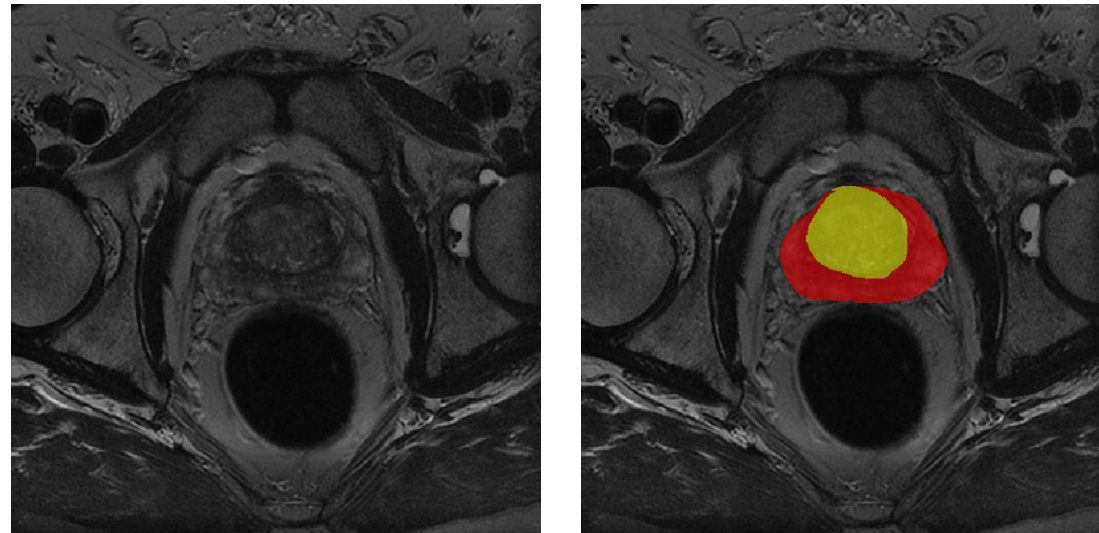
# Image Analysis Problems

- **Classification** refers to a predictive modeling problem where a class label is predicted for a given example of input data. Examples of classification problem in medical imaging: Given a Chest X ray image classify if it as normal or COVID-19.

- **Object localization** aims to locate the main (or most visible) object in an image. Examples of localization problem in medical imaging: Given a prostate MR image find the location of prostate whole gland by specifying a tightly cropped bounding box centered on its instance.

- **Object detection** defines as detecting instances of semantic objects of a certain class in a digital image. Examples of detection problem in medical imaging: Given an abdominal CT image detect renal mass and classify as cyst or tumor.

# Object Localization vs. Object Detection

- **Object localization**: identifying the location of one or more objects in an image and drawing abounding box around their extent.

- **Object detection** combines localization and classification tasks such that localizes and classifies one or more objects in an image.

# Segmentation

**Image segmentation** is partitioning an image into two or more meaningful regions. Examples of segmentation problem in medical imaging: Given a prostate MR image, find the boundary of prostate zones in an image.
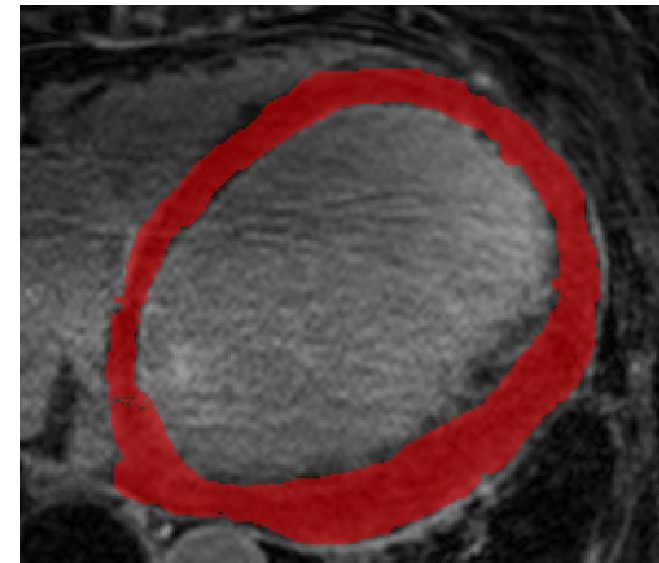
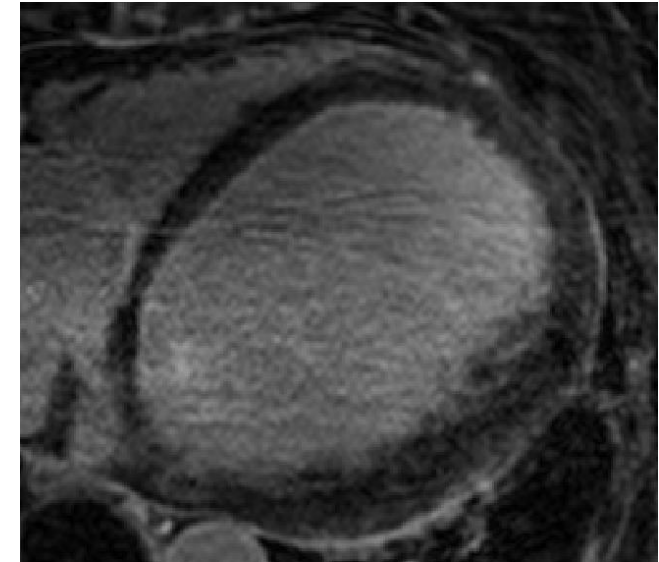# Previous Methods for Image Analysis

- Previous methods developed for classification, object detection and localization, and image segmentation can be categorized into conventional, and Machine learning (ML)-based.

- Conventional techniques are further divided into intensity-based thresholding and functional optimization methods (e.g., active contours, level sets, live wire, graph cuts, convex relaxation techniques, etc.)

# Intensity-based Methods

- Ostu

- Region Growing

- Full Width at Half Maximum (FWHM): identifies the maximum intensity value of ROI and considers its half as a reference (Ir/2) [1].

- Signal Threshold to Reference Mean (STRM): target is considered as any regional signal above a mean intensity value of ROI plus two (STRM2), three (STRM3), four (STRM4), five (STRM5), or six (STRM6) standard deviations [2].

- FWHM and STRM are cardiac specific methods.

[1] Neizel M, Katoh M, Schade E, et al. Rapid and accurate determination of relative infarct size in humans using contrast-enhanced magnetic resonance imaging. Clin Res Cardiol. 2009;98:319–324.
[2] Kolipaka A, Chatzimavroudis GP, White RD, O'Donnell TP, Setser RM. Segmentation of non-viable myocardium in delayed enhancement magnetic resonance images. Int J Cardiovasc Imaging. 2005;21:303–311.

# Limitations of Intensity-based Methods

➢Fail when the image histogram is close to unimodal in binary segmentations.

➢Highly influenced by image noise.

➢Low robustness and accuracy as these types of algorithms solely rely on intensity value as a discriminant feature.

# Optimization-based Segmentation Algorithms

- Energy optimization methods first define a mathematical criterion for the "goodness" of a given segmentation that translates the formulation of the segmentation problem as an optimization problem under certain geometric constraints.

- Practically in functional optimization techniques, an initial contour is defined around or inside the target and evolved by minimizing an energy function.

- The energy function could be a combination of shape, length-based, and region-based energy terms that reaches its minimum value when the contour lies on the boundary of the target.

- In other words, functional optimization methods search for a unique contour that lies on the boundary of the target by minimizing energy function based on optimizing an objective function.

**Limitations of Optimization-based Segmentation Methods:**

➢ Subject to high operator variability.

➢ The performance of energy optimization-based methods plateaus despite the increased number of available images.

# Machine Learning Algorithms

- ML methods are a set of algorithms developed to learn meaningful patterns from example data aiming to minimize human interaction.

- Selecting an appropriate algorithm and proper training allow the machine to learn patterns more efficiently.

- ML algorithms can be divided into two main categories: supervised and unsupervised.

- In supervised learning, we utilize labeled data to train a model whereas in unsupervised learning we allow the model to work on its own to discover information, and we utilize unlabeled data.

- For supervised learning-based methodologies, labeled data is divided into training and testing parts.

- The training samples are utilized to learn from and develop an algorithm to perform a task. Testing data are used to assess the performance of the trained model.

# K-Means Clustering

- K-means clustering is an example of unsupervised algorithms that has been widely used for image segmentation when we have unlabeled data.

- K-Means clusters or partitions the given data into K-clusters or parts based on the K-centroids. The goal is to find certain groups based on some kind of similarity in the data with the number of groups represented by K.
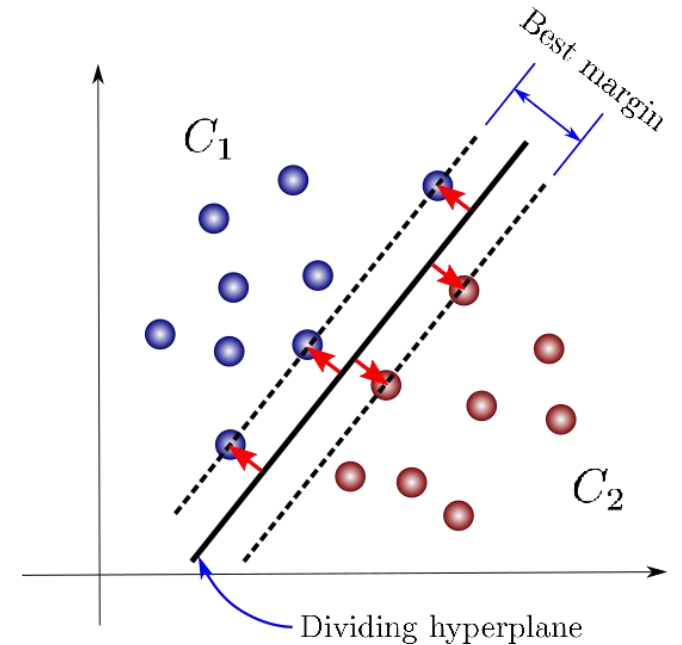
**Limitations of K-Means Clustering:**

➤ In Clustering, we are trying to identify some segments or clusters in the data. When clustering algorithms are used, unexpected things like structures, clusters, and groupings can suddenly pop-up.

➤ Like all clustering methods, the starting locations of the partitions are important for achieving optimal solution since they are susceptible to termination when achieving a local maximum.

➤ Incorporating anatomical and spatial information into segmentation is the main challenge for clustering algorithms in medical image segmentation applications.

# Support Vector Machin (SVM)

- SVM is a type of supervised learning algorithm.

- Like other classifiers, SVM constructs a hyperplane or set of hyperplanes as a decision boundary to separate between different classes.

- A good separation is achieved when the decision boundary has the largest distance (margin) to the nearest training samples of any class.

- This larger margin leads to the lower ***generalization error*** of the classifier, a measure of how accurately an algorithm can predict outcome values for testing data.

**Limitations of SVM:**

➢ Not suitable for large data sets.

➢ Does not perform very well when the target classes are overlapping.

➢ There is no probabilistic explanation for the classification (non interpretable).



Retrieved from https://towardsdatascience.com/support-vector-machines-for-classification-fc7c1565e3

# K Nearest Neighbor (KNN)

- KNN is a non-parametric instance-based supervised learning method that does not construct a general model for classification.

- The input consists of the k closest training examples in a data.

- KNN stores training samples along with the labels in a feature space and classifies an object by a majority vote of its nearest neighbors.

**Limitations of KNN:**

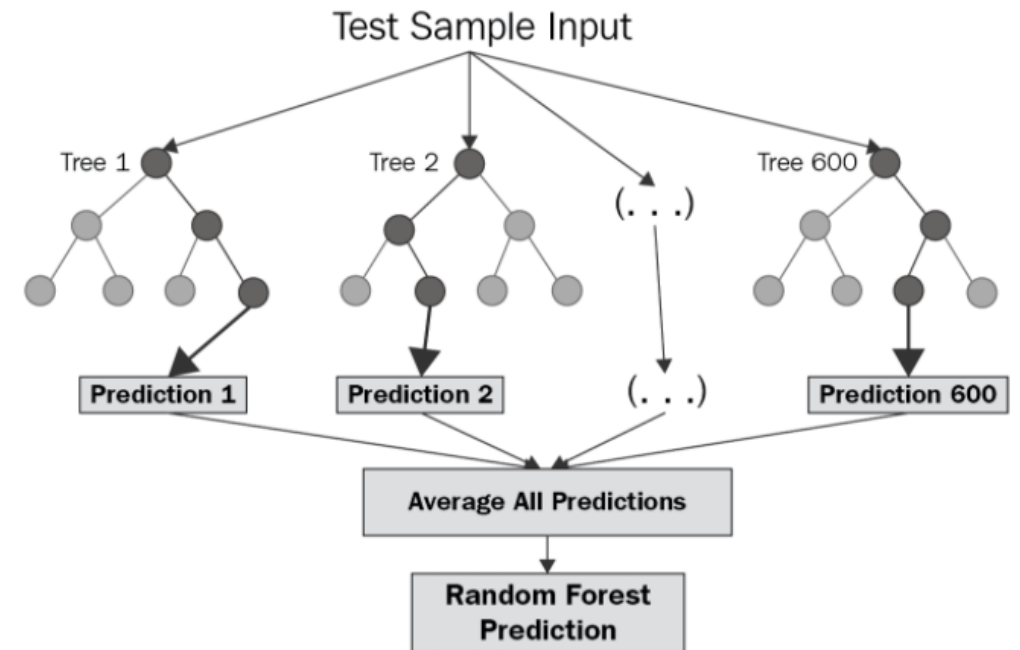➢ Although this algorithm does not require explicit training, KNN is quite sensitive to the local structure of the data.

# Random Forest (RF)

- RF is a type of supervised learning method.

- The RF is an ensemble algorithm, which trains and combines multiple decision trees to produce a highly accurate classifier.

- The prediction for a test sample is then performed by taking a vote from the predictions from all individual trees.

- The ensemble procedure leads to a better performance provided that decision trees are not highly correlated.

- It runs efficiently on large databases and is quite fast to be built and predict.

**Limitation of RF:**

➢ A large number of trees can make the algorithm too slow and ineffective for real-time predictions.



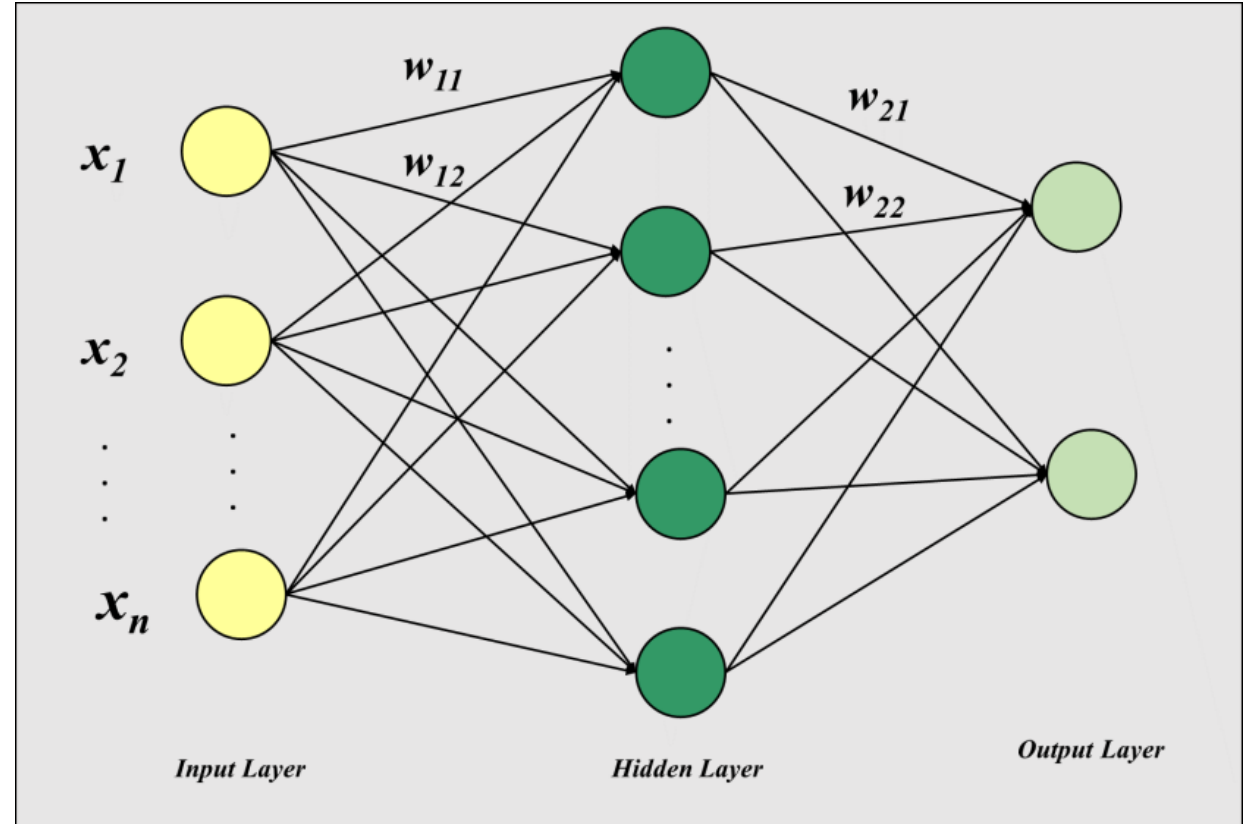Retrieved from https://levelup.gitconnected.com/random-forest-regression-209c0f354c84

# Artificial Neural Network (ANN)

- ANN is another type of supervised learning system that learns how to perform tasks by considering examples.

- ANN is made up of artificial neurons or nodes, where the connection between nodes is modeled as weights and learned from instances.
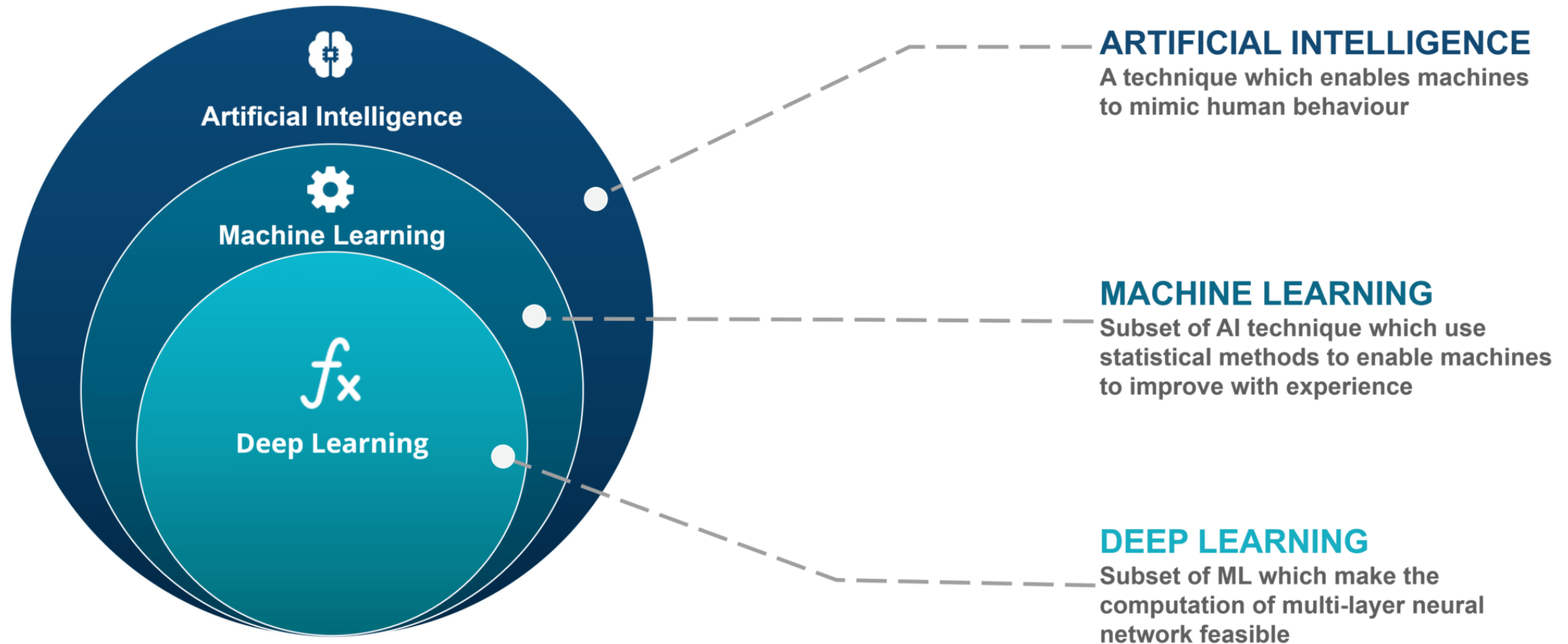


[Retrieved https://www.simplilearn.com/tutorials/deep-learning-tutorial/deep-learning-algorithm]

# ANN Cont.

- In ANN, information ($xi$) flows forward through the network to produce an output (forward propagation).

- The difference between produced output and the true label is then propagated backward through the network to compute the gradient and correct the weights (backpropagation algorithm).

- The backpropagation algorithm updates the weights via iteration technique to improve the network until it can perform the task for which it is being trained.

# Deep Learning



**ARTIFICIAL INTELLIGENCE**
A technique which enables machines to mimic human behaviour

**MACHINE LEARNING**
Subset of AI technique which use statistical methods to enable machines to improve with experience

**DEEP LEARNING**
Subset of ML which make the computation of multi-layer neural network feasible

Retrieved from https://www.edureka.co/blog/ai-vs-machine-learning-vs-deep-learning/

# History of Deep Learning

- 1950s and 1960s: we can learn good features from data

- 1970s and 1980s: the backpropagation algorithm for learning multiple layers of non-linear features was invented.

- Why now?

➢ Large dataset

➢ GPU hardware advance

➢ Improved techniques

# Deep Learning

People are excited about DL because it can solve a lot of AI tasks that are very difficult:


- Image Classification


- Machine translation


- Speech recognition


- Speech synthesis

# Deep Learning Algorithms

1. **Convolutional Neural Networks (CNNs)→ image**
2. Long Short Term Memory Networks (LSTMs)
3. Recurrent Neural Networks (RNNs)→ date with temporal dynamic behavior
4. Generative Adversarial Networks (GANs)
5. Radial Basis Function Networks (RBFNs)
6. Multilayer Perceptron (MLPs)
7. Self Organizing Maps (SOMs)
8. Deep Belief Networks (DBNs)
9. Restricted Boltzmann Machines( RBMs)
10. Autoencoders

# ANN Training

- As the input of ANN is in vector form, the image must be vectorized, where pixel intensities are directly converted into the feature space.

- Since the structural information among neighboring pixels or voxels is a source of information in a given image, this operation destroys spatial information, which is a major drawback of this method for image analysis.

# CNN

- The advantage of CNN compared to ANN is that it releases the constraint of immediate image vectorization by applying convolutional (Conv.) and pooling layers (having neurons arranged in 3 dimensions: width, height, and depth) to the input image.

- The special architecture of CNN helps to extract useful features from the image and use them as the input layer of the classifier (e.g., ANN) that is located at the end layer of the network.
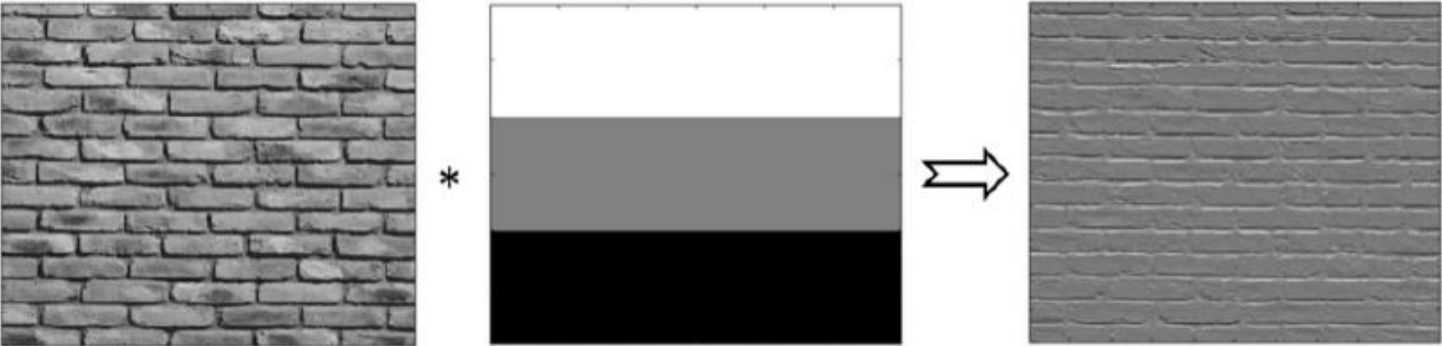


Input Image → Conv. Layer → Detection Layer → Pooling Layer → Artificial Neural Network

Main Components of a Typical CNN

# Convolutional layer Cont.

**Image × Kernel = Feature Map**

$$K = \begin{bmatrix} 1 & 1 & 1 \\ 0 & 0 & 0 \\ -1 & -1 & -1 \end{bmatrix}$$



a)



b)

# Convolutional layer Cont.

# Convolutional layer

- In the convolutional layer, some kernels/filters ($K(i,j)$) convolve with the input image ($I(i,j)$) to produce the feature map ($F(i,j)$).

$$F(i,j) = (I * K)(i,j) = \Sigma\Sigma I(m,n)K(i-m,j-n)nm$$

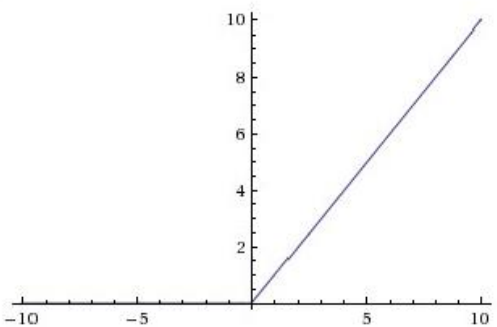$i$ and $j$ denote the position information of each pixel in a digital image of size $m \times n$.

- Conv. layer computes the dot product between the input image and kernel by moving the kernel across the image to detect local features at different positions.

- Mathematically, we get the large value in the feature maps where the template of the filter is found in the input image.

- By doing so, the image features which are similar to the kernel are captured.

- In CNN, filters are considered as weights and are set from the examples during network training via the backpropagation algorithm.

# Convolutional layer Cont.

- Conv. layer improves the ML algorithm through sparse interactions, parameter sharing, and equivariant properties.

- Despite traditional ANN, where every output unit interacts with every input unit (fully connected neural network), CNN has sparse interaction (sparse connectivity/weights) as the size of the kernel is smaller than that of the input image.

- The sparse interaction property not only reduces the memory requirements of the model but also improves its statistical efficiency (a measure of quality of an estimator/an experimental design where fewer observations for statistically efficient is needed compared to less efficient one to achieve a given performance), which leads to fewer operations for output computation.

- In conventional ANN each element in the weight matrix is used once whereas, in CNN, each filter is utilized at every position of the input.

- In other words, in CNN parameter sharing is performed rather than learning a separate set of parameters for every location. Therefore, only one set of parameters is learned.

- This property further reduces the storage requirements of the model.

- Additionally, CNN is equivariance to translation, which means that if the input changes, the output changes in the same way. Put differently, moving an object in the input image will move its representation in the feature map the same amount.
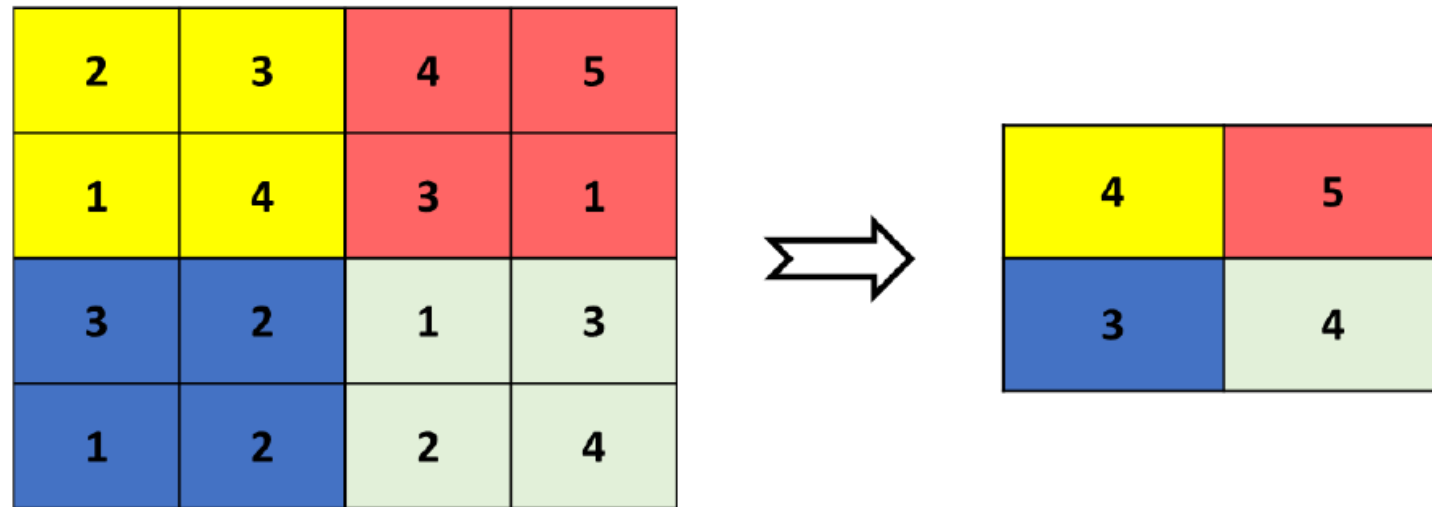
# Detection Layer

- The detection layer (activation function) applies an elementwise activation function to the feature map, in which the size of the feature map does not change.

- Technically, this layer performs non-linear thresholding to the input image.

- This operation increases the nonlinear properties of the decision function as in most of the complex cases, data are not linearly separable.

| | | |
|---|---|---|
| Sigmoid |  | $f(x) = \dfrac{1}{1+exp(-x)}$ |
| tanh |  | $f(x) = \dfrac{e^x - e^{-x}}{e^x + e^{-x}}$ |
| ReLU |  | $f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$ |

# Pooling Layer

- Pooling layers computes a summary statistic (maximum, mean, median, etc.) of some output nodes.

- The advantage of this layer is that it makes the model less sensitive to the local translation.

- Additionally, it yields a simpler model through down-sampling, which reduces the computational complexity of the next set of layers.
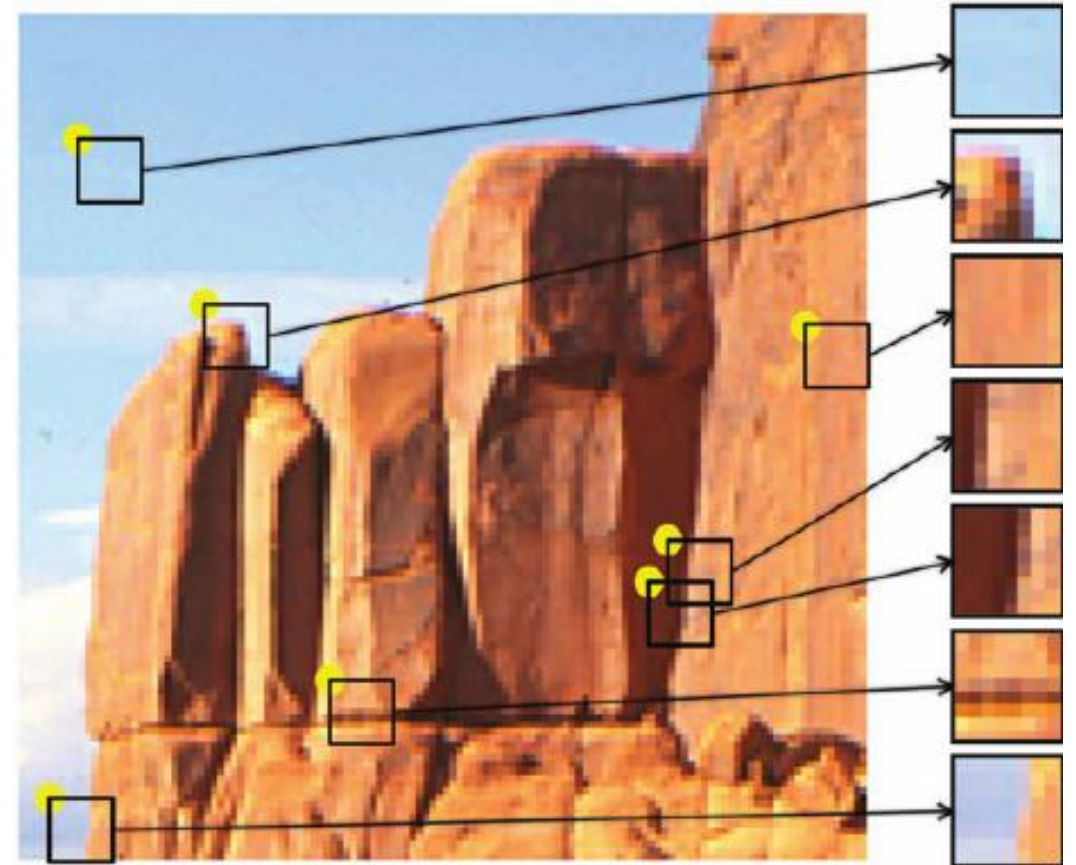


**Feature Map**

# CNN Example

Any combination of the Conv., detection, and pooling layers can be used to design an optimal CNN for a given problem.
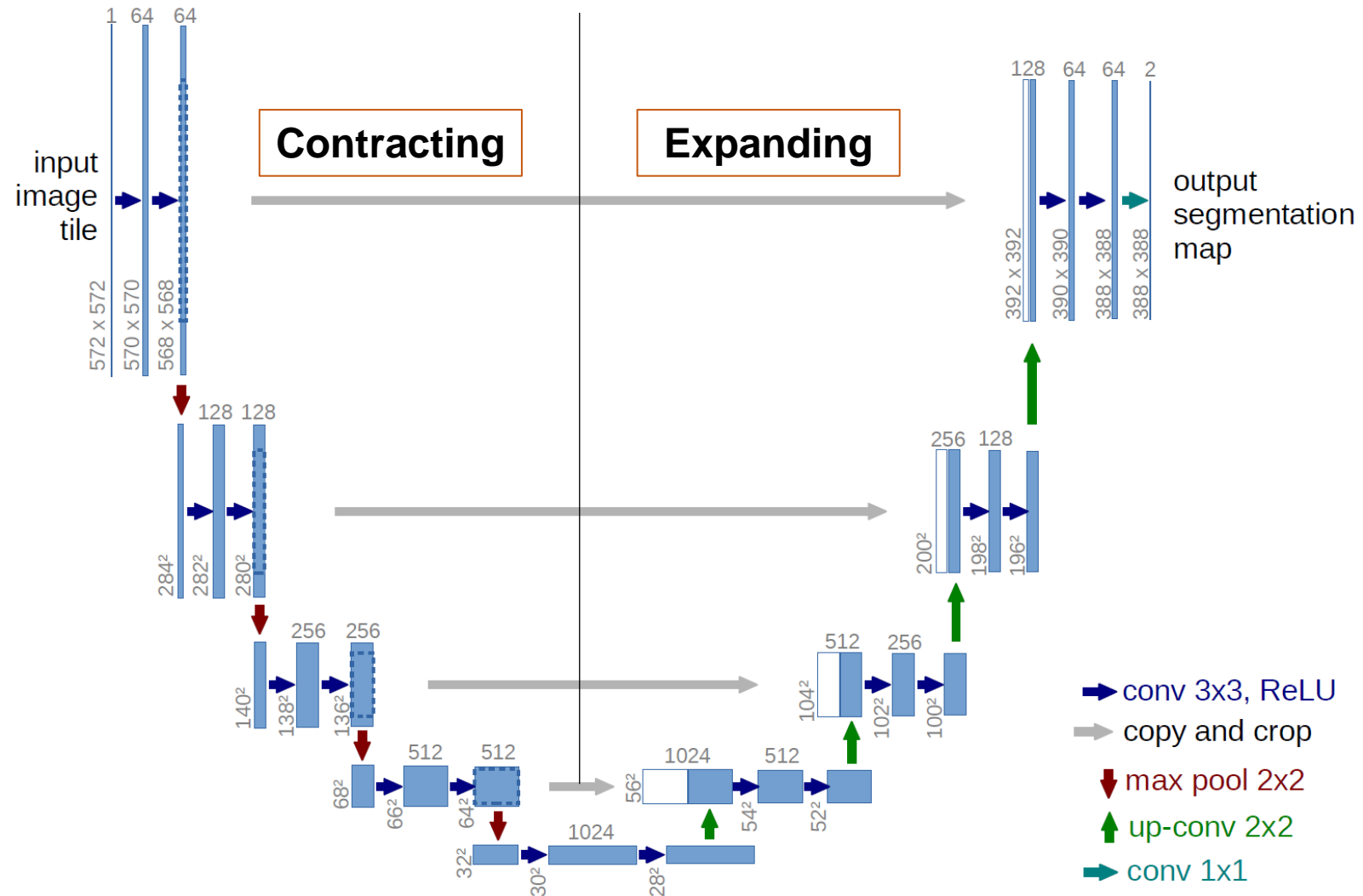
# Limitation of CNN for Image Segmentation

- To employ CNN for segmentation tasks, a local region (patch) around each pixel is extracted in a sliding window fashion and passed through the network to be labeled as foreground or background.

- This method suffers from redundant computation due to the overlap exists between adjacent patches.

- Moreover, global features are not captured when small patches are extracted, and localization error is increased when large patches are used.

- Finally, extracting image patches is time-consuming, requires prior constraint of the target tissue, and overall causes this process to remain slow.

- As the primary goal of applying AI in medical image analysis is to achieve high speed, this approach remains sub-optimal.

- U-Net has been introduced recently to address the limitations of CNN for image segmentation.



Stevens et al., 2013

# U-Net [Ronneberger et al., 2015]

U-Net architecture was introduced explicitly with the segmentation of medical images in mind and used to produce state-of-the-art results on the ISBI challenge for segmentation of neuronal structures in electron microscopic stacks as well as cell tracking in 2015.

# U-Net Cont.

- The U-Net is made up of contracting and expanding paths where pooling and up-sampling layers are used in each way that yields a U-shaped architecture.

- The contraction path is identical to standard CNN in which convolutional layers along with pooling and activation layers are applied to the input data.

- In expanding path, pooling layers are replaced by up-sampling layers to expand the dimension of feature space.

- The output of up-sampling layers is merged with appearance feature representation learned from the corresponding layer in the shrinking path to localize high-resolution features.

- U-Net does not need patch extraction and generates a segmentation map with the corresponding size of the input image that accommodates images of arbitrary sizes.

- U-Net does not have a fully connected neural network at the end.

Note: **Transposed convolution** is used in the Up-sampling Path but not deconvolution, deconvolution attempts to reverse the effects of a convolution. In a transposed convolutional layer, the regular convolution operation is performed but its spatial transformation is reversed.
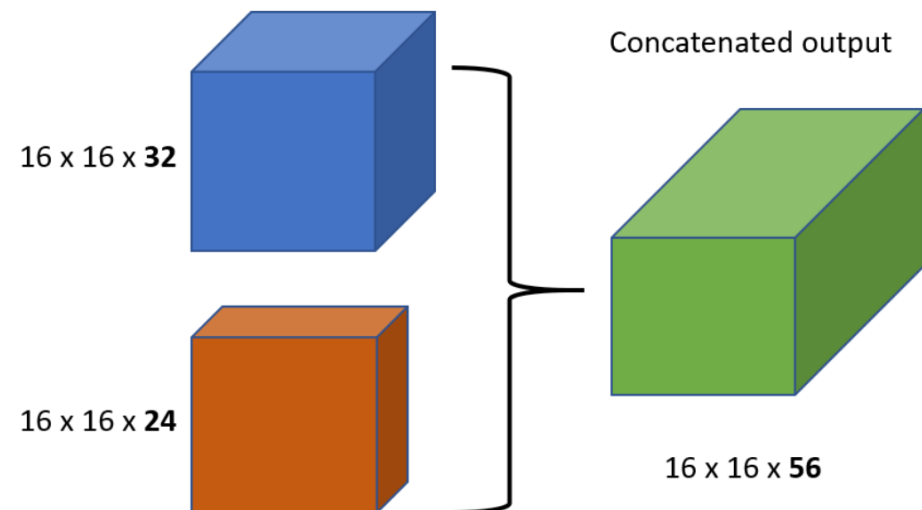
# Why Skip Connection in U-Net?

- One of the main problem in training deep neural networks is Vanishing Gradient problem.

- It happens when the training loss stops decreasing while it is still far away from the desired value.

- Update rule in gradient descent:

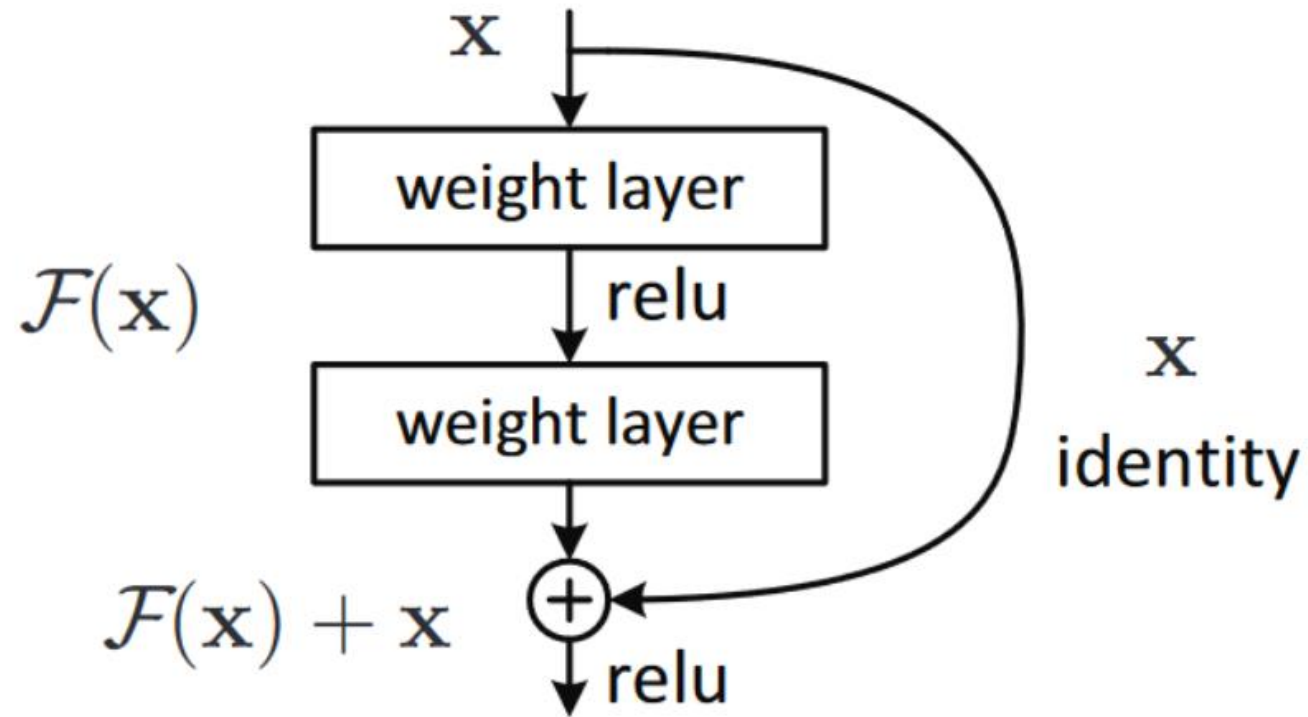$$\theta := \theta - \alpha \frac{d}{d\theta} J(\theta)$$

- The idea of backpropagation is to gradually minimize the loss by updating the parameters of the network.

- When learning rate and gradient of loss functions are very small, we do not observe any change in the model while training our network.

# Skip Connections in Deep Learning

- Skip connections provide an alternative path for the gradient.

- It has been shown that these additional paths are beneficial for the model convergence.

- Skip connections skip some layers in the neural network and feed the output of one layer as the input of next layers (instead of only the next one).

- In general there are two types of skip connections in deep NNs:

➢Concatenation as in densely connected

architecture.

➢Addition as in residual architectures.

Retrieved from https://theaisummer.com/skip-connections/

16 x 16 x 32

16 x 16 x 24

Concatenated output

16 x 16 x 56

# ResNet: Skip Connection via Addition



[He et al., 2015]

# DenseNet: Skip Connections via Concatenation



➢ Transition Layers
➢ Dense Blocks



[Huang et al., 2018]

# Densely Connected Convolutional Networks

- Convolutional networks can be substantially deeper, more accurate, and efficient to train if they contain shorter connections between layers close to the input and those close to the output.

- Dense Convolutional Network (DenseNet) embraces this observation and connects each layer to every other layer in a feed-forward fashion.

- For each layer, the feature-maps of all preceding layers are used as inputs, and its own feature-maps are used as inputs into all subsequent layers.

**DenseNets:**

➤alleviate the vanishing-gradient problem,

➤strengthen feature propagation,

➤encourage feature reuse,

➤and substantially reduce the number of parameters.

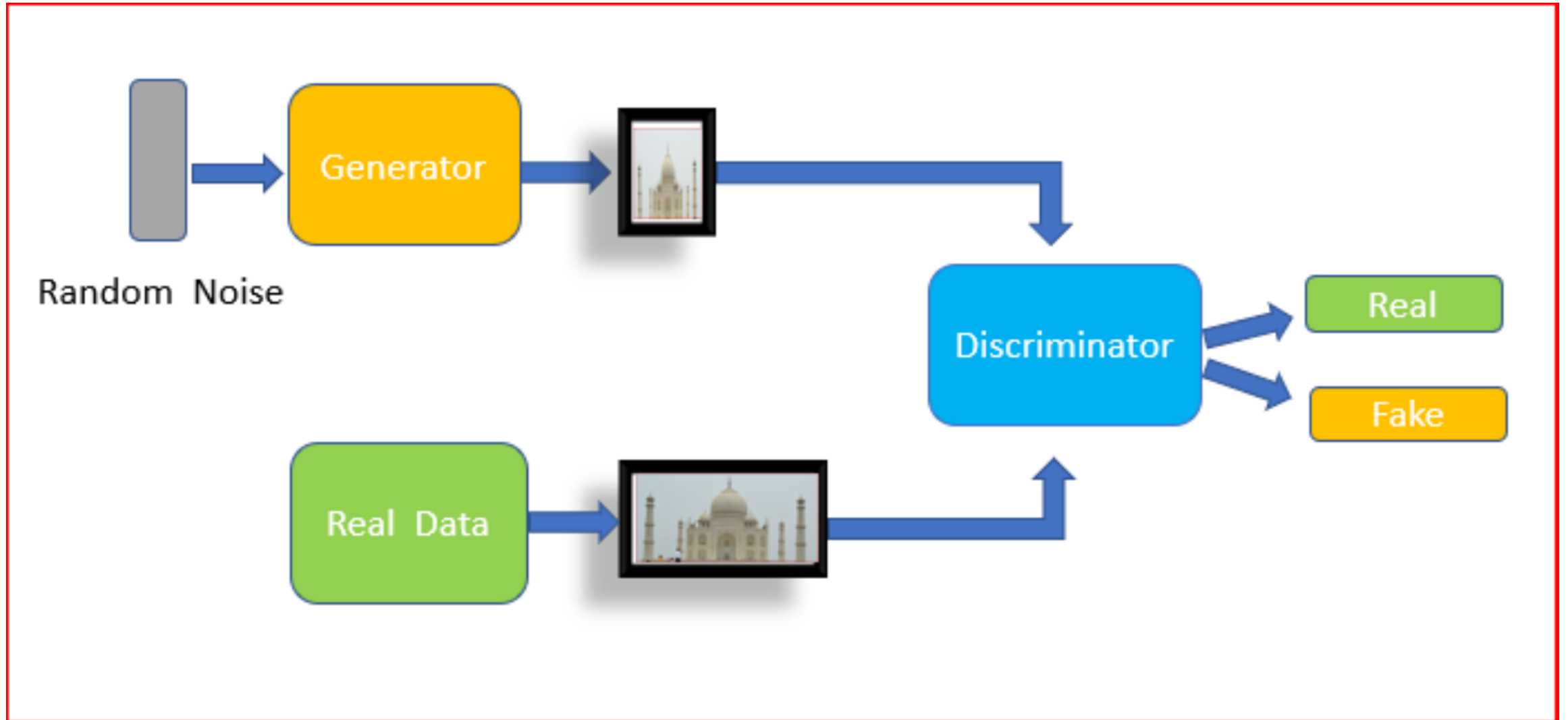# Generative Adversarial Networks (GAN) [M. Kana,April 2020, towardsdatascience.com]

- GAN is an old idea arising from the game theory, introduced by Ian J. Goodfellow and co-authors in 2014.

- "Generative Adversarial Network— the most interesting idea in the last ten years in machine learning" by Yann LeCun.

- We can generate new photo-realistic images with a Variational AutoEncoder. VAEs typically produce blurry and non-photorealistic faces, which was a motivation to built GANs.

- Generating new faces can be expressed by a random variable generation problem, where the face is described by random variables, represented through its RGB values, flatten into a vector of $N$ numbers.

- A GAN generates a new face by generating a new vector following the celebrity face probability distribution over the N-dimensional vector space, which is a very complex one and we don't know how to directly generate complex random variables.

- The complex random variable can be represented by a function applied to a uniform random variable (transform method) in which N uncorrelated uniform random variables is generated. It then applies a very complex function to that simple random variable!

- Very complex functions are naturally approximated by a neural network. Once the network is trained, it will be able to take a simple N-dimensional uniform random variable as input and return another N-dimensional random variable that would follow our celebrity-face probability distribution. This is the core motivation behind generative adversarial networks.

# GAN Cont.

- During each training iteration of the neural network, we need to compare a sample of faces from training set with a sample of generated faces.

- Theoretically, we would compare the true distribution versus the generated distribution based on samples using the Maximum Mean Discrepancy (MMD) approach. To do so the distribution matching error must be computed that is used to update the network via backpropagation, which is practically very complex to implement.

- GANs address this issue by solving a non-discrimination task between true and generated samples Instead of directly comparing both true and generated distributions.

- A GAN has three components: a generator model for generating new data, a discriminator model for classifying whether generated data are real faces, or fake, and the adversarial network that pits them against each other.

- The generative part is responsible for taking N-dimensional uniform random variables as input and generating fake faces.

- The discriminative part, which is a simple classifier, evaluates and distinguished the generated faces from true celebrity faces.

- It is important that both networks learn equally during training and converge together.
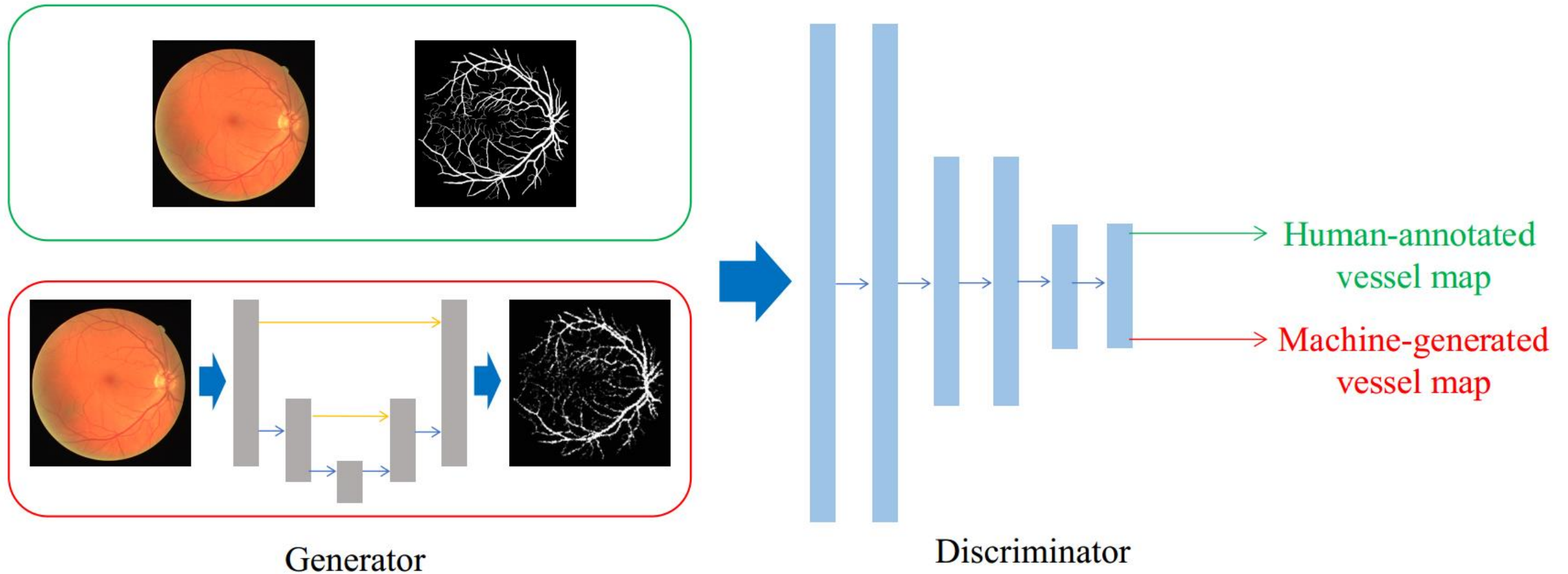
# GAN Cont.

# GAN for Medical Image Segmentation

- The generator can be a U-Net that produce segmentation map.
- The discriminator can be a CNN that take the generated segmentation map and compare it to the ground truth mask.
- The generator is creating segmentation map for an input image and passes it to the discriminator.
- It does so in the hopes that the segmentation mask will be deemed authentic.
- The goal of the generator is to generate passable mask.
- The goal of the discriminator is to identify masks coming from the generator no similar to the ground truth.

# Example: Retinal Vessel Segmentation in Fundoscopic Images with Generative Adversarial Networks [arXiv:1706.09318v1]



Generator

Discriminator

Human-annotated vessel map

Machine-generated vessel map

# Object Detection Algorithms [ Gandhi, 2018, towardsdatascience.com]

- In detection algorithms, the goal is to draw a bounding box around the object of interest to locate it within the image.

- A naïve approach would be to take different regions of interest from the image and use a CNN to classify the presence of the object within that region.

- The problem with this approach is that the objects of interest might have different spatial locations within the image and different aspect ratios. Hence, a huge number of regions must be selected that could be computationally expensive.

- Therefore, algorithms like R-CNN, Fast R-CNN, Faster R-CNN, YOLO, etc. have been developed to find the occurrences of objects in an image fast.

- **Regions with CNN** features (R-CNN) has been proposed in which selective search algorithm is employed to extract just 2000 regions from the image that is called region proposals.

- The CNN acts as a feature extractor and the output dense layer consists of the features extracted from the image and the extracted features are fed into an SVM to classify the presence of the object within that candidate region proposal.

# Problems with R-CNN

- It still takes a huge amount of time to train the network as 2000 region proposals per image must be classified.

- It cannot be implemented real time.

- The selective search algorithm is a fixed algorithm. Therefore, no learning is happening at that stage. This could lead to the generation of bad candidate region proposals.

- Fast R-CNN has been introduced, which is similar to the R-CNN algorithm but, instead of feeding the region proposals to the CNN, the input image is fed to the CNN to generate a convolutional feature map. From the convolutional feature map, the region of proposals are identified.

- Faster R-CNN, YOLO, Mask R-CNN

# Clinical Objective

| Why LV Scar Segmentation? | → | • Treatment Plan<br>• Risk Management for Ventricular Arrhythmia |
|---|---|---|



Left ventricle

Right ventricle

Retrieved from http://www.anatomynote.com



Blocked Lumen in Branch of Left Coronary Artery

Anterior Infarct

Retrieved from http://medimoon.com

# 2D Versus 3D LGE-MRI

- Currently, 2D LGE MRI is used to identify myocardial scar.

- 3D LGE-MRI has emerged, enabling more accurate spatial representation and quantification[1].



An Example of 2D LGE-MRI

[1] Kawaji, K, et al., "3D Late gadolinium enhanced cardiovascular MR with CENTRA-PLUS profile/view ordering: Feasibility of right ventricular myocardial damage assessment using a swine animal model," MRI, Papers 39, 7-14 (2017).

# Problem Statement

**Automated LV scar segmentation is challenging**

**1** High variability of cardiac structures

**2** Complexity in segmenting the apical and basal slice images

**3** Inherent noise associated with cardiac MRI

**4** Dynamic motion of heart

# Limitations of Previous Works

**Intensity-based methods**
Full Width at Half Maximum (FWHM)
Signal Threshold to Reference Mean (STRM)
Region Growing (RG)

**Limitations**
Highly influenced by image noise
Low robustness and accuracy

**Deformable model techniques**
Hierarchical Max Flow (HMF)

**Limitations**
Subject to high operator variability

# Our Proposed Techniques

# Dataset

- Acquired at Robarts Research Institute of Western University (ON, Canada).
- They had chronic MI
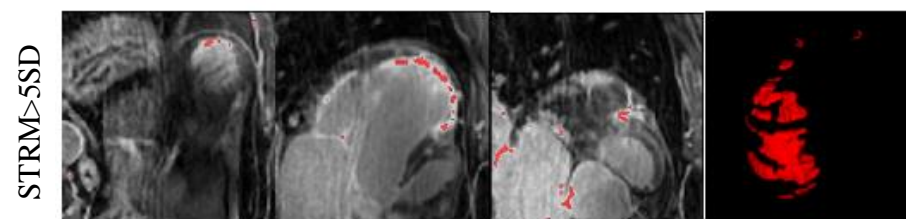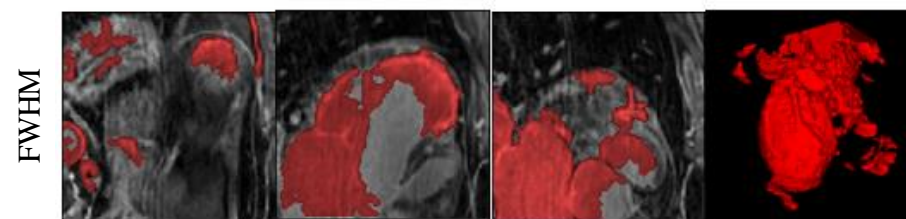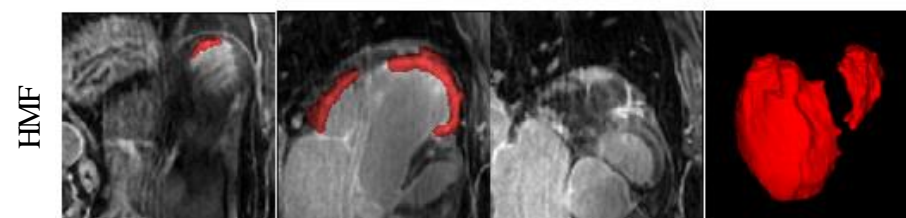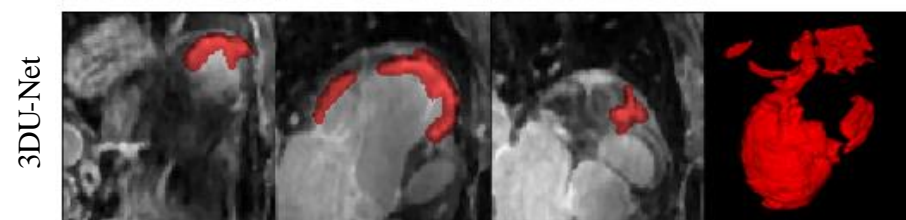- Mean (LVEF) of 32.1 ± 12.7%

Training Set:
N=10

Testing Set:
N=24

# Architecture of the CNN

# Exemplary Results



**Ground Truth**     **CNN**

| Dice Index (%) | AVD (ml) |
|:---:|:---:|
| 93.63 ± 2.61 | 2.08 ± 1.81 |

# Our Proposed Method



Input 3D
LGE-MRI

Extracted Slices from Input
Image in the Axial, Sagittal,
and Coronal Directions

Myo-Net

Segmented Myocardium in the
Axial, Sagittal, and Coronal
Directions

3D View of
Segmented LV
Myocardium

Scar-Net

Segmented Scar in the
Axial, Sagittal, and
Coronal Directions

3D View of
Segmented LV
Scar

# Dataset



- Acquired at Robarts Research Institute of Western University (ON, Canada).
- They had chronic MI
- Mean (LVEF) of 32.1 ± 12.7%

Training Set:
N=18

Testing Set:
N=16

# Results



|  | LV Myocardium | LV Scar |
|---|---|---|
| DSC (%) | 85.14 ± 3.36 | 88.61 ± 2.54 |
| HD (mm) | 19.21 ± 4.74 | 17.04 ± 9.93 |
| AVD (cm$^3$) | 43.72 ± 27.18 | 9.33 ± 7.24 |
| Required Time (s) | 49.96 ± 9.76 | 120.45 ± 23.34 |

# Comparison Results



| Method | DSC (%) |
|--------|---------|
| **Proposed** | **88.61 ± 2.54** |
| **CCU-Nets** | 85.69 ± 4.20 |
| **DMPU-Net** | 85.19 ± 4.59 |
| **DCU-Net** | 83.48 ± 4.81 |
| **3DU-Net** | 72.39 ± 7.02 |
| **HMF** | 77.41 ± 3.05 |
| **FWHM** | 64.16 ± 7.65 |
| **STRM>5SD** | 64.07 ± 7.14 |

# Thank you!
# Any Question?