



# Sparsity in CNNs

Journal Club (11th Dec 2017)



# **Marginal Space Deep Learning: Efficient Architecture for Volumetric Image Parsing**

## 9-Parameter Affine Transform

$$\begin{bmatrix} x' \\ y' \\ z' \\ 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & Tx \\ 0 & 1 & 0 & Ty \\ 0 & 0 & 1 & Tz \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix}$$

4x4 translation  
matrix **T** (3  
parameters)

$$\begin{bmatrix} x' \\ y' \\ z' \\ 1 \end{bmatrix} = \begin{bmatrix} Sx & 0 & 0 & 0 \\ 0 & Sy & 0 & 0 \\ 0 & 0 & Sz & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix}$$

4x4 scale matrix **S**  
(3 parameters)

$$\begin{bmatrix} x' \\ y' \\ z' \\ 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & \cos(\phi_x) & -\sin(\phi_x) & 0 \\ 0 & \sin(\phi_x) & \cos(\phi_x) & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix}$$

rotation matrix **R<sub>x</sub>** for  
rotation about x-axis  
(1 parameter)

$$\begin{bmatrix} x' \\ y' \\ z' \\ 1 \end{bmatrix} = \begin{bmatrix} \cos(\phi_y) & 0 & -\sin(\phi_y) & 0 \\ 0 & 1 & 0 & 0 \\ \sin(\phi_y) & 0 & \cos(\phi_y) & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix}$$

rotation matrix **R<sub>y</sub>** for  
rotation about y-axis  
(1 parameter)

$$\begin{bmatrix} x' \\ y' \\ z' \\ 1 \end{bmatrix} = \begin{bmatrix} \cos(\phi_z) & -\sin(\phi_z) & 0 & 0 \\ \sin(\phi_z) & \cos(\phi_z) & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix}$$

rotation matrix **R<sub>z</sub>** for  
rotation about z-axis  
(1 parameter)

$$M = [S][R_x][R_y][R_z][T]$$

Above ordering does translations first to match origins, then rotations about new origin, and finally scaling of the aligned image. Most software used 9 or 12 parameters (sometimes called degrees of freedom) for registration.

# Problem

Segmenting aortic valve in TEE images

Volumes vary from  $100 \times 100 \times 50$  to  $250 \times 250 \times 150$  voxels,  
and 0.75 to 1mm

2891 volumes from 869 patients

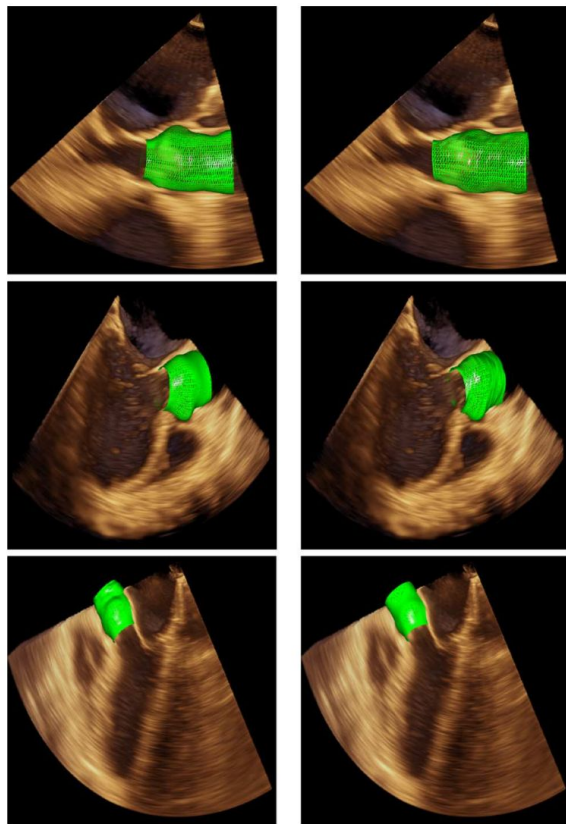


Fig. 9. Example images showing the aortic valve segmentation results for different patients from the test set, using our proposed method. Each row of images corresponds to a different patient, the left image represents the ground truth mesh, while the right image shows the detected mesh.



# Estimate 9D space

Marginal Space Deep Learning

*Learning classifiers in clustered, high probability regions of spaces*

- 1) Roughly initialize the 9 parameters (localization)
- 2) Used deep learning to “guide the shape deformation” (fine-tuning)

# Sparse Adaptive Deep Neural Network (SADNN)

Extracting sufficiently representative training set

Training set =  $m$  patches with a set of class assignment ( $y$ )

Find a sparsity map for the weights ( $w$ )

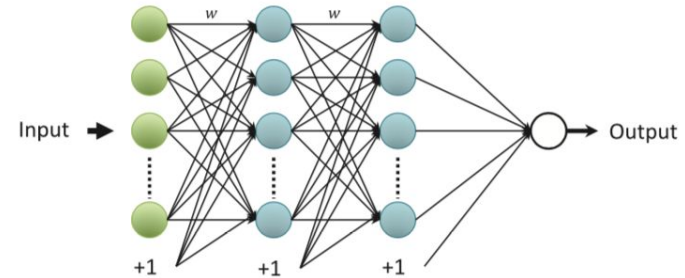


Fig. 2. Example of a fully connected neural network with 3 layers. Every neuron in a layer is connected to all neurons in the previous layer.

---

**Algorithm 1** Learning algorithm with iterative threshold-enforced sparsity

---

- 1: Pre-training using all weights  $\vec{w}^{(0)} \leftarrow \vec{w}$  (small # epochs)
  - 2: Initialize sparsity map  $\vec{s}^{(0)}$  with ones
  - 3:  $t \leftarrow 1$
  - 4: **for** training round  $t \leq T$  **do**
  - 5:     **for all** filters  $i$  with sparsity **do**
  - 6:          $\vec{s}_i^{(t)} \leftarrow \vec{s}_i^{(t-1)}$  + remove smallest active weights
  - 7:          $\vec{w}_i^{(t)} = \vec{w}_i^{(t-1)} \odot \vec{s}_i^{(t)}$
  - 8:         Normalize active coeff. s.t.  $\|\vec{w}_i^{(t)}\|_1 = \|\vec{w}_i^{(t-1)}\|_1$
  - 9:     **end for**
  - 10:      $\vec{b}^{(t)} \leftarrow \vec{b}^{(t-1)}$  (copy current biases)
  - 11:     Train network on active weights (small # epochs)
  - 12:      $t \leftarrow t + 1$
  - 13: **end for**
  - 14: Sparse kernels:  $\vec{w}_s \leftarrow \vec{w}^{(T)}$
  - 15: Biases:  $\vec{b}_s \leftarrow \vec{b}^{(T)}$
-

---

**Algorithm 1** Learning algorithm with iterative threshold-enforced sparsity

---

- 1: Pre-training using all weights  $\vec{w}^{(0)} \leftarrow \vec{w}$  (small # epochs)
  - 2: Initialize sparsity map  $\vec{s}^{(0)}$  with ones
  - 3:  $t \leftarrow 1$
  - 4: **for** training round  $t \leq T$  **do**
  - 5:     **for all** filters  $i$  with sparsity **do**
  - 6:          $\vec{s}_i^{(t)} \leftarrow \vec{s}_i^{(t-1)}$  + remove smallest active weights
  - 7:          $\vec{w}_i^{(t)} = \vec{w}_i^{(t-1)} \odot \vec{s}_i^{(t)}$
  - 8:         Normalize active coeff. s.t.  $\|\vec{w}_i^{(t)}\|_1 = \|\vec{w}_i^{(t-1)}\|_1$
  - 9:     **end for**
  - 10:      $\vec{b}^{(t)} \leftarrow \vec{b}^{(t-1)}$  (copy current biases)
  - 11:     Train network on active weights (small # epochs)
  - 12:      $t \leftarrow t + 1$
  - 13: **end for**
  - 14: Sparse kernels:  $\vec{w}_s \leftarrow \vec{w}^{(T)}$
  - 15: Biases:  $\vec{b}_s \leftarrow \vec{b}^{(T)}$
-



---

**Algorithm 1** Learning algorithm with iterative threshold-enforced sparsity

---

- 1: Pre-training using all weights  $\vec{w}^{(0)} \leftarrow \vec{w}$  (small # epochs)
  - 2: Initialize sparsity map  $\vec{s}^{(0)}$  with ones
  - 3:  $t \leftarrow 1$
  - 4: **for** training round  $t < T$  **do**
  - 5:     **for all filters**  $i$  **with sparsity do**
  - 6:          $\vec{s}_i^{(t)} \leftarrow \vec{s}_i^{(t-1)}$  + remove smallest active weights
  - 7:          $\vec{w}_i^{(t)} = \vec{w}_i^{(t-1)} \odot \vec{s}_i^{(t)}$
  - 8:         Normalize active coeff. s.t.  $\|\vec{w}_i^{(t)}\|_1 = \|\vec{w}_i^{(t-1)}\|_1$
  - 9:     **end for**
  - 10:      $\vec{b}^{(t)} \leftarrow \vec{b}^{(t-1)}$  (copy current biases)
  - 11:     Train network on active weights (small # epochs)
  - 12:      $t \leftarrow t + 1$
  - 13: **end for**
  - 14: Sparse kernels:  $\vec{w}_s \leftarrow \vec{w}^{(T)}$
  - 15: Biases:  $\vec{b}_s \leftarrow \vec{b}^{(T)}$
- 

“...in later stages of the training exponentially less filter weights are set to zero ”



**Dropout:** Similar end result

**Translation Invariance:** No convs. BUT each sparse pattern is invariant to position so we get translation invariance

**Pooling:** None here

**Kernel structure:** No longer square because we have an adaptive sampling method



## Marginal Space Deep Learning

$$\begin{aligned}(\hat{\vec{T}}, \hat{\vec{R}}, \hat{\vec{S}}) &= \arg \max_{\vec{T}, \vec{R}, \vec{S}} p(\vec{T}|I)p(\vec{R}|\vec{T}, I)p(\vec{S}|\vec{T}, \vec{R}, I) \\ &= \arg \max_{\vec{T}, \vec{R}, \vec{S}} \frac{p(\vec{T}, \vec{R}, \vec{S}|I)}{p(\vec{T}|I)p(\vec{R}|\vec{T}, I)},\end{aligned}$$

**Marginal spaces**

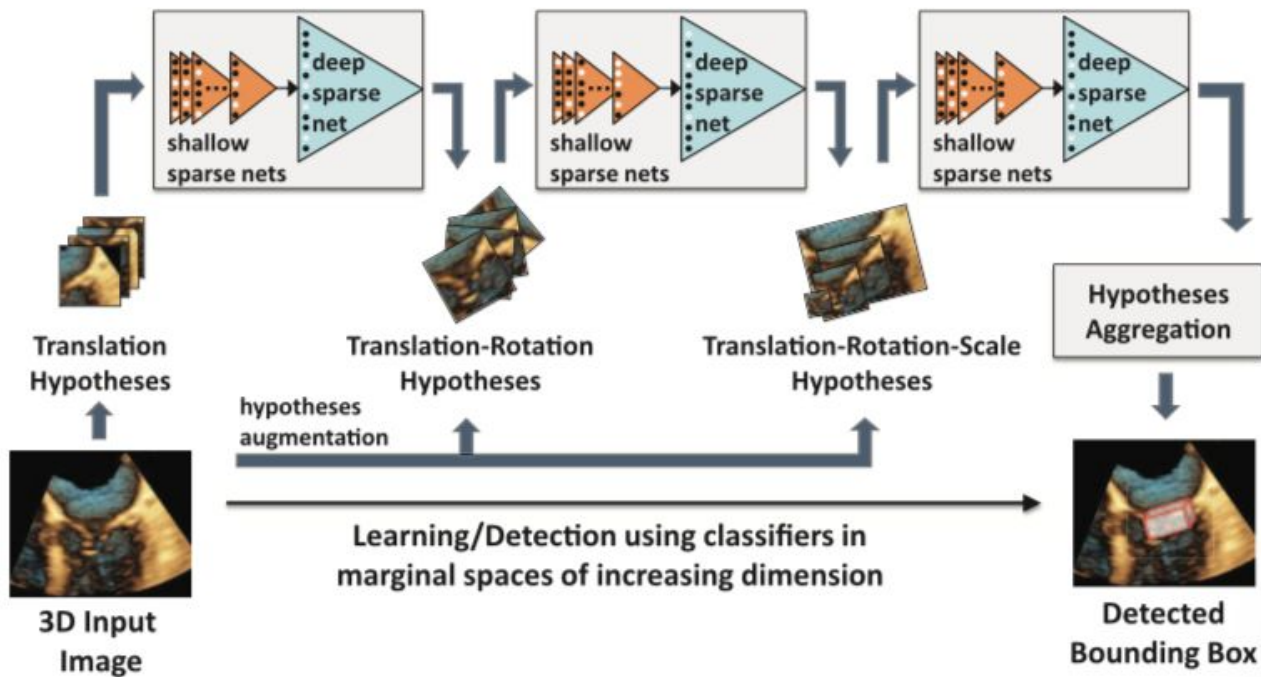


Fig. 4. Schematic visualization of the marginal space deep learning framework. The black/white dots encode the sparse sampling patterns.

# Balancing training set

Due to limited range in position, orientation and space, great class imbalance exists (1:1000)

Series of shallow nets

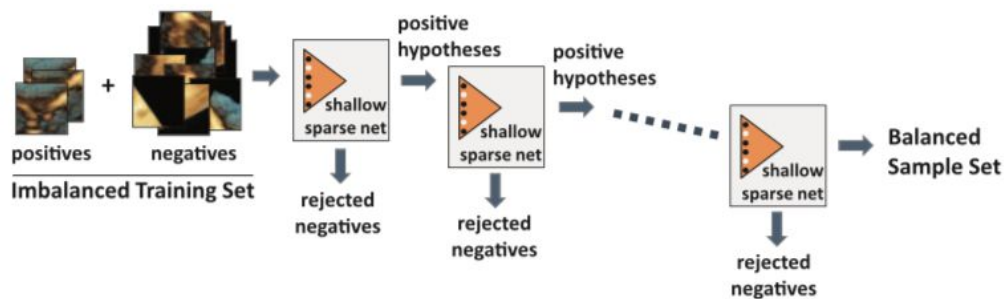


Fig. 5. Schematic visualization of the negative-sample filtering cascade used to balance the training set. The black/white dots encode the sparse sampling patterns of the shallow nets.

# Boundary estimation

Boundary classification problem using translation and orientation along the normal

Training data: samples along the ground truth boundary (positive) and various distances from boundary (negative)

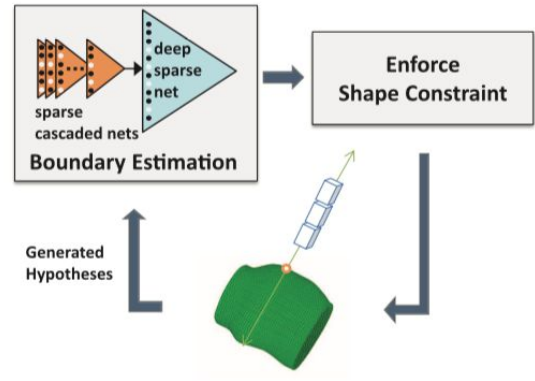
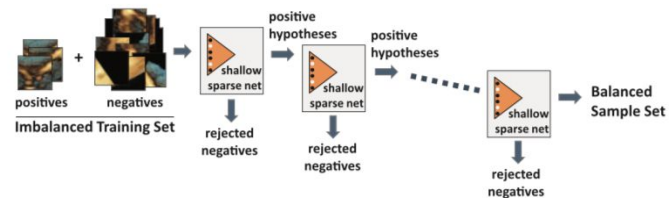


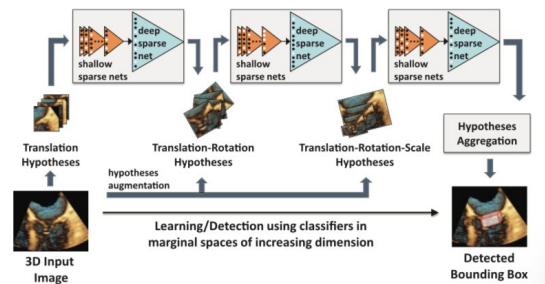
Fig. 6. Schematic visualization of the boundary deformation with SADNN. Starting from the current shape, the SADNN is aligned and applied along the normal for each point of the mesh, the boundary is deformed and projected under the current shape space. The process is iteratively repeated. The black/white dots encode the sparse patterns, learned in the cascaded shallow networks and deep boundary classifier.

# Experiments

Cascade network: 5832 (sparse) x 60 x 1



Main classifier: 5832 (sparse) x 150 x 80 x 50 x 1

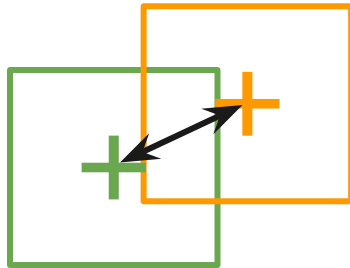


---

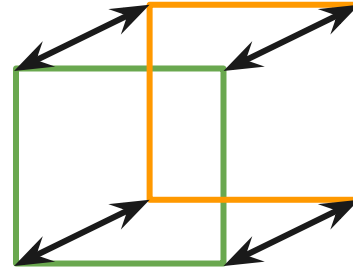
# Experiments

Cascade network: 5832 (sparse) x 60 x 1

Main classifier: 5832 (sparse) x 150 x 80 x 50 x 1



translation



orientation and scale





TABLE I

COMPARISON OF THE PERFORMANCE OF THE STATE-OF-THE-ART MSL [5] AND THE PROPOSED MSDL FRAMEWORK FOR AORTIC VALVE DETECTION. THE MEASURES USED TO QUANTIFY THE QUALITY OF THE RESULTS W.R.T TO THE GROUND TRUTH DATA ARE THE ERROR OF THE POSITION OF THE BOX AND MEAN CORNER DISTANCE (BOTH MEASURED IN MILLIMETRES). THE SUPERIOR RESULTS ARE DISPLAYED IN BOLD

	Position Error [mm]				Corner Error [mm]			
	Training Data		Test Data		Training Data		Test Data	
	MSL	MSDL	MSL	MSDL	MSL	MSDL	MSL	MSDL
Mean	3.12	<b>1.47</b>	3.34	<b>1.83</b>	5.42	<b>2.80</b>	6.16	<b>3.72</b>
Median	2.80	<b>1.27</b>	3.05	<b>1.58</b>	4.98	<b>2.58</b>	5.85	<b>3.34</b>
STD	1.91	<b>0.99</b>	1.85	<b>1.31</b>	2.47	<b>1.23</b>	2.31	<b>1.74</b>

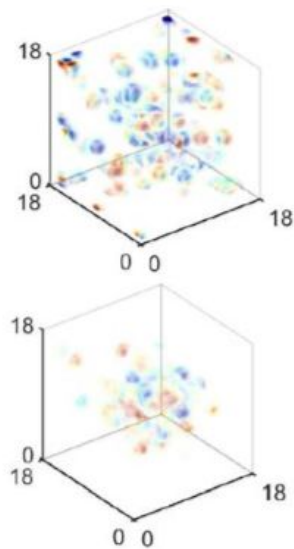
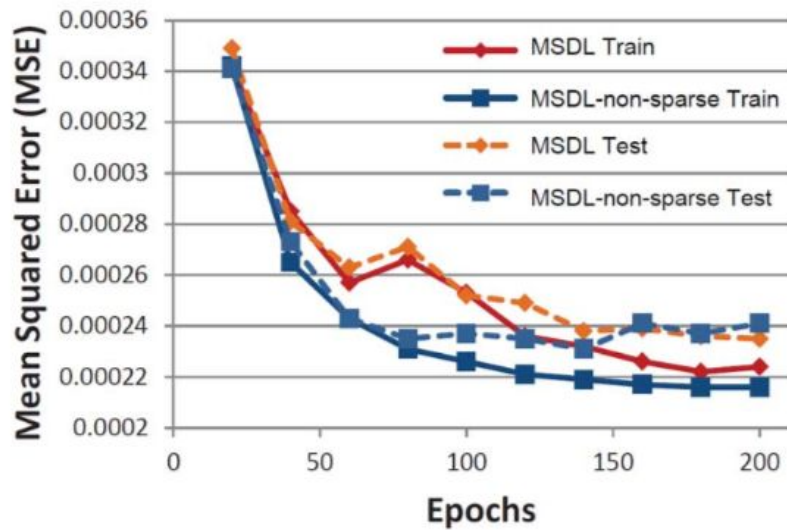


Fig. 7. Left: Training and testing error of the SADNN in the translation stage. The performance oscillation of the sparse model is explained by the enforcement of sparsity in the dense sampling layer of the network. Right: Example color-coded sparse patterns for translation (upper box) and full space (lower box). The latter representation is more compact because of the better data alignment.

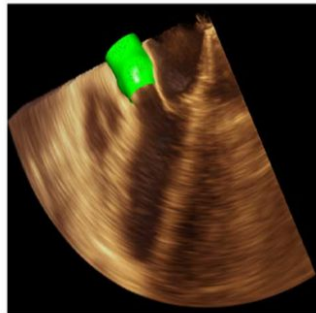
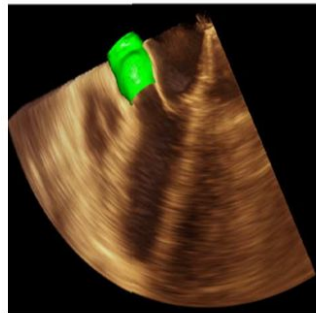
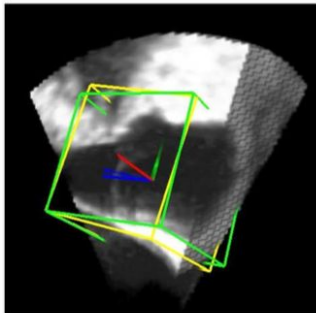
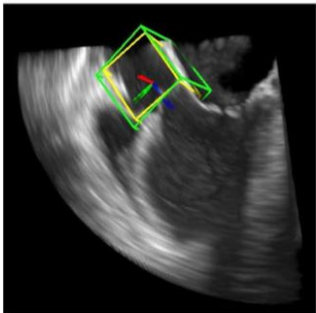
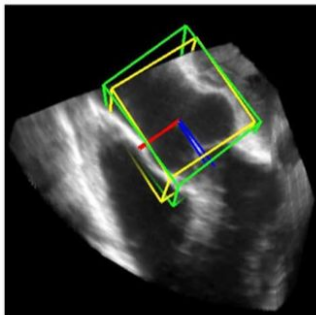
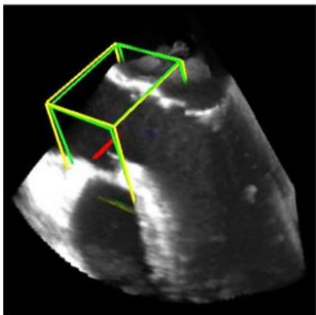
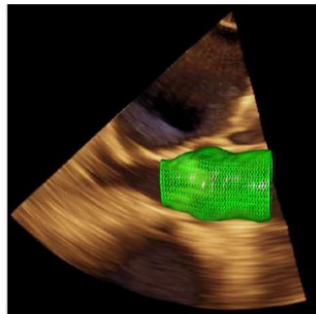
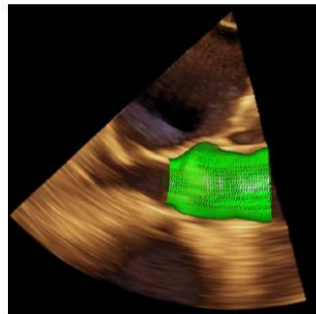
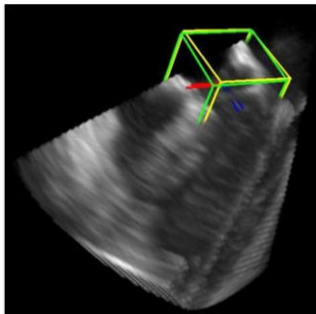
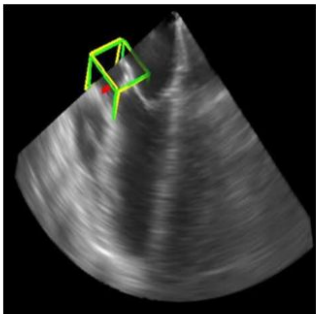




TABLE II

COMPARISON OF THE PERFORMANCE OF THE STATE-OF-THE-ART MSL [5] AND THE PROPOSED MSDL FRAMEWORK FOR AORTIC VALVE SEGMENTATION. THE MEASURE ILLUSTRATES THE DISTANCE OF THE DETECTED MESH TO THE GROUNDTRUTH FOR BOTH THE INITIALIZATION FROM THE DETECTED BOX AS WELL AS THE DISTANCE AFTER BOUNDARY REFINEMENT. THE SUPERIOR RESULTS ARE DISPLAYED IN BOLD

	Training Data, distance to ground truth mesh [mm]				Test Data, distance to ground truth mesh [mm]			
	Initialization		Final		Initialization		Final	
	MSL	MSDL	MSL	MSDL	MSL	MSDL	MSL	MSDL
Mean	2.08	<b>1.16</b>	1.17	<b>0.89</b>	2.06	<b>1.21</b>	1.04	<b>0.90</b>
Median	1.94	<b>1.09</b>	1.05	<b>0.82</b>	1.95	<b>1.10</b>	0.98	<b>0.80</b>
STD	0.83	<b>0.40</b>	0.66	<b>0.35</b>	0.79	<b>0.55</b>	0.50	<b>0.48</b>



TABLE III  
COMPARISON OF THE PERFORMANCE OF OUR DL-BASED FRAMEWORK WITH THE STATE-OF-THE-ART SOLUTION [5] AND THE TWO MIXED VARIANTS: COMBINING THE MSDL BOX DETECTION WITH THE ORIGINAL SEGMENTATION METHOD [5], RESPECTIVELY THE MSL BOX DETECTION [5] WITH OUR DL-BASED SEGMENTATION METHOD. THE SUPERIOR RESULTS ARE DISPLAYED IN BOLD

	Segmentation error [mm]		
	Mean	Median	STD
<b>Ours</b>	<b>0.90</b>	<b>0.80</b>	<b>0.48</b>
Reference [5]	1.04	0.98	0.50
<b>Detection (ours) + Seg. [5]</b>	0.95	0.88	0.48
<b>Detection [5] + Seg. (ours)</b>	1.00	0.90	0.51



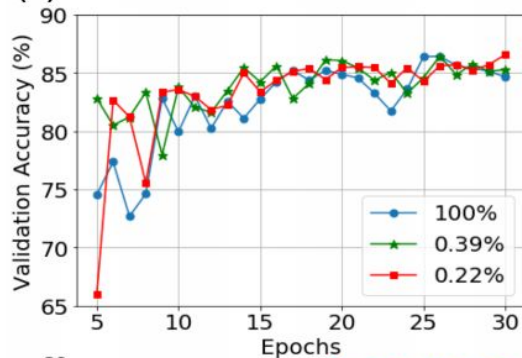
# Characterizing Sparse Connectivity Patterns in Neural Networks

<https://arxiv.org/abs/1711.02131>

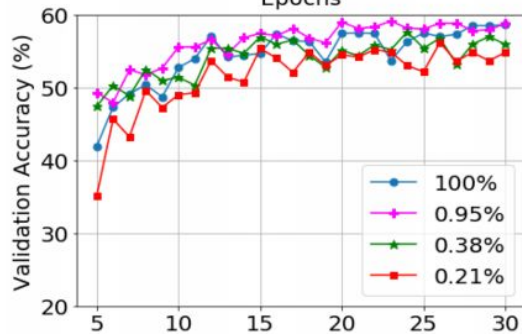
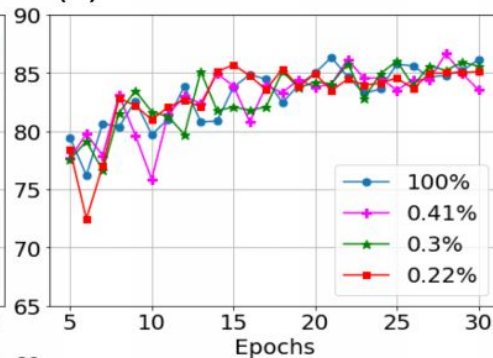
Metric for measuring sparsity in fully  
connected layers in CNN

If you reduced connectivity in fully connected layers, it made little difference to overall performance

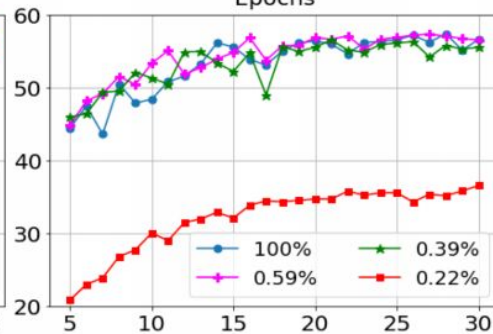
(a) CIFAR10 conv2



(b) CIFAR10 conv3



(d) CIFAR100 conv2



(e) CIFAR100 conv3

# Window Adjacent Matrix ( $A^w$ )

Overall density =  $(8+8)/(32+16) = 33\%$

Junction 1

$N_1 = 8, N_2 = 4$

$fo_1 = 1, fi_1 = 2$

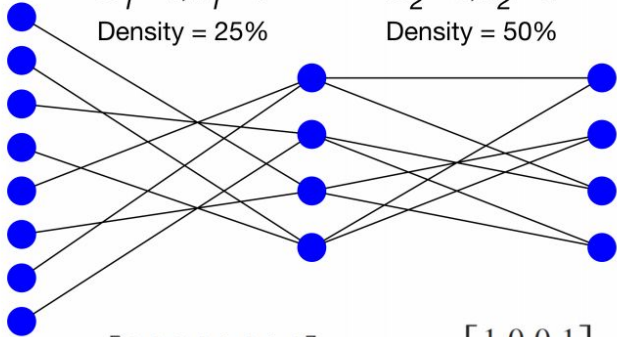
Density = 25%

Junction 2

$N_2 = 4, N_3 = 4$

$fo_2 = 2, fi_2 = 2$

Density = 50%



$$A_1 = \begin{bmatrix} 0 & 0 & 0 & 0 & 1 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$$A_2 = \begin{bmatrix} 1 & 0 & 0 & 1 \\ 0 & 0 & 1 & 1 \\ 1 & 1 & 0 & 0 \\ 0 & 1 & 1 & 0 \end{bmatrix}$$

$$A_1 = \begin{bmatrix} 0 & 0 & 0 & 0 & 1 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$\Sigma$  (red arrow pointing to the 5th column)  
 $\Sigma$  (blue arrow pointing to the 8th row)

$$A_1^{w_{ir}} = \begin{bmatrix} 0 & 2 \\ 1 & 1 \\ 1 & 1 \\ 2 & 0 \end{bmatrix}$$

$$A_1^{w_{ir}} = [1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1]$$

Figure 6: Construction of window adjacency matrices.



# Scatter (S)

= Proportion of values  $> 1$  in window adjacent matrix

We need to consider the number of S values in the network i.e. the number of junction which share the same level of scatter

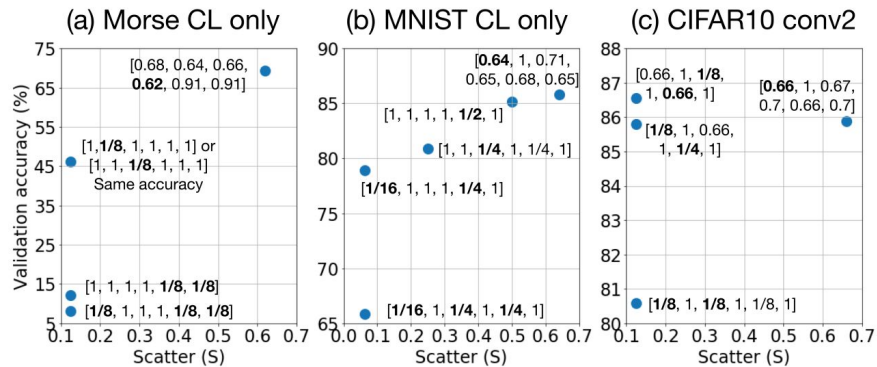


Figure 7: Network performance vs. scatter for CL only networks of (a) Morse (b) MNIST, and convolutional network with 2 CL junctions of (c) CIFAR10. All minimum values that need to be considered to differentiate between connection patterns are bolded.



# Future Work

Extending this concept to conv layers

In the meantime, small Keras script will be made available on [martellab website](http://martellab.com)