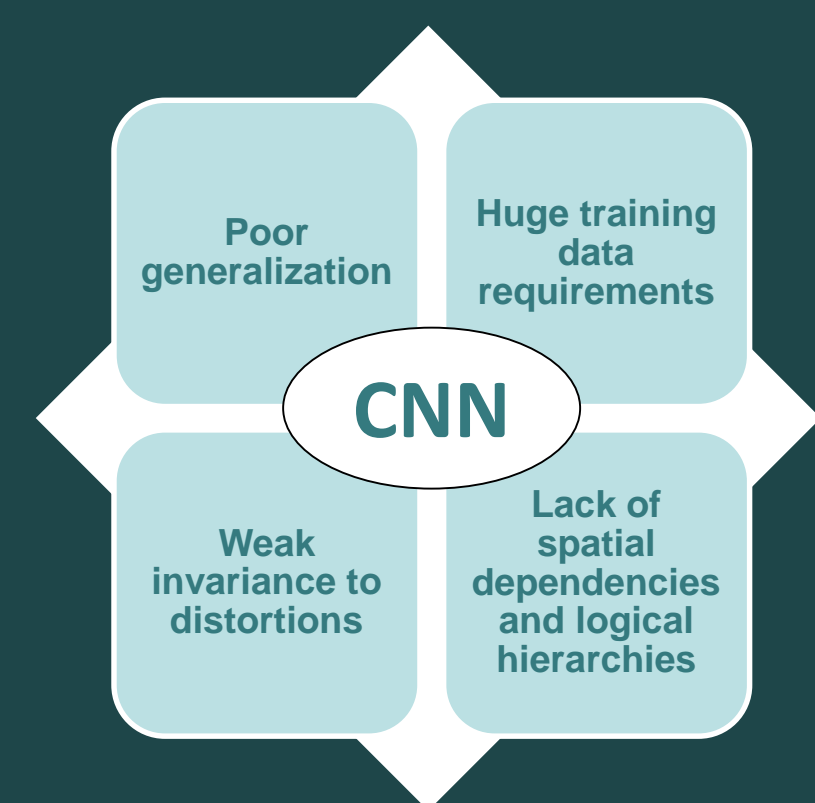


Kolezhuk L., Goatman K. A.

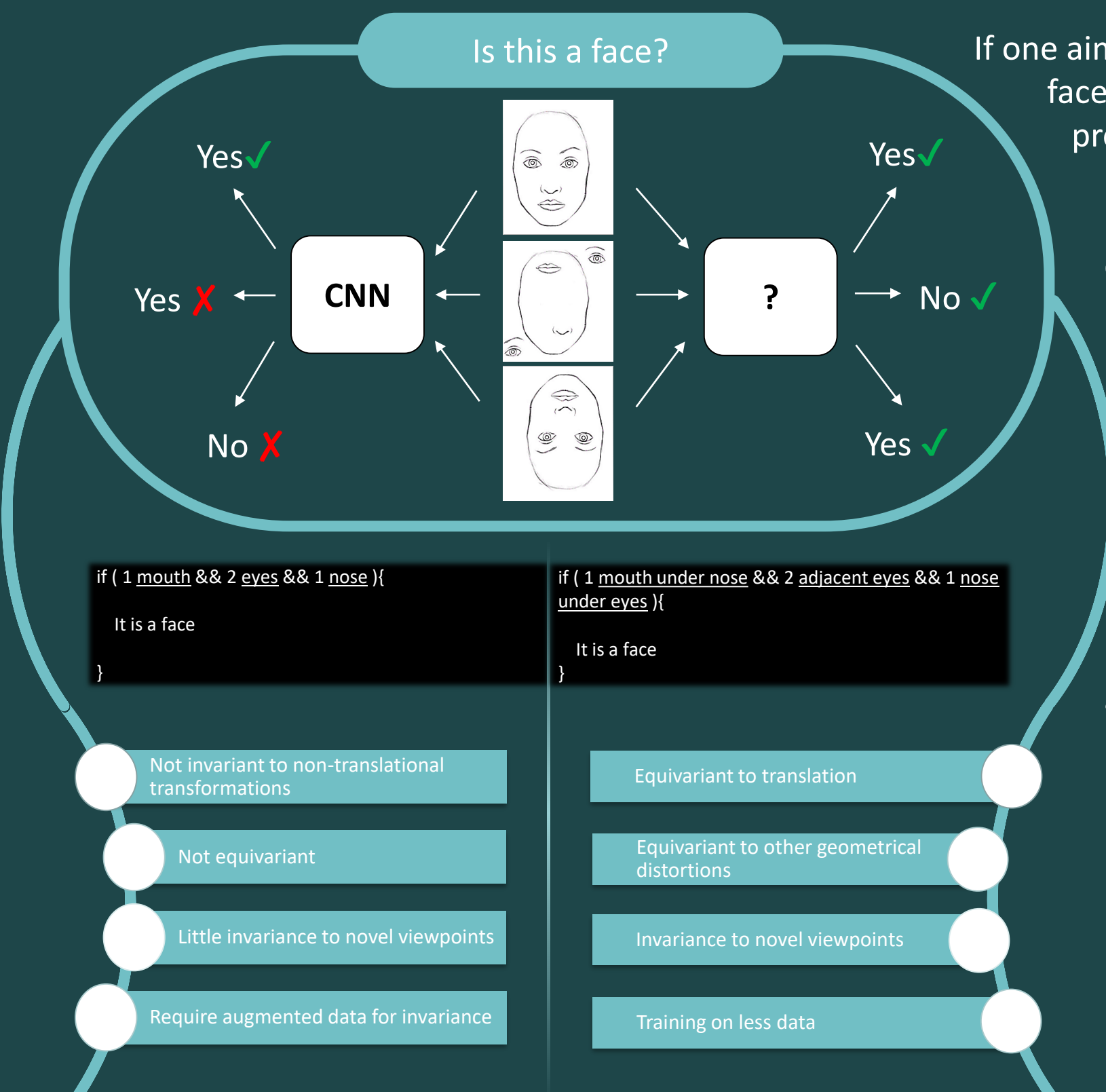
Background



CNNs are the key to obtaining solutions with remarkable results to countless modern computer vision tasks. However, there are some limitations introduced by their concept: absence of invariance to geometrical distortions other than translation; poor generalization; do not consider spatial relationships between features on the image but only their presence; require huge training datasets; low stability to change of viewpoints; no logical structure and neuron hierarchies; generically large networks easily fit random labels to data. Some of these issues are associated with the type of data routing that is used to pass signals to deeper layers of the network. The most common routing algorithm is **max-pooling**.

The pooling operation used in convolutional neural networks is a big mistake and the fact that it works so well is a disaster. If the pools do not overlap, pooling loses valuable information about where things are. We need this information to detect precise relationships between the parts of an object.

G. Hinton



Goal

In this work we investigate a recently proposed approach that claims to tackle these issues by introducing advanced grouping of neurons into capsules, that are designed to encode not only the presence of certain features, but also their instantiation parameters. The advantages and limitations of this architecture are studied. Our work expands the use of capsule networks to the task of false positive reduction for pulmonary nodule detection in lung CT scans while achieving improved classification accuracy when compared to a conventional CNN with max-pooling.

Experiments and results

Handwritten digit classification

In order to verify the correctness of the implementation we have tested CapsNet for the MNIST digit classification task.

- The following experiments were performed:
- MNIST classification performance
 - Robustness to random affine transformations of the input
 - Encoded representation quality investigation

Dataset	Test error rate %		
	CapsNet (Sabour et al. 2017)	CNN (state-of-the-art)	CapsNet (our implementation)
MNIST	0.25	0.39	0.29
affMNIST	21	34	25
multiMNIST	5.20(80% overlap)	5.2(4% overlap)	---

Performance sufficiently similar to the one claimed by Sabour et al. (2017) was achieved by our implementation. Significant robustness to affine transformations of input data has also been verified indicating the conceptual and functional correctness of the implementation.

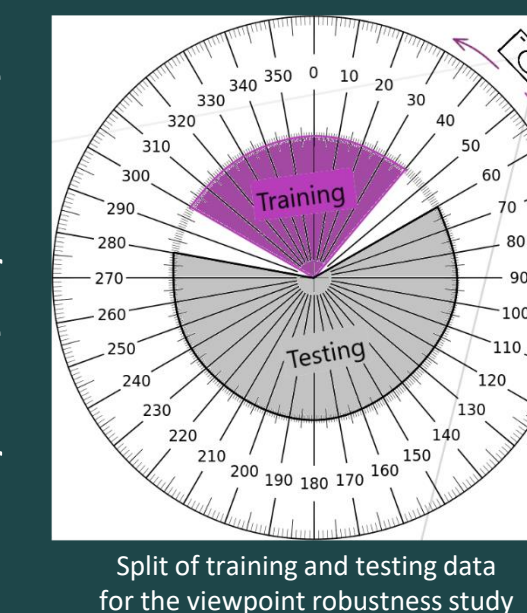
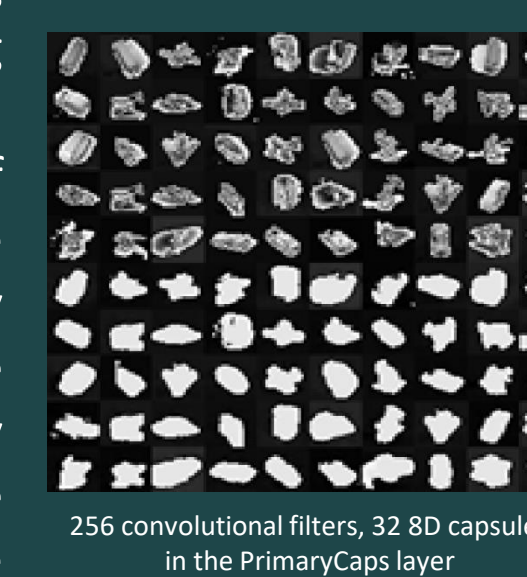
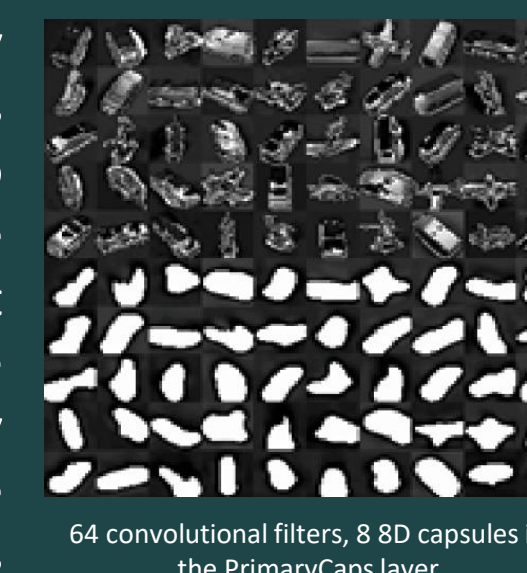
Taking into account that this data is rather simplistic, while having no background and small details in the object bodies, it is of great interest to study how the architecture is able to cope with more complex data.

Classification of images with 3D objects

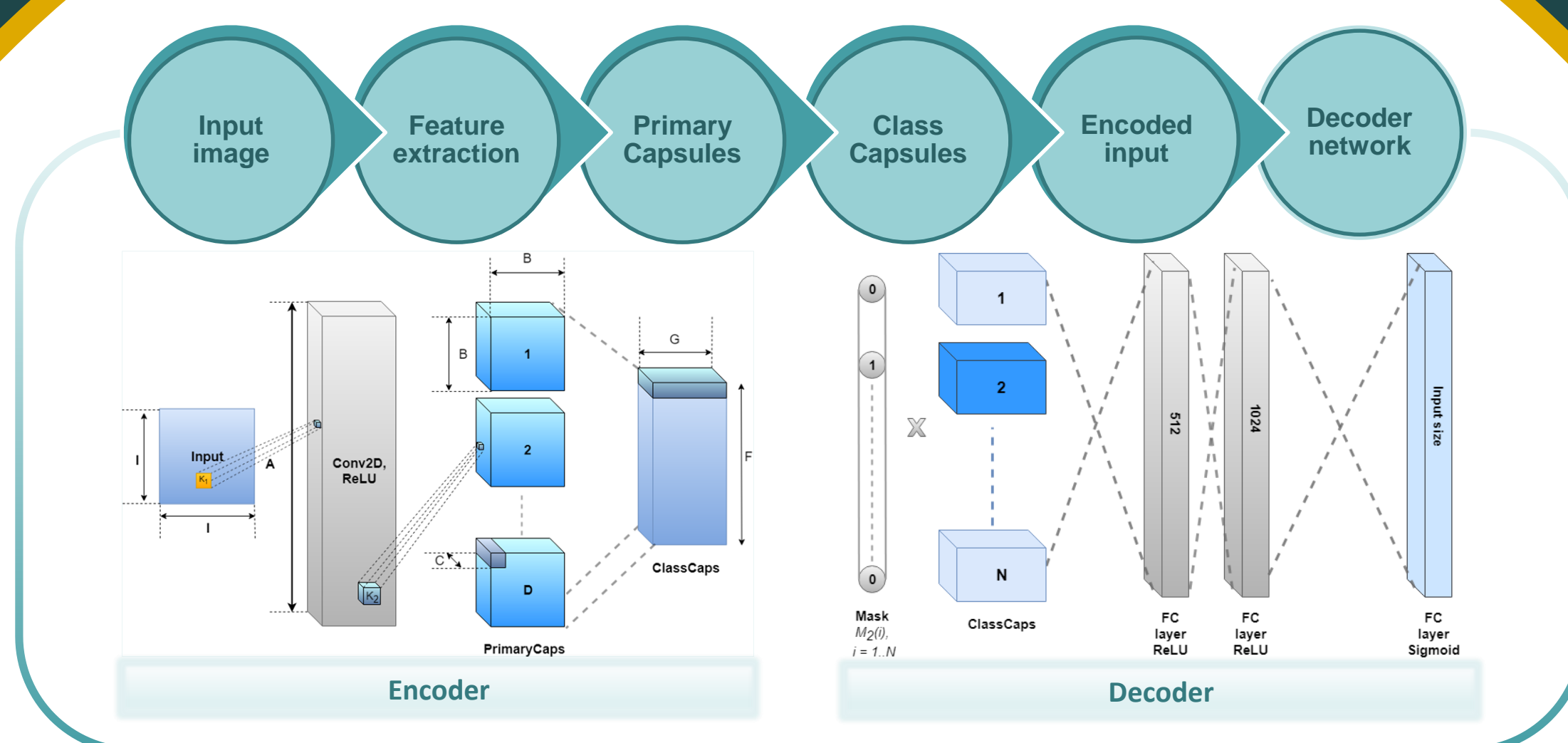
Dataset	Test error rate %	
	CapsNet	CNN
smallNORB	2.7 – Sabour et al.(2017)	2.0
	1.8 – EM Caps	
	10.9 – our approach	
CIFAR10	10.6 – Sabour et al.(2017)	4.5
	11.9 – EM Caps	
	28.5 – Xi E. (2017)	

Once data of greater complexity is processed it becomes troublesome for the network to properly extract and encode sufficient amount of important features from the inputs. The quality of the of this extraction may be judged by the detailization of the reconstructed images. Logically this issue may be tackled by increasing the size of the network by adding more layers, capsule types, extending the dimensions of capsule vectors etc. However, the experiments have proven that any further expansions after the presented ones do not cause any notable raises neither in the reconstruction quality nor the overall classification accuracy.

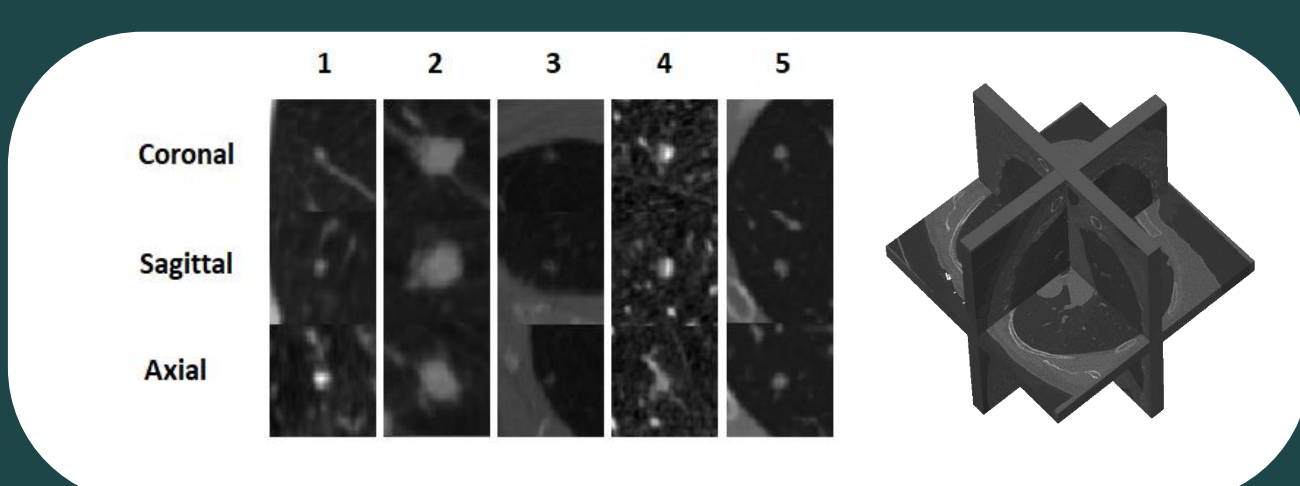
CapsNet is more robust to the change of viewpoints in the data, however it was proven that data augmentation is necessary in order to enforce correct estimation of the objects instantiation parameters. The last allows CapsNet to better deal with previously unseen data than a conventional CNN.



Capsule network

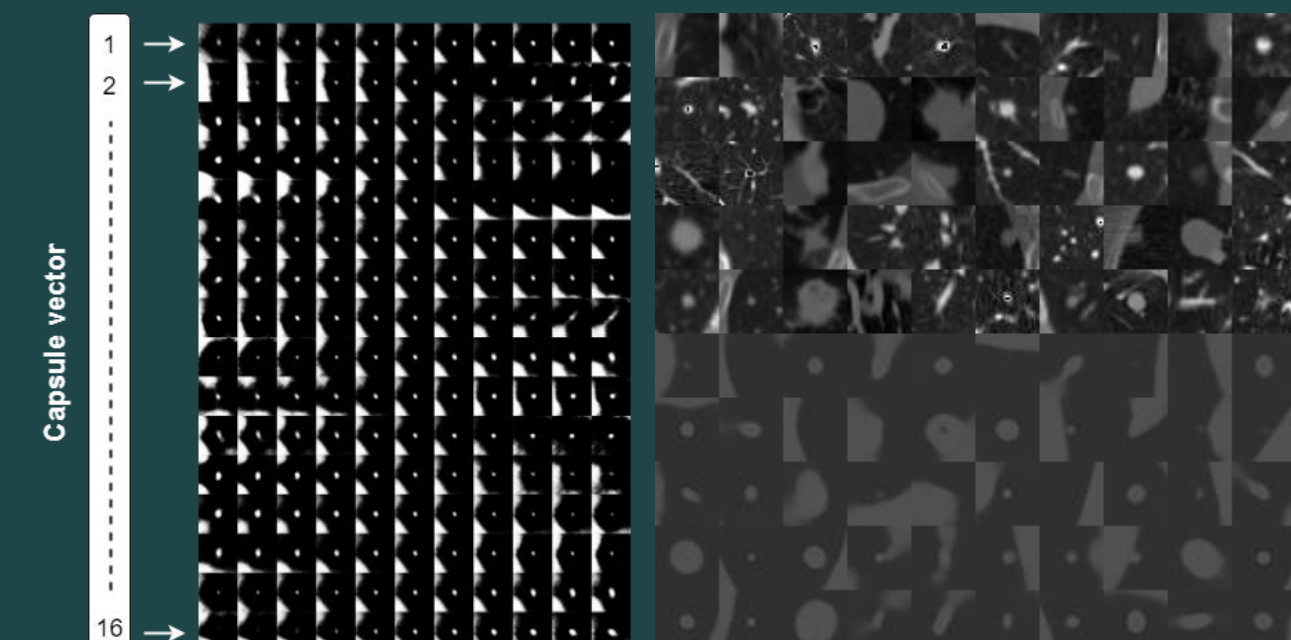


False positive reduction for pulmonary nodule detection



The LUNA16 dataset is used to target the goal. We extract orthogonal patches for each nodule candidate in CT volumes. The classes are balanced by means of generating 125 samples per view for the positive class, and selective reduction of the number of samples in the negative class.

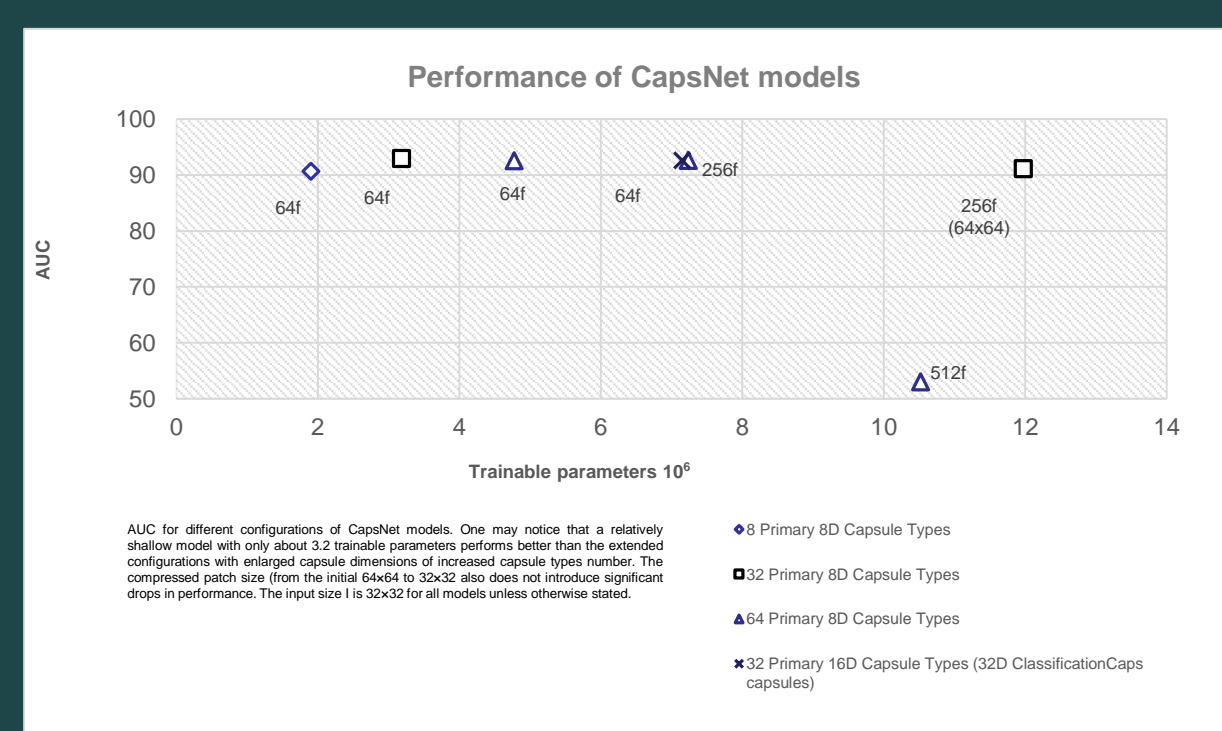
Encoding and reconstruction quality



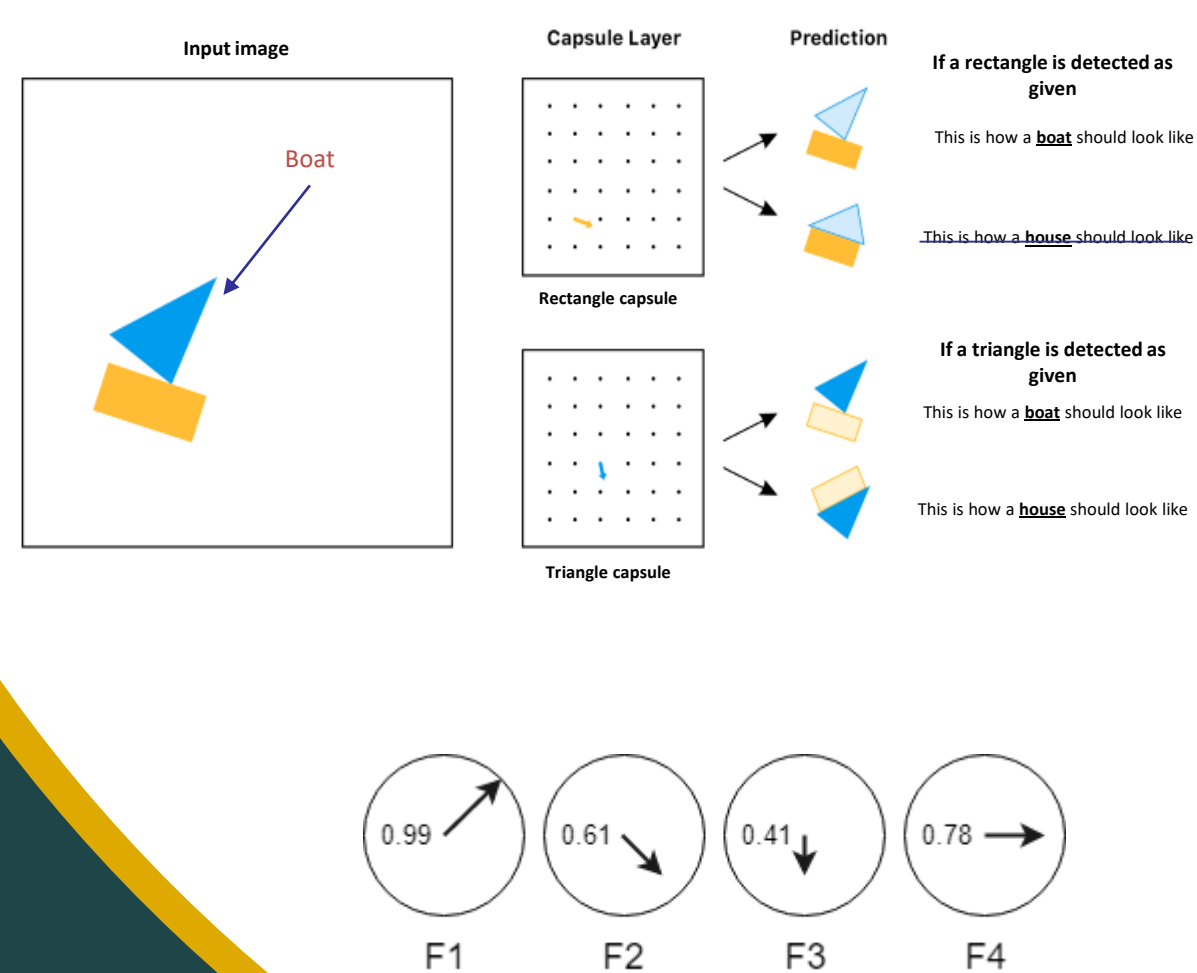
Performance of CapsNet models

Configuration	Test AUC	
	CNN baseline	CapsNet
Separate patches	91.4	93.2
Candidate (patch voting)	93.8	96.0

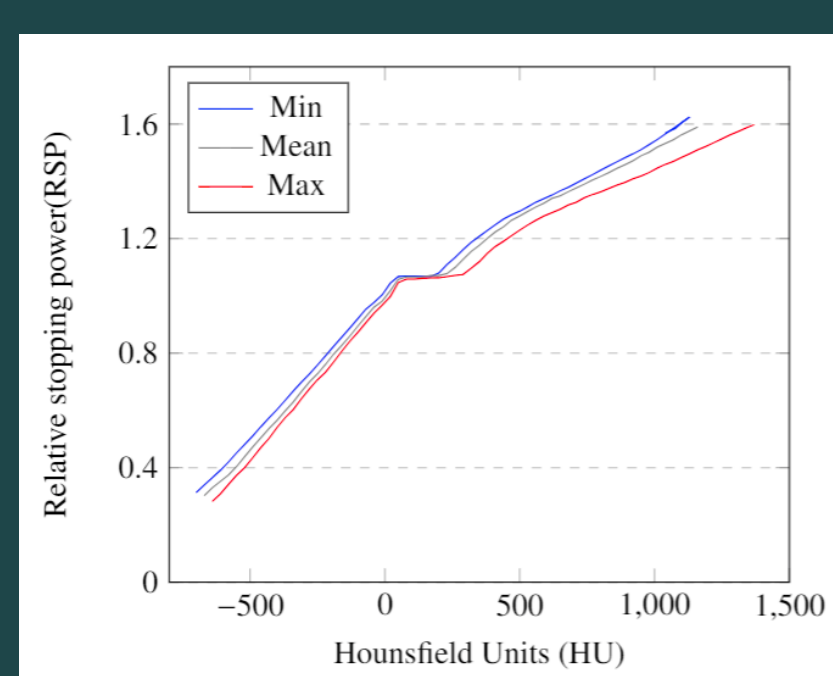
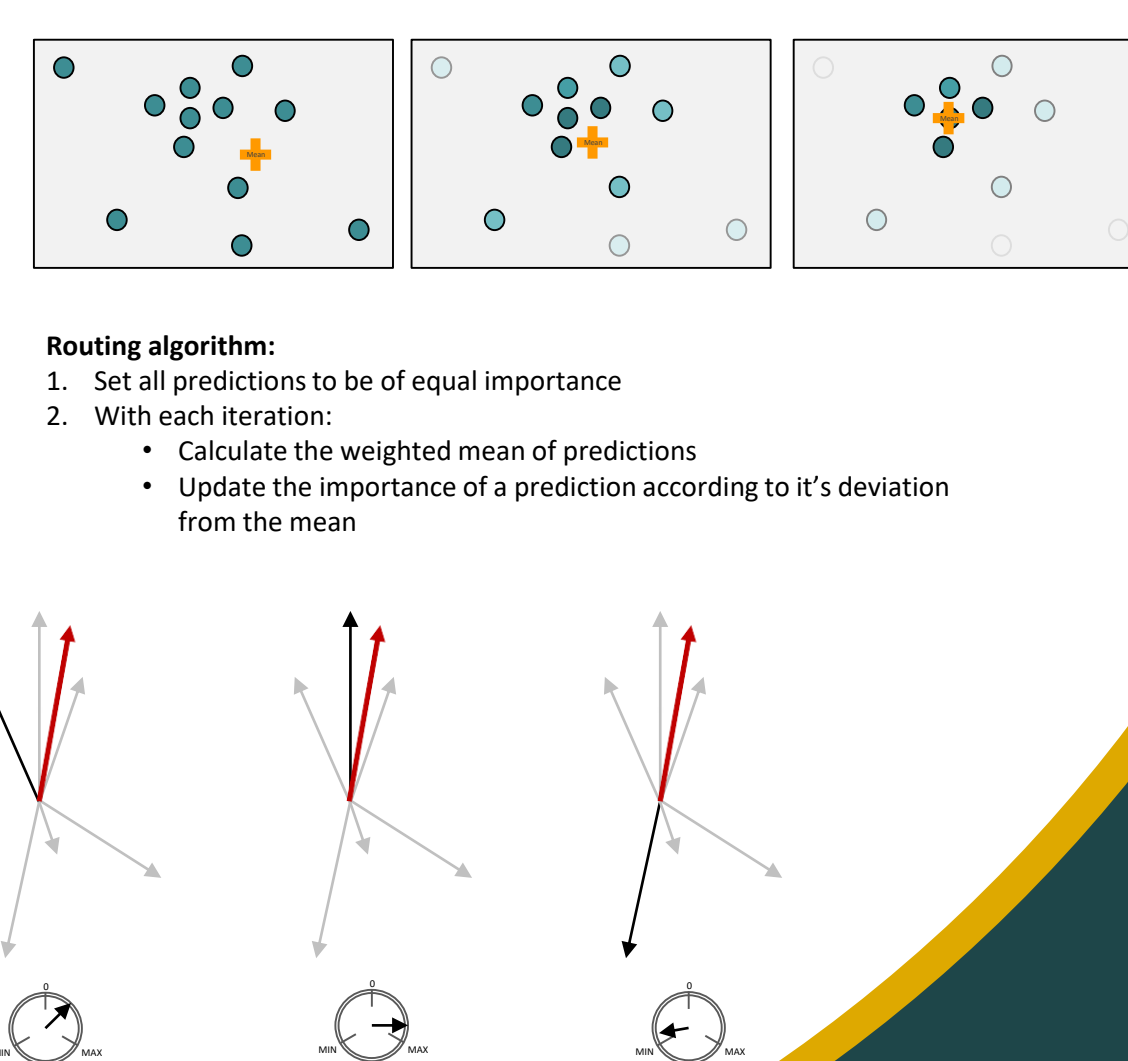
Performance of different CapsNet models was investigated. It was determined that the deeper models do not allow to obtain neither better classification accuracy, nor improved feature encoding quality.



Capsule predictions



Dynamic routing



Computation time per training step	
CapsNet	CNN
158ms	19ms

