

MIXCAPS: A Capsule Network-based Mixture of Experts for Lung Nodule Malignancy Prediction

Parnian Afshar^a, Farnoosh Naderkhani^a, Anastasia Oikonomou^b, Moezedin Javad Rafiee^c, Arash Mohammadi^{a,*}, Konstantinos N. Plataniotis^d

^a*Concordia Institute for Information Systems Engineering, Concordia University, Montreal, QC, Canada*

^b*Department of Medical Imaging, Sunnybrook Health Sciences Centre, University of Toronto, Canada*

^c*Affiliation 1: Department of Medicine and Diagnostic Radiology, McGill University Health Center- Research Institute. Affiliation 2: Babak Imaging Center, Tehran, Iran*

^d*Department of Electrical and Computer Engineering, University of Toronto, Toronto, ON, Canada*

Abstract

Lung diseases including infections such as Pneumonia, Tuberculosis, and novel Coronavirus (COVID-19), together with Lung Cancer are significantly widespread and are, typically, considered life threatening. In particular, lung cancer is among the most common and deadliest cancers with a low 5-year survival rate. Timely diagnosis of lung cancer is, therefore, of paramount importance as it can save countless lives. In this regard, Computed Tomography (CT) scan is widely used for early detection of lung cancer, where human judgment is currently considered as the gold standard approach. Recently, there has been a surge of interest on development of automatic solutions via radiomics, as human-centered diagnosis is subject to inter-observer variability and is highly burdensome. Hand-crafted radiomics, serving as a radiologist assistant, requires fine annotations and pre-defined features. Deep learning radiomics solutions, however, have the promise of extracting the most useful features on their own in an end-to-end fashion without having access to the annotated boundaries. Among different deep learning models, Capsule Networks are proposed to over-

*Corresponding author: Tel.: +1 (514) 848-2424 ext. 2712; This work was partially supported by the Natural Sciences and Engineering Research Council (NSERC) of Canada through the NSERC Discovery Grant RGPIN-2016-04988.

Email address: arash.mohammadi@concordia.ca (Konstantinos N. Plataniotis)

come shortcomings of the Convolutional Neural Networks (CNNs) such as their inability to recognize detailed spatial relations. Capsule networks have so far shown satisfying performance in medical imaging problems. Capitalizing on their success, in this study, we propose a novel capsule network-based mixture of experts, referred to as the MIXCAPS. The proposed MIXCAPS architecture takes advantage of not only the capsule network’s capabilities to handle small datasets, but also automatically splitting dataset through a convolutional gating network. MIXCAPS enables capsule network experts to specialize on different subsets of the data. Our results show that MIXCAPS outperforms a single capsule network and a mixture of CNNs, with an accuracy of 92.88%, sensitivity of 93.2%, specificity of 92.3% and area under the curve of 0.963. Our experiments also show that there is a relation between the gate outputs and a couple of hand-crafted features, illustrating explainable nature of the proposed MIXCAPS. To further evaluate generalization capabilities of the proposed MIXCAPS architecture, additional experiments on a brain tumor dataset are performed showing potentials of MIXCAPS for detection of tumors related to other organs.

Keywords: Tumor type classification, Capsule network, Mixture of experts

1. Introduction

Lung cancer, according to recent statistics [1], is associated with the highest mortality rate, among all different cancer types, and is considered as one of the top three cancers, in terms of incidence. The combined 5-year survival for lung cancer is still low [2], at 18%, because the majority of patients are diagnosed at advanced stages [3]. What makes the early diagnosis of lung cancer significantly challenging is the lack of sufficient visible warning symptoms and signs in early stages of the disease. Computed Tomography (CT) scan [4] is by far one of the most advanced and effective techniques used for lung cancer diagnosis. However, even the CT scans may not reveal convincing signs that can contribute to early diagnosis of lung cancer. In other words, Imaging features of nodule such as size, shape, and attenuation that play an important role

in identifying the cancer may not be immediately accessible to the unaided eye [5]. More importantly, human-centered diagnosis is subject to inter-observer variability, meaning that radiologists can have different judgments, depending on their previous experience. Finally, investigating the test results and coming into an inclusive decision can be extremely time-consuming and burdensome [6].

Radiomics analysis [7, 8, 9], referring to the extraction of several quantitative and semi-quantitative features from the medical images, is one of the most successful approaches towards automatizing the cancer diagnosis/prediction process [10]. Features extracted in the radiomics analysis are aimed at capturing different properties of the nodules, such as their shape and texture. Such features have shown association with the nodule malignancy, its stage, and even the patient’s survival time. Radiomics is often categorized in two groups of hand-crafted [11, 12, 13, 14] and deep learning-based. The former category involves extraction of a set of pre-defined features that are further processed and analyzed by a statistical or Machine Learning (ML) model. Despite showing satisfactory results in different tasks [15, 16], hand-crafted radiomics is limited to the features defined by the radiologists and as such there is no guarantee that the features contribute to the problem at hand. Furthermore, since hand-crafted Radiomics features are extracted from the annotated Region of Interest (ROI), they are still subject to inter-observer variability, and besides being time-consuming, their performance highly depends on the accuracy of the provided annotations [17]. In other words, extra effort is required to enhance the annotations and select features that are more descriptive and robust [18].

Deep learning-based radiomics [19, 20, 21], proposed to overcome the shortcomings of its hand-crafted counterparts, does not require a pre-knowledge about the types of features to be utilized. In other words, deep learning-based techniques are capable of extracting features that can best contribute to the problem at hand in an end-to-end fashion. Furthermore, deep learning-based radiomics does not need to be fed with the annotated ROI, which has the promise of reducing the effect of inter-observer variability as well as the burden of segmenting the images. Among different deep learning techniques, Convolutional

Neural Networks (CNNs) are more popular within the field of radiomics [22], due to their ability to efficiently process and learn meaningful features from medical images [23]. Performance of the CNNs, however, partly depends on the size of the available dataset [24]. More specifically, CNNs, typically, fail to determine the spatial relations between the image instances and identify rotation or transformation of an object. As such, CNNs need to be fed with a large dataset containing all the possible transformations of the objects. Large datasets are, however, not typically available in medical imaging in particular for lung cancer malignancy prediction.

Capsule networks [25], also referred to as the CapsNets, are developed aiming at overcoming the aforementioned drawbacks of the CNNs. CapsNets use capsules, instead of using individual neurons, to represent imaging instances. CapsNets, therefore, can identify the spatial relations via their “Routing by Agreement” process, through which capsules try to come to a mutual agreement about the existence of the objects. In particular, CapsNet’s ability to handle transformations is further investigated in Reference [26] for medical image segmentation. In our recent studies [27, 28, 29], we showed superior performance and capabilities of CapsNets for tumor type classification.

Capitalizing on the success of the CapsNets, in this study we propose a new framework, referred to as the Mixture of Capsule networks (MIXCAPS), for the task of lung nodule malignancy prediction. The proposed MIXCAPS framework is a “Mixture of Experts” type model [30, 31, 32, 33], which has the potential to noticeably improve the classification accuracy by integrating/coupling several experts (individual CapsNets in the context of the proposed MIXCAPS). To be more precise, mixture of experts solves the classification problems by splitting the dataset into similar samples, and each expert specializes in classifying similar instances. To the best of our knowledge, the proposed MIXCAPS is the first CapsNet-based mixture of experts framework. The MIXCAPS model benefits from the following three important properties: (i) The embedded capsule network is capable of classifying the lung nodules without requiring availability of a large dataset; (ii) The mixture of experts approach enables each CapsNet

within the MIXCAPS architecture to focus on a specific subset of the nodules, therefore, improving the overall classification performance of the model, and;

(iii) As shown in our experiments, MIXCAPS is not restricted to the task of lung nodule malignancy prediction. In fact, it can be easily generalized to the prediction of other tumor types such as brain cancer. The following summarizes our contributions:

- CapsNets are utilized, for the first time, as individual experts within a mixture of experts framework.
- A new and modified CapsNet loss function (margin loss) is developed to reflect the loss associated with the experts and gating models.
- Output of the gating model is investigated for potential correlations with nodule hand-crafted features to improve the potential interpretability of the proposed MIXCAPS.
- Generalizability of the proposed MIXCAPS is illustrated via extension and evaluation based on a separate dataset associated with a different prediction task other than the one initially used to design the framework.

The rest of this paper is organized as follows: First, in Section 2 the previous studies on lung nodule malignancy prediction is briefly investigated. In Section 3, the dataset and the pre-processing steps are described, along with the proposed MIXCAPS. Results and discussions are presented in Section 4. Finally, Section 5 concludes the paper.

2. Related Works

Generally speaking, most of the studies based on hand-crafted radiomics follow a pre-defined set of steps [7, 8, 9]:

- (i) The first step is to pre-process the images and segment the nodule;
- (ii) The second and the main step is to extract hundreds of features from the segmented nodule. These features mostly fall into three categories

of intensity-based, shape-based, and texture-based features. The former category captures basic properties of the nodule related to its histogram. While shape-based features quantify shape-related properties such as area, diameter, and volume, texture-based ones capture the heterogeneity of the nodule texture;

- (iii) In the third step of the hand-crafted radiomics analysis, feature reduction techniques are utilized to select the most relevant and robust features;
- (iv) In the final step, extracted features are fed to a statistical or machine learning tool to calculate the desired outcome.

For example, the study performed by authors in Reference [34] is a recent implementation of the above mentioned steps for extracting hand-crafted radiomics for lung nodule malignancy prediction. In this study, a total of 385 features is extracted from the annotated nodules. Consequently, based on a correlation analysis, the non-redundant features are selected and fed to a regression model to output the malignancy probability.

The limitations of the hand-crafted radiomics, including its dependence on the annotated region, have caused a surge of interest in deep learning-based radiomics, especially using CNNs [35, 36]. CNNs are powerful models for analyzing images and extracting features that best contribute to the problem at hand, through trainable filters. Furthermore, filters share weights across the input, which significantly reduces the computational cost, compared to a fully-connected network. CNNs have been recently used for the problem of lung nodule malignancy prediction. While some studies [37, 38] have proposed to adopt previously developed CNN architectures for the radiomics analysis, others [22, 39] have designed and optimized their own specific CNN-based models. Although showing satisfying results, most of these studies had to use a data augmentation or transfer learning strategy to compensate for the lack of large datasets specifically for the problem of lung nodule malignancy prediction. These strategies, however, are associated with more computational costs. Furthermore, there is still no comprehensive study on the effectiveness of these

strategies for the nodule malignancy prediction. Capsule network (CapsNet), briefly described in the following section, is an alternative and attractive modeling paradigm to address the aforementioned issues, i.e., accounting for more variations in the input, without resorting to heavy data augmentation.

2.1. Capsule Networks

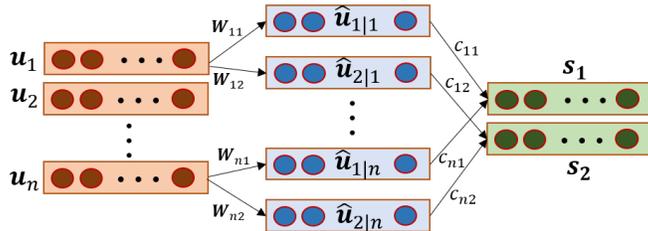


Figure 1: Routing by agreement. For the sake of simplicity number of output capsules is set to two.

Capsule networks are constructed based on capsules, as their main building blocks. A capsule being represented by a vector consists of several neurons representing, collectively, a specific object at a specific location. While neurons capture the instantiation parameters of the object, the length of a capsule determines the existence probability of that object. The most important property of a capsule network, distinguishing it from CNNs, is its routing by agreement process. Generally speaking, each Capsule i , having the instantiation parameter vector \mathbf{u}_i , in a lower layer tries to predict the output of the capsules in the next layer, through a trainable weight matrix \mathbf{W}_{ij} given by

$$\hat{\mathbf{u}}_{j|i} = \mathbf{W}_{ij}\mathbf{u}_i, \quad (1)$$

where $\hat{\mathbf{u}}_{j|i}$ denotes the prediction for parent Capsule j . Through the routing by agreement process, the predictions are evaluated in terms of their similarity to the actual outputs. More weight is then given to the successful predictions,

before calculating the final output \mathbf{s}_j for the capsule j , as follows

$$a_{ij} = \mathbf{s}_j \cdot \hat{\mathbf{u}}_{j|i}, \quad (2)$$

$$b_{ij} = b_{ij} + a_{ij}, \quad (3)$$

$$c_{ij} = \frac{\exp(b_{ij})}{\sum_k \exp(b_{ik})}, \quad (4)$$

$$\text{and } \mathbf{s}_j = \sum_i c_{ij} \hat{\mathbf{u}}_{j|i}, \quad (5)$$

where a_{ij} shows the agreement between actual output \mathbf{s}_j and its prediction $\hat{\mathbf{u}}_{j|i}$, and c_{ij} denotes the score assigned to the prediction based on the obtained agreement. The routing by agreement process, summarized in Fig. 1, enables capsule the network to recognize spatial information between image instances.

Tumor classification based on capsule networks has been investigated in several recent studies, leading to increased performance when compared to CNNs. Lung tumor malignancy prediction is considered in Reference [27], where a multi-scale framework is proposed, outperforming single-scale and multi-scale CNNs. Classifying tumors related to other organs, such as brain, using capsule networks, has also been investigated in several studies [28, 29, 40, 41], leading to satisfying performance. The paper makes a unique contribution in this field by introducing a novel CapsNet architecture based on ‘‘Mixture of Experts’’, which is briefly described below.

2.2. Mixture of Experts

Mixture of experts (MoE) [31] refers to adopting several experts, each of which is specialized on a subset of the data, to collectively perform the final prediction task. As shown in Fig. 2, experts are separately fed with the input data and the final output is a weighted average of all the predictions coming from all the N active experts. The weight g_i assigned to Expert i can be either a pre-determined value, or a trainable one. One simple example of the former case is averaging over all the experts’ predictions [33]. However, more sophisticated approaches such as soft clustering of the input may also be adopted. In the latter case, weights may be trained at the same time with the experts. One

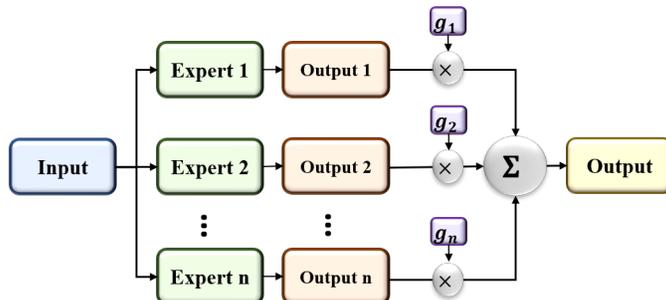


Figure 2: General framework of a mixture of experts approach.

other approach to use trainable gating weights is to concatenate the feature vectors obtained from the individual experts and feed the resulting vector to an external gating model to make the final decision.

The MoE concept has been widely used in medical imaging. The simple averaging scenario is investigated in References [42] and [43] for retinal vessel detection from fundus images and breast cancer detection from histology images, respectively. Trainable gating weights are studied in Reference [44], where hand-crafted and CNN-based features are combined to detect breast cancer from pathology images. The scenario where gating weights are trained at the same time with the experts is investigated in Reference [32] for breast cancer diagnosis. In particular, CNN experts are combined using weights coming from an external gating network. The gating network itself is a CNN, taking the same inputs as the experts, and outputting the probability of each expert being responsible for each particular input. Our proposed MIXCAPS, which is based on the same gating scenario as Reference [32], is explained in the next section, along with its incorporated data pre-processing approach.

3. The Proposed MIXCAPS Framework

In this section, first we present the dataset used to design and develop the proposed MIXCAPS. Afterwards, the pre-processing approach, and the proposed MIXCAPS framework are described.

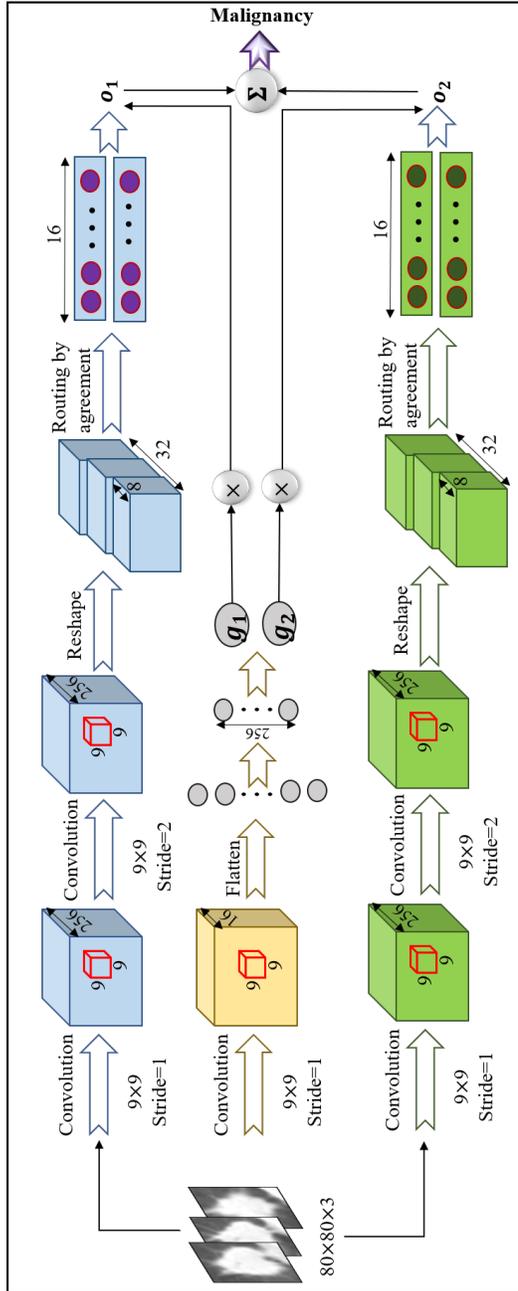


Figure 3: Proposed MIXCAPS.

3.1. Data and Pre-processing Approach

The lung nodule malignancy dataset is adopted from the Lung Image Database Consortium (LIDC) and Image Database Resource Initiative (IDRI) dataset [45, 46, 47]. This dataset consists of CT scans from 1,018 subjects. All the images are labeled and annotated by one to four radiologists. Labels include non-nodule, nodule less than 3 *mm* in size, and nodules with malignancy scores of 1 to 5, where larger numbers denote higher possibility of malignancy. In this study, we discarded all the cases with average malignancy score of 3 which dictates an indeterminate malignancy. Consequently, we regrouped labels 1 and 2 as benign nodules, and labels 4 and 5 as malignant nodules. Therefore, we ended up having a binary classification problem with a total of 2,283 nodules. It is worth mentioning that we included all the annotations provided by all the radiologists as separate nodules. However, the malignancy scores are the average over all the provided scores. For each nodule, we extracted a 3D patch from the center of the nodule (center slice and the two immediate neighbors). Patches are extracted to fit the nodule boundary provided by the radiologists. However, to have fixed size inputs, all patches were zero-padded to 80×80 (the largest possible width and height based on the training data).

3.2. The MIXCAPS Architecture

The proposed capsule network-based mixture of experts for lung nodule malignancy prediction, referred to as the MIXCAPS, is shown in Fig 3. The 3D nodule patches are the inputs to two capsule network experts, as well as the convolutional gating network. The two experts, as shown in Fig 3, consist of two convolutional layers, the last of which is reshaped to form a capsule layer. This capsule layer is followed by a routing by agreement and the final capsule layer. The outputs of the two experts, denoted by \mathbf{o}_1 and \mathbf{o}_2 , represent the class (benign and malignant) probabilities. The gating network, consisting of a convolutional and two fully connected layers, determines the contribution of each expert, denoted by g_1 and g_2 , for a specific input through a Softmax layer,

as follows

$$g_1 = \frac{\exp(G_1)}{\exp(G_1) + \exp(G_2)}, \quad g_2 = \frac{\exp(G_2)}{\exp(G_1) + \exp(G_2)}, \quad (6)$$

where G_1 and G_2 are pre-activation outputs. The Softmax layer ensures that g_1 and g_2 sum to one. These contributions are multiplied by \mathbf{o}_1 and \mathbf{o}_2 to calculate the final prediction \mathbf{o} as follows

$$\mathbf{o} = g_1 \mathbf{o}_1 + g_2 \mathbf{o}_2. \quad (7)$$

Output vector \mathbf{o} encompasses the probability of benign and malignant classes, denoted by $o^{(0)}$ and $o^{(1)}$, respectively. In other words

$$\mathbf{o} = [o^{(0)}, o^{(1)}]^T. \quad (8)$$

where superscript T denotes transpose operator. Originally, margin loss is proposed for the training of the capsule networks. In this study, we adopt the same loss function with the difference that the loss l is calculated over the final output of the MIXCAPS instead of the individual capsule networks, as follows

$$l^{(0)} = T^{(0)} \max(0, m^+ - o^{(0)})^2 + \lambda(1 - T^{(0)}) \max(0, o^{(0)} - m^-)^2, \quad (9)$$

$$l^{(1)} = T^{(1)} \max(0, m^+ - o^{(1)})^2 + \lambda(1 - T^{(1)}) \max(0, o^{(1)} - m^-)^2, \quad (10)$$

$$l = l^{(0)} + l^{(1)}, \quad (11)$$

where $l^{(0)}$ and $l^{(1)}$ denote the losses associated with the benign and malignant classes, respectively. m^+ , λ , and m^- are hyper-parameters. Terms $T^{(0)}$ and $T^{(1)}$ are the ground-truth labels for benign and malignant classes, respectively. According to Reference [31] comparing the desired output with the blend of outputs from the experts, leads to a strong coupling between experts and solutions in which many experts are used for one case. However, in this study, we did not encounter such a problem, and therefore did not adopt non-linear combinations of the outputs.

3.3. CapsNet as a Mixture of Experts

In this subsection, we revisit the idea of the capsule networks and show how they can be viewed within the mixture of experts framework. In other words,

we show that a CapsNet is a series of consecutive MoE layers such that each lower level capsule with instantiation vector \mathbf{u}_i serves as an expert to predict the output of the capsule in the next layer with instantiation vector \mathbf{s}_j .

Recall from Section 2.1 that each capsule (among N_{PrC} number of primary capsules) with instantiation vector \mathbf{u}_i , for $(1 \leq i \leq N_{PrC})$, makes predictions $\hat{\mathbf{u}}_{j|i}$, through Eq. (1). Consequently, each capsule (among N_{PaC} number of parent capsules) with instantiation vector \mathbf{s}_j , for $(1 \leq j \leq N_{PaC})$, receives predictions from all the lower level primary capsules. Each primary Capsule i , therefore, can be considered as an expert making predictions for all the parent (final) capsules. Contribution of each capsule expert i to each final capsule j is represented by c_{ij} , which is basically similar to g_i in an MoE framework, with the difference that in the conventional MoE formulation, each expert contributes equally to all the outputs, whereas capsule experts have different contributions to different final capsules. This is the reason why the notation of c_{ij} is used instead of c_i . The instantiation parameter of each final Capsule j is calculated according to Eq. (5) incorporating predictions from all the experts. Another difference between capsule experts and conventional MoE ones is that the gating model in the latter case is typically a simple or advanced machine learning model, whereas in the former case, routing by agreement serves as the gate to determine contribution through Eq. (2) to (5). It is also worth noting that Eq. (4) ensures that contributions to each final capsule j sum to one, satisfying the requirement of an MoE approach as in Eq. (6).

Having the aforementioned discussion in mind, each CapsNet itself is a series of mixtures of experts. In the proposed MIXCAPS, the CapsNets themselves are utilized as single experts. Therefore, MIXCAPS can be considered as a hierarchical MoE technique. It is also interesting to study how the calculation of c_{ij} s resembles the calculation of experts' weights in an MoE approach. Generally speaking, there are several solutions to an MoE problem [48]. An Expectation Maximization (EM) algorithm is one applicable solution, through which the experts' weights are considered as hidden variables, whose posteriors

are estimated in the E-step, as follows

$$p(z_i^n | \mathbf{t}^n, \mathbf{x}^n) = \frac{p(\mathbf{t}^n | z_i^n = 1, \mathbf{x}^n) p(z_i^n = 1 | \mathbf{x}^n)}{p(\mathbf{t}^n | \mathbf{x}^n)}, \quad (12)$$

where binary variable z_i^n is one when instance n is assigned to expert i , and zero otherwise. Term $p(z_i^n | \mathbf{t}^n, \mathbf{x}^n)$ represents the posterior probability of z_i^n given input vector \mathbf{x}^n and target vector \mathbf{t}^n . Following the Bayes' rule, this posterior is calculated using the likelihood term $p(\mathbf{t}^n | z_i^n = 1, \mathbf{x}^n)$ and the prior over z_i^n , denoted by $p(z_i^n = 1 | \mathbf{x}^n)$. All the terms appearing in Eq. (12) can be calculated through the MIXCAPS framework. The likelihood term can be replaced by the output of the expert capsule networks $o_i^{n(1)}$, which denotes the probability of malignancy for Instance n , based on the i^{th} expert. The prior probability can also be estimated using the output of the gating model g_i^n denoting the probability of assigning Instance n to Expert i . The posterior, therefore, can be defined as

$$p(z_i^n | \mathbf{t}^n, \mathbf{x}^n) = \frac{g_i^n o_i^{n(1)}}{\sum_j^M g_j^n o_j^{n(1)}}, \quad (13)$$

where M is the number of experts.

To further shed light on the MoE view of CapsNets, it would be interesting to note that the EM formulation of the MoE closely resembles the weight update process of a multiple model (MM) [49] approach. In MM formulation, observations are sequentially generated from different models and the goal is to identify the contribution of each single model i given all the observations up to the current time (\mathbf{Y}^k), as follows

$$p(z_i^k | \mathbf{Y}^k) = \frac{p(\mathbf{y}^k | z_i^k = 1, \mathbf{Y}^{k-1}) p(z_i^k = 1 | \mathbf{Y}^{k-1})}{\sum_{j=1}^M p(\mathbf{y}^k | z_j^k = 1, \mathbf{Y}^{k-1}) p(z_j^k = 1 | \mathbf{Y}^{k-1})}, \quad (14)$$

where \mathbf{y}^k is the most recent observation. Comparing Eq. (14) with Eq. (13), it can be seen that while the prior in an MoE approach is determined based on the current input vector, it is calculated based on the previous observations, in the MM case. In other words, in MM, the prior is iteratively replaced with the posterior. The updates of coefficients in the routing by agreement process of the CapsNet is similar to the weight updates in MM. In particular, in each round

of the routing by agreement, the previously calculated c_{ij} serves as the prior to compute the coefficient in the next round.

4. Results and Discussion

In this section, three different experiments on lung cancer malignancy prediction are presented. The main objective is to evaluate performance of the proposed MIXCAPS framework and compare its capabilities with those of its state-of-the-art counterparts. Results are obtained with 200 iterations of bootstrapping, where in each iteration, 80% of the data is sampled (with replacement) from the whole dataset. 20 % of the training dataset is then randomly extracted for validation. A 95% confidence interval (CI) is calculated for all the performance metrics. Adam optimizer with 10 epochs and batch size of 16 is used for training.

Experiment 1: Our first experiment is to compare the performance of the proposed MIXCAPS with a single capsule network and a mixture of CNNs, as shown in Table 1, where performance is measured in terms of sensitivity, specificity, accuracy, and area under the curve (AUC). The architecture of the single capsule network is exactly the same as the CapsNet experts. We tried to keep the complexity as similar as possible to the MIXCAPS, when designing the mixture of CNNs. In particular, the gating network exactly resembles that of the MIXCAPS. The CNN experts consist of two convolutional layers with 256 filters, similar to the experts in the MIXCAPS. The convolutional layers are followed by a dense layer with 32 neurons (the same as the dimension of the last capsule layers), and the final softmax layer for nodule malignancy prediction. As shown in Table 1, MIXCAPS outperforms its two aforementioned counterparts, in terms of sensitivity, specificity, accuracy, and AUC.

Experiment 2: In the second experiment, we compare the proposed MIXCAPS with several well-known studies on the same dataset. Table 2 shows these studies, their methods, and the obtained results. As it can be inferred from Table 2,

Table 1: Performance of the proposed MIXCAPS compared to that of a single capsule network and a mixture of CNNs. Numbers in parenthesis show the 95% confidence intervals.

Model	Sensitivity	Specificity	Accuracy	AUC
Proposed MIXCAPS	89.5 (89.3, 89.7)%	93.4 (93.2, 93.6)%	90.7 (90.6, 90.8)%	0.956 (0.955, 0.956)
Single capsule network	86.1(85.7, 86.4)%	90.8(90.5, 91.1)%	88.6(88.5, 88.7)%	0.938(0.937, 0.939)
Mixture of CNNs	87.5(87.1, 87.8)%	91.3(91.1, 91.6)%	89.5(89.4, 89.7)%	0.948(0.946, 0.948)

Table 2: List of studies that have used LIDC-IDRI to predict lung nodule malignancy based on the ratings provided by radiologists. Note that some of the studies such as References [5] and [50] included hand-crafted features, requiring expert annotations. Numbers in parenthesis show the 95% confidence intervals obtained from 200 iterations of bootstrapping.

Method	Area Under the Curve (AUC)	Accuracy	Specificity	Sensitivity
Proposed MIXCAPS	0.956(0.955, 0.956)	90.7(90.6, 90.8)%	89.5(89.3, 89.7)%	89.5(89.3, 89.7)
CNN [5]	0.938	87.9%	87.9%	87.9%
CNN in combination with hand-crafted features [5]	0.971	93.2%	98.5%	87.9%
Deep residual network [51]	0.9459	89.90%	88.64%	91.07%
Deep belief network [52]	-	81.19%	-	-
CNN in combination with hand-crafted features [50]	-	86.79%	95.42%	60.26%
Multi-crop CNN [53]	0.93	87.14%	93%	77%

the proposed MIXCAPS outperforms all the studies in terms of accuracy and AUC, except Reference [5]. However, it is worth mentioning that the aforementioned study utilizes hand-crafted radiomics, requiring fine annotation of the nodules, from which our proposed approach is independent. Reference [50] has obtained a higher specificity compared to the proposed MIXCAPS. Its low sensitivity, however, is a sign of an unbalanced classification and/or over-fitting. Reference [51] has achieved the highest sensitivity among all the other references. Nevertheless, no confidence interval is provided to ensure the robustness of the result.

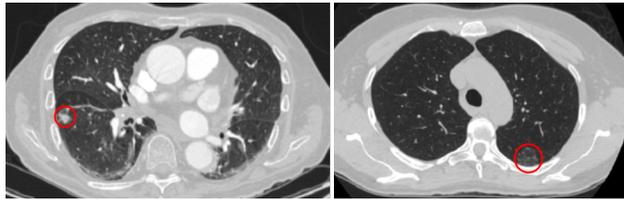


Figure 4: Example of nodules assigned to experts based on their volume and diameter. The nodule on the left, which has a lower probability of belonging to the first expert, is smaller in terms of volume and diameter compared to the nodule on the right.

Experiment 3: Finally, we conduct an experiment to gain an insight on how the data instances are split between the two experts. The LIDC-IDRI dataset is accompanied by a few nodule-related properties, determined by the radiologists. These features include volume, diameter, x center of mass and y center of mass. We calculated the correlation between the output of the gating network and these features. While the correlations with volume and diameter are 0.58 and 0.77, respectively, we observed no correlation with the centers of mass. It should be noted that the inputs to the proposed MIXCAPS are cropped nodule regions. In other words, the model has no access to the location of the nodule. Therefore, the almost zero correlations with the centers of the mass is completely expected. The observed correlations between the gate outputs and the volume and diameter imply that larger nodules have higher probabilities of being assigned to the first expert. Fig. 4 shows two nodules in the test set.

The left nodule, which has a volume of 496.32 and diameter of 9.823, has a low probability of belonging to the first expert, whereas the nodule on the right, with a volume of 6663.44 and diameter of 23.347, has a high probability of being assigned to the first expert. In other words, the first expert tends to handle larger nodules, compared to the second expert.

Although MoE techniques are shown to be able to improve the classification performance, they typically face an objection related to the high computational cost at the test time. This problem, however, can be dealt with by using distillation [54]. Therefore, in our future studies, we will focus on distilling MIXCAPS into a smaller and more time-efficient model.

4.1. MIXCAPS for Brain Tumor Type Classification

Brain tumor is among the deadliest cancers. Determining the type of the tumor, which is a challenging task in terms of accuracy and inter-observer variability, can significantly facilitate the control/treatment process. Therefore, we dedicate this subsection to investigate the generalizability of the proposed MIXCAPS to brain tumor type classification. In a previous study [29], we proposed a capsule network-based framework, which we referred to as the BoxCaps, for brain tumor classification, considering not only raw magnetic resonance imaging (MRI) inputs, but also the coarse tumor boundaries. The motivation behind such framework was that the whole brain image contained valuable information on the location of the tumor with respect to the brain tissue. The CapsNet, however, tends to get distracted from the main tumor region when being fed with all the details from the brain image. As such, we designed a modified architecture where the output capsules were concatenated with the tumor course boundary box. This way, the model had access to both brain tissue and tumor region.

To investigate whether the MIXCAPS can be generalized to brain tumor classification, we replaced the capsule experts in MIXCAPS with the previously designed BoxCaps architecture, as shown in Fig. 5. We then tested the resulting framework on a brain tumor dataset [55], where train, validation, and test splits

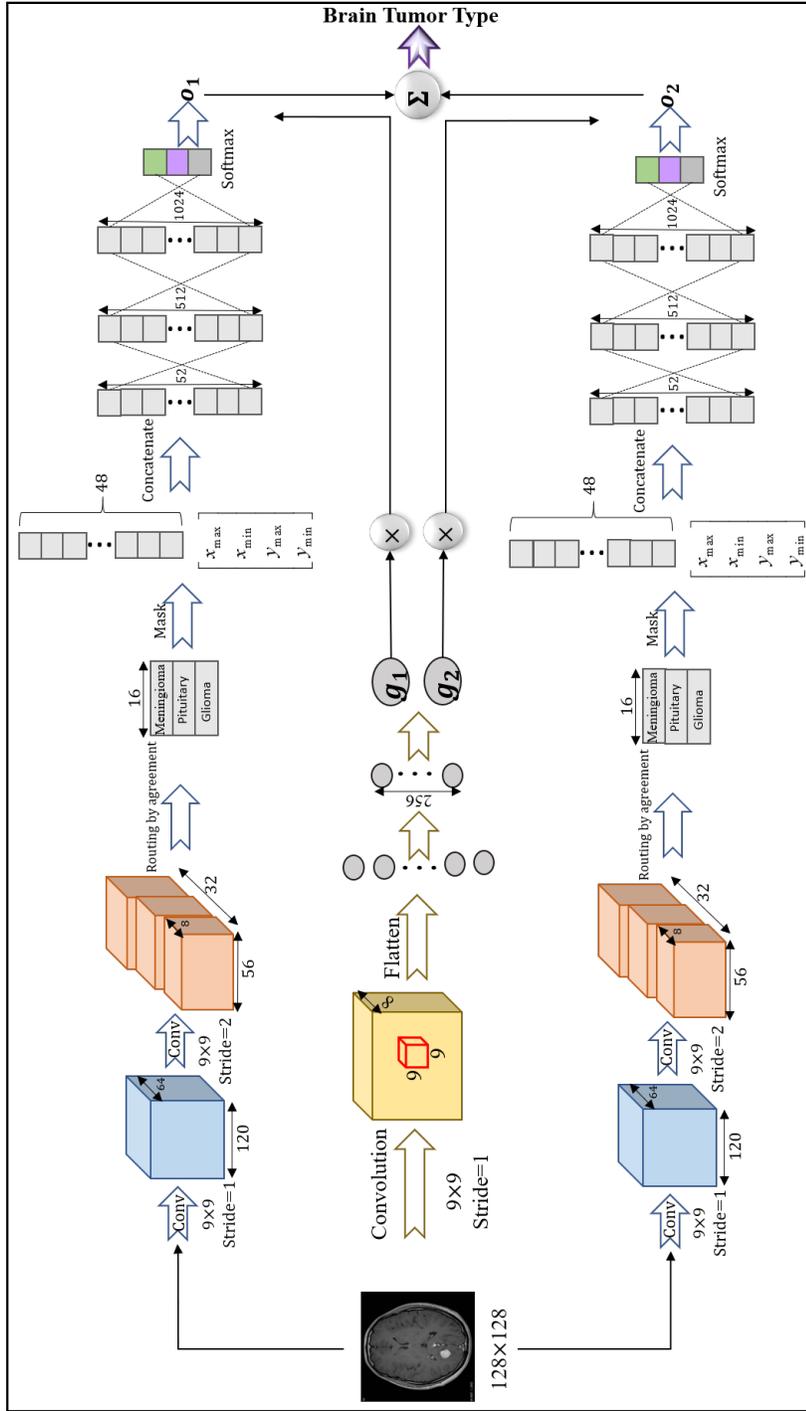


Figure 5: MIXCAPS architecture with BoxCaps as experts for brain tumor type classification.

are obtained from the same bootstrapping approach used for the LIDC-IDRI dataset. The aforementioned dataset consists of 3,064 images from 233 patients, diagnosed with one of the three brain tumor types, i.e., Meningioma, Pituitary, and Glioma. Table 3 presents the obtained results, according to which, the MoE approach leads to higher accuracy compared to a single BoxCaps. Furthermore, the MoE approach leads to higher sensitivity for Glioma and Pituitary, and higher specificity for Meningioma and pituitary tumor types.

Table 3: Performance of the proposed MIXCAPS with BoxCaps as experts. Numbers in parenthesis show the 95% confidence intervals.

	MIXCAPS-BoxCaps	BoxCaps
Accuracy	91.3 (91.1, 91.5) %	90.9 (90.2, 91.5) %
Sensitivity for Meningioma	77.5 (77.1, 77.9) %	80.1 (76.2, 84) %
Sensitivity for Glioma	95.9 (93.2, 98.5) %	92 (90, 94.1) %
Sensitivity for Pituitary	97.7 (97.2, 98.3) %	97.2 (95.6, 98.9) %
Specificity for Meningioma	96.1 (96, 96.1) %	94.1 (92.7, 95.5) %
Specificity for Glioma	88.7 (87.6, 89.8) %	89.8 (88.4, 91.2) %
Specificity for Pituitary	88.7 (86.2, 91.2) %	88.1 (86.9, 89.3) %

Finally, we conduct another experiment to study if the provided boundary box is the only important factor leading to the obtained result. In other words we need to make sure that the input images are not ignored by the model, simply because the boundary box itself can determine the tumor type. To this end, we gradually added zero-mean Gaussian noise to input images and calculated the model’s accuracy. It is observed that while a noise with a standard deviation (STD) of 0.01 does not change the accuracy, increasing STD to 0.1 and 0.5 degrades the accuracy to 84.44% and 76%, respectively. This experiment shows that while the boundary box assists the classification, it does not replace the input images.

5. Conclusion and Future Direction

In this paper, we proposed a capsule network-based mixture of experts framework, referred to as the MIXCAPS, for lung nodule malignancy prediction. The proposed MIXCAPS framework contains two capsule network experts and a convolutional gating network to assign instances to experts. Our obtained results show that MIXCAPS outperforms a single capsule network and a mixture of CNNs. It has also several advantages over the previous methods. First, MIXCAPS utilizes capsule networks and is therefore capable of handling smaller datasets. Second, through the MoE approach, experts get the chance to specialize on a subset of the data. Furthermore, MIXCAPS does not require fine annotations and is independent from pre-defined hand-crafted features. Our future directions include exploring capsule gating networks and optimizing the number of experts, as well as focusing on MIXCAPS knowledge distillation to improve the model's time-efficacy.

References

- [1] F. Bray, J. Ferlay, I. Soerjomataram, R. L. Siegel, L. A. Torre, A. Jemal, Global cancer statistics 2018: Globocan estimates of incidence and mortality worldwide for 36 cancers in 185 countries, *CA: A Cancer Journal for Clinicians* 66 (2018) 7–30.
- [2] R. L. Siegel, K. D. Miller, A. Jemal, *Cancer statistics, 2016*, *CA: A Cancer Journal for Clinicians* 68 (2016) 394–424.
- [3] H. Xie, D. Yang, N. Sun, Z. Chen, Y. Zhang, Automated pulmonary nodule detection in ct images using deep convolutional neural networks, *Pattern Recognition* 85 (2019) 109–119.
- [4] D. R. Aberle, A. M. Adams, C. D. Berg, W. C. Black, J. D. Clapp, R. M. Fagerstrom, I. F. Gareen, C. Gatsonis, P. M. Marcus, J. D. Sicks, Reduced lung-cancer mortality with low-dose computed tomographic screening, *N Engl J Med.* 365 (2011) 395–409.

- [5] J. L. Causey, J. Zhang, S. Ma, B. Jiang, J. A. Qualls, D. G. Politte, F. Prior, S. Zhang, X. Huang, Highly accurate model for prediction of lung nodule malignancy with ct scans., *Scientific Reports* 8.
- [6] Y. Zhang, A. Oikonomou, A. Wong, M. A. Haider, F. Khalvati, Radiomics-based prognosis analysis for non-small cell lung cancer, *Scientific Reports* 7 (2017) 481–487.
- [7] H. J. W. L. Aerts, E. R. Velazquez, R. T. H. Leijenaar, C. Parmar, P. Grossmann, S. Carvalho, J. Bussink, R. Monshouwer, B. Haibe-Kains, D. Rietveld, F. Hoebers, M. M. Rietbergen, C. R. Leemans, A. Dekker, J. Quackenbush, R. J. Gillies, P. Lambin, Decoding tumour phenotype by noninvasive imaging using a quantitative radiomics approach, *Nature Communications* 5.
- [8] A. Oikonomou, F. Khalvati, P. N. Tyrrell, M. A. Haider, U. Tarique, L. Jimenez-Juan, M. C. Tjong, I. Poon, A. Eilaghi, L. Ehrlich, P. Cheung, Radiomics analysis at pet/ct contributes to prognosis of recurrence and survival in lung cancer treated with stereotactic body radiotherapy, *Scientific Reports* 8.
- [9] A. Afshar, A. Mohammadi, N. P. Konstantinos, A. Oikonomou, H. Benali, From hand-crafted to deep learning-based cancer radiomics: Challenges and opportunities, *IEEE Signal Processing Magazine* 36 (2019) 132–160.
- [10] P. Lambin, E. Rios-Velazquez, R. Leijenaar, S. Carvalho, R. G. van Stiphout, P. Granton, C. M. Zegers, R. Gillies, R. Boellard, A. Dekker, H. J. Aerts, Radiomics: extracting more information from medical images using advanced feature analysis., *Eur J Cancer* 48 (2012) 441–446.
- [11] C. Chen, C. Chang, C. Tu, W. Liao, B. Wu, K. Chou, Y. Chiou, S. Yang, G. Zhang, T. Huang, Radiomic features analysis in computed tomography images of lung nodule classification., *PLoS One* 13.

- [12] C. Parmar, R. T. H. Leijenaar, P. Grossmann, E. R. Velazquez, J. Bussink, D. Rietveld, M. M. Rietbergen, B. Haibe-Kains, P. Lambin, H. J. W. L. Aerts, Radiomic feature clusters and prognostic signatures specific for lung and head and neck cancer, *Scientific Reports* 5.
- [13] T. P. Coroller, V. Agrawal, V. Narayan, Y. Hou, P. Grossmann, S. W. Lee, R. H. Mak, H. J. W. L. Aerts, Multiview convolutional neural networks for lung nodule classification, *Radiotherapy and Oncology* 119 (2016) 480–486.
- [14] E. Huynh, T. P. Coroller, V. Narayan, V. Agrawal, Y. Hou, J. Romano, I. Franco, R. H. Mak, H. J. W. L. Aerts, Ct-based radiomic analysis of stereotactic body radiation therapy patients with lung cancer, *Radiotherapy and Oncology* 120 (2016) 258–266.
- [15] R. J. Gillies, P. E. Kinahan, H. Hricak, Radiomics: Images are more than pictures, they are data, *Radiology*.
- [16] A. Oikonomou, P. Salazar, Y. Zhang, D. Hwang, A. Petersen, A. Dmytriw, N. Paul, E. Nguyen, Histogram-based models on non-thin section chest ct predict invasiveness of primary lung adenocarcinoma subsolid nodules, *Scientific Reports* 9.
- [17] S. S. F. Yip, H. J. W. L. Aerts, Applications and limitations of radiomics, *Physics in Medicine and Biology* 61.
- [18] J. E. Park, S. Y. Park, H. J. Kim, H. S. Kim, Reproducibility and generalizability in radiomics modeling: Possible strategies in radiologic and statistical perspectives, *Korean Journal of Radiology* 20 (2019) 1124–1137.
- [19] Z. Li, Y. Wang, J. Yu, Y. Guo, W. Cao, Deep learning based radiomics (dlr) and its usage in noninvasive idh1 prediction for low grade glioma, *Scientific Reports* 7.
- [20] L. Oakden-Rayner, G. Carneiro, T. Bessen, J. C. Nascimento, A. P. Bradley, L. J. Palmer, Precision radiology: Predicting longevity using fea-

ture engineering and deep learning methods in a radiomics framework, Scientific Reports 7.

- [21] K. H. Cha, L. Hadjiiski, H. Chan, A. Z. Weizer, A. Alva, R. H. Cohan, E. M. Caoili, C. Paramagul, R. K. Samala, Bladder cancer treatment response assessment in ct using radiomics with deep-learning, Scientific Reports 7.
- [22] D. Kumar, A. G. Chung, M. J. Shaifee, F. Khalvati, M. A. Haider, A. Wong, Discovery radiomics for pathologically-proven computed tomography lung cancer prediction, Karray F., Campilho A., Cheriet F. (eds) Image Analysis and Recognition. ICIAR 2017. Lecture Notes in Computer Science, Springer, Cham 10317.
- [23] A. Krizhevsky, I. Sutskever, G. E. Hinton, Imagenet classification with deep convolutional neural networks, Neural Information Processing Systems (NIPS).
- [24] R. Yamashita, M. Nishio, R. K. G. Do, K. Togashi, Convolutional neural networks: An overview and application in radiology, Insights into Imaging 9 (2018) 611629.
- [25] S. Sabour, N. Frosst, G. E. Hinton, Dynamic routing between capsules, Neural Information Processing Systems (NIPS).
- [26] R. LaLondea, Z. Xub, S. Jainc, U. Bagcia, Capsules for biomedical image segmentation, arXiv:2004.04736v1.
- [27] P. Afshar, A. Oikonomou, F. Naderkhani, P. N. Tyrrell, K. N. Plataniotis, K. Farahani, A. Mohammadi, A 3d multi- scale capsule network for lung nodule malignancy classification, Scientific Reports 10.
- [28] P. Afshar, A. Mohammadi, K. N. Plataniotis, Brain tumor type classification via capsule networks, 25th IEEE International Conference on Image Processing (ICIP) (2018) 3129–3133.

- [29] P. Afshar, K. N. Plataniotis, A. Mohammadi, Capsule networks for brain tumor classification based on mri images and coarse tumor boundaries, 44th IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP) (2019) 1368–1372.
- [30] A. Mohammadi, K. N. Plataniotis, Improper complex-valued multiple-model adaptive estimation, IEEE Transactions on Signal Processing 63 (2015) 1528–1542.
- [31] R. A. Jacobs, M. I. Jordan, S. J. Nowlan, G. E. Hinton, Adaptive mixtures of local experts, Neural Computation 3 (1991) 79–87.
- [32] R. Rasti, M. Teshnehlab, S. M. Phung, Breast cancer diagnosis in dce-mri using mixture ensemble of convolutional neural networks, Pattern Recognition 72 (2017) 381–390.
- [33] J. Guo, S. Gould, Deep cnn ensemble with data augmentation for object detection, arXiv:1506.07224.
- [34] L. Mao, H. Chen, M. Liang, K. Li, J. Gao, P. Qin, X. Ding, X. Li, X. Liu, Quantitative radiomic model for predicting malignancy of small solid pulmonary nodules detected by low-dose ct screening, Quant Imaging Med Surg. 9 (2019) 263272.
- [35] A. de Carvalho Filho, A. Correa Silva, A. de Paiva, R. Nunes, M. Gattass, Classification of patterns of benignity and malignancy based on ct using topology-based phylogenetic diversity index and convolutional neural network, Pattern Recognition 81 (2018) 200–212.
- [36] X. Liu, F. Hou, H. Qin, A. Hao, Multi-view multi-scale cnns for lung nodule type classification from ct images, Pattern Recognition 77 (2018) 262–275.
- [37] W. Sun, B. Zheng, W. Qian, Automatic feature learning using multichannel roi based on deep structured algorithms for computerized lung cancer diagnosis, Computers in Biology and Medicine 89 (2017) 530–539.

- [38] H. Wang, Z. Zhou, Y. Li, Z. Chen, P. Lu, W. Wang, W. Liu, L. Yu, Comparison of machine learning methods for classifying mediastinal lymph node metastasis of non-small cell lung cancer from 18f-fdg pet/ct images, *EJNMMI Research* 7.
- [39] L. Fu, J. Ma, Y. Ren, Y. S. Han, J. Zhao, Automatic detection of lung nodules: false positive reduction using convolution neural networks and handcrafted features, *Proc.SPIE* 10134.
- [40] K. Adu, Y. Yu, J. Cai, N. Tashi, Dilated capsule network for brain tumor type classification via mri segmented tumor region, *IEEE International Conference on Robotics and Biomimetics (ROBIO)*.
- [41] Y. Cheng, G. Qin, R. Zhao, Y. Liang, M. Sun, Convcaps: Multi-input capsule network for brain tumor classification, *International Conference on Neural Information (2019)* 524–534.
- [42] D. Maji, A. Santara, P. Mitra, D. Sheet, Ensemble of deep convolutional neural networks for learning to detect retinal vessels in fundus images, *arXiv:1603.04833*.
- [43] H. Chen, Q. Dou, X. Wang, J. Qin, P. Heng, Mitosis detection in breast cancer histology images via deep cascaded networks, *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence (2016)* 1160–1166.
- [44] H. Wang, A. Cruz-Roa, A. Basavanhally, H. Gilmore, N. Shih, M. Feldman, J. Tomaszewski, F. Gonzalez, A. Madabhushi, Mitosis detection in breast cancer pathology images by combining handcrafted and convolutional neural network features, *Journal of Medical Imaging* 1.
- [45] S. Armato III, G. McLennan, L. Bidaut, M. McNitt-Gray, C. Meyer, A. Reeves, B. Zhao, D. Aberle, C. Henschke, E. Hoffman, Eric A. Kazerooni, H. MacMahon, E. van Beek, D. Yankelevitz, A. Biancardi, P. Bland, M. Brown, R. Engelmann, G. Laderach, D. Max, R. Pais, D. Qing, R. Roberts, A. Smith, A. Starkey, P. Batra, P. Caligiuri, A. Farooqi,

- G. Gladish, C. Jude, R. Munden, I. Petkowska, L. Quint, L. Schwartz, B. Sundaram, L. Dodd, C. Fenimore, D. Gur, N. Petrick, J. Freymann, J. Kirby, B. Hughes, A. Castele, S. Gupte, M. Sallam, M. Heath, M. Kuhn, E. Dharaiya, R. Burns, D. Fryd, M. Salganicoff, V. Anand, U. Shreter, S. Vastagh, B. Croft, L. Clarke, Data from lidc-idri., The Cancer Imaging Archive.
- [46] S. Armato III, G. McLennan, L. Bidaut, M. McNitt-Gray, C. Meyer, A. Reeves, B. Zhao, D. Aberle, C. Henschke, E. Hoffman, Eric A. Kazerooni, H. MacMahon, E. van Beek, D. Yankelevitz, A. Biancardi, P. Bland, M. Brown, R. Engelmann, G. Laderach, D. Max, R. Pais, D. Qing, R. Roberts, A. Smith, A. Starkey, P. Batra, P. Caligiuri, A. Farooqi, G. Gladish, C. Jude, R. Munden, I. Petkowska, L. Quint, L. Schwartz, B. Sundaram, L. Dodd, C. Fenimore, D. Gur, N. Petrick, J. Freymann, J. Kirby, B. Hughes, A. Castele, S. Gupte, M. Sallam, M. Heath, M. Kuhn, E. Dharaiya, R. Burns, D. Fryd, M. Salganicoff, V. Anand, U. Shreter, S. Vastagh, B. Croft, The lung image database consortium (lidc) and image database resource initiative (idri): A completed reference database of lung nodules on ct scans., *Medical Physics* 38 (2011) 915–931.
- [47] K. Clark, B. Vendt, K. Smith, J. Freymann, J. Kirby, P. Koppel, S. Moore, S. Phillips, D. Maffitt, M. Pringle, L. Tarbox, F. Prior, The cancer imaging archive (tcia): Maintaining and operating a public information repository, *Journal of Digital Imaging* 26 (2013) 1045–1057.
- [48] M. Jordan, R. Jacobes, Hierarchical mixtures of experts and the em algorithm, *Proceedings of 1993 International Joint Conference on Neural Networks* (1993) 1339–1344.
- [49] A. Mohammadi, K. Plataniotis, Improper complex-valued multiple-model adaptive estimation, *IEEE TRANSACTIONS ON SIGNAL PROCESSING* 63 (2015) 1528–1542.
- [50] Y. Xie, J. Zhang, S. Liu, W. Cai, Y. Xia, Lung nodule classification by

jointly using visual descriptors and deep features, *Medical Computer Vision and Bayesian and Graphical Models for Biomedical Imaging. BAMBI 2016, MCV 2016. Lecture Notes in Computer Science*, Springer, Cham 10081.

- [51] A. Nibali, H. Zhen, D. Wollersheim, Pulmonary nodule classification with deep residual networks, *International Journal of Computer Assisted Radiology and Surgery* 12 (2017) 1799–1808.
- [52] W. Sun, B. Zheng, W. Qian, Computer aided lung cancer diagnosis with deep learning algorithms, *Proceedings of SPIE* 9785.
- [53] W. Shen, M. Zhou, F. Yang, D. Yu, D. Dong, C. Yang, Y. Zang, J. Tian, Multi-crop convolutional neural networks for lung nodule malignancy suspiciousness classification, *Pattern Recognition* 61 (2017) 663–673.
- [54] G. Hinton, V. Oriol, J. Dean, Distilling the knowledge in a neural network, *NIPS Deep Learning and Representation Learning Workshop*.
- [55] J. Cheng, W. Yang, M. Huang, W. Huang, J. Jiang, Y. Zhou, R. Yang, J. Zhao, Y. Feng, Q. Feng, W. Chen, Retrieval of brain tumors by adaptive spatial pooling and fisher vector representation, *PloS one*.