



# Reasoning with Uncertainty in Deep Learning for Safer Medical Image Computing

Ryutaro Tanno

Department of Computer Science  
University College London  
Gower Street, London, United Kingdom

THESIS

Submitted for the degree of  
**Doctor of Philosophy, University College London**

July 19, 2021

I, Ryutaro Tanno, confirm that the work presented in this thesis is my own. Where information has been derived from other sources, I confirm that this has been indicated in the thesis.

# Abstract

Deep learning is now ubiquitous in the research field of medical image computing. As such technologies progress towards clinical translation, the question of safety becomes critical. Once deployed, machine learning systems unavoidably face situations where the correct decision or prediction is ambiguous. However, the current methods disproportionately rely on deterministic algorithms, lacking a mechanism to represent and manipulate uncertainty. In safety-critical applications such as medical imaging, reasoning under uncertainty is crucial for developing a reliable decision making system.

Probabilistic machine learning provides a natural framework to quantify the degree of uncertainty over different variables of interest, be it the prediction, the model parameters and structures, or the underlying data (images and labels). Probability distributions are used to represent all the uncertain unobserved quantities in a model and how they relate to the data, and probability theory is used as a language to compute and manipulate these distributions. In this thesis, we explore probabilistic modelling as a framework to integrate uncertainty information into deep learning models, and demonstrate its utility in various high-dimensional medical imaging applications. In the process, we make several fundamental enhancements to current methods.

We categorise our contributions into three groups according to the types of uncertainties being modelled: (i) predictive; (ii) structural and (iii) human uncertainty. *Firstly*, we discuss the importance of quantifying predictive uncertainty and understanding its sources for developing a risk-averse and transparent medical image enhancement application. We demonstrate how a measure of predictive uncertainty can be used as a proxy for the predictive accuracy in the absence of ground-truths. Furthermore, assuming the structure of the model is flexible enough for the task, we introduce a way to decompose the predictive uncertainty into its orthogonal sources i.e. aleatoric and parameter uncertainty. We show the potential utility of such decoupling in providing a quantitative “explanations” into the model performance. *Secondly*, we introduce our recent attempts at learning model structures directly from data. One work proposes a method based on variational inference to learn a posterior distribution over connectivity structures within a neural network architecture for multi-task learning, and share some preliminary results in the MR-only radiotherapy planning application. Another work explores how the training algorithm of decision trees could be extended to grow the architecture of a neural network to adapt to the given availability of data and the complexity of the task. *Lastly*, we develop methods to model the “measurement noise” (e.g., biases and skill levels) of human annotators, and integrate this information into the learning process of the neural network classifier. In particular, we show that explicitly modelling the uncertainty involved in the annotation process not only leads to an improvement in robustness to label noise, but also yields useful insights into the patterns of errors that characterise individual experts.



# Impact Statement

This thesis introduces new algorithms for modelling different types of uncertainty in deep learning models and explores utility thereof in high-dimensional, challenging medical imaging applications where safety is critical. We have demonstrated that appropriate modelling of uncertainty not only improves the predictive performance of the machine learning model in an array of settings (e.g., within-distribution and out-of-distribution generalisation, presence of outliers/noisy labels in the training data, scenarios in which multiple tasks are jointly learned, etc), but also enables quantification of risks associated with the prediction, which can be used to manage the safety of downstream decision-making. While most of the algorithmic advances were motivated by specific medical applications, the resultant solutions are general and thus their applicability extends to other problems in medical imaging and beyond (e.g., computer vision).



# Acknowledgements



# Contents

<b>1</b>	<b>Introduction</b>	<b>13</b>
1.1	Taxonomy of Uncertainty and Why Should We Care? . . . . .	14
1.2	“Classics” in Uncertainty Modelling for Medical Imaging . . . . .	16
1.3	Deep Learning in Medical Imaging and What’s Lacking . . . . .	20
<b>2</b>	<b>Quantifying Predictive Uncertainty and Its Sources for Super-Resolution</b>	<b>25</b>
2.1	Introduction . . . . .	25
2.2	Background: Image Quality Transfer . . . . .	26
2.2.1	Image Enhancement as Patch-wise Regression . . . . .	26
2.2.2	Baseline Super-Resolution Model . . . . .	27
2.3	Methods . . . . .	28
2.3.1	Sources of Predictive Uncertainty . . . . .	28
2.3.2	Aleatoric Uncertainty and Heteroscedastic Noise Model . . . . .	29
2.3.3	Parameter Uncertainty and Variational Dropout . . . . .	30
2.3.4	Joint Modelling of Aleatoric and Parameter Uncertainty . . . . .	32
2.3.5	Uncertainty Decomposition and Propagation . . . . .	32
2.4	Related works . . . . .	34
2.5	Data preprocessing and implementation details . . . . .	35
2.5.1	Datasets . . . . .	35
2.5.2	Network Architectures and Training . . . . .	37
2.6	Results . . . . .	38
2.6.1	Benefits on Super-resolution Performance . . . . .	38
2.6.2	Reliability Assessment of Model Predictions . . . . .	40
2.7	Discussion and Conclusion . . . . .	45
<b>3</b>	<b>Uncertainty in Multitask Learning (I): Spatially Adaptive Weighting of Task Loss Functions</b>	<b>49</b>
3.1	Introduction . . . . .	49
3.2	Related work . . . . .	50
3.3	Methods . . . . .	50
3.3.1	Bi-task architecture . . . . .	50
3.3.2	Task weighting with heteroscedastic uncertainty. . . . .	51
3.3.3	Parameter uncertainty with approximate Bayesian inference. . . . .	52
3.4	Data Preprocessing and Implementation Details . . . . .	53
3.4.1	Data . . . . .	53
3.4.2	Network architectures and training . . . . .	53
3.5	Results . . . . .	53
3.5.1	Experimental set-up . . . . .	53

---

3.5.2	Model performance . . . . .	53
3.5.3	Uncertainty estimation for radiotherapy . . . . .	55
3.6	Conclusions . . . . .	55
<b>4</b>	<b>Uncertainty in Multitask Learning (II): Stochastic Filter Groups for Learning Structured Sparsity</b>	<b>57</b>
4.1	Introduction . . . . .	57
4.2	Related works . . . . .	58
4.3	Methods . . . . .	59
4.3.1	Stochastic Filter Groups . . . . .	59
4.3.2	Optimisation: T+1 Way “Drop-Out” . . . . .	61
4.4	Experiments . . . . .	63
4.4.1	Baselines . . . . .	63
4.5	Data preprocessing and implementation details . . . . .	64
4.5.1	Optimisation, regularisation and initialisation . . . . .	64
4.5.2	UTKFace . . . . .	64
4.5.3	Medical imaging dataset . . . . .	64
4.5.4	Implementation details . . . . .	65
4.5.5	CNN architecture details . . . . .	65
4.6	Results . . . . .	66
4.6.1	Age regression and gender prediction . . . . .	66
4.6.2	Image regression and semantic segmentation . . . . .	66
4.6.3	Learned architectures . . . . .	68
4.6.4	Learned grouping probability plots . . . . .	69
4.6.5	Effect of $p$ initialisation . . . . .	69
4.6.6	Learned filter groups on duplicate tasks . . . . .	69
4.7	Discussion . . . . .	74
<b>5</b>	<b>How to Learn Network Architecture like a Decision Tree</b>	<b>75</b>
5.1	Introduction . . . . .	75
5.2	Related work . . . . .	77
5.3	Adaptive Neural Trees . . . . .	78
5.3.1	Model Topology and Operations . . . . .	78
5.3.2	Probabilistic Model and Inference . . . . .	80
5.4	Optimisation . . . . .	81
5.4.1	Loss function: optimising parameters of $\mathbb{O}$ . . . . .	81
5.4.2	Growth phase: learning architecture $T$ . . . . .	81
5.4.3	Refinement phase: global tuning of $\mathbb{O}$ . . . . .	82
5.5	Experiments and Results . . . . .	82
5.5.1	Set-up details . . . . .	82
5.5.2	Model Performance . . . . .	83
5.5.3	Interpretability . . . . .	86
5.5.4	Effect of global refinement . . . . .	88
5.5.5	Adaptive model complexity . . . . .	88
5.6	Discussion and Conclusion . . . . .	92
<b>6</b>	<b>Modelling Human Uncertainty (I): Classification</b>	<b>93</b>
6.1	Introduction . . . . .	93
6.1.1	Other Related Works. . . . .	94

6.2	Methods . . . . .	95
6.2.1	Noisy Observation Model . . . . .	95
6.2.2	Joint Estimation of Confusion and True labels . . . . .	96
6.2.3	Motivation for Trace Regularization . . . . .	99
6.3	Experiments and Results . . . . .	100
6.3.1	Datasets, training and architectures . . . . .	100
6.3.2	Set-Up . . . . .	100
6.3.3	Comparing with EM-based Approaches . . . . .	101
6.3.4	Value of Modelling Individual Annotators . . . . .	105
6.3.5	Experiments on Cardiac View Classification . . . . .	106
6.4	Discussion and Conclusion . . . . .	108
<b>7</b>	<b>Modelling Human Uncertainty (II): Semantic Segmentation</b>	<b>109</b>
7.1	Introduction . . . . .	109
7.2	Related Works . . . . .	110
7.3	Method . . . . .	112
7.3.1	Problem Set-up . . . . .	112
7.3.2	Probabilistic Model and Proposed Architecture . . . . .	112
7.3.3	Learning Confusion Matrices and True Segmentation . . . . .	112
7.3.4	Justification for the Trace Norm . . . . .	113
7.4	Experiments . . . . .	114
7.4.1	MNIST and MS lesion segmentation datasets . . . . .	115
7.4.2	BraTS Dataset and LIDC-IDRI Dataset . . . . .	118
7.5	Discussion and Conclusion . . . . .	122
<b>8</b>	<b>Conclusions and Future Work</b>	<b>123</b>
<b>Bibliography</b>		<b>127</b>



# Chapter 1

## Introduction

Deep learning [1] has become ubiquitous in the field of medical image processing. With the combination of the increasing volume of digitised imaging data and advances in both hardware and software, these so-called “black-box” machine learning (ML) algorithms have enabled a substantial leap in performance in a plethora of applications from image analysis to surgical assistance [2, 3]. For certain well-defined tasks with access to homogeneous and carefully annotated data, such machine learning systems have begun to surpass the performance of clinical experts [4, 5, 6, 7]. However, translation of these research innovations into clinical practice requires care. In applications where the algorithm’s outputs inform scientific conclusions in research, and diagnostic, prognostic or interventional decisions in clinics, we need principled protocols to ensure safety.

In practice, ML systems often face situations where the correct decision or prediction is ambiguous. We therefore need a mechanism to quantify the confidence of the model output (e.g. error bounds) and act upon it to prevent catastrophic failures. We would also like to be able to reason about the sources of such uncertainty, and further improve the performance. Does the training data need to be more diverse? Were the images or the annotations too noisy? It might have been that the choice of the model was not adequate. Or perhaps we may have been just unlucky and the particular instance of failures is an inherently challenging case where the input image did not contain enough information. Implementing systematic approaches to answering these questions is important not only to improve the predictive performance, but also to build trust with the practitioners. However, to date, the majority of deep learning techniques used in the medical imaging context rely on deterministic methods, and lack a mechanism to communicate uncertainty, a key ingredient to address such problems.

In contrast, *probabilistic machine learning* provides a natural framework to quantify the degree of uncertainty over different variables of interest, be it the prediction, the model parameters and structures, or the underlying data (images and labels) [8]. Probability distributions are used to represent all the uncertain unobserved quantities in a model and how they relate to the data, and probability theory is used as a language to compute and manipulate these distributions. The main goal of this thesis is to develop a practical framework for medical imaging applications to model and reason with different types of uncertainty in deep learning models by translating ideas from the paradigm of probabilistic machine learning. In the process, several fundamental enhancements to current methods arise.

## 1.1 Taxonomy of Uncertainty and Why Should We Care?

There are many different types of uncertainty that are of practical importance in medical imaging applications. Here we provide a taxonomy of such uncertainty “species” as illustrated in Fig. 1.1 and explain their relations and differences.

Imagine you have a machine learning model,  $F_\theta$ , which takes some input  $\mathbf{x}$  (e.g. a MR image) and makes a prediction or decision,  $\hat{\mathbf{y}} = F_\theta(\mathbf{x})$ , about the quantity of interest,  $\mathbf{y}$  (e.g. presence of pathology). Here we use different notations,  $\mathbf{y}$  and  $\hat{\mathbf{y}}$ , to emphasize the difference between the target output variable and its estimate from the model. In this case,  $\theta$  denotes the parameters of the model, and has been optimised based on training data, consisting of  $N$  pairs of inputs and labels of interest  $\mathcal{D} = \{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^N$ . You now have thoroughly evaluated the performance on a held-out test dataset (e.g. a large public imaging dataset), and are about to deploy it in the “wild” (e.g. your local hospital).

*Predictive uncertainty* describes the degree of ambiguity (or confidence) in the model’s prediction on a given input. For example, error bounds are a common measure of predictive uncertainty, and can be used to assess the reliability of prediction for the particular data instance at hands. You might be interested in knowing how well your model performs in a new environment, but you may not have access to a sufficient volume of ground-truth labels for validation. However, quantification of predictive uncertainty provides a proxy measure of performance on new data point, which can be used to manage failure risks in a principled way—if the model encounters input examples with very high predictive uncertainty, then we should not trust the model’s output.

One is often interested not only in quantifying the predictive uncertainty, but also understanding its *sources*. Having quantitative answers to questions such as “is the model uncertain on this particular image because the observed feature is not represented in the training data or the image quality is too low to make definitive decisions? Perhaps the current model does not best explain the data and we should use a different one?” has clear benefits for building a more reliable system. Commonly, such *sources* of predictive uncertainty is divided into two types, *aleatoric* and *epistemic* uncertainty [9, 10].

*Aleatoric uncertainty* (also known as intrinsic uncertainty) — from the Latin word, *alea*, meaning a “die” — refers to uncertainty inherent to a problem that in principle cannot be reduced by additional physical or experimental knowledge. The mapping from input  $\mathbf{x}$  to target  $\mathbf{y}$  which one wants to approximate from data may be intrinsically stochastic. For example, an acquired image  $\mathbf{x}$  does not contain enough information to conclude whether an observed feature is pathological or not, and thus both predictions “pathology is present” and “pathology is absent” are equally probable. Such inherent ambiguity is only worsened by the fact that the training data  $\mathcal{D}$  may be corrupted by observational noise (such as measurement errors and annotation noise).

*Epistemic uncertainty* — from the Greek word *episteme*, meaning “knowledge” — refers to uncertainty arising from lack of knowledge. In the context of modelling, epistemic uncertainty is often subdivided into *parameter uncertainty*, in which one believes that the form of the model reflects reality well, but one is uncertain about which values of the parameters  $\theta$  in the model to use, and *structural uncertainty*, in which one has significant doubts that the model  $F_\theta(\cdot)$  is even ‘structurally correct’ (e.g. is linear regression appropriate or a neural network, if the latter, how many layers, etc.).

We should, however, note that the distinction between epistemic and aleatoric uncer-

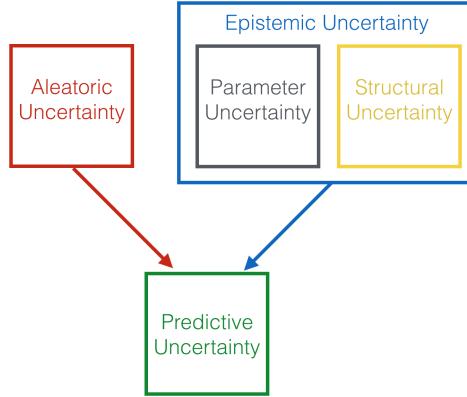


Figure 1.1: Illustration of different types of uncertainty. The combined effects of both aleatoric and epistemic uncertainties induce *predictive uncertainty*.

tainty is a rather ambiguous one and depends somewhat on the interpretation of what probability means [11, 12]. For example, an orthodox frequentist, who believes in the definition of a probability based on the objective relative frequency of events, would say that the results of a coin flip are random and represents aleatoric uncertainty. On the other hand, a devout Bayesian, who believes probabilities represent degrees of subjective ignorance, might say that the coin flip is characterised by epistemic uncertainty, arguing one is simply ignorant about the set of influential parameters e.g. initial conditions of the coin such as the angle and height of the drop, the material and shape of the coin, and the dynamics of the wind which may affect the trajectory, etc. In any case, putting aside the details of philosophical debates surrounding such categorisation of uncertainty, the language of probability theory provides a powerful tool in describing different forms of uncertainty [13]. This thesis, in particular, focuses on the ways to describe such uncertainty with(in) deep learning models, and aims to demonstrate their practical benefits in medical imaging applications.

In the probabilistic framework of machine learning, all uncertain quantities are treated as random variables, and probability distributions are used to express their associated information of uncertainty (see Table 1.1). For example, the aleatoric uncertainty of the target mapping  $\mathbf{x} \rightarrow \mathbf{y}$  is summarised by the underlying conditional distribution of the task  $P(\mathbf{y}|\mathbf{x})$ . This distribution describes the inherent stochasticity in predicting  $\mathbf{y}$  for the given input  $\mathbf{x}$ . On the other hand, the parameter uncertainty is described by  $P(\theta|\mathcal{D}, m)$ , the posterior distribution over the unknown parameters  $\theta$  of the model class  $m = \{F_\theta(\cdot); \theta \in \Theta\}$  given the training observations  $\mathcal{D}$ . Similarly, structural uncertainty is described by the distribution  $P(m|\mathcal{D})$  which quantifies how probable the model class is given the observations. Finally, the predictive uncertainty is represented in the conditional distribution  $P(\hat{\mathbf{y}}|\mathbf{x})$  over the model's output  $\hat{\mathbf{y}} = F_\theta(\mathbf{x})$ , and its standard deviation is commonly used as a confidence interval. Many technical problems tackled in this thesis boil down to the estimation of these distributions in different settings. As will be explained later, for complex models such as deep neural networks and decision trees, many of these distributions are not tractable, necessitating efficient and effective approximations.

To make matters even worse, in practical applications such as medical imaging, the ground truths for such distributions of interest,  $P(\mathbf{y}|\mathbf{x})$ ,  $P(\hat{\mathbf{y}}|\mathbf{x})$ ,  $P(\theta|\mathcal{D}, m)$  and  $P(m|\mathcal{D})$  are not available, thus rendering the direct evaluation of uncertainty estimation unfeasible. In this thesis, I take a pragmatic position and focus on evaluating the utility—rather than the fidelity—of the derived uncertainty estimates via surrogate measures, such as

Table 1.1: Uncertainty types and their distributional forms.

Uncertainty Type	Distributional Form	Ambiguity in
Predictive	$P(\hat{\mathbf{y}} \mathbf{x})$	the model's output
Aleatoric	$P(\mathbf{y} \mathbf{x})$	the data formation process
Epistemic - Parametric	$P(\theta \mathcal{D}, m)$	the estimation of the model parameters
Epistemic - Structural	$P(m \mathcal{D})$	the model specification

the generalisation to out-of-distribution data and domain, robustness to noise and biases, detection of unseen structures, certification of performance with confidence intervals, etc.

## 1.2 “Classics” in Uncertainty Modelling for Medical Imaging

The importance of representing uncertainty information has long been recognised in the medical imaging community before the recent surge of deep learning methods. There is a large body of prior research on uncertainty quantification, based on traditional probabilistic machine learning and statistical techniques in a variety of medical image processing tasks such as classification, registration, segmentation and synthesis. Here I provide a brief survey of such prior research.

### Classification

Many medical decisions can be viewed as classification tasks where the information about patients are used as inputs to infer some discrete states of their health (e.g. diagnosis) and the subsequent course of treatment (e.g. prognosis). It is widely accepted that the information available to the physician about her patient and about medical relationships in general is inherently uncertain. Since the advent of computers, academics have attempted to formalise and emulate the process of such medical reasoning performed under uncertainty, in an attempt to support clinicians through systemisation and standardisation of knowledge. Ledley et al. [14] proposed in the 1950s the first decision theoretic framework to aggregate medical evidence of varying degrees of uncertainty, and derive the most likely diagnosis in a restricted setting. They introduced scoring cards for doctors to represent the presence of discriminative symptoms, and designed an automatic system to process them to compute the conditional probability of the disease of interest. During the subsequent twenty years from then, more advanced computerised diagnostic systems were developed for different problems [15, 16, 17], however, most were still manually designed and deterministic rule-based systems without a mechanism to account for uncertainty. Adlassing et al. [18, 19] employed ‘fuzzy set theory’ to quantify and reason with uncertainty in each local decision of the if-then-else statement within such decision trees, which improved accuracy and robustness in applications in internal medicine<sup>1</sup>. Specifically, both the objective uncertainty of each decision measured in frequency of past occurrences of events, and the subjective uncertainty based on the confidence of practitioners were integrated to infer the plausibility of different diagnostic outcomes. More recently, the above approach was augmented with a variety of pre-processing steps to learn fuzzy rules automatically from data [21, 22, 23, 24, 25].

<sup>1</sup>Incidentally, a recent study [20] showed that over-confidence of doctors is also a big cause of diagnostic errors in medicine!

Another popular alternative to the above rule-based approaches is based on Bayesian networks (BNs) [26]. A typical BN-based decision model consists of a set of nodes and edges with directions, where the nodes are random variables which represent a set of symptoms and diseases, and the directed edges define the probabilistic relationships (conditional dependencies) between them, which can be more complex than the simple if-then-else conditional statements. In addition to such generality, BNs are attractive because they can naturally communicate the uncertainty in possible diagnostic alternatives as well as its underlying compositional reasoning even when data is incomplete or partially correct [27]. Some notable applications of BNs include the diagnosis of cardiovascular diseases from echocardiography [28], prediction of mental retardation in newborn babies [29], and detection of breast cancer in mammography [30, 31, 32]. A more extensive review of the existing research, including attempts to mine unknown causal structures between symptoms and diseases, are provided in Sec. 2 in Luiz et al. [33]. Some of such development have also gone beyond the realm of academic research; several computerised diagnosis (also known as computer-aided diagnosis (CAD) [34]) systems have been approved by FDA for mammographic screening and used in clinical practice [35].

## Registration

*Image registration* is the problem of determining a geometric transformation to establish spatial correspondences between images. This task plays a foundational role in numerous image-guided medical tasks [36, 37, 38] performed in both research and clinics. Firstly, it allows the spatial normalisation of multiple subjects into a common reference frame—often referred to as inter-subject registration—which is an important pre-processing step for population modeling of anatomical variability. Secondly, registration can also be used to align images of the same subject—referred to as intra-subject registration—acquired at (i) different time points or (ii) different imaging devices. The former is crucial for longitudinal studies where temporal structural or anatomical changes are examined, while the latter aims to fuse information from multiple imaging modalities to facilitate diagnosis and treatment planning.

However, specification of the optimal alignment suffers from significant uncertainty due to the ill-posed nature of the problem, presence of imaging artifacts, and the inherent variability of human anatomy. The most established approach to combat such ambiguity is probabilistic registration techniques [39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50] where the process of registration is modelled as a hierarchical generative model, and the inference yields an estimation of the posterior distributions over the parameters of both regularisation and transformation models. These techniques not only enhance the performance of registration by automatically tuning the level of regularisation on the given data, but also enable ones to derive measures of registration uncertainty from the estimated distribution of possible transformations, which can be accounted for in downstream tasks.

Probabilistic registration methods can be broadly categorised into two classes, depending on whether the space of transformation is modelled as discrete [41, 45, 49] and continuous [39, 43, 44, 46, 47, 48, 50] random variables. The associated registration uncertainty are typically quantified by summary statistics of such transformation distributions (e.g., the Shannon entropy in the discrete case [42], and the variance [43, 50], standard deviation[48] or inter-quantile range [44, 40] in the continuous case), and are visualised on top of the registered image as a surrogate measure of registration accuracy, which is difficult to measure in practice. Some works have also demonstrated the utility of

uncertainty measures in downstream tasks such as dose estimation in radiotherapy [51], segmentation of human brains [52] and classification tasks where features are derived from spatially normalised images [53]. Luo et al. [54], however, recently has pointed to a glitch in the current treatment of registration uncertainty in the existing research, highlighting the discrepancy between the uncertainty over transformations and the uncertainty over the output registered image, which is of real practical interest.

## Segmentation

*Image segmentation* has been, along with image registration, one of the main challenges in modern medical image analysis, and describes the process of assigning each pixel or voxel in images with biologically meaningful discrete labels, such as anatomical structures and tissue types (e.g. pathology and healthy tissues). The task is required in many clinical and research applications, including surgical planning [55, 56], and the study of disease progression, aging or healthy development [57, 58, 59]. However, there are often cases in practice where the correct delineation of structures is truly uncertain; this is also reflected in the well-known presence of high inter- and intra-reader variability in segmentation labels obtained from trained experts [60, 61, 62].

The vast majority of probabilistic segmentation methods model such uncertainty in pixel-wise manner, estimating the probability vector over classes in each pixel of the output. Such approaches often fall into the category of either *discriminative* or *generative approaches*, although some hybrid models also exist [63, 64, 65, 66, 67] to combine the benefits of the both worlds. In the former, the segmentation probabilities given the input image(s) is directly modelled from labelled data. Different types of classification algorithms, such as the naive Bayes classifiers, AdaBoost, probabilistic boosting trees and variants of random forests have been employed to perform probabilistic segmentation of various structures, including tumours [68, 69], pulmonary emphysema [70], MS lesions [71], organs in whole-body scans [72, 73] and the human brains [74, 75]. On the other hand, the latter approach models how the images are generated based on the underlying class information, commonly by fitting mixture distributions [76, 77, 78, 79, 80, 81, 82], with different components being associated to different classes such as tissue types. Such methods typically consist of the *prior* term that encodes the prior knowledge about where anatomical structures typically occur throughout the image, such as Markov random field models or probabilistic atlases, and the *likelihood* term that defines the distribution of image intensities given the segmentation labels. Once these terms are specified, the (*posterior*) distribution over the possible segmentations can be inferred using Bayes' rule.

One notable limitation of these methods, however, is their reliance on point estimates of parameters, which increases the risk of overfitting. In the discriminative category, many decision tree based approaches ameliorate this issue by using “ensembles” of weak classifiers [65, 66, 72, 73, 69] which implicitly accounts for parameter uncertainty and enhances the predictive performance. However, its effects on the quality of uncertainty estimate over the segmentation map have not been studied properly. In the generative category, the methods introduced above all resort to maximum likelihood (ML) or maximum a posteriori (MAP) estimation of the model parameters, without exploiting the full potential of Bayesian inference. In particular, some methods include deformable registration component that warps the probabilistic atlases into the target image domain, introducing thousands of free parameters [83, 82, 84, 85, 86]. Because many plausible ways to warp the atlases may exist, and the employed numerical optimizer may constrain the quality of reached solution, computing segmentations based on a single estimate of warp may lead

to biased results. To address such drawback, several works have attempted to model the distributions over the model parameters to obtain a better approximation of the posterior distribution. For example, Woolrich *et al.* [87] and Tian *et al.* [88] employed variants of variational inference (VI) to approximate the posterior with a more tractable family of distributions. Blaiotta *et al.* [89] recently considered a more complex, hierarchical generative model, so that the appropriate model complexity (the number of components) can also be inferred from data. Another promising approach is Markov chain Monte Carlo (MCMC) sampling; in particular, Iglesias *et al.* [90, 91] employed the Hamiltonian Monte Carlo variant to sample from the segmentation posterior and showed the derived uncertainty measure provided more meaningful error bars over volume estimates of sub-structures than the one computed from the pixel-wise segmentation probabilities. This work has been, to my knowledge, the only work amongst many other lines of Bayesian segmentation research, that considered propagating segmentation uncertainty into such downstream measurements. However, as consistent with the well-known slow convergence rate of MCMC sampling in high dimensional posterior, they report a running time of 3 hours for the hippocampus, a very small anatomical structure, indicating a need of improved efficiency for wider adoption.

## Synthesis

Medical imaging enjoys a multitude of imaging modalities (e.g. CT, T1- and T2-weighted MRI, FLAIR, DWI, etc.), which provide complementary information about the underlying anatomy. For example, CT images describe the local tissue densities, whilst diffusion weighted images quantify the directionality of tissue structures. *Image synthesis* describes the task of generating an image of a *target modality* (e.g. CT scan), without an actual acquisition, from a given image from the *source modality* (e.g. MRI). This process has been exploited for an array of different purposes such as improving multimodal image registration and tissue segmentation of MRI [92], synthesis of DTI-FA images from structured MRI [93], super-resolution [94, 95], simulation of ultrasound images from PET scans [96] for image-guided intervention, and CT-synthesis from MRI for PET attenuation-map reconstruction [97].

However, in comparison with other established problems in medical imaging such as registration and segmentation, uncertainty quantification in image synthesis has received limited attention. Most methods assume that the process of image synthesis is deterministic or at best its uncertainty is constant spatially, which not only is erroneous but also provides no meaningful measure of how reliable the synthesised image is for downstream processing. Some exceptions include [98, 99, 100], all of which have tackled the problem in different applications. Cardoso *et al.* [98] proposed a template-based generative model for the synthesis of CT images from T1 MRI data, building upon the deterministic framework of modality propagation [93]. They showed that the estimated posterior distribution over the synthetic CT given a new input T1 image and training pairs of templates, captures the presence of multiple solutions in the sinus and skull regions, where the one-to-many nature of the T1-to-CT mapping is expected from the physical properties of the modalities [97]. Cordier *et al.* [99] employed a similar generative model, but instead focused on the process of synthesising multi-modal MR contrasts (T1, T2C, T2 and FLAIR) of pathological brains given their segmentation labels. Lastly, one of my early works [100] instead took a discriminative approach and studied the importance of quantifying predictive uncertainty in the context of diffusion MR super-resolution where high-resolution images are synthesised from the low-resolution counterparts. To be specific, we intro-

duced a locally Bayesian variant of decision trees, and showed that the derived estimate of predictive uncertainty not only correlates well with the reconstruction accuracy, but also is capable of highlighting structures under-represented in the training data such as pathology.

## 1.3 Deep Learning in Medical Imaging and What's Lacking

In the last few years, deep learning techniques have continually updated the state-of-art performance in an array of medical image processing tasks, and enabled people to envision novel applications to complex problems which were previously regarded unfeasible [3, 2]. However, despite the importance of representing uncertainty in medical imaging problems and the community's long-standing recognition of it as surveyed in the previous section (Sec. 1.2), treatment of such information seems to have been largely forgotten in the present performance-centric trend in deep learning research. The more performant deep learning methods become and the more of their applications are integrated into high-stake decision making systems, the more crucial it becomes to quantify and reason with what the model does not know. In this thesis, we explore probabilistic modelling as a framework to describe uncertainty information in deep learning methods, and demonstrate utility in various medical image processing applications.

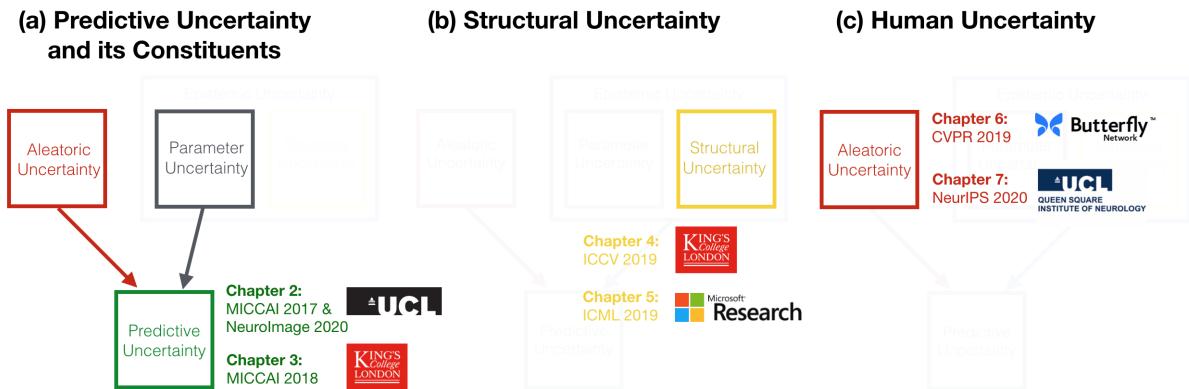


Figure 1.2: Grouping of the thesis chapters according to the types of uncertainty under study. The affiliation of the primary collaborators are also indicated.

## Overview

The thesis consists of 6 papers that touch on different types of uncertainties as summarised in Figure. 1.1. I would like to note that this thesis solely focuses on supervised discriminative models, leaving aside for simplicity other modelling problems and learning scenarios, such as generative models, unsupervised learning, and reinforcement learning for which the deep learning applications in medical imaging still remain nascent. Moreover, while the healthcare domain offers much richer modalities of data beyond images such as genomics, sensors, electronic healthcare records, etc [101], here I mostly focus on tasks of processing high-dimensional imaging data.

The thesis chapters can be broadly categorised into the following three groups as illustrated in Figure. 1.2 according to the types of uncertainties being modelled.

**Predictive Uncertainty and its Constituents:** In Chapter 2, we investigate the importance of quantifying predictive uncertainty and understanding its sources for developing a risk-averse and transparent image enhancement application for diffusion MRI. We demonstrate how a measure of predictive uncertainty can be used as a proxy for the predictive accuracy in the absence of ground-truths. Furthermore, assuming the structure of the CNN model is flexible enough for the task, we introduce a way to decompose the predictive uncertainty into its orthogonal sources i.e. aleatoric and parameter uncertainty. We show the potential utility of such decoupling in providing quantitative “explanations” into the model performance. This chapter is based on the publication [102] and its recent extension [103]:

- **R. Tanno**, D. E. Worrall, A. Ghosh, E. Kaden, S. N. Sotropoulos, A. Criminisi, and D. C. Alexander, “Bayesian image quality transfer with CNNs: Exploring uncertainty in dMRI super-resolution”. (2017) **MICCAI**
- **R. Tanno**, D. E. Worrall, E. Kaden, A. Ghosh, F. Grussu, A. Bizzi, S. N. Sotropoulos, A. Criminisi, and D. C. Alexander, “Uncertainty Quantification in Deep Learning for Safer Neuroimage Enhancement”. (2020) **Neuroimage**

In Chapter 3, based on the publication [104] below, we show that the same concept could be naturally adapted to the multi-task learning paradigm. We demonstrate the benefits in the MR-only radiotherapy treatment planning application where the synthetic CT image and the segmentation of organs at risk are simultaneously predicted from the input MR image.

- F.J.S. Bragman, **R. Tanno**, Z. Eaton-Rosen, W. Li, D. J. Hawkes, S. Ourselin, D. C. Alexander, J. R. McClelland, M. J. Cardoso, “Uncertainty in multitask learning: joint representations for probabilistic MR-only radiotherapy planning”. (2018) **MICCAI**

**Structural Uncertainty:** Motivated by the same multi-task learning application, Chapter 4 introduces a Bayesian method to learn the connectivity structures in a convolutional neural network. The methods, originally proposed in [105], aims to separate task-specific and shared features across different tasks, and thereby increase the benefit from jointly learning from multiple related tasks.

- F.J.S. Bragman\*, **R. Tanno\***, S. Ourselin, D. C. Alexander, M. J. Cardoso, “Stochastic Filter Groups for Multi-Task CNNs: Learning Specialist and Generalist Convolution Kernels”. (2019) **ICCV** (\* equal contributions)

In Chapter 5, based on the publication [106], we explore how the training algorithm of decision trees could be extended to grow the structure of a neural network architecture from simple building blocks to adapt to the given availability of data and the complexity of the task. Both chapters explore different approaches to learning meaningful sparsity in a neural network architecture.

- **R. Tanno**, K. Arulkumaran, D. C. Alexander, A. Criminisi and A. Nori, “Adaptive Neural Trees”. (2019) **ICML**

**Human Uncertainty:** In Chapter 6, we introduce a method for modelling the “measurement noise” (e.g., biases and skill levels) in the human annotation process, and integrate this information into the learning process of the neural network classifier [107]. We propose a well-grounded and practical optimisation method to learn such noise model,

and demonstrate in the classification of ultrasound cardiac images where the annotations are very noisy and sparse. Specifically, we show that the method not only improves the robustness of the model to label noise, but also yields insights into the performance of different human annotators. In Chapter 7, we extend this idea, based on our recent work [108], to the more challenging task of semantic segmentation where every pixel in the input image is classified, and demonstrate similar benefits.

- **R. Tanno**, A. Saheedi, S. Sankaranarayanan, D. C. Alexander, N. Silberman, “Learning From Noisy Labels By Regularized Estimation Of Annotator Confusion”. (2019) **CVPR**
- L. Zhang\*, **R. Tanno\***, M. Xu, C. Jin, J. Jacob, O. Ciccarelli, F. Barkhof, D. C. Alexander, “Disentangling Human Error from the Ground Truth in Segmentation of Medical Images”. (2020) **NeurIPS** (\*equal contributions)

Finally, in Chapter 8, we make overall conclusions, discuss the application of this technology and suggest directions for future research. For each chapter, we review the relevant prior works.

## Other Publications

1. **R. Tanno**, A. Ghosh, F. Grussu, E. Kaden, A. Criminisi, and D. C. Alexander, “Bayesian image quality transfer”. (2016) **MICCAI**
2. D. C. Alexander, D. Zikic, A. Ghosh, **R. Tanno**, V. Wottschel, J. Zhang, E. Kaden, T. B. Dyrby, S. N. Sotiropoulos et al., “Image quality transfer and applications in diffusion MRI”. (2017) **Neuroimage**
3. **R. Tanno**, A. Makropoulos, S. Arslan, O. Oktay, S. Mischkewitz, F. Al-Noor1, J. Oppenheimer, R. Mandegaran, B. Kainz, M. Heinrich. “AutoDVT: Joint Real-time Classification for Vein Compressibility Analysis in Deep Vein Thrombosis Ultrasound Diagnostics”. (2018) **MICCAI**
4. S. B. Blumberg, **R. Tanno**, I. Kokkinos, D. C Alexander. “Deeper Image Quality Transfer: Training Low-Memory Neural Networks for 3D Images”. (2018) **MICCAI**
5. K. Kamnitsas, D. Castro, L. Folgoc, **R. Tanno**, D. Rueckert, B. Glocker, A. Criminisi, A. Nori. “Semi-Supervised Learning via Compact Latent Space Clustering”. (2018) **ICML**
6. F.J.S. Bragman\*, **R. Tanno\***, S. Ourselin, D. C. Alexander, M. J. Cardoso, “Learning task-specific and shared representations in medical imaging”. (2019) **MICCAI** (\*equal contributions)
7. C. Sudre, B.G. Anson, S. Ingala, D. Jimenez, C. Lane, L. Haider, T. Varsavsky, **R. Tanno**, L. Smith, S. Ourselin, R. Jager, M. J. Cardoso, “Let’s agree to disagree: learning highly debatable multirater labelling”. (2019) **MICCAI**
8. S. B. Blumberg, M. Palombo, C. S. Khoo, C. Tax, **R. Tanno**, D. C Alexander. “Multi-Stage Prediction Networks for Data Harmonization”. (2019) **MICCAI**
9. K. Quan, **R. Tanno**, M. Duong, R. Shipley, M. Jones, C. Bereton, J. Hurst, D. Hawkes, J. Jacobs, “Modelling Airway Geometry as Stock Market Data using Bayesian Changepoint Detection”, (2019) **MICCAI** 10<sup>th</sup>Machine Learning in Medical Imaging Workshop
10. C. Jin, **R. Tanno**, M. Xu, T. Mertzanidou, D. C. Alexander, “Foveation for Segmentation of Mega-Pixel Histology Images”. (2020) **MICCAI**
11. L. Zhang\*, **R. Tanno\***, C. Jin, J. Jacob, O. Ciccarelli, F. Barkhof, D. C. Alexander, “Learning to Segment When Experts Disagree”. (2020) **MICCAI** (\*equal contributions)

## Patents

1. **R. Tanno**, K. Arulkmaran, A. Nori, A. Criminisi, “Neural Trees”, G.B. Microsoft Technology Licensing LLC. (2018). Patent No. GB201810736D0. (Filed in Aug 2018).
2. F. A. Noor, S. Mischkewitz, A. Makropoulos, **R. Tanno**, B. Kainz, O. Oktay, “Blood vessel obstruction diagnosis method, apparatus & system” Patent No.WO2018162888A1. (Published in Sep 2018).



# Chapter 2

## Quantifying Predictive Uncertainty and Its Sources for Super-Resolution

**Abstract:** Deep learning (DL) has shown great potential in medical image enhancement problems, such as super-resolution or image synthesis. However, to date, most existing approaches are based on deterministic models, neglecting the presence of different sources of uncertainty in such problems. Here we introduce methods to characterise different components of uncertainty, and demonstrate the ideas using diffusion MRI super-resolution. Specifically, we propose to account for *aleatoric uncertainty* through a heteroscedastic noise model and for *parameter uncertainty* through approximate Bayesian inference, and integrate the two to quantify *predictive uncertainty* over the output image. Moreover, we introduce a method to propagate the predictive uncertainty on a multi-channelled image to derived scalar parameters, and separately quantify the effects of aleatoric and parameter uncertainty therein. The methods are evaluated for super-resolution of two different signal representations of diffusion MR images—Diffusion Tensor images and Mean Apparent Propagator MRI—and their derived quantities such as mean diffusivity and fractional anisotropy, on multiple datasets of both healthy and pathological human brains. Results highlight three key potential benefits of modelling uncertainty for improving the safety of DL-based image enhancement systems. Firstly, modelling uncertainty improves the predictive performance even when test data departs from training data (“out-of-distribution” datasets). Secondly, the predictive uncertainty highly correlates with reconstruction errors, and is therefore capable of detecting predictive “failures”. Results on both healthy subjects and patients with brain glioma or multiple sclerosis demonstrate that such an uncertainty measure enables subject-specific and voxel-wise risk assessment of the super-resolved images that can be accounted for in subsequent analysis. Thirdly, we show that the method for decomposing predictive uncertainty into its independent sources provides high-level “explanations” for the model performance by separately quantifying how much uncertainty arises from the inherent difficulty of the task or the limited training examples. The introduced concepts of uncertainty modelling extend naturally to many other imaging modalities and data enhancement applications. This chapter is based on the publications [102, 103].

### 2.1 Introduction

In the last few years, deep learning techniques have permeated the field of medical image processing [3, 2]. Beyond the automation of existing radiological tasks— e.g. segmentation [109], detection [110], disease grading and classification [111]—deep learning has been successfully applied to so-called “data enhancement” problems. Data enhancement aims to improve the quality, the information content, or the quantity of medical images available for research and clinics by transforming images from one domain to another [112]. Previous research has shown the efficacy of data enhancement in different forms such as super-resolution [113, 114, 115], image synthesis [116, 117], denoising [118, 119], data harmonisation [120, 121] across scanners and protocols, reconstruction [122, 123, 124, 125, 126, 127, 128], registration [129, 130] and quality control [131, 132]. These advances have the potential not only to enhance the quality and efficiency of radiological care, but also facilitate scientific discoveries in medical research through increased volume and content of usable data.

However, most efforts in the development of data enhancement techniques have focused on improving the accuracy of deep learning algorithms, with little consideration of risk management. Blindly trusting the output of a given machine learning tool risks undetected failures e.g. spurious features and removal of structures [133]. In medical imaging applications where ultimately images can inform life-and-death decisions, it is crucial to have mechanisms for quantifying reliability of prediction (i.e. predictive uncertainty), and explaining its sources [134].

In this chapter, we introduce methods for capturing different components of uncertainty, and show its benefits in building a deep-learning based image enhancement system that is more robust, risk-averse and transparent. In particular, we focus on two types of source uncertainties introduced in Chapter 1, namely *aleatoric* and *parameter* uncertainty, which we characterise through a input-dependent (heteroscedastic) noise model [135] and variational dropout [136], respectively (see Figure 2.3). In this chapter, however, on the assumption that the model structure is sufficiently flexible, *structural* uncertainty is discounted. We then combine and propagate these two source uncertainties into a spatial map of *predictive uncertainty* over the output image, which can be used to assess the output reliability on subject-specific and voxel-wise basis. Lastly, we propose a method to propagate the predictive uncertainty to arbitrary derived quantities of the output images, such as scalar indices that are commonly used for subsequent analysis, and decompose it into distinct components which separately quantify the contributions of aleatoric and parameter uncertainty.

We demonstrate these ideas through the super-resolution application of Image Quality Transfer (IQT) [137, 100, 138, 139], a data-enhancement framework for propagating information from rare or expensive high quality images to lower quality but more readily available images. In particular, we evaluate the utility of uncertainty quantification in terms of three aspects; i) performance on unseen datasets; ii) safety assessment of system output; iii) explainability of failures. For two different types of diffusion signal representations, we evaluate the effects of uncertainty modeling on generalisation by measuring the predictive accuracy on unseen test subjects in the Human Connectome Project (HCP) dataset [140] and the Lifespan dataset [141]. We additionally test the value of improved predictive performance in a downstream tractography application. We then test the capability of the predictive uncertainty map to indicate predictive errors and thus to detect potential failures on images of both healthy subjects and those in which pathologies unseen in the training data arise, specifically from glioma and multiple-sclerosis (MS) patients. Lastly, we perform the decomposition of predictive uncertainty on HCP subjects with benign abnormalities, and assess its potential value in gaining high-level interpretations of predictive performance.

## 2.2 Background: Image Quality Transfer

Alexander *et al.* [137] proposed Image Quality Transfer (IQT), the first supervised learning based framework for data enhancement of medical images. We use this framework as the core application for evaluating the practical benefits of the proposed methods of uncertainty quantification. We first provide the general formulation of IQT. We then describe the baseline neural network architecture upon which we build to account for different components of uncertainty.

### 2.2.1 Image Enhancement as Patch-wise Regression

IQT performs data enhancement via regression of low quality against high quality image content. In order to overcome the memory demands of processing 3-dimensional medical images, along with other subsequent work such as [142, 113, 143, 144], IQT assumes factorisability over local neighbourhoods (also called patches) and models the conditional distribution of high-quality image  $I_{High}$  given the corresponding low-quality input  $I_{Low}$  as:

$$p(I_{High}|I_{Low}) = \prod_{i \in \mathcal{S}} p(\mathbf{y}_i|\mathbf{x}_i) \quad (2.1)$$

where  $\{\mathbf{y}_i\}_{i \in \mathcal{S}}$  is a set of disjoint high-quality subvolumes with  $\mathcal{S}$  denoting the set of their indices, which together constitute the whole image  $I_{High}$ , while  $\{\mathbf{x}_i\}_{i \in \mathcal{S}}$  is a set of potentially overlapping low-quality subvolumes, each of which contains and is spatially larger than the corresponding  $\mathbf{y}_i$ , as illustrated in Fig. 2.1. Here we assume that each local neighbourhood is a cubic sub-volume. The locality assumption

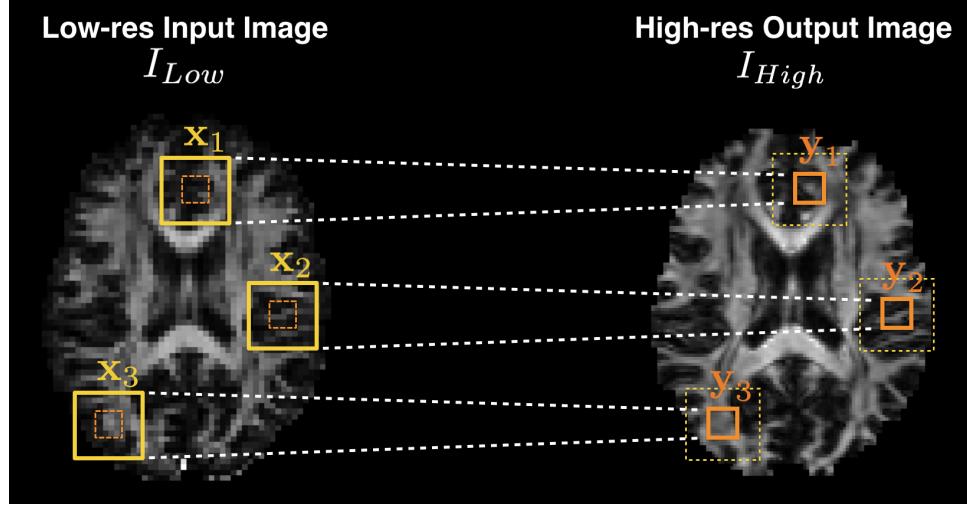


Figure 2.1: Illustration of the patch-wise regression in super-resolution application. The conditional distribution over the high quality image  $p(I_{High}|I_{Low})$  is assumed to factorise over local neighbourhoods  $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i \in \mathcal{S}}$ . In this case, for each input subvolume  $\mathbf{x}_i$  (in yellow), the high resolution version of the smaller centrally located neighbourhood,  $\mathbf{y}_i$  (in orange) is regressed.

reduces the problem of learning  $p(I_{High}|I_{Low})$  to the much less memory intensive problem of learning  $p(\mathbf{y}|\mathbf{x})$ . In other words, IQT formulates the data enhancement task as a patch-wise regression where an input low-quality image  $I_{Low}$  is split into smaller overlapping sub-volumes  $\{\mathbf{x}_i\}_{i \in \mathcal{S}}$  and the corresponding non-overlapping high-quality sub-volumes  $\{\mathbf{y}_i\}_{i \in \mathcal{S}}$  are independently predicted according to the patch regressor  $p(\mathbf{y}|\mathbf{x})$ . The final prediction for the 3D high-quality volume  $I_{high}$  is constructed by tessellating the output patches  $\{\mathbf{y}_i\}_{i \in \mathcal{S}}$ .

The original implementation of IQT [137, 138, 100] employed a variant of random forests (RFs) to model  $p(\mathbf{y}|\mathbf{x})$  while more recent [142, 113, 143, 144] approaches use variants of convolutional neural networks (CNNs). Either way, the machine learning algorithm is trained on pairs of high-quality and low-quality patches  $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$  extracted from a set of image volumes, and is used to perform the data-enhancement task of interest. Typically, such patch pairs  $\mathcal{D}$  are synthesised by down-sampling a collection of high quality images to approximate their counterparts in a particular low-quality scenario [137, 113]. In this work, we focus on the task of super-resolution (SR) where the spatial resolution of  $I_{high}$  is higher than the input image  $I_{low}$ .

## 2.2.2 Baseline Super-Resolution Model

As the baseline architecture for modelling  $p(\mathbf{y}|\mathbf{x})$ , we use a variant of convolutional neural networks (CNNs), a highly specialized form of neural network for learning image representations [1, 145]. In particular, we adapt efficient subpixel-shifted convolutional network (ESPCN) [146] to 3D data. ESPCN is a recently proposed architecture with the capacity to perform real-time per-frame super-resolution of videos while retaining high accuracy on 2D natural images. We have chosen to base on this architecture for its simplicity and computational performance.

ESPCN is a fully convolutional network, with a special *shuffling operation* on the output, which identifies individual feature channel dimensions with spatial locations in the high-resolution output. Fig. 2.2 shows a 2D illustration of an example ESPCN when the fully convolutional part of the network consists of 3 convolutional layers, each followed by a ReLU, and the final layer has  $cr^2$  feature maps where  $r$  is the upsampling rate and  $c$  is the number of channels in the output image (e.g. 6 in the case of DT images). The shuffling operation takes the feature maps of shape  $h \times w \times cr^2$  and remaps pixels from different channels into different spatial locations in the high-resolution output, producing a  $rh \times rw \times c$  image, where  $h$  and  $w$  denote height and width of the pre-shuffling feature maps. This shuffling operation in 3D is given by  $\mathcal{S}(F)_{i,j,k,c} = F_{[i/r],[j/r],[k/r],(r^3-1)c+\text{mod}(i,r)+r\cdot\text{mod}(j,r)+r^3\cdot\text{mod}(k,r)}$  where  $F$  is the pre-shuffled feature maps. The combined effects of the last convolution and shuffling is effectively a learned interpolation, and an efficient implementation of deconvolution layer [?] where the kernel size is divisible by the size of the stride [146]. Therefore, it is less susceptible to checker-board like artifacts commonly

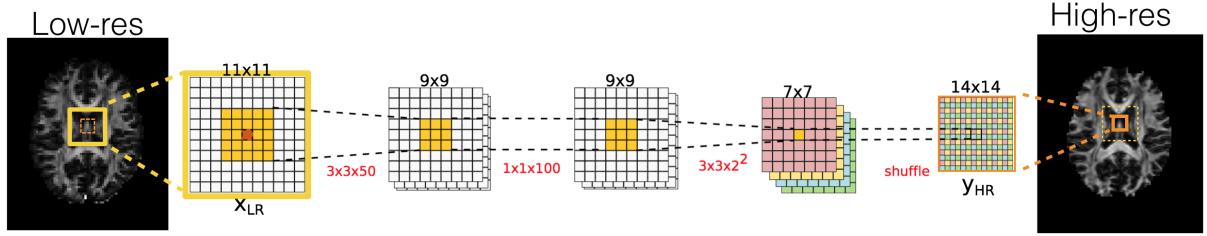


Figure 2.2: 2D illustration of an example baseline network (ESPCN [146]) with 3 convolution layers and upsampling rate,  $r = 2$ . The receptive field of the central  $2^2$  pixels in the output patch is  $5^2$  pixels in the input patch and is shown in yellow. The shuffling operation at the end periodically rearranges the final feature maps from the low-resolution space into the high-resolution space.

observed with deconvolution operations [147].

At test time, the prediction of higher resolution volume is performed through *shift-and-stitch* operation. The network takes each subvolume  $\mathbf{x}$  in a low-resolution image, and predicts the corresponding high-resolution sub-volume  $\mathbf{y}$ . By tessellating the predictions from appropriately shifted inputs  $\mathbf{x}$ , the whole high-resolution volume is reconstructed. With convolutions being local operations, each output voxel is only inferred from a local region in the input volume, and the spatial extent of this local connectivity is referred to as the *receptive field*. For a given input subvolume, the network increases the resolution of the central voxel of each receptive field e.g. the central  $2^3$  output voxels are estimated from the corresponding  $5^3$  receptive field in the input volume, as coloured yellow in Fig. 2.2.

Given training pairs of high-resolution and low-resolution patches  $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$ , we optimise the network parameters by minimising the sum of per-pixel mean-squared-error (MSE) between the ground truth  $\mathbf{y}$  and the predicted high-resolution patch  $\mu_\theta(\mathbf{x})$  over the training set. Here  $\theta$  denotes all network parameters. This is equivalent to minimising the negative log likelihood (NLL) under the Gaussian noise model  $p(\mathbf{y}|\mathbf{x}, \theta) = \mathcal{N}(\mathbf{y}; \mu_\theta(\mathbf{x}), \sigma^2 I)$  with fixed isotropic variance  $\sigma^2$ .

## 2.3 Methods

We now build on the baseline ESPCN architecture to model different types of uncertainty that arise in super-resolution. In particular, we introduce methods to model *aleatoric* and *parameter uncertainty*. We then combine the two approaches and estimate the overall uncertainty over prediction (i.e. *predictive uncertainty*) by approximating the variance of the predictive distribution (eq. (2.9)). Lastly, we propose a method for decomposing predictive uncertainty into its sources—*aleatoric* and parameter uncertainty—in an attempt to provide quantifiable explanations for the confidence on model output (eq. (2.13)).

### 2.3.1 Sources of Predictive Uncertainty

Predictive failures of deep learning systems, by and large, occur due to two reasons: i) the task itself is inherently ambiguous or ii) the learned model is not adequate to describe the data [9, 10, 102, 150], as illustrated in Fig. 2.3. The former stems from *aleatoric uncertainty* [148], which describes ambiguity in the underlying data generating process (e.g. presence of stochasticity such as measurement noise), and cannot be alleviated by increasing available training data or model complexity<sup>1</sup>. The latter is characterised by *model uncertainty*[149], which describes ambiguity in model specification<sup>2</sup>. Model uncertainty arises from a) *parameter uncertainty*: ambiguity in fitting the model to the target mapping due to limited training data, or b) *structural uncertainty*: errors due to insufficient flexibility of the model class (e.g. fitting a linear model to a sinusoidal process). These types of uncertainty can be reduced by collecting more data or specifying a different class of models. In this chapter, we make the assumption that the selected neural network architecture is expressive enough to capture well the target mapping of super-resolution given enough data, and we discount entirely the presence of structural uncertainty. Under this assumption, *aleatoric* and parameter uncertainty (Fig. 2.3) fully characterise the predictive failures

<sup>1</sup>Aleatoric uncertainty is also known as intrinsic or statistical uncertainty.

<sup>2</sup>Model uncertainty refers to *epistemic uncertainty* in the context of modelling [9].

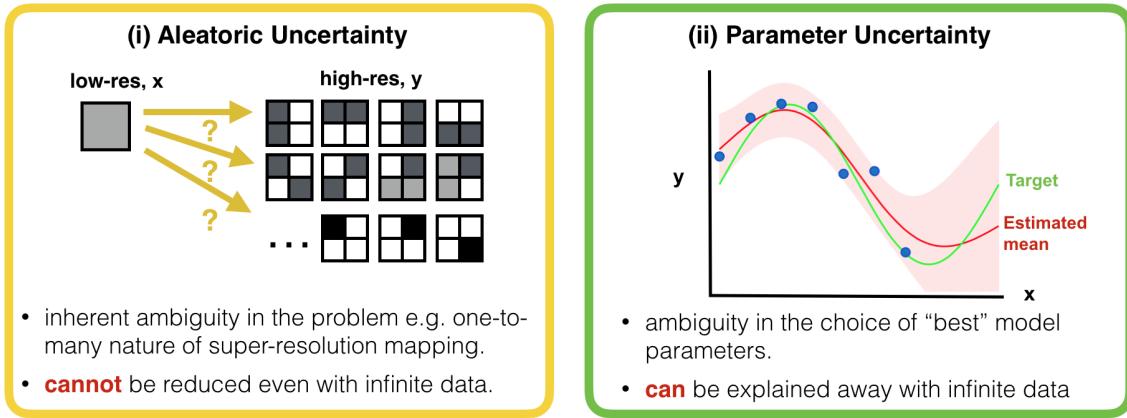


Figure 2.3: Illustration of two different types of uncertainty [9]. **Aleatoric uncertainty** [148] quantifies the degree of inherent ambiguity in the underlying problem. For example, in the case of super-resolution, there exist many possible high-resolution images  $y$  that would get mapped onto the same low-resolution input  $x$ . **Aleatoric uncertainty** is irreducible with training data. On the other hand, the parameter uncertainty [149] (a subtype of model uncertainty) arises from the finite training set. There exist more than one model that can explain the given training data equally well, and the parameter uncertainty quantifies the ambiguity in selecting the model parameters that best captures the target data-generating process. As illustrated in the figure on the right, parameter uncertainty decreases with more data; the green line shows the target function, the red line is the estimated mean, and the shaded region signifies the associated parameter uncertainty (standard deviation), which is higher in regions where we have fewer observations.

of deep learning models. Therefore, accurate estimation of these uncertainties are needed and would potentially allow practitioners to understand better the limits of the models, flag doubtful predictions, and highlight test cases that are not well represented in the training data. In the remaining part of this section, we describe our approaches to estimating these uncertainty components.

### 2.3.2 Aleatoric Uncertainty and Heteroscedastic Noise Model

*Aleatoric uncertainty* quantifies the inherent ambiguity of the underlying problem that is irreducible with data as illustrated in Fig. 2.3(i). Here we capture **aleatoric** uncertainty by estimating the variance of the target conditional distribution  $p(y|x, \theta)$ . In medical images, **aleatoric** uncertainty is often spatially and channel-wise varying. For example, super-resolution could be fundamentally harder on some anatomical structures than others due to signal variability as shown in [100]. It may also be the case that some channels of the image volume might contain more complex, non-linear and noisy signals than other channels e.g. higher order terms in diffusion signal representations. To capture such potential variation of **aleatoric** uncertainty, we model  $p(y|x, \theta)$  as a Gaussian distribution with input-dependent varying variance:

$$p(y|x, \theta_1, \theta_2) = \mathcal{N}(y; \mu(x; \theta_1), \Sigma(x; \theta_2)) = \frac{\exp\left((y - \mu(x; \theta_1))^T \Sigma^{-1}(x; \theta_2) (y - \mu(x; \theta_1))\right)}{\sqrt{(2\pi)^k \det \Sigma(x; \theta_2)}} \quad (2.2)$$

where the mean  $\mu(x; \theta_1)$  and the covariance  $\Sigma(x; \theta_2)$  are functions of input  $x$  and modelled by two separate 3D-ESPCNs (as shown in Fig. 2.4), which we refer to as “mean network” and “covariance network”, and are parametrised by  $\theta_1$  and  $\theta_2$ , respectively. We note that the input patch  $x$  varies spatially, which makes the estimated variance spatially varying and different for respective channels. Fig. 2.4 shows a 2D illustration of our 3D architecture. For each low-resolution input patch  $x$ , we use the output of the mean network  $\mu(x; \theta_1)$  at the top as the final estimate of the high-resolution ground truth  $y$  whilst the diagonal elements of the covariance  $\Sigma(x; \theta_2)$  quantify the corresponding **aleatoric** uncertainty over individual components in  $\mu(x; \theta_1)$  and over different channels. Lastly, we note that this is a specific instance of a broad class of models, called *heteroscedastic noise models* [151, 135] where the variance is a function of the value of the input. In contrast, the baseline 3D-ESPCN can be viewed as an example

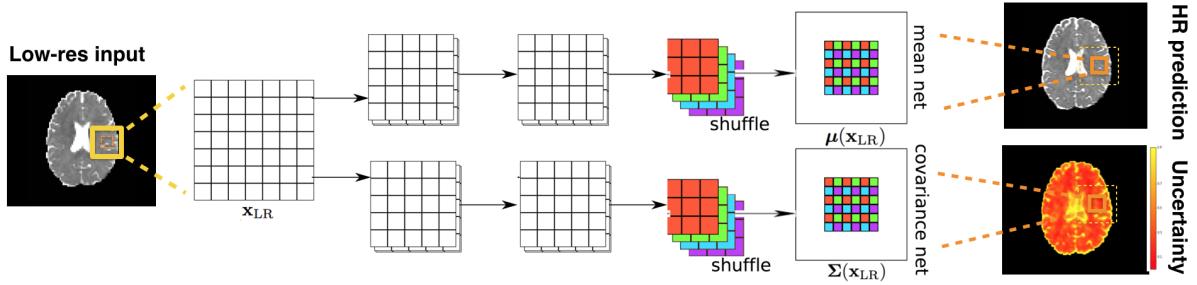


Figure 2.4: 2D illustration of the proposed dual-path architecture which estimates the mean and diagonal covariance of the Gaussian conditional distributions as functions of the input low-resolution subvolume  $\mathbf{x}$ . The “mean network”  $\mu(\cdot)$  at the top generates the high-resolution prediction, while the “covariance network”  $\Sigma(\cdot)$  at the bottom estimates the corresponding covariance matrix at the selected location in the volume. The diagonal entries of the covariance are used to quantify the **aleatoric** uncertainty. The parameters of both networks are learned by minimising the common loss function (eq. (2.5)).

of *homoscedastic noise models* with  $\mathbf{y} = \mu_\theta(\mathbf{x}) + \sigma\epsilon$ ,  $\epsilon \sim \mathcal{N}(0, I)$  with constant variance  $\sigma^2$  across all spatial locations and image channels, which is highly unrealistic in most medical images.

We jointly optimise the parameters  $\theta = \{\theta_1, \theta_2\}$  of the mean network and the covariance network by minimising the negative loglikelihood (NLL):

$$\mathcal{L}_\theta(\mathcal{D}) = \sum_{(\mathbf{x}_i, \mathbf{y}_i) \in \mathcal{D}} -\log p(\mathbf{y}|\mathbf{x}, \theta_1, \theta_2) \quad (2.3)$$

$$= \sum_{(\mathbf{x}_i, \mathbf{y}_i) \in \mathcal{D}} -\log \mathcal{N}(\mathbf{y}_i; \mu(\mathbf{x}_i; \theta_1), \Sigma(\mathbf{x}_i; \theta_2)) \quad (2.4)$$

$$= \mathcal{M}_\theta(\mathcal{D}) + \mathcal{H}_\theta(\mathcal{D}) + c \quad (2.5)$$

where  $c$  is a constant and the remaining terms are given by

$$\mathcal{M}_\theta(\mathcal{D}) = \frac{1}{N} \sum_{i=1}^N (\mathbf{y}_i - \mu(\mathbf{x}_i; \theta_1))^T \Sigma^{-1}(\mathbf{x}_i; \theta_2) (\mathbf{y}_i - \mu(\mathbf{x}_i; \theta_1)), \quad \mathcal{H}_\theta(\mathcal{D}) = \frac{1}{N} \sum_{i=1}^N \log \det \Sigma(\mathbf{x}_i; \theta_2).$$

Here  $\mathcal{M}_\theta(\mathcal{D})$  denotes the mean squared Mahalanobis distance with respect to the predictive distribution  $p(\mathbf{y}|\mathbf{x}, \theta)$ . For simplicity, in this work we assume diagonality of the covariance matrix  $\Sigma(\mathbf{x}; \theta_2)$ . This means that the Mahalanobis distance term  $\mathcal{M}_\theta(\mathcal{D})$  equates to the sum of MSEs across all pixels and channels in the output, weighted by the inverse of the corresponding variance (estimated **aleatoric** uncertainty)<sup>3</sup>. This term naturally encourages assigning low uncertainty to regions with higher MSEs, robustifying the training to noisy labels and outliers. On other other hand,  $\mathcal{H}_\theta(\mathcal{D})$  represents the mean differential entropy and discourages the spread of  $\Sigma(\mathbf{x}; \theta_2)$  from growing too large. We note that the covariance network is used to modulate the training of the mean network and quantify **aleatoric** uncertainty during inference while only the mean network generates the final prediction, requiring a single 3D-ESPCN to perform super-resolution.

### 2.3.3 Parameter Uncertainty and Variational Dropout

*Parameter uncertainty* signifies the ambiguity in selecting the parameters of the model that best describes the training data as illustrated in Fig. 2.3.(ii). The limitation of the previously introduced 3D-ESPCN baseline (Sec. 2.2.2) and its heteroscedastic extension (Sec. 2.3.2) is their reliance on a single estimate of network parameters. In many medical imaging problems, the amount of training data is modest; in such cases, this point estimate approach increases the risk of overfitting [152].

We combat this problem with a Bayesian approach. Specifically, instead of resorting to a single network of fixed parameters, we consider the (posterior) distribution over all the possible settings of network parameters given training data  $p(\theta|\mathcal{D})$ . This probability density encapsulates the parameter

<sup>3</sup>In the case of full covariance,  $\mathcal{M}_\theta(\mathcal{D})$  becomes the MSE in the basis of principle components, weighted by the corresponding eigenvalues.

uncertainty, with its spread of mass describing the ambiguity in selecting most appropriate models to explain the training data  $\mathcal{D}$ . However, in practice, the posterior  $p(\theta|\mathcal{D})$  is intractable due to the difficulty in computing the normalisation constant. We, therefore, propose to approximate  $p(\theta|\mathcal{D})$  with a simpler distribution  $q_\phi(\theta)$  [153]. Specifically, we adapt a technique called *variational dropout* [136] to convolution operations from its original version introduced for feedforward NNs.

Binary dropout [154] is a popular choice of method for approximating posterior distributions [152] with demonstrated utility in medical imaging applications [155, 142, 156, 157, 158, 159, 160]. However, typically hyper-parameters (dropout rates) need to be pre-set before the training, requiring inefficient cross-validation and thus substantially constraining the flexibility of approximate distribution family  $q_\phi(\cdot)$  (often a fixed dropout rate per layer). This limitation motivates us to use variational dropout [136] that extends such approach with a way to learn the dropout rate from data for every single weight in the network and theoretically enables a more effective approximation of the posterior distribution. Another established class of methods is stochastic gradient Markov chain Monte Carlo (SG-MCMC) method [161, 162, 163, 164]. However, in this work, we do not consider SG-MCMC methods because they remain, although unbiased, computationally inefficient due to the requirement of evaluating an ensemble of models for posterior computation, and are slow to converge for high-dimensional problems.

Variational dropout [136] employs a form of variational inference to approximate the posterior  $p(\theta|\mathcal{D})$  by a member of tractable family of distributions  $q_\phi(\theta) = \prod_{ij} \mathcal{N}(\theta_{ij}; \eta_{ij}, \alpha_{ij}\eta_{ij}^2)$  parametrised by  $\phi = \{\eta_{ij}, \alpha_{ij}\}_{ij}$ , such that Kullback-Leibler (KL) divergence  $\text{KL}(q_\phi(\theta)||p(\theta|\mathcal{D}))$  is minimised. Here,  $\theta_{ij}$  denotes an individual element in the convolution filters of CNNs as a random variable with parameters  $\alpha_{ij}$  (dropout rate) and  $\eta_{ij}$  (mean), and the posterior over the set of all weights is effectively approximated with a product of univariate Gaussian distributions. In practice, introducing a prior  $p(\theta)$  and applying Bayes' rule allow us to rewrite the minimization of the KL divergence as maximization of the quantity known as the evidence lower bound (ELBO) [153]. Here during training, we learn the variational parameters  $\phi = \{\eta_{ij}, \alpha_{ij}\}_{ij}$  by minimizing the negative ELBO (to be consistent with the NLL cost function in eq.(3)):

$$\mathcal{L}_\phi(\mathcal{D}) = \sum_{(\mathbf{x}_i, \mathbf{y}_i) \in \mathcal{D}} \left( \mathbb{E}_{q_\phi(\theta)}[-\log p(\mathbf{y}_i|\mathbf{x}_i, \theta)] + \text{KL}(q_\phi(\theta)||p(\theta)) \right) \quad (2.6)$$

An accurate approximation for the KL term for log-uniform prior  $p(\theta)$  is proposed in [165], which is employed here. On the other hand, the first term (referred to as the reconstruction term) cannot be computed exactly, thus we employ the following MC approximation by sampling  $S$  samples of network parameters from the posterior:

$$\mathbb{E}_{q_\phi(\theta)}[-\log p(\mathbf{y}|\mathbf{x}, \theta)] \approx \frac{1}{S} \sum_{s=1}^S -\log p(\mathbf{y}|\mathbf{x}, \theta^{(s)}), \quad \theta^{(s)} \sim q_\phi(\theta) \quad (2.7)$$

Adapting the local reparametrisation trick presented in [136] to a convolution operation, we derive the implementation of posterior sampling  $\theta^{(s)} \sim q_\phi(\theta)$  such that the variance of gradients over each mini-batch is low <sup>4</sup>. In practice, this amounts to replacing each standard convolution kernel with a “Bayesian” convolution, which proceeds as follows. Firstly, we define two separate convolution kernels:  $\eta \in \mathbb{R}^{c \times k^2}$  (“mean” kernels) and  $\alpha \odot \eta^2 \in \mathbb{R}^{c \times k^2}$  (“variance” kernels) where  $\odot$  denotes the element-wise multiplication,  $c$  is the number of input channel and  $k$  is the kernel width. Input feature maps  $F_{\text{in}}$  and its elementwise squared values are convolved by respective kernels to compute the “mean” and “variance” of the output feature maps  $\mu_Y \triangleq F_{\text{in}} \star \eta$  and  $\sigma_Y^2 \triangleq F_{\text{in}}^2 \star (\alpha \odot \eta^2)$ . Lastly, the final output feature maps  $F_{\text{out}}$  are computed by drawing a sample from  $\mathcal{N}(\mu_Y, \sigma_Y^2)$  i.e. computing the following quantity:

$$F_{\text{out}} \triangleq \mu_Y + \sigma_Y \odot \epsilon, \quad \epsilon \sim \mathcal{N}(0, I). \quad (2.8)$$

Every forward pass (i.e. computation of each  $p(\mathbf{y}|\mathbf{x}, \theta^{(s)})$ ) with variational dropout is thus performed via a sequence of Bayesian convolutions. Since the injected Gaussian noise  $\epsilon$  is independent of the variational parameters  $\phi = \{\eta_{ij}, \alpha_{ij}\}_{ij}$ , the approximate reconstruction term in eq. 2.7 is differentiable with respect to them [166].

---

<sup>4</sup>See the proof for feedforward networks given in [136] which generalises to convolutions

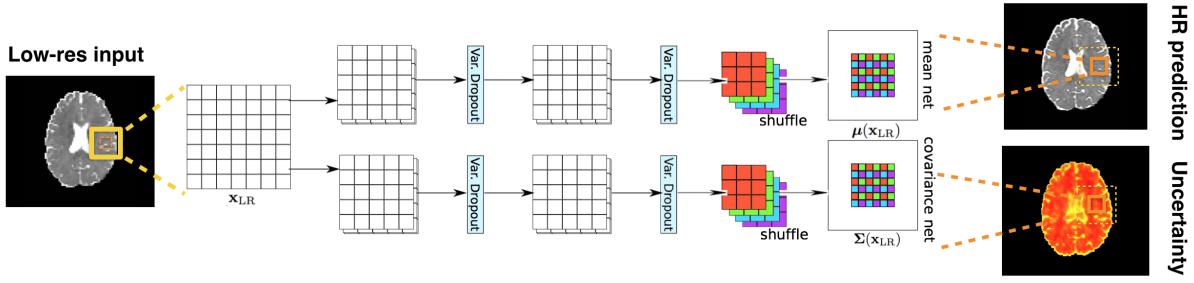


Figure 2.5: 2D illustration of a heteroscedastic network with variational dropout. Diagonal covariance is again assumed. The top 3D-ESPCN estimates the mean and the bottom one estimates the covariance matrix of the likelihood. Variational dropout is applied to feature maps after every convolution where Gaussian noise is injected into feature maps  $F_{\text{out}} = \mu_Y + \sigma_Y \odot \epsilon$  where  $\epsilon \sim \mathcal{N}(0, I)$  (see eq. 2.8).

### 2.3.4 Joint Modelling of Aleatoric and Parameter Uncertainty

We now describe how to combine the methods for modelling aleatoric and parameter uncertainty. Operationally, we take the dual architecture (Fig. 2.4) used to model aleatoric uncertainty, and apply variational dropout to every convolution layer in it. The aleatoric uncertainty is modelled in the heteroscedastic Gaussian model  $p(\mathbf{y}|\mathbf{x}, \theta_1, \theta_2) = \mathcal{N}(\mathbf{y}; \mu(\mathbf{x}; \theta_1), \Sigma(\mathbf{x}; \theta_2))$  while the parameter uncertainty is captured in the approximate posterior  $q_\phi(\theta_1, \theta_2) \approx p(\theta_1, \theta_2 | \mathcal{D})$  obtained from variational dropout.

At test time, for each low-resolution input subvolume  $\mathbf{x}$ , we would like to compute the predictive distribution  $p(\mathbf{y}|\mathbf{x}, \mathcal{D})$  over the high-resolution output  $\mathbf{y}$ . We approximate this quantity by  $q_\phi^*(\mathbf{y}|\mathbf{x})$  by taking the ‘‘average’’ of all possible network predictions  $p(\mathbf{y}|\mathbf{x}, \theta) = \mathcal{N}(\mathbf{y}; \mu(\mathbf{x}; \theta_1), \Sigma(\mathbf{x}; \theta_2))$  from all settings of the parameters  $\theta_1, \theta_2$ , weighted by the associated approximate posterior distribution  $q_\phi(\theta_1, \theta_2)$ . More formally, we need to compute the integral below:

$$q_\phi^*(\mathbf{y}|\mathbf{x}) \triangleq \underbrace{\int \mathcal{N}(\mathbf{y}; \mu(\mathbf{x}; \theta_1), \Sigma(\mathbf{x}; \theta_2))}_{\text{Network prediction}} \cdot \underbrace{q_\phi(\theta_1, \theta_2)}_{\text{Approx. posterior}} d\theta_1 d\theta_2 \quad (2.9)$$

$$\approx \int p(\mathbf{y}|\mathbf{x}, \theta_1, \theta_2) \cdot p(\theta_1, \theta_2 | \mathcal{D}) d\theta_1 d\theta_2 = p(\mathbf{y}|\mathbf{x}, \mathcal{D}) \quad (2.10)$$

where the last line represents the true predictive distribution  $p(\mathbf{y}|\mathbf{x}, \mathcal{D})$  which is estimated by our model  $q_\phi^*(\mathbf{y}|\mathbf{x})$ . However, in practice, the integral  $q_\phi^*(\mathbf{y}|\mathbf{x})$  cannot be evaluated in closed form because the likelihood  $\mathcal{N}(\mathbf{y}; \mu(\mathbf{x}; \theta_1), \Sigma(\mathbf{x}; \theta_2))$  is a highly non-linear function of input  $\mathbf{x}$  as given in eq. 2.2. At test time, we therefore estimate, for each input  $\mathbf{x}$ , the mean and covariance of the approximate predictive distribution  $q_\phi^*(\mathbf{y}|\mathbf{x})$  with the unbiased Monte Carlo estimators:

$$\hat{\mu}_{\mathbf{y}|\mathbf{x}} \triangleq \frac{1}{T} \sum_{t=1}^T \mu(\mathbf{x}; \theta_1^t) \xrightarrow{T \rightarrow \infty} \mathbb{E}_{q_\phi^*(\mathbf{y}|\mathbf{x})}[\mathbf{y}] \quad (2.11)$$

$$\hat{\Sigma}_{\mathbf{y}|\mathbf{x}} \triangleq \frac{1}{T} \sum_{t=1}^T \left( \Sigma(\mathbf{x}; \theta_2^t) + \mu(\mathbf{x}; \theta_1^t) \mu(\mathbf{x}; \theta_1^t)^T \right) - \hat{\mu}_{\mathbf{y}|\mathbf{x}} \hat{\mu}_{\mathbf{y}|\mathbf{x}}^T \xrightarrow{T \rightarrow \infty} \text{cov}_{q_\phi^*(\mathbf{y}|\mathbf{x})}[\mathbf{y}, \mathbf{y}] \quad (2.12)$$

where  $\{(\theta_1^t, \theta_2^t)\}_{t=1}^T$  are samples of the network parameters (i.e. convolution kernels) drawn from the approximate posterior  $q_\phi(\theta_1, \theta_2)$ . In other words, the inference performs  $T$  stochastic forward passes at test time by injecting noise into features according to eq. 2.8, and amalgamates the corresponding network outputs to compute the sample mean  $\hat{\mu}_{\mathbf{y}|\mathbf{x}}$  and sample covariance  $\hat{\Sigma}_{\mathbf{y}|\mathbf{x}}$ . We use the sample mean  $\hat{\mu}_{\mathbf{y}|\mathbf{x}}$  as the final prediction of an high-resolution output patch  $\mathbf{y}$  and use the diagonal elements of the sample covariance  $\hat{\Sigma}_{\mathbf{y}|\mathbf{x}}$  to quantify the corresponding uncertainty, which we refer to as *predictive mean* and *predictive uncertainty*, respectively.

### 2.3.5 Uncertainty Decomposition and Propagation

Predictive uncertainty arises from the combination of two source effects, namely aleatoric and parameter uncertainty, for which we have previously introduced methods for estimation. Lastly, we introduce a

method based on variance decomposition for disentangling these effects and quantifying their contributions separately in predictive uncertainty. We consider such decomposition problem in the presence of an arbitrary transformation of the output variable  $\mathbf{y}$ .

The users of super-resolution algorithms are often interested in the quantities that are derived from the predicted high-resolution images, rather than the images themselves. For example, quantities such as the principal direction (first eigenvalue of the DT), mean diffusivity (MD) and fractional anisotropy (FA) are typically calculated from diffusion tensor images (DTIs) and used in the downstream analysis (see [138] for the details of the equations). We therefore consider a generic function<sup>5</sup>  $g : \mathcal{Y} \rightarrow \mathbb{R}^m$  which transforms the high-resolution multi-channel data  $\mathbf{y}$  to a quantity of interest e.g. MD and FA maps, and propose a way to propagate the predictive uncertainty over  $\mathbf{y}$  to the transformed domain (i.e. compute the variance of  $p(g(\mathbf{y})|\mathcal{D}, \mathbf{x})$ ) and decompose it into the “aleatoric” and “parameter” components. Specifically, by using the law of total variance [167], we perform the following decomposition:

$$\mathbb{V}_{p(\mathbf{y}|\mathbf{x}, \mathcal{D})}[g(\mathbf{y})] = \Delta_m(g(\mathbf{y})) + \Delta_i(g(\mathbf{y})) \quad (2.13)$$

where the respective component terms are given by:

$$\Delta_m(g(\mathbf{y})) = \mathbb{E}_{p(\theta|\mathcal{D})}[\mathbb{V}_{p(g(\mathbf{y})|\theta, \mathbf{x}, \mathcal{D})}[g(\mathbf{y})] - \mathbb{V}_{p(g(\mathbf{y})|\theta, \mathbf{x}, \mathcal{D})}[g(\mathbf{y})|\theta]] \quad (2.14)$$

$$= \underbrace{\mathbb{V}_{p(\theta|\mathcal{D})}[\mathbb{E}_{p(g(\mathbf{y})|\theta, \mathbf{x}, \mathcal{D})}[g(\mathbf{y})|\theta]]}_{\text{propagated parameter uncertainty}} \quad (2.15)$$

$$\Delta_i(g(\mathbf{y})) = \underbrace{\mathbb{E}_{p(\theta|\mathcal{D})}[\mathbb{V}_{p(g(\mathbf{y})|\theta, \mathbf{x}, \mathcal{D})}[g(\mathbf{y})|\theta]]}_{\text{propagated aleatoric uncertainty}} \quad (2.16)$$

We refer to the components  $\Delta_m(g(\mathbf{y}))$  and  $\Delta_i(g(\mathbf{y}))$  as “propagated” parameter and aleatoric uncertainty. Intuitively, the first term quantifies the difference in variance between the cases where we have variable parameters and fixed parameters. In other words, this quantifies how much predictive uncertainty on the derived quantity arises, on average, from the variability in parameters. The second term on the other hand quantifies the average variance of the model prediction when the parameters are fixed, which signifies the model-independent uncertainty due to data i.e. aleatoric uncertainty. Assuming that the considered neural network is identifiable<sup>6</sup> and sufficiently complex to capture the underlying data generating process, as the amount of training data increases, the posterior  $p(\theta|\mathcal{D})$  tends to a Dirac delta function and thus the first term diminishes to zero while the second term remains. A similar variance decomposition technique was employed in [168] to understand how the variation in cell signals of interest (e.g. gene expression) in a bio-chemical network is caused by the fluctuations of other environmental variables (e.g. transcription rate and biological noise). In our case, we employ the variance decomposition technique to separate the effects of network parameters from the aleatoric uncertainty in the prediction of  $g(\mathbf{y})$ .

We first consider a special case where the transform  $g$  is an identify map i.e.  $g(\mathbf{y}) = \mathbf{y}$ . Assuming the likelihood is modelled by a Gaussian distribution with heteroscedastic noise i.e.  $p(\mathbf{y}|\theta_1, \theta_2, \mathbf{x}, \mathcal{D}) = \mathcal{N}(\mathbf{y}; \mu(\mathbf{x}; \theta_1), \Sigma(\mathbf{x}; \theta_2))$ , then we can show that the parameter and aleatoric uncertainty are given by

$$\Delta_m(\mathbf{y}) = \mathbb{V}_{p(\theta_1|\mathcal{D})}[\mu_{\theta_1}(\mathbf{x})], \quad \Delta_i(\mathbf{y}) = \mathbb{E}_{p(\theta_2|\mathcal{D})}[\Sigma_{\theta_2}(\mathbf{x})] \quad (2.17)$$

which can be approximated by the components of the MC variance estimator in eq. (2.12) :

$$\widehat{\Delta}_m(\mathbf{y}) = \frac{1}{T} \sum_{t=1}^T \mu(\mathbf{x}; \theta_1^t) \mu(\mathbf{x}; \theta_1^t)^T - \hat{\mu}_{\mathbf{y}|\mathbf{x}} \hat{\mu}_{\mathbf{y}|\mathbf{x}}^T \quad (2.18)$$

$$\widehat{\Delta}_i(\mathbf{y}) = \frac{1}{T} \sum_{t=1}^T \Sigma(\mathbf{x}; \theta_2^t) \quad (2.19)$$

---

<sup>5</sup>We assume here that the transform  $g$  is a measurable function with well-defined expectation and variance.

<sup>6</sup>We note that a neural network is, in general, not identifiable i.e. there exist more than a single set of parameters that capture the same target distribution  $p(g(\mathbf{y})|\mathbf{x})$ . In such cases, the posterior distribution  $p(\theta|\mathcal{D})$  does not collapse to a single Dirac Delta function with infinite amount of observations—it rather converges to a mixture of all sets of network parameters  $\Theta$  such that  $p(g(\mathbf{y})|\theta^*, \mathbf{x}) = p(g(\mathbf{y})|\mathbf{x}) \forall \theta^* \in \Theta$ . However, the expectation  $\mathbb{E}_{p(g(\mathbf{y})|\theta, \mathbf{x}, \mathcal{D})}[g(\mathbf{y})|\theta]$  is the same for all  $\theta \in \Theta$  and thus the propagated parameter uncertainty  $\Delta_m(g(\mathbf{y}))$  converges to zero.

where  $\{(\theta_1^t, \theta_2^t)\}_{t=1}^T$  are drawn from the approximate posterior  $q_\phi(\theta_1, \theta_2)$ .

More generally, when the transform  $g$  is complicated, MC sampling provides an alternative implementation. Given samples of model parameters  $\{\theta_t\}_{t=1}^T \sim q(\theta|\mathcal{D})$  and  $\{g_j^t\}_{j=1}^J \sim p(g(\mathbf{y})|\theta_t, \mathbf{x}, \mathcal{D})$  for  $t = 1, \dots, T$ , we estimate both the propagated parameter and aleatoric uncertainty as follows:

$$\hat{\Delta}_m(g(\mathbf{y})) \triangleq \frac{1}{T} \sum_t (\hat{\mu}^t)^2 - \left( \frac{1}{(J-1)T} \sum_{j,t} (g_j^t) \right)^2 \quad (2.20)$$

$$\hat{\Delta}_i(g(\mathbf{y})) \triangleq \frac{1}{(J-1)T} \sum_{j,t} (g_j^t)^2 - \frac{1}{T} \sum_t (\hat{\mu}^t)^2 \quad (2.21)$$

$$\hat{\mu}^t = \frac{1}{J} \sum_j g_j^t. \quad (2.22)$$

These estimators are, although unbiased, higher in variance than the case where  $g$  is the identity (eq. (2.18) and eq. (2.19)), due to two sources of sampling, thus requiring more samples for reliable estimation of respective uncertainty components.

## 2.4 Related works

Here we provide a review of related works under several different themes. We start by reviewing the development of learning-based image enhancement methods in medical imaging applications. We then discuss the recent advances made to model and quantify uncertainty in such image enhancement applications based on deep learning techniques. Next we describe the existing strands of research in uncertainty modelling for other medical imaging problems and fields of applications. Lastly, we briefly survey the theoretical development of probabilistic deep learning research.

**Learning-based Image Enhancement:** Various forms of image enhancement can be cast as image transformation problems where the input image from one domain is mapped to an output image from another domain. Numerous recent methods have proposed to perform image transformation tasks as supervised regression of low quality against high quality image content. Alexander *et al.* [137] proposed Image Quality Transfer (IQT), a general framework for supervised quality enhancement of medical images. They demonstrated the efficacy of their method through a random forest (RF) implementation of super-resolution (SR) of brain diffusion tensor images and estimation of advanced microstructure parameter maps from sparse measurements. More recently, deep learning, typically in the form of convolutional neural networks (CNNs), has shown additional promise in this kind of task. For example, Oktay *et al.* [113] proposed a CNN model to upsample a stack of 2D MRI cardiac volumes in the through-plane direction, where the SR mapping is learnt from 3D cardiac volumes of nearly isotropic voxels. This work was later extended by [144] with the addition of global anatomical prior based on auto-encoder. Zhao *et al.* [169] proposed a solution to the same SR problem for brains that utilises the high frequency information in in-plane slices to super-resolve in the through-plane direction without requiring external training data. In addition, a range of different architectures of CNNs have been considered for SR of other modalities and anatomical structures such as structural MRI [114] of brains, retinal fundus images [170] and computer tomography (CT) scans of chest [171]. Another problem of growing interest is image synthesis, which aims to synthesise an image of a different modality given the input image. Nie *et al.* [172] employed a conditional generative adversarial network to synthesise CT from MRI with fine texture details whilst Wolterink *et al.* [173] extended this idea using a CycleGAN [174] to leverage the abundance of unpaired training sets of CT and MR scans. In [143], a variant of CNN was applied to predict 7T images from 3T MRI, where both contrast and resolution are enhanced. Another notable application is the harmonisation of diffusion MRIs [120, 121, 139, 175] where images acquired at different scanners or magnetic field strengths are mapped to the common reference image space to allow for joint analysis.

**Uncertainty Quantification in Image Enhancement:** Despite this advancement, all of these methods commit to a single prediction and lack a mechanism to communicate uncertainty in the output image. In medical applications where images can ultimately inform life-and-death decisions, quantifying reliability of output is crucial. Tanno *et al.* [100] aimed to address this problem for supervised image enhancement for the first time by proposing a Bayesian variant of random forests to quantify uncertainty over predicted high-resolution MRI. They showed that the uncertainty measure correlates well with the accuracy and can highlight abnormality not represented in the training data. In our preliminary work

[102], we made an initial attempt to extend this approach with probabilistic deep-learning formulation, and showed that modelling different components of uncertainty—**aleatoric** and parameter uncertainty—allows one to build a more generalisable model and quantify predictive confidence. Kendall *et al.* [150] concurrently investigated the same problem in computer vision, suggesting its utility for safety-critical applications such as self-driving cars. More recently, Hu *et al.* [176] extended these works in the context of medical image segmentation and proposed a mechanism to learn the **aleatoric** uncertainty in a supervised manner, when multiple labels are available. Dalca *et al.* [177] proposed a CNN-based probabilistic model for diffeomorphic image registration with a learning algorithm based on variational inference, and demonstrated the state-of-the-art registration accuracy on established benchmarks while providing estimates of registration uncertainty. An alternative approach is ensembling where the variance of the predictions of multiple networks is used to quantify the predictive uncertainty [178]. Schlemper *et al.* [179] proposed a novel combination of the cascaded CNN architecture and compressive sensing, equipped with a variant of ensemble techniques, which enabled robust reconstruction of highly undersampled cardiovascular diffusion MR images, and quantification of reconstruction uncertainty. Bragman *et al.* [160] studied the value of uncertainty modelling for multi-task learning in the context of MR-only radiotherapy treatment planning where the synthetic CT image and the segmentation of organs at risk are simultaneously predicted from the input MRI image.

However, these lines of research performed only limited validation of the quality and utility of uncertainty modelling. In this chapter, we formalise and extend the preliminary ideas in Tanno *et al.* [102] and provide a comprehensive set of experiments to evaluate the proposed uncertainty modelling techniques in a diverse set of datasets, which vary in demographics, scanner types, acquisition protocols or pathology. Moreover, with the exception of [102], none of the previous methods model different components of uncertainty, namely **aleatoric** and parameter uncertainty. Our method accounts for both, and provides conclusive evidence that this improves performance thanks to different regularisation effects. In addition, we propose a method to decompose predictive uncertainty over an arbitrary function of the output image (e.g. morphological measurements) into its sources, in order to provide a high-level explanation of model performance on the given input.

**Uncertainty Quantification in Other Applications:** We should also note that, although not the focus of this work, research on uncertainty modelling in deep learning techniques extend to other medical image processing tasks beyond data enhancement, such as segmentation, detection and classification. For example, Nair *et al.*, [159] demonstrated for lesion segmentation of multiple sclerosis that the voxel-wise uncertainty metrics can be used for quality control; by filtering out predictions with high uncertainty, the model could achieve higher lesion detection accuracy. A concurrent work by Eaton-Rosen *et al.* [158] showed for the task of brain tumour segmentation that the Monte Carlo (MC) sample variance from dropout [152] can be calibrated to provide meaningful error bars over estimates of tumour volumes. Similarly, [157] introduced ways to turn voxel-wise uncertainty score into structure-wise uncertainty metrics for brain parcellation task, and showed their values in performing more reliable group analysis. The uncertainty metric based on MC dropout has also shown promise in disease grading of retinal fundal images [155, 156], and more recently an extension based on test-time augmentation was introduced by [180]. An alternative approach is to train a model to predict uncertainty score directly; [181] showed that this approach is more effective when opinions from multiple experts are available for each image. Koh *et al.* [182] and Baumgartner *et al.* [183] proposed methods to generate a set of diverse and plausible segmentation proposals on a given image, capturing more realistically the high inter-reader annotation variability, which is commonly observed in medical image segmentation tasks. Lastly, [184, 107] demonstrated for the classification of mammograms and cardiac ultrasound images, respectively that modelling uncertainty and biases of individual annotators enables robust learning from noisy labels in the presence of large disagreement.

## 2.5 Data preprocessing and implementation details

### 2.5.1 Datasets

We make use of the following four diffusion MRI datasets to evaluate different benefits of the proposed technique:

- **Human Connectome Project dataset:** we use the diffusion MRI data from the WU-Minn HCP (release Q3) [185] as the source of the training datasets. The dataset enjoys very high image

resolution, signal levels and coverage of the measurement space, enabled by the combination of custom imaging, reconstruction innovations and a lengthy acquisition protocol [140]. Each subject’s data set contains 288 diffusion weighted images (DWIs) of voxel size  $1.25^3 \text{ mm}^3$  of which 18 have nominal  $b = 0$  and the three high-angular-resolution-diffusion-imaging (HARDI) shells of 90 directions have nominal b-values of 1000, 2000, and  $3000 \text{ s mm}^{-2}$  (see [140] for the full acquisition details). The data are preprocessed by correcting distortions including susceptibility-induced, eddy currents and motion as outlined in [186].

- **Lifespan dataset:** this dataset (available online at <http://lifespan.humanconnectome.org>) contains 26 subjects of much wider age range (8 – 75 years) than the main HCP cohorts (22 – 36 years), and is acquired with a shortened version of the main HCP protocol with lower resolution ( $1.5 \text{ mm}$  isotropic voxels) and only two HARDI shells, with  $b = 1000$  and  $2500 \text{ s mm}^{-2}$ . However, we also note that the protocol still leverages the special features of the HCP scanners, providing images of substantially better quality than standard sequences. We utilise this out-of-training-distribution dataset to assess the robustness of our techniques to domain shifts.
- **Prisma dataset:** two healthy male adults (29 and 33 years old respectively) were scanned twice at different image resolutions using the clinical 3T Siemens Prisma scanner in FMRIB, Oxford. Both datasets contain diffusion MRI data with 21  $b = 0$  images and three 90-direction HARDI shells, b-values of 1000, 2000, and  $3000 \text{ s mm}^{-2}$ , each for two resolutions,  $2.50 \text{ mm}$  and  $1.35 \text{ mm}$  isotropic voxels (see [138] for full acquisition details). In addition, each of these datasets also includes a standard 3D T1-weighted MPRAGE ( $1 \text{ mm}$  isotropic resolution). The Prisma scanner is less powerful than the bespoke HCP scanner and cannot achieve sufficient signal at  $1.25 \text{ mm}$  resolution, but the  $1.35 \text{ mm}$  data provides a pseudo ground-truth for IQT resolution enhancement of the  $2.5 \text{ mm}$  data.
- **Pathology dataset:** we use two separate datasets which consist of images of brain tumour (Glioma) [187] and multiple sclerosis (MS) patients, respectively. The data of each wubject with glioma contains DWIs with  $b = 700 \text{ s/mm}^2$  while the measurement of each MS patient is of  $b = 1200 \text{ s/mm}^2$ . Both datasets have isotropic voxel size  $2^3 \text{ mm}^3$ , which is closer to the image resolution of commonplace clinical scanners. We use these datasets to assess the behaviour of predictive uncertainty on images with pathological features that are not represented in the training data set.

In all the experiments, super-resolution are performed on diffusion parameter maps derived from the DWIs in the above datasets. In particular, we consider two diffusion MRI models, namely the diffusion tensor (DT) model [188] and Mean Apparent Propagator (MAP) MRI [189], where the former is the simplest and most standard diffusion parameter map, and the latter is a high-order generalisation of the former with the capacity to characterise signals from more complex tissue structures (e.g. fibre crossing regions), a requirement for successful tractography applications. We compute both of these diffusion parameter maps using the implementation from [138], which is available at <https://github.com/ucl-mig/iqt>.

We fit the DT model to the combination of  $b = 0$  images and  $b = 1000 \text{ s/mm}^2$  HARDI shell for the HCP and Lifespan datasets, and  $b = 700 \text{ s/mm}^2$  shell for the brain tumour dataset. In all cases, weighted linear least squares are employed for the fitting, taking into account the spatially varying b-values and gradient directions in the HCP dataset. On the other hand, in the case of MAP-MRI, 22 coefficients of basis functions up to order 4 are estimated via (unweighted) least squares to all three shells of the HCP, Lifespan and Prisma datasets. As noted in [138], the choice of scale parameters (see [189])  $\mu_x = \mu_y = \mu_z = 1.2 \times 10^{-3} \text{ mm}$  empirically minimises the fitting error in the HCP dataset, and is used for all datasets.

Training datasets in all experiments are constructed by artificially downsampling very high-resolution images in the HCP dataset. In particular, we employ the following downsampling procedure: (i) the raw DWIs of selected subjects are blurred by applying the mean filter of size  $r \times r \times r$  independently over channels with  $r$  denoting the upsampling rate; (ii) the DT or MAP parameters are computed for every voxel; (iii) the spatial resolution of the resultant parameter maps are reduced by taking every  $r$  pixels. A coupled library of low-resolution and high-resolution patches is then constructed by associating each patch in the downsampled DTI/MAP-MRI with the corresponding patch in the ground truth DTI or MAP-MRI. In this case, we ensure the low-resolution patch to be centrally and entirely contained within the corresponding high-resolution patch (as illustrated by the yellow and orange squares in Fig. 2.2). We then randomly select a pre-set number of patches from each subject in the training pool to create a training

Table 2.1: Details of training data for two diffusion MR signal representations, DTIs and MAP-MRIs. The first two columns from the right denote the size of the input  $\mathbf{x}$  and output patches  $\mathbf{y}$  of dimension [width, height, depth, channels] while the third and the fourth columns show the number of patch pairs ( $\mathbf{x}, \mathbf{y}$ ) extracted from each subject, and the total number of training subjects used, respectively.

Data	Size of input $\mathbf{x}$	Size of output $\mathbf{y}$	No. pairs ( $\mathbf{x}, \mathbf{y}$ ) per subject	No. subjects
DTIs	$11 \times 11 \times 11 \times 6$	$14 \times 14 \times 14 \times 6$	8000	16
MAP-MRIs	$21 \times 21 \times 21 \times 22$	$14 \times 14 \times 14 \times 22$	4000	16

dataset as detailed in Table 2.1. In addition to the 8 subjects used in the prior work [137, 100, 102], we randomly select additional 8 subjects from the HCP cohort and include them in the training subject pool. Patches are standardized channel-wise by subtracting the mean of foreground pixel intensities of the corresponding subject and dividing by its standard deviation. Moreover, since MAP-MRI datasets contain outliers due to model fitting, in large enough quantity to influence the training of the baseline 3D-ESPCN model, we remove them by clipping the voxel intensity values of the respective 22 channels separately at 0.1% and 99.9% percentiles computed over all the foreground voxels in the whole training dataset.

## 2.5.2 Network Architectures and Training

For the training of all CNN models, we minimised the associated loss function using Adam [190] for 200 epochs with initial learning rate of  $10^{-3}$  and  $\beta = [0.9, 0.999]$ , with minibatches of size 12. We hold out 50% of training patch pairs as a validation set. The best performing model was selected based on the mean-squared-error (MSE) on the validation set.

For the super-resolution of DTIs, as in [146], we use a minimal architecture for the baseline 3D-ESPCN, consisting of three 3D convolutional layers with filters  $(3^3, 50) \rightarrow (1^3, 100) \rightarrow (3^3, 6r^3)$  where  $r$  is upsampling rate and 6 is the number of channels in DTIs. As illustrated in Fig. 2.2, the dimensions of convolution filters are chosen, so each  $5^3 \cdot 6$  low-resolution receptive field patch maps to a  $r^3 \cdot 6$  high-resolution patch, which mirrors competing random forest based methods [137, 100] for a fair comparison. On the other hand, for MAP-MRI, which is a more complex image modality with 21 channels, we employ a deeper model with 6 convolution layers  $(5^3, 256) \rightarrow (3^3, 256) \rightarrow (3^3, 128) \rightarrow (3^3, 128) \rightarrow (3^3, 64) \rightarrow (3^3, 21r^3)$  prior to the shuffling operation, which expands the receptive field on each  $r^3 \cdot 21$  high-resolution patch to  $15^3 \cdot 21$  input low-resolution patch. Every convolution layer is followed by a ReLU non-linearity except the last one in the architecture, and batch-normalization [191] is additionally employed for MAP-MRI super-resolution between convolution layer and ReLU non-linearity.

The mean and variance networks in the heteroscedastic noise model introduced in Sec. 2.3.2 are implemented as two separate baseline 3D-ESPCNs of the architectures, specified above for DTIs and MAP-MRIs. Positivity of the variance is enforced by passing the output through a softplus function  $f(x) = \ln(1 + e^x)$  as in [178].

For variational dropout, we considered two flavours: Var.(I) optimises per-weight dropout rates, and Var.(II) optimises per-filter dropout rates. More formally, the “drop-out rate”  $\alpha_{ij}$  in the approximate posterior  $q_\phi(\theta_{ij}) = \mathcal{N}(\theta_{ij}; \eta_{ij}, \alpha_{ij}\eta_{ij}^2)$  is different for every element in each convolution kernel in the former while the latter has common  $\alpha_{ij}$  shared across each kernel. In preliminary analysis, we found that the number of samples per data point for estimating reconstruction term (eq. 2.7) can be set to  $S = 1$  so long as the batch size is sensibly large ( $M = 12$ ).

We also note the default training with binary and Gaussian dropout also employs  $S = 1$  [154] along with other MC variational inference methods for neural networks such as [166, 136, 192]. Variational dropout is applied to both the baseline and heteroscedastic models without changing the architectures. For both binary and Gaussian dropout modes, we incorporate the dropout operations of fixed rate  $p$  in every convolution layer of the baseline 3D-ESPCN architecture.

All models are trained on simulated datasets generated from 16 HCP subjects as detailed in Sec. 2.5.1. We also retrained the random forest models employed in [100, 138] on equivalent datasets. It takes under 60/360 mins to train a single network on DTI/MAP-MRI data on a single TITAN X GPU. All models are implemented in the TensorFlow framework [193].

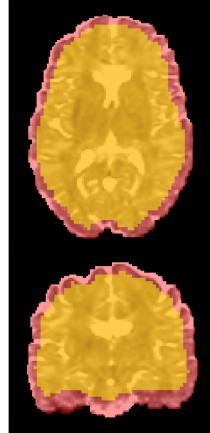


Figure 2.6: Visualisation of “interior” (yellow) and “exterior” regions (red). The interior region consists of a set of patches contained entirely within the brain while the exterior region consists of partial patches that contain mixtures of brain and background voxels

## 2.6 Results

In this section, we evaluate the proposed uncertainty modelling techniques for super-resolution of diffusion MR images. We first compare quantitatively the reconstruction performance of our probabilistic CNN models against the relevant baselines in two different types of diffusion signal representations. Secondly, we study the real-world utility of the technique in downstream tractography applications. Thirdly, we evaluate the value of predictive uncertainty as a reliability metric of output images on multiple datasets of both healthy subjects and those with unseen pathological structures such as brain tumour (Glioma) and multiple sclerosis (MS).

### 2.6.1 Benefits on Super-resolution Performance

We evaluate the prediction performance of our models for super-resolution of DTI and MAP-MRI on two datasets—HCP and Lifespan as detailed in Sec. 2.5.1. The first dataset contains 16 unseen subjects from the same HCP cohort used for training, while the second one consists of 10 subjects from the HCP Lifespan dataset. The latter tests generalisability, as they are acquired with a different protocol at lower resolution (1.5 mm isotropic), and contain subjects of a different age range (45–75 years) to the original HCP data (22–36 years). We perform  $\times 2$  upsampling in all spatial directions. The reconstruction quality is measured with root-mean-squared-error (RMSE), peak-signal-to-noise-ratio (PSNR) and mean-structural-similarity (MSSIM) [194] on two separate regions: i) “interior”; set of patches contained entirely within the brain mask; ii) “exterior”; set of patches containing some brain and some background voxels, as shown in Fig. 2.6. This is because the current state-of-the-art methods based on random forests (RFs) such IQT-RF [138] and BIQT-RF [100] are only trained on patches from the interior region and requires a separate procedure on the brain boundary. In addition, the estimation problem is quite different in boundary regions, but remains valuable particularly for applications such as tractography where seed or target regions are often in the cortical surface of the brain. We only present the RMSE results, but the derived conclusions remain the same for the other two metrics. Aside from the interpolation techniques, for each method an ensemble of 10 models are trained on different trainings set (generated by randomly extracting patch pairs from the common 16 HCP training subjects) and for each model, the average error metric over the test subjects are first calculated. The mean and standard deviations of such average errors are computed across the model ensemble and reported in Table 2.2 and Table 2.3.

Table 2.2 shows that our baseline achieves 8.5%/39.8% reduction in RMSE for the super-resolution of DTIs on the HCP dataset on the interior/exterior regions with respect to the best published method, BIQT-RF[100]. While the standard deviations are higher, the improvements are more pronounced in MAP-MRI super-resolution, reducing the average RMSEs by 49.6% and 63.5% on the interior and exterior regions. We note that that IQT-RF and BIQT-RF are only trained on interior patches, and super-resolution on boundary patches requires a separate *ad hoc* procedure. Despite including exterior

Models	HCP (interior)	HCP (exterior)	Life (interior)	Life (exterior)
CSpline-interpolation	$10.069 \pm n/a$	$31.738 \pm n/a$	$32.483 \pm n/a$	$49.066 \pm n/a$
$\beta$ -Spline interpolation	$9.578 \pm n/a$	$98.169 \pm n/a$	$33.429 \pm n/a$	$186.049 \pm n/a$
IQT-RF	$6.974 \pm 0.024$	$23.139 \pm 0.351$	$10.038 \pm 0.019$	$25.166 \pm 0.328$
BIQT-RF	$6.972 \pm 0.069$	$23.110 \pm 0.362$	$9.926 \pm 0.055$	$25.208 \pm 0.290$
3D-ESPCN(baseline)	$6.212 \pm 0.017$	$13.609 \pm 0.084$	$8.902 \pm 0.020$	$16.389 \pm 0.114$
+ Binary Dropout ( $p = 0.1$ )	$6.319 \pm 0.015$	$13.738 \pm 0.048$	$9.093 \pm 0.024$	$16.489 \pm 0.099$
+ Gaussian Dropout ( $p = 0.05$ )	$6.463 \pm 0.034$	$14.168 \pm 0.051$	$9.184 \pm 0.048$	$16.653 \pm 0.092$
+ Variational Dropout (I)	$6.194 \pm 0.013$	<b><math>13.412 \pm 0.041</math></b>	$8.874 \pm 0.027$	<b><math>16.147 \pm 0.051</math></b>
+ Variational Dropout (II)	$6.201 \pm 0.015$	<b><math>13.479 \pm 0.047</math></b>	$8.878 \pm 0.031$	<b><math>16.230 \pm 0.075</math></b>
+ Hetero.	$6.135 \pm 0.029$	$15.469 \pm 0.231$	$8.885 \pm 0.041$	$17.208 \pm 0.211$
+ Hetero. + Variational Dropout (I)	<b><math>6.121 \pm 0.015</math></b>	$13.591 \pm 0.051$	<b><math>8.837 \pm 0.043</math></b>	$16.261 \pm 0.053$
+ Hetero. + Variational Dropout (II)	<b><math>6.116 \pm 0.013</math></b>	$13.622 \pm 0.099$	<b><math>8.861 \pm 0.031</math></b>	$16.387 \pm 0.098$

Table 2.2: Super-resolution results on diffusion tensor images (DTIs) of HCP and Lifespan datasets for different upsampling methods. For each method, an ensemble of 10 models are trained on different training sets generated by randomly extracting a set of patch pairs from the common 16 HCP subjects. For each model, the average RMSE ( $\times 10^{-4} \text{mm}^2/\text{s}$ ) over subjects in respective datasets is first computed and the mean/std of such average RMSE over the ensemble are then reported. Best results in red, and the second best in blue. **In addition, the best results are shown in bold when they are statistically better than the second best ( $p < 0.05$ )**

patches in training our model, which complicates the learning task, the baseline CNN out-performs the RF methods on both regions. We see similar improvements in the out-of-distribution Lifespan dataset.

Reconstruction is faster than the RF baselines; the 3D-ESPCN is capable of estimating the whole high-resolution DTI/MAP-MRI under 10/60 seconds on a CPU and 1/10 second(s) on a GPU. On the other hand, BIQT-RF takes  $\sim 10$  mins with 8 trees on both DTIs and MAP-MRIs. The fully convolutional architecture of the model enables to process input patches of different size from that of training inputs, and we achieve faster reconstruction by using larger input patches of dimension  $25^3 \cdot c$  where  $c$  is the number of channels. We also note that the reconstruction time of the variational dropout based models increases by a factor of the number of MC samples used at test time, although it is possible, with more memory, to leverage GPU parallelisation by making multiple copies of each input patch and treating them as a mini-batch. On the other hand, the heteroscedastic CNN enjoys the same inference speed of the baseline since only the mean network is used for reconstruction (the covariance network is only employed to quantify the estimated **aleatoric** uncertainty).

Table 2.2 shows that, on both HCP and Lifespan data, modelling both **aleatoric** and parameter uncertainty (i.e. Hetero. + Variational Dropout (I), (II)) achieves the best reconstruction accuracy in DTI super-resolution. We observe that modelling **aleatoric** uncertainty with the heteroscedastic network on its own further reduces the average RMSE of the baseline 3D-ESPCN on the interior region with high statistical significance ( $p < 10^{-3}$ ). However, poorer performance is observed on the exterior than the baseline. On the other hand, using 200 MC weight samples, we see modelling parameter uncertainty with variational dropout (see Variational Dropout.(I)-CNN) performs best on both datasets on the exterior region. Combination of heteroscedastic model and variational dropout (i.e. Hetero. + Variational Dropout (I) or (II)) leads to the top 2 performance on both datasets on the interior region and reduces errors on the exterior to the level comparable or better than the baseline.

Similarly, Table 2.3 shows that the best performance in MAP-MRI super-resolution comes from the combined models (i.e. Hetero.+Variational Dropout.(I) and (II)). We observe that as with the DTI case, modelling **aleatoric** uncertainty through the heteroscedastic network improves the reconstruction accuracy on the interior region, whilst the errors on the exterior are increased with respect to the baseline 3D-ESPCN. Moreover, the improvement is pronounced when the outliers due to model fitting errors are not removed in the training data. In this case, we see that the reconstruction accuracy of 3D-ESPCN dramatically decreases, whilst in contrast it is only marginally compromised when equipped with the heteroscedastic noise model, displaying robustness to outliers. Lastly, we note that the top-2 accuracy are consistently achieved by the joint modelling of **aleatoric** and parameter uncertainty (i.e. Hetero.+Variational Dropout.(I) and (II)) on both the interior and exterior regions on both HCP and Lifespan datasets.

The performance difference of heteroscedastic network between the interior and the exterior region

Models	HCP (interior)	HCP (exterior)	Life (interior)	Life (exterior)
CSpline interpolation	$5.234 \pm \text{n/a}$	$30.362 \pm \text{n/a}$	$7.135 \pm \text{n/a}$	$29.232 \pm \text{n/a}$
$\beta$ -Spline interpolation	$4.852 \pm \text{n/a}$	$63.446 \pm \text{n/a}$	$6.523 \pm \text{n/a}$	$56.937 \pm \text{n/a}$
IQT-RF [138]	$4.538 \pm 0.113$	$25.541 \pm 0.131$	$5.882 \pm 0.121$	$26.137 \pm 0.279$
BIQT-RF [100]	$4.838 \pm 0.129$	$25.523 \pm 0.175$	$5.949 \pm 0.131$	$27.509 \pm 0.233$
3D-ESPCN(baseline)	$2.285 \pm 0.126$	$9.316 \pm 0.127$	$4.195 \pm 0.163$	$11.922 \pm 0.192$
+ Binary Dropout ( $p = 0.1$ )	$2.283 \pm 0.154$	$9.272 \pm 0.132$	$4.120 \pm 0.178$	$11.652 \pm 0.204$
+ Gaussian Dropout ( $p = 0.1$ )	$2.370 \pm 0.155$	$9.335 \pm 0.144$	$4.327 \pm 0.157$	$11.907 \pm 0.211$
+ Variational Dropout (I)	$2.155 \pm 0.122$	$9.205 \pm 0.193$	$3.997 \pm 0.153$	$11.547 \pm 0.177$
+ Variational Dropout (II)	$2.172 \pm 0.128$	$9.112 \pm 0.173$	$3.972 \pm 0.132$	$11.511 \pm 0.172$
+ Hetero.	$1.998 \pm 0.132$	$11.294 \pm 0.216$	$3.872 \pm 0.140$	$12.084 \pm 0.129$
+ Hetero + Variational Dropout (I)	$1.951 \pm 0.122$	$9.102 \pm 0.181$	$3.572 \pm 0.171$	$11.037 \pm 0.192$
+ Hetero + Variational Dropout (II)	$1.969 \pm 0.119$	$9.052 \pm 0.162$	$3.606 \pm 0.141$	$11.311 \pm 0.195$
3D-ESPCN(without outlier removal)	$3.425 \pm 0.163$	$13.284 \pm 0.239$	$6.032 \pm 0.229$	$15.513 \pm 0.273$
+ Hetero.	$2.264 \pm 0.153$	$11.306 \pm 0.172$	$3.919 \pm 0.140$	$12.821 \pm 0.150$
+ Hetero + Variational Dropout (I)	$2.138 \pm 0.159$	$10.022 \pm 0.187$	$3.681 \pm 0.193$	$12.133 \pm 0.205$
+ Hetero + Variational Dropout (II)	$2.133 \pm 0.188$	$9.988 \pm 0.209$	$3.690 \pm 0.184$	$12.052 \pm 0.212$

Table 2.3: Super-resolution results on MAP-MRIs of HCP and Lifespan datasets for different upsampling methods. For each method, an ensemble of 5 models are trained on different training sets generated by randomly extracting a set of patch pairs from the common 16 HCP subjects. For each model, the average RMSE over subjects in respective datasets is first computed and the mean/std of such average RMSEs over the ensemble are then reported. Best results in red, and the second best in blue. In addition, the performance of 3D-ESPCN and its probabilistic variants trained on data without outlier removal are also included. **The best results are also shown in bold when they are statistically better than the second best ( $p < 0.05$ ).**

roots from the loss function. The Mahalanobis term  $\mathcal{M}_\theta(\mathcal{D})$  in eq.(2.5) imposes a larger penalty on the regions with smaller aleatoric uncertainty. The network therefore allocates less of its resources towards the regions with higher uncertainty (e.g. boundary regions) where the statistical mapping from the low-resolution to high-resolution space is more ambiguous, and biases the model to fit the regions with lower uncertainty. However, we note that the performance of the heteroscedastic network is still considerably better than the standard interpolation and RF-based methods. By augmenting the model with variational dropout, the exterior error of the heteroscedastic model is dramatically reduced, indicating its regularisation effect against overfitting to low-uncertainty areas. We also observe concomitant performance improvement on the interior regions on both datasets, which additionally shows the benefits of such regularisation even in low-uncertainty areas.

Both Table 2.2 and Table 2.3 show that the use of variational dropout attains lower errors than the models with fixed dropout probabilities  $p$ , namely, Binary and Gaussian dropout [154]. Different instances of both dropout models are trained for a range of  $p$  by linearly increasing on the interval  $[0.05, 0.3]$  with increment 0.05, and the test errors for the configurations with smallest RMSE on the validation set are reported in Table 2.2 and Table 2.3. As with variational dropout models, 200 MC samples are used for inference. In all cases, two variants of variational dropout (I) and (II) outperform the networks with the best binary or Gaussian dropout models, showing the benefits of learning dropout probabilities  $p$  rather than fixing them in advance.

## 2.6.2 Reliability Assessment of Model Predictions

In this section, we investigate the value of uncertainty modelling in enhancing the safety of super-resolution system beyond reduced reconstruction errors. Firstly, we study the utility of predictive uncertainty map as a proxy measure of reconstruction accuracy on healthy test subjects from both HCP and Lifespan datasets. Secondly, we look into the behaviour of uncertainty maps in the presence of abnormal features that are not present in the training data.

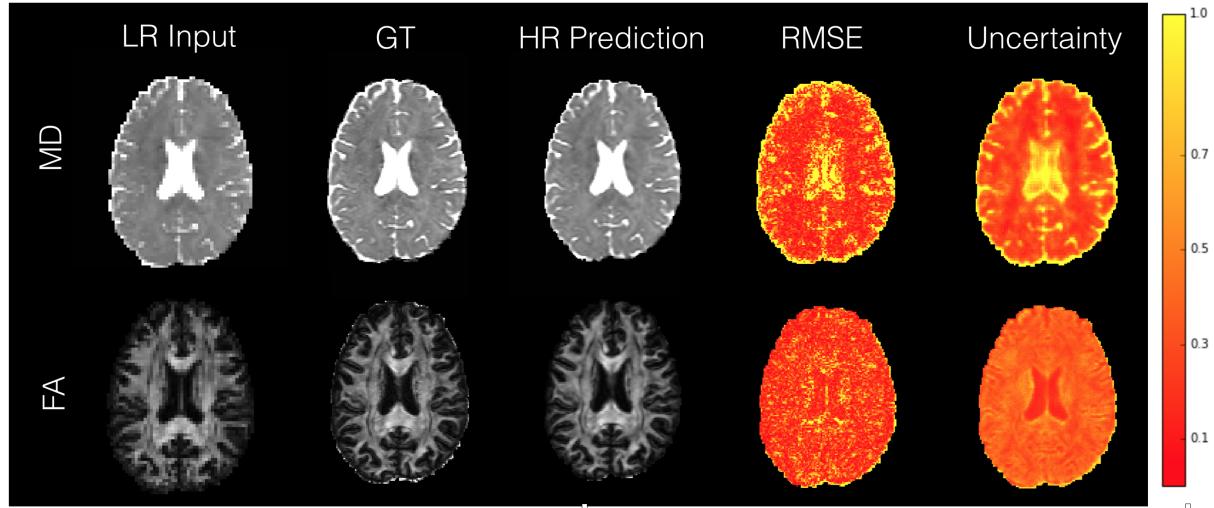


Figure 2.7: Comparison between voxel-wise RMSE and predictive uncertainty maps for FA and MD computed on a HCP test subject (min-max normalised for MD and FA separately). Low-res input, ground truth and the mean of high-resolution predictions are also shown.

### Healthy Test Subjects

We employ the most performant CNN model (3D-ESPCN + Hetero. + Variational Dropout(I)) to generate the high-resolution predictions of *mean diffusivity* (MD) and *fractional anisotropy* (FA), and their associated predictive uncertainty maps. Here we draw 200 samples of high-resolution DTI predictions for each subject from the predictive distribution  $q_\phi^*(\mathbf{y}|\mathbf{x})$ , and then the FA and MD maps of each prediction are computed. The sample mean and standard deviation are then calculated from these samples to generate the final estimates of high-resolution MD/FA maps and their corresponding predictive uncertainty.

Fig. 2.7 displays high correspondence between the error (RMSE) maps and the predictive uncertainty on both FA and MD of a HCP test subject. This demonstrates the potential utility of uncertainty map as a surrogate measure of prediction accuracy. In particular, the MD uncertainty map captures subtle variations within the white matter and the cerebrospinal fluid (CSF) at the centre. Also, in accordance with the low reconstruction accuracy, high predictive uncertainty is observed in the CSF in MD. This is expected since the CSF is essentially free water with low signal-to-noise-ratio (SNR) and is also affected by biological noise such as cardiac pulsations. The reconstruction errors are high in FA prediction on the bottom-right quarter of the brain boundary, close to the skull, which is also reflected in the uncertainty map.

Fig. 2.8 tests the utility of predictive uncertainty map in discriminating potential predictive failures in the predicted high-resolution MD map. We define ground truth “safe” voxels as the ones with reconstruction error (RMSE) smaller than a fixed value, and the task is to separate them from the remaining ground-truth “risky” voxels by thresholding on their predictive uncertainty values. The threshold for defining safe voxels is set to  $1.5 \times 10^{-4} \text{ s/mm}^2$ , such that the risky voxels mostly concentrate on the outer-boundary and the CSF regions (which account for 17.5% of all voxels under consideration). Here the positive class is defined as “safe” while the negative class is defined as “risky”. Fig. 2.8 (a) shows the corresponding receiver operating characteristic (ROC) curve of such binary classification task, which plots the true-positive-rate (TPR) against the false-positive-rate (FPR) computed based on all the voxels in the 16 HCP training subjects. In this case, TPR describes the percentage of correctly detected safe voxels out of all the safe ones, while FPR is defined as the percentage of risky voxels that are wrongly classified as safe out of all the risky voxels. We then select the best threshold by maximising the F1 score, and use this to classify the voxels in each predicted high-resolution MD into “safe” and “risky” ones for all subjects in the test HCP dataset and the Lifespan dataset. Fig. 2.8 (b) shows the inter-subject average of the TPR and FPR on both datasets. While on average TPR slightly worsens compared to the results on the training subjects, FPR improves in both cases—notably, this uncertainty-based classification is able to correctly identify 96% of risky predictions on unseen subjects from out-of-training-distribution dataset, namely Lifespan, which differs in demographics and underlying acquisition. Fig. 2.8 (c) visualises

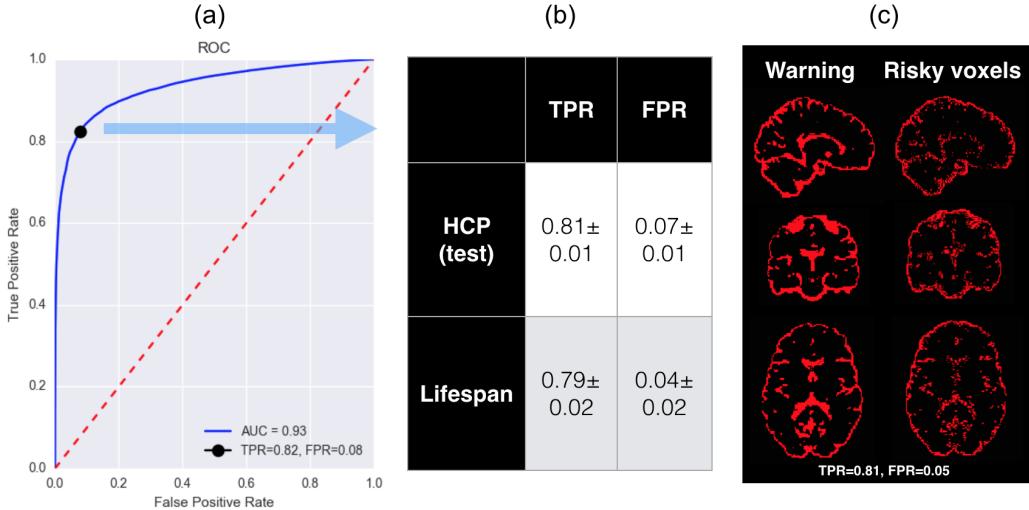


Figure 2.8: Discrimination of “safe” voxels in the predicted high-resolution MD map by thresholding on predictive uncertainty. Here a single 3D-ESPCN + Hetro. + Variational Dropout (I) model is used to quantify the predictive uncertainty over each image volume. (a) the ROC curve plots the true positive rate (TPR) against false positive rate (FPR) computed for a range of threshold values on the foreground voxels in the training subjects. Best threshold (black dot) was selected such that F1 score is maximised and is employed to separate “safe” voxels from “risky” ones; (b) the average TPR and FPR over the 16 test HCP subjects and the 16 Lifespan subjects are shown; (c) an example visualisation of the “ground truth” safe (black) and risky (red) voxels on a Lifespan subject along with the corresponding classification results denoted as “warning”.

the classification results to the pre-defined “ground truth” on one of the Lifespan subjects, which illustrates that the generated “warning” aggressively flags potentially risky voxels at the cost of thresholding out the safe ones.

## Unseen Abnormalities and Uncertainty Decomposition

We separately visualise the propagated **aleatoric** and parameter uncertainty over the predicted high-resolution MD map on images of subjects with a variety of different unseen abnormal structures, such as benign cysts, tumours (Glioma) and focal lesions caused by multiple sclerosis (MS). We emphasise here that the all these images have been acquired with different protocols. Specifically, benign cysts in the HCP datasets represent abnormalities in images acquired with the same protocol as the training data, while tumours and MS lesions are examples of pathologies present in out-of-distribution imaging protocols. In all cases, we use the SR network, Hetero.+Variational Dropout (I), trained on healthy subjects from HCP dataset. For each of 200 different sets of parameters  $\{\theta_t\}_{t=1}^{200}$  sampled from the posterior distribution  $q(\theta|\mathcal{D})$ , we draw 10 samples of high-resolution DTIs from the likelihood,  $\{\mathbf{y}_j^t\}_{j=1}^{10} \sim p(\mathbf{y}|\theta_t, \mathbf{x}, \mathcal{D})$ , compute the corresponding MD, and approximate the two constituents of predictive uncertainty with the MC estimators given in eq.(2.20) and (2.21).

Fig. 2.9 shows the reconstruction accuracy along with the components of predictive uncertainty over the high-resolution MD map of a HCP test subject, which contains a benign abnormality (a small posterior midline arachnoid cyst). The error (RMSE) and propagated **aleatoric** uncertainty are plotted on the same scale whereas the propagated model uncertainty is plotted on 1/5 of the scale for clear visualisation. In this case, the predictive uncertainty is dominated by the **aleatoric** component. In particular, low propagated **aleatoric** uncertainty is observed in the interior of the cyst relative to its boundary in accordance with the high accuracy in the region. This is expected as the interior structure of a cyst is highly homogeneous with low variance in signals and the super-resolution task should therefore be relatively straightforward. On the other hand, the component of parameter uncertainty is high on the interior structure which also makes sense as such homogeneous features are underrepresented in the training data of healthy subjects. This example illustrates how decoupling the effects of **aleatoric** and parameter uncertainty potentially allows one to make sense of the predictive performance.

Fig.2.10 visualises the uncertainty components generated by the same CNN model trained on datasets

of varying size. We see that the propagated parameter uncertainty diminishes as the training set size increases, while the propagated aleatoric uncertainty stays more or less constant. This result is indeed what is expected as described in Fig. 2.3; the specification of network weights becomes more confident i.e. the variance of the posterior distribution decreases as the amount of training data increases, while the effect of aleatoric uncertainty is irreducible with the amount of data. On the other hand, when the standard binary or Gaussian dropout was employed instead of variational dropout, we observed that the effect of parameter uncertainty stayed more or less constant with the size of training data. This may be a consequence of the posterior variance, largely determined by the prespecified drop-out rates, which in turn results in more static variance of predictive distribution.

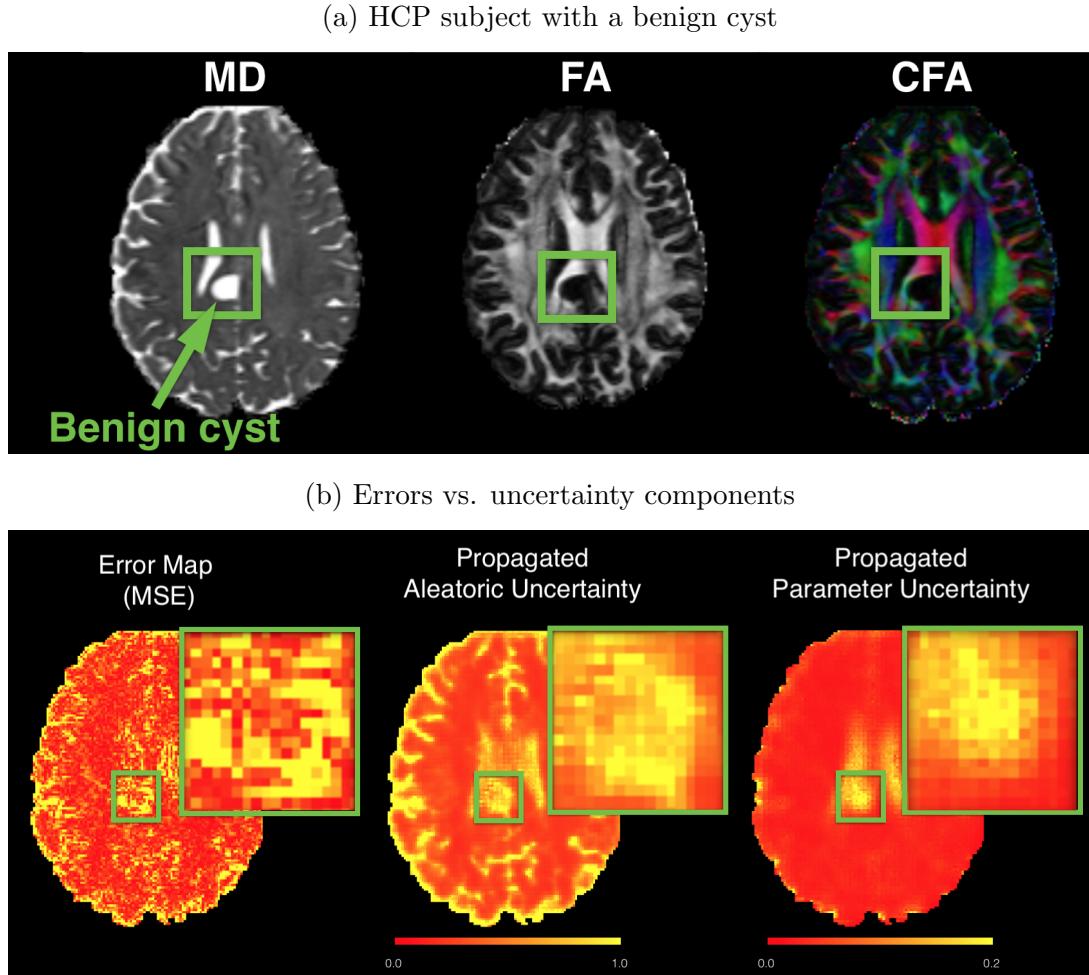


Figure 2.9: Visualisation of (a) MD, FA and colour FA maps computed from the DTI of a HCP subject with a small posterior midline arachnoid cyst in the central part of the brain. (b) the corresponding reconstruction accuracy (RMSE) in MD and the corresponding components of predicted uncertainty.

We further validate our method on clinical images with previously unseen pathologies. We note that the pathology data contain images acquired with standard clinical protocols with voxel size slightly smaller than that of the training low-resolution images and lower signal-to-noise ratio.

Fig. 2.11 shows that pathological areas not represented in the training set are flagged as highly uncertain. Although the ground truth is not available in this case, the uncertainty can be quantified instead to flag potential low accuracy areas. Fig. 2.11 (a) shows that the propagated parameter uncertainty highlights the tumour core, and speckly artefacts in the input image, which are not represented in the training data. On the other hand, the aleatoric uncertainty component is high on the whole region of pathology covering both the tumour core and its surrounding edema. Fig. 2.11 (b) shows that high parameter uncertainty is assigned to a large part of focal lesions in MS, while the aleatoric uncertainty is mostly prevalent around the boundaries between anatomical structures and CSF. We also observe that the super-resolution sharpens the original image without introducing noticeable artifacts; in particular, for the brain tumour image, some of the partial volume effects are cleared.

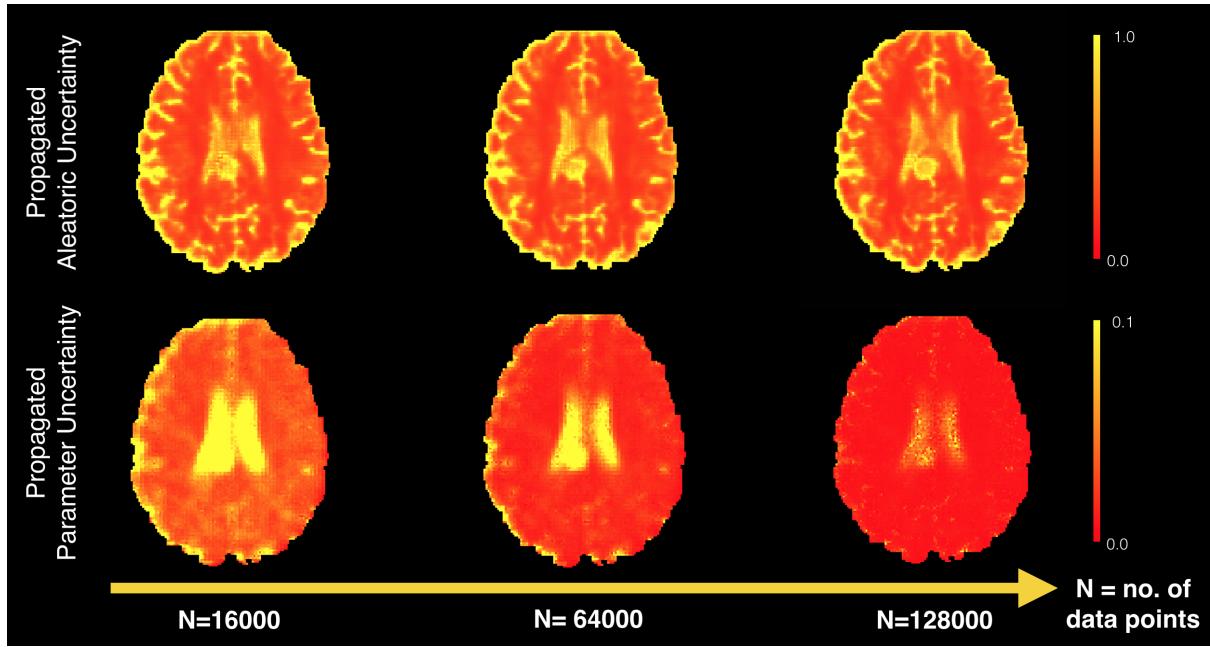


Figure 2.10: Training set size vs propagated aleatoric/parameter uncertainty on the MD map of an unseen HCP subject with a benign cyst. The uncertainty maps are normalised across all the figures.

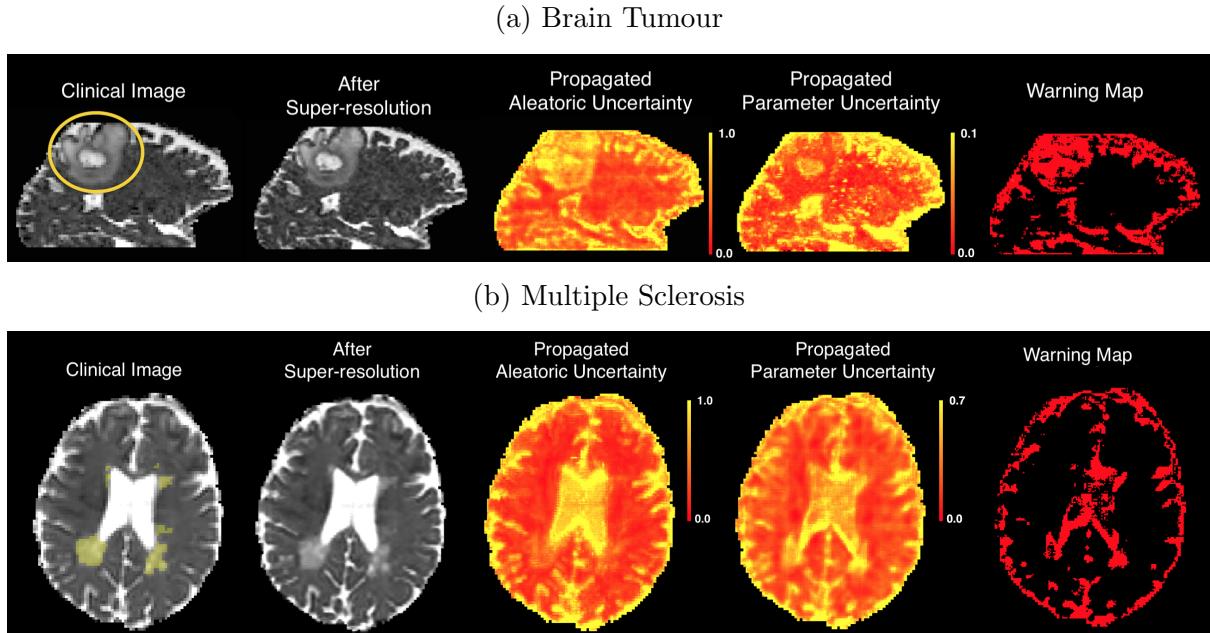


Figure 2.11: Visualisation of propagated uncertainty components on clinical images with pathology that was not present in the training data. The super-resolution is performed on the clinical images due to low-resolution, and thus the ground truths are not available in both cases. (a) shows the results on the data of a Glioma patient, and the yellow circle indicates the region of tumour. (b) shows the same set of results on a MS patient with labels of focal lesions obtained from a neurologist indicated in yellow. Each row shows from left to right: (i) MD map computed from the original DTI; (ii) MD map computed from the output of super-resolution; (iii), (iv) maps of the estimated propagated aleatoric and parameter uncertainty; (v) “warning map” obtained from the same threshold value used in Sec. 2.6.2, which flag large parts of the pathological features in both cases.

## 2.7 Discussion and Conclusion

We introduce a probabilistic deep learning (DL) framework for quantifying three types of uncertainties that arise in data-enhancement applications, and demonstrate its potential benefits in improving the safety of such systems towards practical deployment. The framework models *aleatoric uncertainty* through heteroscedastic noise model and *parameter uncertainty* through approximate Bayesian inference in the form of variational dropout, and finally integrates the two to quantify *predictive uncertainty* over the system output. Experiments focus on the super-resolution application of image quality transfer (IQT)[138] and study several desirable properties of such framework, which lack in the existing body of data enhancement methods based on deterministic DL models.

Firstly, results on a range of applications and datasets show that modelling uncertainty improves overall prediction performance. Table 2.2 and 2.3 show that modelling the combination of both *aleatoric* and *parameter* uncertainty achieves the state-of-the-art accuracy on super-resolution of DTIs and MAP-MRI coefficients in both of the HCP test dataset and the Lifespan dataset, improving on the present best methods based on random-forests (RF-IQT[138] and RF-BIQT[100]) and interpolation—the standard method to estimate sub-voxel information used in clinical visualisation software. In particular, results on the Lifespan dataset, which differs from the training data in age range and acquisition protocol, indicates the better generalizability of our method. Such improvement in the predictive performance arises from the regularisation effects imparted by the modelling of respective uncertainty components. Specifically, modelling *aleatoric* uncertainty through the heteroscedastic network improves robustness to outliers, while modelling parameter uncertainty via variational dropout defends against overfitting. For example, Table 2.3 shows that the predictive performance of the 3D-ESPCN + Hetero. model is only marginally compromised even when the outliers are not removed from training data, while the baseline 3D-ESPCN results in much poorer performance. This can be ascribed to the ability of the variance network  $\Sigma_{\theta_2}(\cdot)$  in the 3D-ESPCN + Hetero. architecture to attenuate the effects of outliers by assigning small weights (i.e. high uncertainty) in the weighted MSE loss function as shown in eq. (2.21). However, this loss attenuation mechanism can also encourage the network to overfit to low-uncertainty regions, potentially focusing less on ambiguous yet important parts of the data—we indeed observe in Table 2.3 that the heteroscedastic network performs considerably worse than the baseline 3D-ESPCN on the exterior regions while the reverse is observed on the interior part. Such overfitting to low-uncertainty interior regions is alleviated by modelling parameter uncertainty with variational dropout [136], as evidenced by the dramatic error reduction in the exterior region on both HCP and Lifespan datasets.

Secondly, experiments on the images of healthy and pathological brains have demonstrated the utility of *predictive uncertainty* as a reliability metric of output images. Fig. 2.7 illustrates the strong correspondence between the maps of predictive uncertainty and the reconstruction quality (voxel-wise RMSE) in the downstream derived quantities such as FA and MD maps. In addition, Fig. 2.11 shows that such uncertainty measure also highlights pathological structures not observed in the training data. We have also tested the utility of predictive uncertainty in discriminating voxels with sufficiently low RMSEs in the predicted high-resolution MD maps. As shown in Fig. 2.8, the optimal threshold selected on the HCP training dataset is capable to detecting over 90% of non-reliable predictions—voxels with RMSE above a certain threshold—not only on the unseen subjects in the same HCP cohort but also on subjects from the out-of-sample Lifespan dataset, that are statistically disparate from the training distribution (e.g. different age range and acquisition protocol). These results combined demonstrate the utility of predictive uncertainty map as a means to quantify output safety, and provides a subject-specific alternative to standard population-group reliability metrics (e.g. mean reconstruction accuracy in a held-out cohort of subjects). Such conventional group statistics can be misleading in practice; for instance, the information that a super-resolution algorithm is reliable 99% of the time on a dataset of 1000 subjects may not accurately represent the performance on a new unseen individual if the person is not well-represented in the cohort (e.g. pathology, different scanners, etc). In contrast, predictive uncertainty provides a metric of reliability, tailored to each individual at hand.

Thirdly, our preliminary experiments show that decomposition of the effects of *aleatoric* and parameter uncertainty in the predictive uncertainty provides a layer of explanations into the performance of the considered deep learning methods. Fig. 2.9 shows that the low reconstruction error in the centre of the benign cyst can be explained by the dominant *aleatoric* uncertainty, which indicates the inherent simplicity of super-resolution task in such homogeneous region, whilst the unfamiliarity of such structure in the healthy training dataset is reflected in the high parameter uncertainty. Assuming that the esti-

mates of decomposed uncertainty components are sufficiently accurate, we could act on them to further improve the overall safety of the system. Imagine a scenario where reconstruction error is consistently high on certain image structures, if the parameter uncertainty is high but aleatoric uncertainty is low, this indicates that collecting more training data would be beneficial. On the other hand, if the parameter uncertainty is low and aleatoric uncertainty is high, this would mean that we need to regard such errors as inevitability, and abstain from predictions to ensure safety or account for them appropriately in subsequent analysis.

The proposed methods for estimating aleatoric and parameter uncertainty, however, make several simplifying assumptions in the forms of likelihood model  $p(\mathbf{y}|\theta, \mathbf{x})$  and posterior distributions over network parameters  $p(\theta|\mathcal{D})$ . Firstly, the likelihood model takes the form of a Gaussian distribution with a diagonal covariance matrix. This means that the likelihood model is not able to capture multi-modality of the predictive distribution i.e. the presence of multiple different solutions. While the full predictive distribution (eq. (2.9)) is not necessarily unimodal in theory due to the integration with the posterior distribution, we observe in practice that the drawn samples are not very diverse. Future work should explore the benefits of employing more complex forms of likelihood functions such as mixture models [195, 182], diversity losses [196, 197, 198] and more powerful density estimators [199, 200, 201, 202, 182]. Also, the diagonality of covariance matrices means that the output pixels are assumed statistically independent given the input. Although the predicted images display high inter-pixel consistency, modelling the correlations between neighbouring pixels [203] may further improve the reconstruction quality. Analogous to the likelihood function, variational dropout [136], which is used in this work, approximates the posteriors  $p(\theta|\mathcal{D})$  by Gaussian distributions with diagonal covariance, imposing restrictive assumptions of unimodality and statistical independence between neural network weights. More recent advances in the Bayesian deep learning research [204, 205, 206, 207, 208, 209] could be used to enhance the quality of parameter uncertainty estimation by allowing the model to capture multi-modality and statistical dependencies between parameters. We also refer the readers to a recent review paper by [210] on this topic for a balanced perspective on possible approaches. We should note that both the mean and variance MC estimators of very high dimensional posterior distribution converge with only a few hundred samples in our case, because of this simplistic choice of the variational distributions. However, it is likely that, in order to approximate the posterior with a more complex family of distributions, a larger number of samples would be necessary.

An important future challenge is the clinical validation of predictive uncertainty as a reliability metric of output images. To this end, we need to design a more clinically meaningful definition of success and failure of the data enhancement algorithm at hand. Despite the high accuracy in distinguishing between predictive failures and successes attained with our method (Fig. 2.8), our definition of reconstruction quality, namely voxel-wise RMSE, does not necessarily represent the real utility of the output image. One possible approach would be to have clinical experts to label the potential failures in the super-resolved images, be it for a targeted application (e.g. diagnosis of some neurological conditions) or for general usage in clinical practice. A more economical alternative, which does not require extra label acquisition, is to define the prediction success in downstream measurements of interest i.e. functions of the output images  $g(\cdot)$ , such as morphometric measurements of anatomical or pathological structures (e.g. volumes). The propagation method (eq. (2.13)) introduced in Sec. 2.3.5 can be utilised to quantify uncertainty components in the space of target measurement  $g(\cdot)$ . Measuring the correlation between such propagated uncertainty estimates and the errors would be a useful indicator of how well the uncertainty measure reflects the accuracy of the chosen measurement  $g(\cdot)$ . Lastly, our initial results on the brain tumour dataset motivate a larger-scale quantitative validation of uncertainty estimates in the presence of pathology. Future work must examine the effect of including patients' images in the training data on the estimate of uncertainty components.

There are many ways in which uncertainty information could be utilised by radiologists or other users of data enhancement algorithms. First, predictive uncertainty can be used to decide when to abstain from predictions in high-risk regions of images (e.g. anomalies, out-of-distribution examples or inherently ambiguous features). For example, the original input low-resolution image can be augmented by overlaying the high-resolution prediction only in locations with sufficiently low uncertainty, before presenting to clinicians. As demonstrated by Fig. 2.8 in the context of super-resolution, such uncertainty-based quality control of predictions is potentially an effective means to maintain high accuracy of output images and also to safeguard against hallucination or removal of structures [133]. Second, the uncertainty information could be used for active learning [211] to decide which images should be labelled and included in the training set to maximally improve the model performance. Prior work [212, 213] define the acquisition function so as to select examples with high parameter uncertainty, and achieve promising

results in classification and segmentation tasks. In particular, these methods are able to construct a compact and effective training dataset, and consequently improve the prediction accuracy while reducing the training time. The same idea could be naturally extended to data enhancement problems, that are typically formulated as multivariate regression tasks. For example, in the case of IQT, we could simulate a library of low-resolution and high-resolution image pairs from a large public dataset (e.g. HCP), and incrementally expand the training data by adding more examples from such a library. We should note, however, that in many data enhancement applications, obtaining a new “label” may require an extra acquisition possibly with a different scanner or modality, which may be logically challenging. Third, another important application is transfer learning [214] where uncertainty information could be used to leverage knowledge from different but related domains or tasks. In many data enhancement applications, the test distribution can considerably deviate from the training distribution. For example, the algorithm might be trained on a synthetic dataset or images acquired from a scanner that is very different from the one used in the hospital where one plans to deploy the model. Therefore, a mechanism to adapt performance within a specific environment (e.g., based on the local patient population) [215], possibly in an online fashion [216, 217], is in demand. Recent work have shown that the Bayesian formalism provides a natural framework to use uncertainty in order to account for the difference and commonality between distributions to guide information transfer in continual learning [218, 219] or few-shot learning [220, 221] settings. Exploring the benefits of these ideas in the context of medical image enhancement remains future work.

The proposed framework for uncertainty quantification is formulated for multivariate regression in the general form, and thus is naturally applicable to many other image enhancement challenges such as: rapid image acquisition techniques e.g., compressed sensing [122], MR fingerprinting [222, 223] or sparse reconstruction [125, 124]; denoising [118] and dealiasing [127, 224]; image synthesis tasks e.g., estimating T2-weighted images from T1 [94, 93, 225], estimating CT images from MRI [226, 160, 172], and generating a high-field scan from a low-field scan [143]; data harmonisation [227, 120, 121] which aims to learn mappings among imaging protocols to reduce confounds in multicentre studies. Our results on image quality transfer [138] illustrate the potential of the uncertainty modelling techniques to improve the safety of these applications by not only improving the predictive accuracy, but also providing a mechanism to quantify risks and safeguard against potential malfunction.



# Chapter 3

## Uncertainty in Multitask Learning (I): Spatially Adaptive Weighting of Task Loss Functions

**Abstract:** In this chapter, we extend the methods of uncertainty modelling introduced in Chapter 2 to the multi-task learning setting. Such adaptation naturally yields a mechanism to automatically determine, in a spatially adaptive fashion, the relative weighting between the task losses, which is a key determinant in the efficacy of multi-task learning. Focusing on the task of structured predictions, we evaluate the benefits of this idea in the context of MR-only radiotherapy planning by considering the multi-task learning problem of simultaneously regressing a synthetic CT (synCT) scan and segmenting organs at risk (OAR) from the input MRI. We test our method on prostate cancer scans and show that it achieves state-of-the-art performance in the regression and segmentation of prostate cancer scans. We further show that the estimates of uncertainty correlate strongly in areas prone to errors across both tasks, which can be used as mechanism for quality control in radiotherapy treatment planning. This chapter is based on a joint work [160] with Felix Bragman. My primary contributions are method development and experiment design.

### 3.1 Introduction

Radiotherapy is a common and effective treatment modality for cancer. The state-of-the-art protocol in such treatment requires acquiring a magnetic resonance (MR) scan to accurately segment the target and surrounding organs at risk (OARs), and a registered computed tomography (CT) scan for dose calculation. However, this approach has seen limited translation to clinical practice because the acquisition of both MRI and CT is time-consuming, and the registration step for spatially aligning the modalities may introduce unacceptable errors that propagate in the planning process. MR-only treatment planning has recently attracted a lot of attention as a potential solution to these issues [228, 229, 230, 231], and involves translating MR image data into CT like image, so called synthetic CTs (synCT) [232, 233]. This synthesis process, when combined with a hand drawn region of interest and a set of safety margins, enable clinicians to devise a radiotherapy plan. In this chapter, we employ a convolutional neural network to jointly generate the synCT and the segmentation of the OARs for a given MR input. This multi-task learning formulation [234] provides an end-to-end system that operates on a single MR scan and provides the outputs necessary for radio-therapy planning, while improving the prediction quality by fusing information across the tasks.

However, the efficacy of such approach largely depends on the quality of the training data. While synthesising CT from MR can be fundamentally ill-posed in some regions [98], the variability in physicians' delineations of OARs and pathology has been empirically shown to be high [235, 236]. This means that the training data suffers from such sources of inherent variability. Considering the dire consequences of failures in radio-therapy treatment [236, 237, 238], it is important to quantify the uncertainty in the predicted SynCT and segmentations of OARs, and account for it in the subsequent treatment planning.

In this chapter, we therefore take the methods introduced in Chapter 2, where intrinsic uncertainty

is estimated through a heteroscedastic noise model and parameter uncertainty is modelled using approximate Bayesian inference, and adapt them to the relevant multi-task learning setting. This provides a data-driven adaptation of task losses on a voxel-wise basis and importantly, a measure of uncertainty over the prediction of both tasks, which can be potentially used for quality assurance in radiotherapy treatment planning (e.g. quantification of dose delivery uncertainty).

## 3.2 Related work

**MR-only radiotherapy planning:** Methods for simulating a synCT and segmenting corresponding MR scans have originated from multi-atlas propagation [239]. Recently, applications of convolutional neural networks (CNNs) to CT synthesis from MRI have become a topic of growing interest due to their reconstruction performance. To alleviate the problem of missing high-frequency information in synCT due to mean-squared reconstruction loss, Nie et al. [240] employed a conditional generative adversarial network to capture fine texture details whilst Wolterink et al. [241] extended this idea using a CycleGAN to leverage the abundance of unpaired training sets of CT and MR scans. However, despite this advancement, all of these methods commit to a single prediction with no measure of predictive uncertainty, limiting the utility in view of current and future probabilistic dose delivery systems. Moreover, none of the CNN-based methods segment OARs. Here we employ a probabilistic multi-task architecture to simultaneously estimate OAR segmentations and a synCT that are anatomically consistent, along with their associated uncertainty.

**Uncertainty in Multi-task learning:** Past approaches to multi-task learning have relied on uniform or hand-tuned weighting of task losses [242, 243]. Recently, Kendall et al. [244] proposed to model the uncertainty of each task, and thereby adjust each task’s relative weight in the cost function automatically. They demonstrated on various structured prediction tasks that the method outperformed separate models trained individually on each task. However, the work assumes that the uncertainty is constant for each task and spatially, which is unrealistic for imaging data. For example, Figure 1 in Asman and Landman [245] illustrates in the context of brain segmentation that even human experts have a clear inclination to mislabel boundary pixels and other ambiguous regions. Here we enrich their probabilistic multi-task learning method by modelling the spatial variation of intrinsic uncertainty via a so-called heteroscedastic noise model, and integrating parameter uncertainty via dropout. We show later that this confers a mechanism to select the relative weighting of task losses in a pixel-wise fashion.

## 3.3 Methods

We introduce a probabilistic bi-task CNN architecture which takes an MR image, and simultaneously estimates the distributions over the corresponding CT image and the segmentation probability of the OARs. Analogous to the methods introduced in Chapter 2, we use a Gaussian heteroscedastic noise model [135] and binary dropout [154] to account for *intrinsic* and *parameter* uncertainty, respectively, and show that we obtain not only a measure of uncertainty over prediction, but also a mechanism for data-driven adaptation of weightings of task losses, which is integral for benefiting from the multi-task learning framework. As before, to address the memory burden of 3D medical images, we employ a patch-based approach to perform both tasks, in which the input MR image is split into smaller overlapping patches that are processed independently. For each input patch  $\mathbf{x}$ , our dual-task model estimates the conditional distributions  $p(\mathbf{y}_i|\mathbf{x})$  for tasks  $i = 1, 2$  where  $\mathbf{y}_1$  and  $\mathbf{y}_2$  denote the Hounsfield Unit and class probabilities of OARs at the center of the input patch. At inference, the probability maps over the synCT and OARs are obtained by stitching together outputs from appropriately shifted versions of the input patches.

### 3.3.1 Bi-task architecture

Here we use a standard multi-task learning architecture with hard-parameter sharing [246] where the model shares initial few layers across the two tasks and branches out into four task-specific components with separate parameters as illustrated in Fig.3.1. There are two components per task, where one

aims to performs CT synthesis (regression) or OAR segmentation, and the remaining models *intrinsic* uncertainty associated to the data and the task.

The rationale behind shared layers is to learn a joint representation between two tasks to regularise the learning of features for one task by using cues from the other. We used a high-resolution network architecture (HighResNet) [247] as the shared trunk of the model for its compactness and accuracy shown in brain parcellation. HighResNet is a fully convolutional architecture that utilises dilated convolutions and residual connections to produce an end-to-end mapping from an input patch ( $\mathbf{x}$ ) to voxel-wise predictions ( $\mathbf{y}$ ). The combination of dilated convolution and the residual connections enable the model to integrate efficiently both the short-range (e.g., textures, intensities) and long-range information (e.g., semantic features), both of which are important for the synthesis and anatomical segmentation tasks.

The final layer of the shared representation is split into two task-specific compartments (Fig. 3.1). Each compartment consists of two fully convolutional networks which operate on the output of representation network and together learn task-specific representation and define likelihood function  $p(\mathbf{y}_i|\mathbf{W}, \mathbf{x})$  for each task  $i = 1, 2$  where  $\mathbf{W}$  denotes the set of all parameters of the model.

### 3.3.2 Task weighting with heteroscedastic uncertainty.

Previous probabilistic multitask learning methods based on deep learning [244] assumed constant intrinsic uncertainty in respective tasks. In our context, this means that the inherent ambiguity present in synthesis or segmentation tasks do not depend on the spatial locations within an image volume. This is a highly unrealistic assumption since these tasks can be more challenging on some anatomical structures (e.g. tissue boundaries) than others as evidenced by previous work [245]. In order to capture potential spatial variation in intrinsic uncertainty, we adapt the *heteroscedastic* (data-dependent) noise model to our multitask learning problem.

In particular, for the CT synthesis task, we define our likelihood as a normal distribution  $p(\mathbf{y}_1|\mathbf{W}, \mathbf{x}) = \mathcal{N}(f_1^{\mathbf{W}}(\mathbf{x}), \sigma_1^{\mathbf{W}}(\mathbf{x})^2)$  where mean  $f_1^{\mathbf{W}}(\mathbf{x})$  and variance  $\sigma_1^{\mathbf{W}}(\mathbf{x})^2$  are modelled by the regression output and uncertainty branch as functions of the input patch  $\mathbf{x}$  (see Fig.3.1). We define the task loss for CT synthesis to be the negative log-likelihood  $\mathcal{L}_1(\mathbf{y}_1, \mathbf{x}; \mathbf{W}) = \frac{1}{2\sigma_1^{\mathbf{W}}(\mathbf{x})^2} \|\mathbf{y}_1 - f_1^{\mathbf{W}}(\mathbf{x})\|^2 + \log \sigma_1^{\mathbf{W}}(\mathbf{x})^2$ . This loss encourages assigning high-uncertainty to regions of high errors, enhancing the robustness of the network against noisy labels and outliers, which are prevalent at organ boundaries especially close to the bone.

For the segmentation, we define the classification likelihood as softmax function of scaled logits i.e.  $p(\mathbf{y}_2|\mathbf{W}, \mathbf{x}) = \text{Softmax}(f_2^{\mathbf{W}}(\mathbf{x})/\sigma_2^{\mathbf{W}}(\mathbf{x})^2)$  where the segmentation output  $f_2^{\mathbf{W}}(\mathbf{x})$  is scaled by the uncertainty term  $\sigma_2^{\mathbf{W}}(\mathbf{x})^2$  before softmax (Fig.3.1). As the uncertainty term  $\sigma_2^{\mathbf{W}}(\mathbf{x})$  increases, the Softmax output approaches a uniform distribution, which corresponds to the maximum entropy discrete distribution. We simplify the scaled Softmax likelihood by considering an approximation used in [244]

$$\frac{1}{\sigma_2^{\mathbf{W}}(\mathbf{x})^2} \sum_{c'} \exp\left(\frac{1}{\sigma_2^{\mathbf{W}}(\mathbf{x})^2} f_{2,c'}^{\mathbf{W}}(\mathbf{x})\right) \approx \left( \sum_{c'} \exp(f_{2,c'}^{\mathbf{W}}(\mathbf{x})) \right)^{1/\sigma_2^{\mathbf{W}}(\mathbf{x})^2}$$

where  $c'$  is denotes a segmentation class. This approximation becomes an equality when  $\sigma_2^{\mathbf{W}} \rightarrow 1$ . This has the advantage of simplifying the optimisation objective, as well as empirically improving results. This yields the NLL task-loss of the form  $\mathcal{L}_2(\mathbf{y}_2 = c, \mathbf{x}; \mathbf{W}) \approx \frac{1}{\sigma_2^{\mathbf{W}}(\mathbf{x})^2} \text{CE}(f_2^{\mathbf{W}}(\mathbf{x}), \mathbf{y}_2 = c) + \log \sigma_2^{\mathbf{W}}(\mathbf{x})^2$ , where CE denotes cross-entropy. Finally, assuming that the two tasks  $\mathbf{y}_1, \mathbf{y}_2$  are statistically independent given

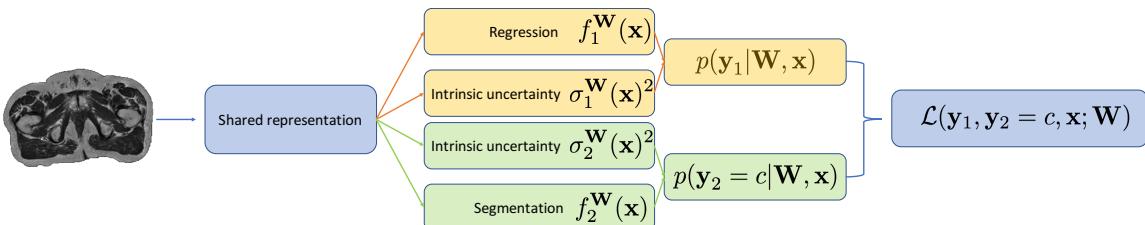


Figure 3.1: Multi-task learning architecture. The predictive mean and variance  $[f_i^{\mathbf{W}}(\mathbf{x}), \sigma_i^{\mathbf{W}}(\mathbf{x})^2]$  are estimated for the regression and segmentation. The task-specific likelihoods  $p(\mathbf{y}_i|\mathbf{W}, \mathbf{x})$  are combined to yield the multi-task likelihood  $p(\mathbf{y}_1, \mathbf{y}_2|\mathbf{W}, \mathbf{x})$ .

the input image  $\mathbf{x}$ , the joint likelihood factorises over tasks  $p(\mathbf{y}_1, \mathbf{y}_2 | \mathbf{W}, \mathbf{x}) = p(\mathbf{y}_1 | \mathbf{W}, \mathbf{x}) p(\mathbf{y}_2 | \mathbf{W}, \mathbf{x})$ , and thus we can derive the NLL loss for the dual-task model as

$$\mathcal{L}(\mathbf{y}_1, \mathbf{y}_2 = c, \mathbf{x}; \mathbf{W}) = -\log p(\mathbf{y}_1, \mathbf{y}_2 | \mathbf{W}, \mathbf{x}) \quad (3.1)$$

$$= \frac{\|\mathbf{y}_1 - f_1^{\mathbf{W}}(\mathbf{x})\|^2}{2\sigma_1^{\mathbf{W}}(\mathbf{x})^2} + \frac{\text{CE}(f_2^{\mathbf{W}}(\mathbf{x}), \mathbf{y}_2 = c)}{\sigma_2^{\mathbf{W}}(\mathbf{x})^2} + \log(\sigma_1^{\mathbf{W}}(\mathbf{x})^2 \sigma_2^{\mathbf{W}}(\mathbf{x})^2) \quad (3.2)$$

where the MSE and CE terms are weighted by the inverse of heteroscedastic intrinsic uncertainty terms  $\sigma_i^{\mathbf{W}}(\mathbf{x})^2$ , that enables automatic weighting of task losses on a per-sample basis. The log-term controls the spread. By contrast, the previous approaches resort to spatially constant [244] or manually specified weighting of task losses [242, 243].

### 3.3.3 Parameter uncertainty with approximate Bayesian inference.

Analogous to Chapter 2, we approximate the posterior distribution over the network weights using variational inference, and assess the benefit of modelling parameter uncertainty in the context of our multitask learning problem. However, here we instead use the standard binary dropout [154] following its Bayesian interpretation introduced by Gal et al.[248]. Let  $q(\mathbf{W})$  denote the variational distribution we use to approximate the true posterior  $p(\mathbf{W} | \mathbf{X}, \mathbf{Y}_1, \mathbf{Y}_2)$  where  $\mathbf{X} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}$ ,  $\mathbf{Y}_1 = \{\mathbf{y}_1^{(1)}, \dots, \mathbf{y}_1^{(N)}\}$ ,  $\mathbf{Y}_2 = \{\mathbf{y}_2^{(1)}, \dots, \mathbf{y}_2^{(N)}\}$  is the training data. During training, we minimise the following variational objective similarly to eq. (2.6):

$$\mathcal{L}(\mathcal{D}; \mathbf{W}) = \sum_{(\mathbf{x}^{(i)}, \mathbf{y}_1^{(i)}, \mathbf{y}_2^{(i)}) \in \mathcal{D}} \left( \mathbb{E}_{q_{\phi}(\mathbf{W})}[-\log p(\mathbf{y}_1^{(i)}, \mathbf{y}_2^{(i)} | \mathbf{x}^{(i)}, \mathbf{W})] + \text{KL}(q(\mathbf{W}) || p(\mathbf{W})) \right) \quad (3.3)$$

The second term (referred to as the prior term) is simply given by the L2 weight decay [152]. On the other hand, the first term (referred to as the reconstruction term) cannot be computed exactly, thus we employ the following MC approximation by drawing  $S$  samples of network parameters from the approximate posterior  $\mathbf{W}^{(s)} \sim q_{\phi}(\mathbf{W})$ :

$$\begin{aligned} \mathbb{E}_{q_{\phi}(\mathbf{W})}[-\log p(\mathbf{y}_1, \mathbf{y}_2 = c | \mathbf{x}, \mathbf{W}^{(s)})] &\approx \frac{1}{S} \sum_{s=1}^S -\log p(\mathbf{y}_1, \mathbf{y}_2 = c | \mathbf{x}, \mathbf{W}^{(s)}) \\ &\propto \frac{\sum_{s=1}^S \sigma_1^{\mathbf{W}}(\mathbf{x})^2 \|\mathbf{y}_1 - f_1^{\mathbf{W}}(\mathbf{x})\|^2}{2 \sum_{s=1}^S \sigma_1^{\mathbf{W}}(\mathbf{x})^2} \\ &+ \frac{\sum_{s=1}^S \sigma_2^{\mathbf{W}}(\mathbf{x})^2 \text{CE}(f_2^{\mathbf{W}^{(s)}}(\mathbf{x}), \mathbf{y}_2 = c)}{\sum_{s=1}^S \sigma_2^{\mathbf{W}^{(s)}}(\mathbf{x})^2} \\ &+ \sum_{s=1}^S \log(\sigma_1^{\mathbf{W}^{(s)}}(\mathbf{x})^2 \sigma_2^{\mathbf{W}^{(s)}}(\mathbf{x})^2) \end{aligned}$$

where the first two terms correspond to the weighted average of the respective task loss functions weighted by the corresponding intrinsic uncertainty. We should also note the slight abuse of notation here; we use  $\propto$  to denote the equality up to additive constants. In our experiments we found that the number of samples  $S$  per datapoint can be set to 1 as long as the minibatch size was large enough.

At test time, for each input patch  $\mathbf{x}$  in a MR scan, we collect output samples (the mean and variance terms)  $\{\mathbf{f}^{w^{(t)}}(\mathbf{x})\}_{t=1}^T$  where  $\mathbf{f}^{w'}(\mathbf{x}) := [f_1^{w'}(\mathbf{x}), f_2^{w'}(\mathbf{x}), \sigma_1^{w'}(\mathbf{x})^2, \sigma_2^{w'}(\mathbf{x})^2]$  by performing  $T$  stochastic forward-passes with  $\{w^{(t)}\}_{t=1}^T \sim q(\mathbf{W})$ . To compute the final estimates of synCT and segmentation of OARs, we calculate the mean predictions over the  $T$  samples for the respective tasks:

$$\hat{\mu}_{\mathbf{y}_j | \mathbf{x}} = \frac{1}{T} \sum_{t=1}^T f_j^{w^{(t)}}(\mathbf{x}), \quad j \in \{1, 2\}$$

Finally, we estimate the intrinsic and parameter components of the predictive uncertainty in both tasks by considering the approximations given by eq. (2.18) and eq. (2.19):

$$\hat{\Delta}_m(\mathbf{y}_j) = \frac{1}{T} \sum_{t=1}^T f_j^{w^{(t)}}(\mathbf{x}) f_j^{w^{(t)}}(\mathbf{x})^\top - \hat{\mu}_{\mathbf{y}_j|\mathbf{x}} \hat{\mu}_{\mathbf{y}_j|\mathbf{x}}^\top \quad (3.4)$$

$$\hat{\Delta}_i(\mathbf{y}_j) = \frac{1}{T} \sum_{t=1}^T \sigma_j^{w^{(t)}}(\mathbf{x})^2 \quad (3.5)$$

where  $j \in \{1, 2\}$  denotes the task index. As before, the total predictive uncertainty is the sum of the above intrinsic and parameter components.

## 3.4 Data Preprocessing and Implementation Details

### 3.4.1 Data

We validated on 15 prostate cancer patients, who each had a T2-weighted MR image (3T,  $1.46 \times 1.46 \times 5\text{mm}^3$ ) and a CT image (140kVp,  $0.98 \times 0.98 \times 1.5\text{mm}^3$ ) acquired on the day. Organ delineation was performed by a clinician with labels for the left and right femur head, bone, prostate, rectum and bladder. All images were resampled to isotropic resolution. The CT scans were spatially aligned with the T2 scans using the method of Burgos et al. [239]. In the segmentation, we predicted labels for the background, left/right femur head, prostate, rectum and bladder. The bone region was used for quantifying the synCT.

### 3.4.2 Network architectures and training

We trained our model on randomly selected 2D axial slices and reconstructed the 3D volume at test time. The representation network was composed of a convolutional layer followed by 3 sets of twice repeated dilated convolutions [249] with dilation factors  $[1, 2, 4]$  and a final convolutional layer. Each layer ( $l$ ) used a  $3 \times 3$  kernel with features  $f_R = [64, 64, 128, 256, 2048]$ . Each task-specific branch was a set of 5 convolutional layers of size  $[256_{l=1,2,3,4}, n_{i,l=5}]$  where  $n_{i,l=5}$  is equal to 1 for regression and  $\sigma$  and equal to the number of segmentation classes. The first two layers were  $3 \times 3$  kernels whilst the final convolutional layers were fully connected. A Bernoulli drop-out mask with probability  $p = 0.5$  was applied on the final layer of the representation network. We minimised the loss using the ADAM optimiser [190] with a learning rate  $10^{-3}$  and trained for 19,000 iterations. For the stochastic sampling, we performed model inference 10 times at iterations 18000 and 19000 leading to a set of  $T = 20$  samples.

## 3.5 Results

### 3.5.1 Experimental set-up

We performed a 3-fold cross-validation. Statistics over all hold-out sets are reported. We considered four separate models; (M1) baseline networks for regression/segmentation, (M2a) baseline network with drop-out for regression/segmentation, (M2b) the baseline with drop-out and heteroscedastic noise, (M3) multi-task network using homoscedastic task weighting [244] and (M4) multi-task network using task-specific heteroscedastic noise and drop-out. The baseline networks used only the representation network with  $1/2f_R$  and a fully-connected layer for the final output. We also compared our results against the current state of the art in atlas propagation (AP) [239], which was validated on the same dataset.

### 3.5.2 Model performance

An example of the model output is shown in Fig. 3.2. We calculated the Mean Absolute Error (MAE) between the predicted and reference scans across the body and at each organ (Tab. 3.1). The fuzzy

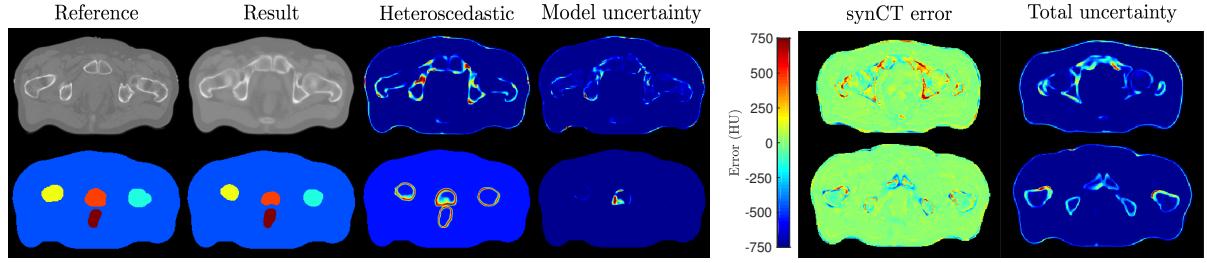


Figure 3.2: Example outputs from the proposed model. Right: the propagated *intrinsic* (“Heteroscedastic”) and *parameter* (“Model”) uncertainty are shown on an example image for both tasks (CT synthesis and anatomical segmentation) along with the ground truth (“Reference”) and the predictions (“Result”). The uncertainty maps correlate with regions of high contrast (bone in the regression, organ boundary for segmentation). Left: The errors in CT-synthesis (“synCT error”) are compared against the predictive uncertainty (“Total”—the sum of the heteroscedastic and model uncertainty—on two additional images. Note the correlation between model error and the predicted uncertainty.

DICE score between the probabilistic segmentation and the reference was calculated for the segmentation (Tab. 3.1). Best performance was in our presented method (M4) for the regression across all segmentation masks except at the bladder. Application of the multi-task heteroscedastic network with drop-out (M4) produced the most consistent synCT across all models with the lowest average MAE and the lowest variation across patients ( $43.3 \pm 2.9$  versus  $45.7 \pm 4.6$  [239] and  $44.3 \pm 3.1$  [244]). This was significantly lower when compared to M1 ( $p < 0.001$ ) and M2 ( $p < 0.001$ ). This was also observed at the bone and prostate ( $p < 0.001$ ). Whilst differences at  $p < 0.05$  levels of significance was not observed versus M2b and M3, the consistently lower MAE and standard deviation across patients in M4 demonstrates the added benefit of modelling heteroscedastic noise and the inductive transfer from the segmentation task. Moreover, we performed better than the current state of the art in atlas propagation [239]. Despite equivalence with the state of the art (Tab. 3.1), we did not observe any significant differences between our model and the baselines despite an improvement in mean DICE at the smaller organs such as prostate and rectum ( $0.70 \pm 0.06$  and  $0.74 \pm 0.12$ ) versus the baseline M1 ( $0.67 \pm 0.12$ ,  $0.70 \pm 0.15$ ). The *intrinsic uncertainty* (Fig. 3.2) captures the uncertainty specific to the data and thus penalises regions of high error, leading to an under-segmentation yet with higher confidence in the result.

Table 3.1: Model comparison. Bold values indicate where a model was significantly worse than M4  $p < 0.05$ . No data was available for significance testing with AP. M2b was statistically better  $p < 0.05$  than M4 in the prostate segmentation.

Models	All	Bone	<i>L</i> femur	<i>R</i> femur	Prostate	Rectum	Bladder
Regression - synCT - Mean Absolute Error (HU)							
M1	<b>48.1(4.2)</b>	<b>131(14.0)</b>	78.6(19.2)	<b>80.1(19.6)</b>	<b>37.1(10.4)</b>	63.3(47.3)	<b>24.3(5.2)</b>
M2a	<b>47.4(3.0)</b>	<b>130(12.1)</b>	78.0(14.8)	77.0(13.0)	<b>36.5(7.8)</b>	67(44.6)	<b>24.1(7.5)</b>
M2b [150]	44.5(3.6)	128(17.1)	75.8(20.1)	74.2(17.4)	31.2(7.0)	56.1(45.5)	17.8(4.7)
M3 [244]	44.3(3.1)	126(14.4)	74.0(19.5)	73.7(17.1)	29.4(4.7)	58.4(48.0)	18.2(3.5)
AP [239]	45.7(4.6)	125(10.3)	-	-	-	-	-
M4 (ours)	43.3(2.9)	121(12.6)	69.7(13.7)	67.8(13.2)	28.9(2.9)	55.1(48.1)	18.3(6.1)
Segmentation - OAR - Fuzzy DICE score							
M1	-	-	0.91(0.02)	0.90(0.04)	0.67(0.12)	0.70(0.15)	0.92(0.05)
M2a	-	-	0.85(0.03)	0.90(0.04)	0.66(0.12)	0.69(0.13)	0.90(0.07)
M2b [150]	-	-	0.92(0.02)	0.92(0.01)	0.77(0.07)	0.74(0.13)	0.92(0.03)
M3 [244]	-	-	0.92(0.02)	0.92(0.02)	0.73(0.07)	0.76(0.10)	0.93(0.02)
AP [239]	-	-	0.89(0.02)	0.90(0.01)	0.73(0.06)	0.77(0.06)	0.90(0.03)
M4 (ours)	-	-	0.91(0.02)	0.91(0.02)	0.70(0.06)	0.74(0.12)	0.93(0.04)

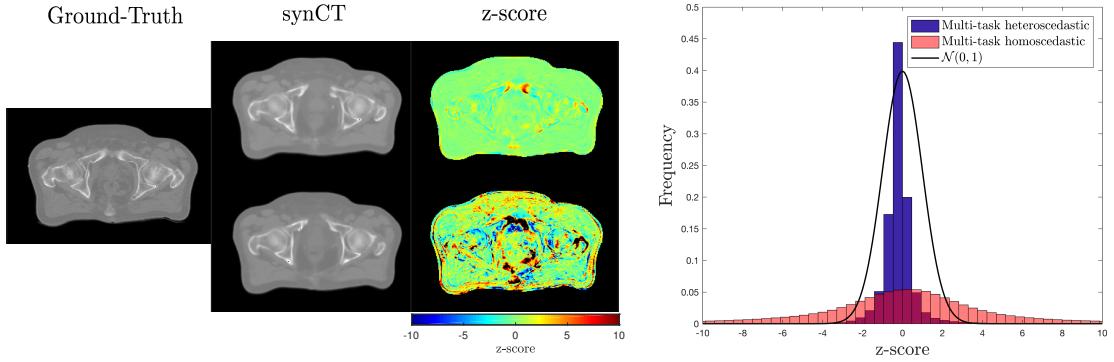


Figure 3.3: Analysis of uncertainty estimation. a) synCTs and z-scores for the a subject between M4 (top) and M3 (bottom) models. b) z-score distribution of all patients (15) between both models.

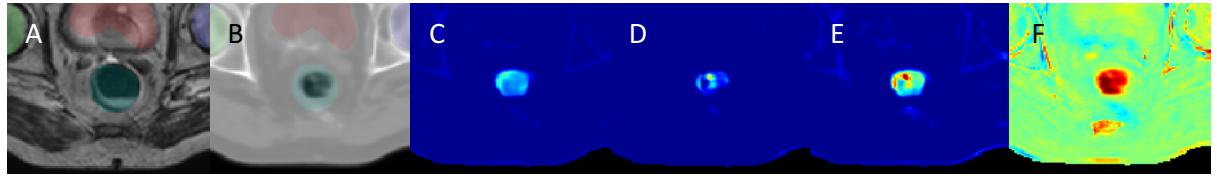


Figure 3.4: Uncertainty in problematic areas. a) T2 with reference segmentation, b) synCT with localised errors at the rectum, c) propagated intrinsic uncertainty in synCT, d) propagated parameter uncertainty in synCT, e) total predictive uncertainty and f) error in HU (range [-750HU, 750HU]).

### 3.5.3 Uncertainty estimation for radiotherapy

We tested the ability of the multi-task heteroscedastic network to better predict associated uncertainties in the synCT error. To verify that our network produces clinically viable samples for treatment planning, we quantified the distribution of regression z-scores for the multi-task heteroscedastic and homoscedastic models. In the former, the total predictive uncertainty is the sum of the propagated *intrinsic* and *parameter* uncertainties, which is used to normalise the error between the synCT and the reference. This should lead to a better approximation of the variance in the model. In contrast, the total uncertainty in the latter reduces to the variance of the stochastic test-time samples. This is likely to lead to a mis-calibrated variance. A  $\chi^2$  goodness of fit test was performed, showing that the homoscedastic z-score distribution is not normally distributed ( $0.82 \pm 0.54$ ,  $p < 0.01$ ) in contrast to the heteroscedastic model ( $0.04 \pm 0.84$ ,  $p > 0.05$ ). This is apparent in Fig. 3.3 where there is greater confidence in the synCT produced by our model in contrast the homoscedastic case.

The predictive uncertainty can be exploited for quality assurance (Fig. 3.4). There may be issues whereupon time differences have caused variations in bladder and rectum filling across MR and CT scans causing patient variability in the training data. This is exemplified by large errors in the synCT at the rectum (Fig. 3.4) and quantified by large localised z-scores (Fig. 3.4g), which correlate strongly with the propagated *intrinsic* and *parameter* uncertainty across tasks.

## 3.6 Conclusions

We have adapted the methods of uncertainty modelling introduced in Chapter 2 to the multi-task learning setting. Our network extends prior work in multi-task learning by integrating heteroscedastic uncertainty modelling to naturally weight task losses and thus facilitate inductive transfer between tasks. We have demonstrated the applicability of our network in the context of MR-only radiotherapy treatment planning where the synthetic CT scan (SynCT) and the segmentation of OARs are simultaneously generated from the input MR image. We have shown that accounting for uncertainty information leads to more accurate and consistent synCTs with a constraint on anatomical consistency with the segmentations. Importantly, we have demonstrated that the output of our network leads to consistent anatomically correct stochastic synCT samples that can potentially be effective in treatment planning. Furthermore, we have also shown that the estimates of predictive uncertainty with our method is more calibrated than the equivalent model

with homoscedastic noise model. In the future work, we will evaluate the downstream utility of such uncertainty information in the predicted synCT and OAR segmentation in designing a safer treatment planning of radiotherapy e.g., [250, 251].

# Chapter 4

## Uncertainty in Multitask Learning (II): Stochastic Filter Groups for Learning Structured Sparsity

**Abstract:** The performance of multi-task learning in Convolutional Neural Networks (CNNs) hinges on the design of feature sharing between tasks within the architecture. The number of possible sharing patterns are combinatorial in the depth of the network and the number of tasks, and thus hand-crafting an architecture, purely based on the human intuitions of task relationships can be time-consuming and suboptimal. In this chapter, we present a probabilistic approach to learning task-specific and shared representations in CNNs for multi-task learning. Specifically, we propose “stochastic filter groups” (SFG), a mechanism to assign convolution kernels in each layer to “specialist” or “generalist” groups, which are specific to or shared across different tasks, respectively. The SFG modules determine the connectivity between layers and the structures of task-specific and shared representations in the network. We employ variational inference to learn the posterior distribution over the possible grouping of kernels and network parameters. We demonstrate the utility of learning such structured sparsity in the architecture in the context of MR-only radiotherapy planning application (also considered in Chapter 3) as well as the age and gender prediction task on the UTKFace dataset. This is based on the joint work [105] in collaboration with Felix Bragman and Jorge Cardoso.

### 4.1 Introduction

Multi-task learning (MTL) aims to enhance learning efficiency and predictive performance by simultaneously solving multiple related tasks [234]. Recently, applications of convolutional neural networks (CNNs) in MTL have demonstrated promising results in a wide-range of computer vision applications, ranging from visual scene understanding [252, 253, 254, 255, 256, 257] to medical image computing [242, 258, 104, 243].

A key factor for successful MTL neural network models is the ability to learn shared and task-specific representations [254]. A mechanism to understand the commonalities and differences between tasks allows the model to transfer information between tasks while tailoring the predictive model to describe the distinct characteristics of the individual tasks. The quality of such representations is determined by the architectural design of where model components such as features [259] and weights [260] are shared and separated between tasks. However, the space of possible architectures is combinatorially large, and the manual exploration of this space is inefficient and subject to human biases. For example, Fig. 4.1 shows a typical CNN architecture for MTL comprised of a shared “trunk” feature extractor and task-specific “branch” networks [243, 261, 262, 244, 256, 104]. The desired amount of shared and task-specific representations, and their interactions within the architecture are dependent on the difficulty of the individual tasks and the relation between them, neither of which are a priori known in most cases [263]. This illustrates the challenge of handcrafting an appropriate architecture, and the need for an effective automatic method to learn it from data.

In this paper, we propose *Stochastic Filter Groups* (SFGs); a probabilistic mechanism to learn the

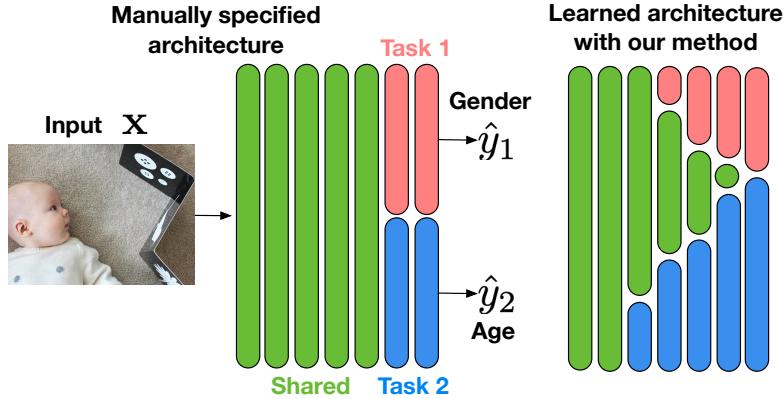


Figure 4.1: Figure on the left illustrates a typical multi-task architecture, while the figure on the right shows an example architecture that can be learned with our method. We propose *Stochastic Filter Groups*, a principled way to learn the assignment of convolution kernels to task-specific and shared groups.

amount of task-specific and shared representations needed in each layer of MTL architectures (Fig. 4.1). Specifically, the SFGs learn to allocate kernels in each convolution layer into either ‘‘specialist’’ groups or a ‘‘shared’’ trunk, which are specific to or shared across different tasks, respectively (Fig. 4.2). The SFG equips the network with a mechanism to learn inter-layer connectivity and thus the structures of task-specific and shared representations. We cast the learning of SFG modules as a variational inference problem.

We evaluate the efficacy of SFGs on a variety of tasks. In particular, we focus on two multi-task learning problems: 1) age regression and gender classification from face images on UTKFace dataset [264] and 2) semantic regression (i.e. image synthesis) and semantic segmentation on a real-world medical imaging dataset, both of which require predictions over all pixels. Experiments show that our method achieves considerably higher prediction accuracy than baselines with no mechanism to learn connectivity structures, and either higher or comparable performance than a cross-stitch network [254], while being able to learn meaningful architectures automatically.

## 4.2 Related works

Our work is concerned with the goal of learning where to share neural network components across different tasks to maximise the benefit of MTL. The main challenge of such methods lies in designing a mechanism that determines how and where to share weights within the network. There are broadly two categories of methods that determine the nature of weight sharing in MTL networks.

The first category is composed of methods that directly optimise the sharing of weights in order to maximise task-wise performance. These methods set out to learn a set of vectors that control which features are shared within a layer and how these are distributed across [265, 260, 254, 259]. They start with a baseline CNN architecture where they learn additional connections and pathways that define the final MTL model. For instance, Cross-Stitch networks [254] control the degree of weight sharing at each convolution layer whilst Soft-Layer Ordering [260] goes beyond the assumption of parallel ordering of feature hierarchies to allow features to mix at different layers depending on the task. In contrast, the second group of MTL methods focuses on weight clustering based on task-similarity [266, 267, 268, 269, 270]. For example, [269] employed a greedy, iterative algorithm to grow a tree-like deep architecture that clusters similar tasks hierarchically or [270] which determines the degree of weight sharing based on statistical dependency between tasks.

Our method falls into first category, and differentiates itself by performing ‘‘hard’’ partitioning of task-specific and shared features. By contrast, prior methods are based on ‘‘soft’’ sharing of features [254, 259] or weights [265, 260]. These methods generally learn a set of mixing coefficients that determine the weighted sum of features throughout the network, which does not impose connectivity structures on the architecture. On the other hand, our method learns a distribution over the connectivity of layers by grouping kernels. This allows our model to learn meaningful grouping of task-specific and shared features as illustrated in Fig. 4.7.

## 4.3 Methods

We introduce a new approach for determining where to learn task-specific and shared representation in multi-task CNN architectures. We propose *stochastic filter groups* (SFG), a probabilistic mechanism to partition kernels in each convolution layer into “specialist” groups or a “shared” group, which are specific to or shared across different tasks, respectively. We employ variational inference to learn the distributions over the possible grouping of kernels and network parameters that determines the connectivity between layers and the shared and task-specific features. This naturally results in a learning algorithm that optimally allocate representation capacity across multi-tasks via gradient-based stochastic optimization, e.g. stochastic gradient descent.

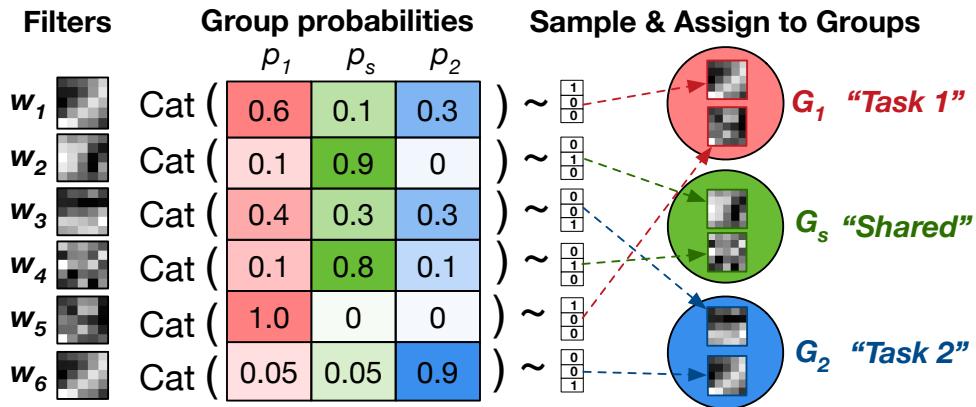


Figure 4.2: Illustration of filter assignment in a SFG module. Each kernel  $\{w_k\}$  in the given convolution layer is probabilistically assigned to one of the filter groups  $G_1, G_s, G_2$  according to the sample drawn from the associated categorical distribution  $\text{Cat}(p_1, p_s, p_2)$ .

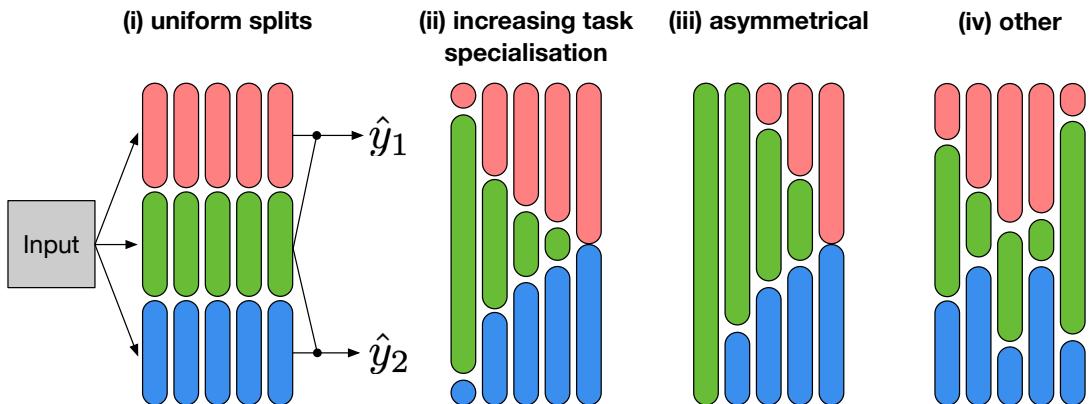


Figure 4.3: Illustration of possible grouping patterns learnable with the proposed method. Each set of green, pink and yellow blocks represent the ratio of filter groups  $G_1$  (red),  $G_s$  (green) and  $G_2$  (blue). (i) denotes the case where all kernels are uniformly split. (ii) & (iii) are the cases where the convolution kernels become more task-specific at deeper layers. (iv) shows an example with more heterogeneous splits across tasks.

### 4.3.1 Stochastic Filter Groups

SFGs introduce a sparse connection structure into the architecture of CNN for multi-task learning in order to separate features into task-specific and shared components. Ioannou et al. [271] introduced

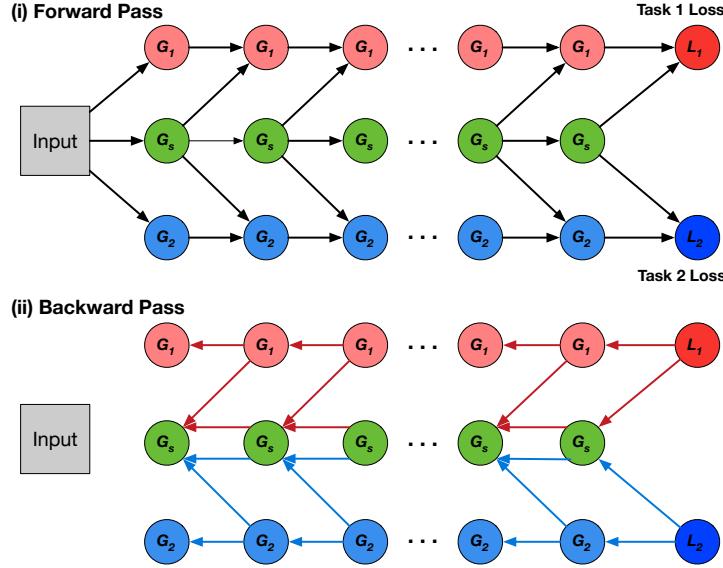


Figure 4.4: Illustration of feature routing. The circles  $G_1, G_s, G_2$  denote the task-specific and shared filter groups in each layer. (i) shows the directions of routing of activations between different filter groups while (ii) shows the directions of the gradient flow from the task losses  $L_1$  and  $L_2$ . The red and blue arrows denote the gradients that step from  $L_1$  and  $L_2$ , respectively. The task-specific groups  $G_1, G_2$  are only updated based on the associated losses, while the shared group  $G_s$  is updated based on both.

*filter groups* to partition kernels in each convolution layer into groups, each of which acts only on a subset of the preceding features, and demonstrated that such sparsity reduces computational cost and number of parameters without compromising accuracy. Here we adapt the concept of filter groups to the multi-task learning paradigm and propose an extension with an additional mechanism for learning an optimal kernel grouping rather than pre-specifying them.

For simplicity, we describe SFGs for the case of multitask learning with two tasks, but can be trivially extended to a larger number of tasks. At the  $l^{\text{th}}$  convolution layer in a CNN architecture with  $K_l$  kernels  $\{\mathbf{w}^{(l),k}\}_{k=1}^{K_l}$ , the associated SFG performs two operations:

1. **Filter Assignment:** each kernel  $\mathbf{w}_k^{(l)}$  is stochastically assigned to either: i) the “task-1 specific group”  $G_1^{(l)}$ , ii) “shared group”  $G_s^{(l)}$  or iii) “task-2 specific group”  $G_2^{(l)}$  with respective probabilities  $\mathbf{p}^{(l),k} = [p_1^{(l),k}, p_s^{(l),k}, p_2^{(l),k}] \in [0, 1]^3$ . Convolving with the respective filter groups yields distinct sets of features  $F_1^{(l)}, F_s^{(l)}, F_2^{(l)}$ . Fig. 4.2 illustrates this operation and Fig. 4.3 shows different learnable patterns.
2. **Feature Routing:** as shown in Fig. 4.4 (i), the features  $F_1^{(l)}, F_s^{(l)}, F_2^{(l)}$  are routed to the filter groups  $G_1^{(l+1)}, G_s^{(l+1)}, G_2^{(l+1)}$  in the subsequent  $(l+1)^{\text{th}}$  layer in such a way to respect the task-specificity and sharedness of filter groups in the  $l^{\text{th}}$  layer. Specifically, we perform the following routing for  $l > 0$ :

$$\begin{aligned} F_1^{(l+1)} &= h^{(l+1)}([F_1^{(l)} | F_s^{(l)}] * G_1^{(l+1)}) \\ F_s^{(l+1)} &= h^{(l+1)}(F_s^{(l)} * G_s^{(l+1)}) \\ F_2^{(l+1)} &= h^{(l+1)}([F_2^{(l)} | F_s^{(l)}] * G_2^{(l+1)}) \end{aligned}$$

where each  $h^{(l+1)}$  defines the choice of non-linear function, “ $*$ ” denotes convolution operation and “ $|$ ” denotes a merging operation of arrays (e.g. concatenation). At  $l = 0$ , input image  $\mathbf{x}$  is simply convolved with the first set of filter groups to yield  $F_i^{(1)} = h^{(1)}(\mathbf{x} * G_i^{(1)}), i \in \{1, 2, s\}$ . Fig. 4.4(ii) shows that such sparse connectivity ensures the parameters of  $G_1^{(l)}$  and  $G_2^{(l)}$  are only learned based on the respective task losses, while  $G_s^{(l)}$  is optimised based on both tasks.

Fig. 4.5 provides a schematic of our overall architecture, in which each SFG module stochastically generates filter groups in each convolution layer and the resultant features are sparsely routed as described

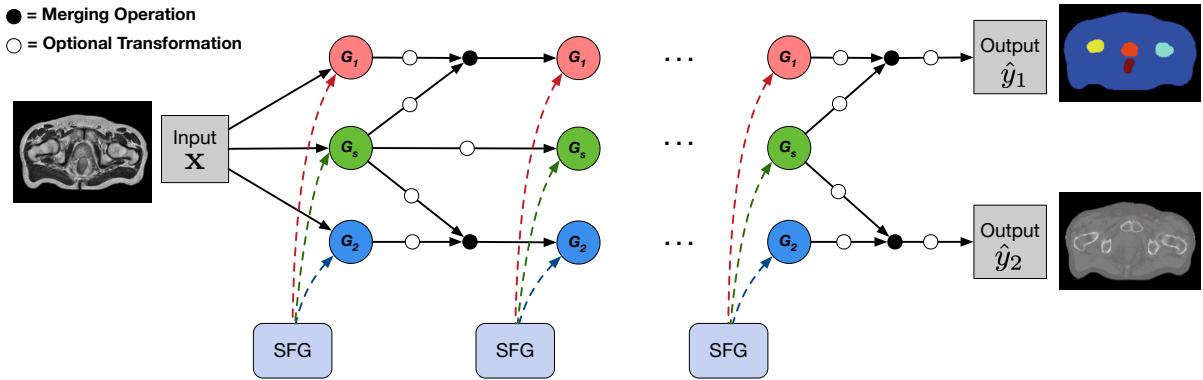


Figure 4.5: Schematic of the proposed multi-task architecture based on a series of SFG modules in the presence of two tasks. At each convolution layer, kernels are stochastically assigned to task-specific and shared filter groups  $G_1, G_s, G_2$ . Each input image is first convolved with the respective filter groups to yield three distinct sets of output activations, which are routed sparsely to the filter groups in the second layer layer. This process repeats in the remaining SFG modules in the architecture until the last layer where the outputs of the final SFG module are combined into task-specific predictions  $\hat{y}_1$  and  $\hat{y}_2$ . Each small white circle denotes an optional transformation (e.g. extra convolutions) and black circle merges the incoming inputs (e.g. concatenation).

above. The merging modules, denoted as black circles, combine the task-specific and shared features appropriately, i.e.  $[F_i^{(l)}|F_s^{(l)}], i = 1, 2$  and pass them to the filter groups in the next layer. Each white circle denotes the presence of additional transformations (e.g. convolutions or fully connected layers) in each  $h^{(l+1)}$ , performed on top of the standard non-linearity (e.g. ReLU).

The proposed sparse connectivity is integral to ensure task performance and structured representations. In particular, one might argue that the routing of “shared” features  $F_s^{(l)}$  to the respective “task-specific” filter groups  $G_1^{(l+1)}$  and  $G_2^{(l+1)}$  is not necessary to ensure the separation of gradients across the task losses. However, this connection allows for learning more complex task-specific features at deeper layers in the network. For example, without this routing, having a large proportion of “shared” filter group  $G_s$  at the first layer (Fig. 4.3 (ii)) substantially reduces the amount of features available for learning task-specific kernels in the subsequent layers—in the extreme case in which all kernels in one layer are assigned to  $G_s$ , the task-specific filter groups in the subsequent layers are effectively unused.

Another important aspect that needs to be highlighted is the varying dimensionality of feature maps. Specifically, the number of kernels in the respective filter groups  $G_1^{(l)}, G_s^{(l)}, G_2^{(l)}$  can vary at each iteration of the training, and thus, so does the depth of the resultant feature maps  $F_1^{(l)}, F_s^{(l)}, F_2^{(l)}$ . Instead of directly working with feature maps of varying size, we implement the proposed architecture by defining  $F_1^{(l)}, F_s^{(l)}, F_2^{(l)}$  as sparse tensors. At each SFG module, we first convolve the input features with all kernels, and generate the output features from each filter group by zeroing out the channels that root from the kernels in the other groups, resulting in  $F_1^{(l)}, F_s^{(l)}, F_2^{(l)}$  that are sparse at non-overlapping channel indices. In the simplest form with no additional transformation (i.e. the grey circles in Fig. 6.1 are identity functions), we define the merging operation  $[F_i^{(l)}|F_s^{(l)}], i = 1, 2$  as pixel-wise summation. In the presence of more complex transforms (e.g. residual blocks), we concatenate the output features in the channel-axis and perform a 1x1 convolution to ensure the number of channels in  $[F_i^{(l)}|F_s^{(l)}]$  is the same as in  $F_s^{(l)}$ .

### 4.3.2 Optimisation: T+1 Way “Drop-Out”

Here we derive the method for simultaneously optimising the CNN parameters and grouping probabilities. We achieve this by extending the variational interpretation of binary dropout [272, 192] to the  $(T + 1)$ -way assignment of each convolution kernel to the filter groups where  $T$  is the number of tasks. As before, we consider the case  $T = 2$ .

Suppose that the architecture consists of  $L$  SFG modules, each with  $K_l$  kernels where  $l$  is the index and let us denote the parameters of the kernels in all layers by  $\mathcal{W} = \{\mathbf{W}^{(l),k}\}_{k=1,\dots,K_l,l=1,\dots,L}$ . As the

posterior distribution over the convolution kernels in SFG modules  $p(\mathcal{W}|\mathbf{X}, \mathbf{Y}^{(1)}, \mathbf{Y}^{(2)})$  is intractable, we approximate it with a simpler distribution  $q_\phi(\mathcal{W})$  parametrised by  $\phi$  which we describe below. Assuming that the posterior distribution factorizes over layers and kernels up to group assignment, we define the variational distribution as:

$$\begin{aligned} q_\phi(\mathcal{W}) &:= \prod_{l=1}^L \prod_{k=1}^{K_l} q_{\phi_{lk}}(\mathbf{W}^{(l),k}) \\ &= \prod_{l=1}^L \prod_{k=1}^{K_l} q_{\phi_{lk}}(\mathbf{W}_1^{(l),k}, \mathbf{W}_s^{(l),k}, \mathbf{W}_2^{(l),k}) \end{aligned}$$

where  $\{\mathbf{W}_1^{(l),k}, \mathbf{W}_s^{(l),k}, \mathbf{W}_2^{(l),k}\}$  denotes the  $k^{\text{th}}$  kernel in  $l^{\text{th}}$  convolution layer after being routed into task-specific  $G_1^{(l)}, G_2^{(l)}$  and shared group  $G_s^{(l)}$ . We define each  $q_{\phi_{lk}}(\mathbf{W}_1^{(l),k}, \mathbf{W}_2^{(l),k}, \mathbf{W}_s^{(l),k})$  as:

$$\mathbf{W}_i^{(l),k} = z_i^{(l),k} \cdot \mathbf{M}^{(l),k} \quad \text{for } i \in \{1, s, 2\} \quad (4.1)$$

$$\mathbf{z}^{(l),k} = [z_1^{(l),k}, z_2^{(l),k}, z_s^{(l),k}] \sim \text{Cat}(\mathbf{p}^{(l),k}) \quad (4.2)$$

where  $\mathbf{z}^{(l),k}$  is the one-hot encoding of a sample from the categorical distribution over filter group assignments, and  $\mathbf{M}^{(l),k}$  denotes the parameters of the pre-grouping convolution kernel. In summary, the set of variational parameters for each kernel in each layer is thus given by  $\phi_{lk} = \{\mathbf{M}^{(l),k}, \mathbf{p}^{(l),k} = [p_1^{(l),k}, p_s^{(l),k}, p_2^{(l),k}]\}$ .

We minimize the KL divergence between the approximate posterior  $q_\phi(\mathcal{W})$  and  $p(\mathcal{W}|\mathbf{X}, \mathbf{Y}^{(1)}, \mathbf{Y}^{(2)})$ . Assuming that the joint likelihood over the two tasks factorizes, we have the following optimization objective:

$$\mathcal{L}_{\text{MC}}(\phi) = -\frac{N}{M} \sum_{i=1}^M \left[ \log p(y_i^{(1)}|\mathbf{x}_i, \mathcal{W}_i) + \log p(y_i^{(2)}|\mathbf{x}_i, \mathcal{W}_i) \right] + \sum_{l=1}^L \sum_{k=1}^{K_l} \text{KL}(q_{\phi_{lk}}(\mathbf{W}^{(l),k})||p(\mathbf{W}^{(l),k})) \quad (4.3)$$

where  $M$  is the size of the mini-batch,  $N$  is the total number of training data points, and  $\mathcal{W}_i$  denotes a set of model parameters sampled from  $q_\phi(\mathcal{W})$ . The first two terms are approximations of the expected log-likelihoods for the respective tasks. The last KL term regularizes the deviation of the approximate posterior from the prior  $p(\mathbf{W}^{(l),k}) = \mathcal{N}(0, \mathbf{I}/l^2)$  where  $l > 0$ . Adapting the approximation presented in [272] to our scenario, we obtain:

$$\text{KL}(q_{\phi_{lk}}(\mathbf{W}^{(l),k})||p(\mathbf{W}^{(l),k})) \propto \frac{l^2}{2} \|\mathbf{M}^{(l),k}\|_2^2 - \mathcal{H}(\mathbf{p}^{(l),k}) \quad (4.4)$$

where  $\mathcal{H}(\mathbf{p}^{(l),k}) = -\sum_{i \in \{1, 2, s\}} p_i^{(l),k} \log p_i^{(l),k}$  is the entropy of the grouping probabilities. While the first term performs the L2-weight norm, the second term pulls the grouping probabilities towards the uniform distribution. Plugging eq.(4.4) into eq.(4.3) yields the overall loss:

$$\mathcal{L}_{\text{MC}}(\phi) = -\frac{N}{M} \sum_{i=1}^M \left[ \log p(y_i^{(1)}|\mathbf{x}_i, \mathcal{W}_i) + \log p(y_i^{(2)}|\mathbf{x}_i, \mathcal{W}_i) \right] + \lambda_1 \cdot \sum_{l=1}^L \sum_{k=1}^{K_l} \|\mathbf{M}^{(l),k}\|^2 - \lambda_2 \cdot \sum_{l=1}^L \sum_{k=1}^{K_l} \mathcal{H}(\mathbf{p}^{(l),k}) \quad (4.5)$$

where  $\lambda_1 > 0, \lambda_2 > 0$  are regularization coefficients. During training, we perform the SGD based on the above loss function to learn both the network weights  $\mathcal{W}$  and the filter assignment probabilities  $\{\mathbf{p}^{(l),k}\}_{l,k}$ . By default, the assignment probabilities are initialised as  $\mathbf{p}^{(l),k} = [0.2, 0.6, 0.2]$  although different options are explored e.g., see Fig. 4.13.

We note that the discrete sampling operation during filter group assignment (eq. (4.2)) creates discontinuities, giving the first term in the objective function (eq. 4.5) zero gradient with respect to the grouping probabilities  $\{\mathbf{p}^{(l),k}\}$ . We therefore, as employed in [244] for the binary case, approximate each of the categorical variables  $\text{Cat}(\mathbf{p}^{(l),k})$  by the Gumbel-Softmax distribution,  $\text{GSM}(\mathbf{p}^{(l),k}, \tau)$  [273, 274], a continuous relaxation which allows for sampling, differentiable with respect to the parameters  $\mathbf{p}^{(l),k}$  through a reparametrisation trick. The temperature term  $\tau$  adjusts the bias-variance tradeoff of gradient approximation; as the value of  $\tau$  approaches 0, samples from the GSM distribution become one-hot vectors (vectors whose one element is 1 and the rest are 0s) i.e. lower bias while the variance of the gradients increases. In practice, we start at a high  $\tau$  and anneal to a small but non-zero value as in [274, 192]. The details are available in supplementary materials of [105].

## 4.4 Experiments

We tested *stochastic filter groups* (SFG) on two multi-task learning (MTL) problems: 1) age regression and gender classification from face images on UTKFace dataset [264] and 2) semantic image regression (synthesis) and segmentation on a medical imaging dataset.

**UTKFace dataset:** We tested our method on UTKFace [264], which consists of 23,703 cropped faced images in the wild with labels for age and gender. We created a dataset with a 70/15/15% split. We created a secondary separate dataset containing only 10% of images from the initial set, so as to simulate a data-starved scenario.

**Medical imaging dataset:** We used a medical imaging dataset to evaluate our method in a real-world, multi-task problem where paucity of data is common and hard to mitigate. The goal of radiotherapy treatment planning is to maximise radiation dose to the tumour whilst minimising dose to the organs. To plan dose delivery, a Computed Tomography (CT) scan is needed as CT voxel intensity scales with tissue density, thus allowing dose propagation simulations. An MRI scan is needed to segment the surrounding organs. Instead of acquiring both an MRI and a CT, algorithms can be used to synthesise a CT scan (task 1) and segment organs (task 2) given a single input MRI scan. For this experiment, we acquired 15 3D prostate cancer scans with respective CT and MRI scans with semantic 3D labels for organs (prostate, bladder, rectum and left/right femur heads) obtained from a trained radiologist. We created a training set of 10 patients, with the remaining 5 used for testing. We trained our networks on 2D subimages of size 128x128 randomly sampled from axial slices, and reconstructed the 3D volumes of size 288x288x62 at test time by stitching together the subimage-wise predictions.

### 4.4.1 Baselines

We compared our model against four baselines in addition to Cross-Stitch networks [254] trained end-to-end rather than sequentially for fair comparison. The four baselines considered are: 1) single-task networks, 2) hard-parameter sharing multi-task network (MT-hard sharing), 3) SFG-networks with constant  $1/3$  allocated grouping (MT-constant mask) as *per* Fig. 4.3(i), and 4) SFG-networks with constant grouping probabilities (MT-constant  $\mathbf{p}$ ). We train all the baselines in an end-to-end fashion for all the experiments.

We note that all four baselines can be considered special cases of an SFG-network. Two *single-task networks* can be learned when the shared grouping probability of kernels is set to zero. Considering Fig. 6.1, this would remove the diagonal connections and the shared network. This may be important when faced with two unrelated tasks which share no contextual information. A *hard-parameter sharing network* exists when all shared grouping probabilities are maximised to one leading to a scenario where all features are shared within the network up until the task-specific layers. The *MT-constant mask network* is illustrated in Fig. 4.3(i), where  $1/3$  of kernels are allocated to the task 1, task 2 and shared groups, yielding uniform splits across layers. This occurs when an equal number of kernels in each layer obtain probabilities of  $\mathbf{p}^{(l),k} = [1, 0, 0], [0, 1, 0]$  and  $[0, 0, 1]$ . Lastly, the *MT-constant  $\mathbf{p}$*  model represents the situation where the grouping is non-informative and each kernel has equal probability of being specific or shared with probability  $\mathbf{p}^{(l),k} = [1/3, 1/3, 1/3]$ .

**UTKFace network:** We used VGG-11 CNN architecture [275] for age and gender prediction. The network consists of a series of 3x3 convolutional layers interleaved with max pooling layers. In contrast to the original architecture, we replaced the final max pooling and fully connected layers with global average pooling (GAP) followed by a fully connected layers for prediction. Our model’s version of VGG (SFG-VGG) replaces each convolutional layer in VGG-11 with a SFG layer with max pooling applied to each feature map  $F_1^{(l)}, F_2^{(l)}, F_s^{(l)}$ . We applied GAP to each final feature map before the final merging operation and two fully connected layers for each task.

**Medical imaging network:** We used a high-resolution network architecture (HighResNet) [276] for CT synthesis and organ segmentation. This network has been successfully developed for semantic segmentation in medical imaging and has been used in a variety of medical applications such as CT

synthesis [104, 277], brain segmentation [276] and tumour segmentation [278]. It consists of a series of residual blocks, which group two  $3 \times 3$  convolutional layers with dilated convolutions. The baseline network is composed of a  $3 \times 3$  convolutional layer followed by three sets of twice repeated residual blocks with dilated convolutions using factors  $d = [1, 2, 4]$ . There is a  $3 \times 3$  convolutional layer between each set of repeated residual blocks. The network ends with two final  $3 \times 3$  layers and either one or two  $1 \times 1$  convolutional layers for single and multi-task predictions. In our model, we replace each convolutional layer with an SFG module. After the first SFG layer, three distinct repeated residual blocks are applied to  $F_1^{(l=0)}$ ,  $F_2^{(l=0)}$ ,  $F_s^{(l=0)}$ . These are then merged according the feature routing methodology followed by a new SFG-layer and subsequent residual layers. Our model concludes with 2 successive SFG-layers followed by  $1 \times 1$  convolutional layers applied to the merged features  $F_1^{(l=L)}$  and  $F_2^{(l=L)}$ .

## 4.5 Data preprocessing and implementation details

### 4.5.1 Optimisation, regularisation and initialisation

All networks were trained with ADAM optimiser [279] with an initial learning rate of  $10^{-3}$  and  $\beta = [0.9, 0.999]$ . We used values of  $\lambda_1 = 10^{-6}$  and  $\lambda_2 = 10^{-5}$  for the weight and entropy regularisation factors in Equation (5) in Section 3.2. All *stochastic filter group* (SFG) modules were initialised with grouping probabilities  $\mathbf{p} = [0.2, 0.6, 0.2]$  for every convolution kernel. **Positivity of the grouping probabilities  $\mathbf{p}$  is enforced by passing the output through a *softplus* function  $f(x) = \ln(1 + e^x)$  as in [178], and they are then normalised so they sum to 1.** The scheduler  $\tau = \max(0.10, \exp(-rt))$  recommended in [274] was used to anneal the Gumbel-Softmax temperature  $\tau$  where  $r$  is the annealing rate and  $t$  is the current training iteration. We used  $r = 10^{-5}$  for our models.

Hyper-parameters for the annealing rate and the entropy regularisation weight were obtained by analysis of the network performance on a secondary randomly split on the UTK dataset (70/15/15). They were then applied to all trained models (large and small dataset for UTKFace and medical imaging dataset).

### 4.5.2 UTKFace

For training the VGG networks (Section 4.1 - UTKFace network), we used the root-mean-squared-error (RMSE) for age regression and the cross entropy loss for gender classification. The labels for age were divided by 100 prior to training. The input RGB images ( $200 \times 200 \times 3$ ) were all normalised channel wise to have unit variance and zero mean prior to training and testing. A batch-size of 10 was used. No augmentation was applied. We monitored performance during training using the validation set ( $n = 3554$ ) and trained up to 330 epochs. We performed 150 validation iterations every 1000 iterations, leading to 1500 predictions per validation iteration. Performance on the validation set was analysed and the iteration where Mean Absolute Error (MAE) was minimised and classification Accuracy was maximised was chosen for the test set.

### 4.5.3 Medical imaging dataset

We used T2-weighted Magnetic Resonance Imaging (MRI) scans (3T, 2D spin echo, TE/TR: 80/2500ms, voxel size  $1.46 \times 1.46 \times 5 \text{ mm}^3$ ) and Computed Tomography (CT) scans (140 kVp, voxel size  $0.98 \times 0.98 \times 1.5 \text{ mm}^3$ ). The MR and CT scans were resampled to isotropic resolution ( $1.46 \text{ mm}^3$ ). We performed intensity non-uniformity correction on the MR scans [280].

In the HighResNet networks (Section 4.1 - Medical imaging network), we used the RMSE loss for the regression task and the Dice + Cross-Entropy loss [281] for the segmentation task. The CT scans were normalised using the transformation  $\text{CT}/1024 + 1$ . The original range of the CT voxel intensity was  $[-1024, 2500]$  with the background set to  $-1024$ . The input MRI scans were first normalised using histogram normalisation based on the  $1^{\text{st}}$  and  $99^{\text{th}}$  percentile [282]. The MRI scans were then normalised to zero mean and unit variance. At test time, input MRI scans were normalised using the histogram normalisation transformation obtained from the training set then normalised to have zero mean and unit variance.

All scans were of size 288x288x62. We sub-sampled random patches from random axial slices of size 128x128. We sampled from all axial slices in the volume ( $n = 62$ ). We trained up to 200,000 iterations using a batch-size of 10. We applied augmentation to the randomly sampled patches using random scaling factors in the range  $[-10\%, 10\%]$  and random rotation angles in the range  $[-10^\circ, 10^\circ]$ . The trained patches were zero-padded to increase their size to 136x136. However, the loss during training was only calculated in non-padded regions.

The inference iteration for the test set was determined when the performance metrics on the training set (Mean Absolute Error and Accuracy) first started to converge for at least 10,000 iterations. In our model where the grouping probabilities were learned, the iteration when convergence in the update of the grouping probabilities was first observed was selected since performance generally increased as the grouping probabilities were updated.

#### 4.5.4 Implementation details

We used Tensorflow and implemented our models within the NiftyNet framework [283]. Models were trained on NVIDIA Titan Xp, P6000 and V100. All networks were trained in the Stochastic Filter Group paradigm. Single-task networks were trained by hard-coding the allocation of kernels to task 1 and task 2 i.e. 50% of kernels per layer were allocated to task 1 and 50% were allocated to task 2 with constant probabilities  $\mathbf{p}=[1,0,0]$  and  $\mathbf{p}=[0,0,1]$  respectively. The multi-task hard parameter sharing (MT hard-sharing) network was trained by hard-coding the allocation of kernels to the shared group i.e. 100% of kernel per layer were allocated to the shared group with constant probability  $\mathbf{p}=[0, 1, 0]$ . The cross-stitch (CS) [254] networks were implemented in a similar fashion to the single-task networks, with CS modules applied to the output of the task-specific convolutional layers. The other baselines (MT-constant mask and MT-constant  $\mathbf{p}=[1/3, 1/3, 1/3]$ ) were trained similarly.

We used Batch-Normalisation [191] to help stabilise training. We observed that the deviation between population statistics and batch statistics can be high, and thus we did not use population statistic at test time. Rather, we normalised using batch-statistics instead, and this consistently lead to better predictive performance. We also used the Gumbel-Softmax approximation [274] at test-time using the temperature value  $\tau$  that corresponded to the iteration in  $\tau$  annealing schedule.

#### 4.5.5 CNN architecture details

We include schematics and details of the single-task VGG11 [275] and HighResNet [276] networks in Fig. 4.6. In this work, we constructed multi-task architectures by augmenting these networks with the proposed SFG modules. We used the PReLU activation function [284] in all networks. For the residual blocks used in the HighResNet networks in Fig. 4.6 (ii), we applied PReLU and batch-norm as pre-activation [285] to the convolutional layers. The SFG module was used to cluster the kernels in every coloured layer in Fig. 4.6, and distinct sets of additional transformations (pooling operations for VGG and high-res blocks for HighResNet) were applied to the outputs of the respective filter groups  $G_1, G_2, G_s$ . For a fair comparison, the CS units [254] were added to the same set of layers.

For clarification, the SFG layer number  $n$  (e.g. SFG layer 2) corresponds to the  $n^{th}$  layer with an SFG module. In the case of SFG-VGG11, each convolutional layer uses SFGs. The SFG layer number thus corresponds with layer number in the network. In the case of SFG-HighResNet, not every convolutional layer uses SFGs such as those within residual blocks. Consequently, SFG layer 1 corresponds to layer 1, SFG layer 2 is layer 6, SFG layer 3 is layer 11, SFG layer 4 is layer 16 and SFG layer 5 is layer 17.

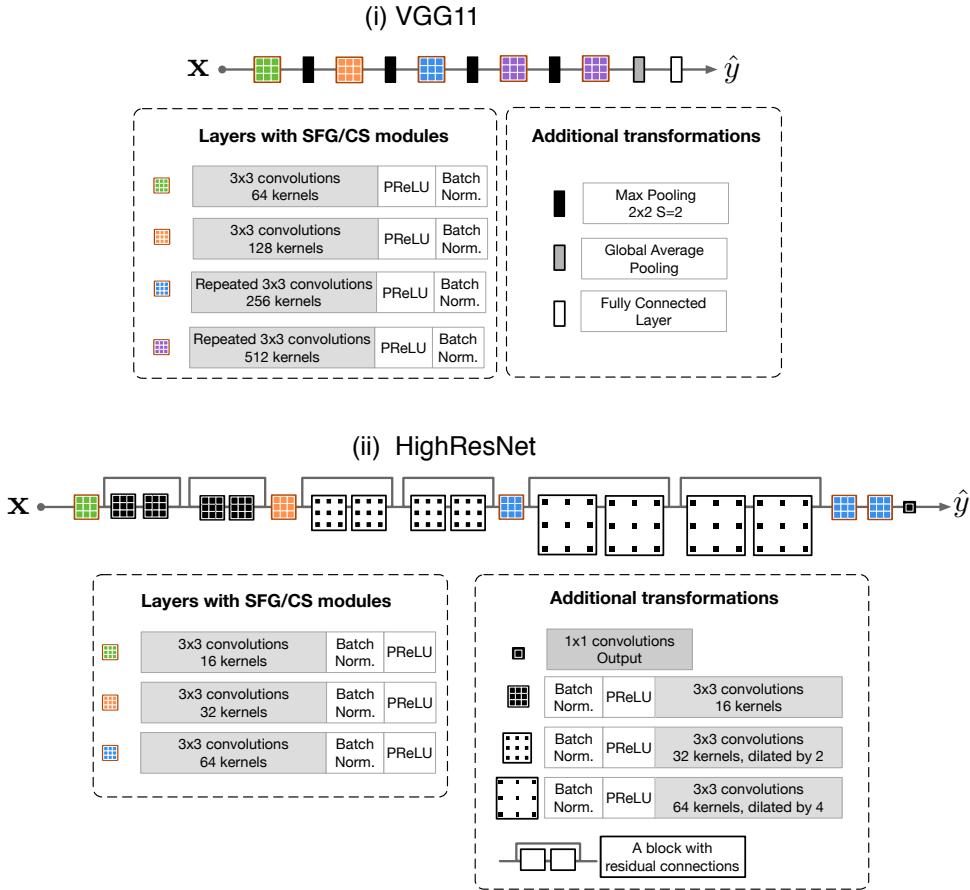


Figure 4.6: Illustration of the single-task architectures, (i) VGG11 and (ii) HighResNet used for UTKFace and medical imaging dataset, respectively. In each architecture, the coloured components indicate the layers to which SFG or cross-stitch (CS) modules are applied when extended to the multi-task learning scenario, whilst the components in black denote the additional transformations applied to the outputs of respective filter groups or CS operations (see the description of black circles in the schematic provided in Fig. 5 of the main text)

## 4.6 Results

### 4.6.1 Age regression and gender prediction

Results on age prediction and gender classification on both datasets are presented in Tab. 4.1a and 4.1b. Our model (MT-SFG) achieved the best performance in comparison to the baselines in both data regimes. In both sets of experiments, our model outperformed the hard-parameter sharing (*MT-hard sharing*) and constant allocation (*MT-constant mask*). This demonstrates the advantage of learning to allocate kernels. In the *MT-constant mask* model, kernels are equally allocated across groups. In contrast, our model is able to allocate kernels in varying proportions across different layers in the network (Fig. 4.8 - SFG-VGG11) to maximise inductive transfer. Moreover, our methods performed better than a model with constant, non-informative grouping probabilities (*MT-constant p= [1/3, 1/3, 1/3]*), displaying the importance of learning structured representations and connectivity across layers to yield good predictions.

### 4.6.2 Image regression and semantic segmentation

Results on CT image synthesis and organ segmentation from input MRI scans is detailed in Tab. 4.2. Our method obtains equivalent (non-statistically significant different) results to the Cross-Stitch network [254] on both tasks. We have, however, observed best synthesis performance in the bone regions (femur

(a) Full training data			(b) Small training data		
Method	Age (MAE)	Gender (Accuracy)	Method	Age (MAE)	Gender (Accuracy)
One-task (VGG11) [275]	7.32	90.70	One-task (VGG11) [275]	8.79	85.54
MT-hard sharing	7.92	90.60	MT-hard sharing	9.19	85.83
MT-constant mask	7.67	89.41	MT-constant mask	9.02	85.98
MT-constant $\mathbf{p}=[\frac{1}{3}, \frac{1}{3}, \frac{1}{3}]$	6.34	92.10	MT-constant $\mathbf{p}=[\frac{1}{3}, \frac{1}{3}, \frac{1}{3}]$	9.15	86.01
VGG11 Cross Stitch [254]	6.78	90.30	VGG11 Cross Stitch [254]	8.85	83.72
MT-SFG (ours)	<b>6.00</b>	<b>92.46</b>	MT-SFG (ours)	<b>8.54</b>	<b>87.01</b>

Table 4.1: Age regression and gender classification results on UTKFace [264] with (a) the full and (b) limited training set. The best and the second best results are shown in red and blue. The mean absolute error (MAE) is reported for the age prediction and classification accuracy for gender prediction. For our model, we performed 50 stochastic forward passes at test time by sampling the kernels from the approximate posterior  $q_\phi(\mathcal{W})$ . We calculated the average age per subject and obtained gender prediction using the mode of the test-time predictions. We initialised our model with grouping probabilities  $\mathbf{p}=[0.2, 0.6, 0.2]$  for all convolution kernels.

(a) CT Synthesis (PSNR)						
Method	Overall	Bones	Organs	Prostate	Bladder	Rectum
One-task [276]	25.76 (0.80)	30.35 (0.58)	38.04 (0.94)	51.38 (0.79)	33.34 (0.83)	34.19 (0.31)
MT-hard sharing	26.31 (0.76)	31.25 (0.61)	39.19 (0.98)	52.93 (0.95)	34.12 (0.82)	34.15 (0.30)
MT-constant mask	24.43(0.57)	29.10(0.46)	37.24(0.86)	50.48(0.73)	32.29(1.01)	33.44(2.88)
MT-constant $\mathbf{p}=[\frac{1}{3}, \frac{1}{3}, \frac{1}{3}]$	26.64(0.54)	31.05 (0.55)	39.11 (1.00)	<b>53.20 (0.86)</b>	34.34 (1.35)	35.61 (0.35)
Cross Stitch [254]	<b>27.86 (1.05)</b>	<b>32.27 (0.55)</b>	40.45 (1.27)	<b>54.51 (1.01)</b>	<b>36.81 (0.92)</b>	36.35 (0.38)
MT-SFG (ours)	<b>27.74 (0.96)</b>	<b>32.29 (0.59)</b>	<b>39.93 (1.09)</b>	53.01 (1.06)	35.65 (0.44)	35.65 (0.37)

(b) Segmentation (DICE)						
Method	Overall	Left Femur Head	Right Femur Head	Prostate	Bladder	Rectum
One-task [276]	0.848(0.024)	0.931 (0.012)	<b>0.917 (0.013)</b>	0.913 (0.013)	0.739 (0.060)	0.741 (0.011)
MT-hard sharing	0.829(0.023)	<b>0.933 (0.009)</b>	0.889 (0.044)	0.904 (0.016)	0.685 (0.036)	0.732 (0.014)
MT-constant mask	0.774(0.065)	0.908 (0.012)	0.911 (0.015)	0.806 (0.0541)	0.583 (0.178)	0.662 (0.019)
MT-constant $\mathbf{p}=[\frac{1}{3}, \frac{1}{3}, \frac{1}{3}]$	0.752(0.056)	0.917 (0.004)	<b>0.917 (0.01)</b>	0.729 (0.086)	0.560 (0.180)	0.639 (0.012)
Cross Stitch [254]	<b>0.854 (0.036)</b>	0.923 (0.008)	0.915 (0.013)	<b>0.933 (0.009)</b>	<b>0.761 (0.053)</b>	<b>0.737 (0.015)</b>
MT-SFG (ours)	<b>0.852(0.047)</b>	<b>0.935 (0.007)</b>	0.912 (0.013)	<b>0.923 (0.016)</b>	0.750 (0.062)	<b>0.758 (0.011)</b>

Table 4.2: Performance on the medical imaging dataset with best results in red, and the second best results in blue. The PSNR is reported for the CT-synthesis (synCT) across the whole volume (overall), at the bone regions, across all organ labels and individually at the prostate, bladder and rectum. For the segmentation, the average DICE score per patient across all semantic labels is computed. The standard deviations are computed over the test subject cohort. For our model, we perform 50 stochastic forward passes at test-time by sampling the kernels from the approximated posterior distribution  $q_\phi(\mathcal{W})$ . We compute the average of all passes to obtain the synCT and calculate the mode of the segmentation labels for the final segmentation. We initialised our model with grouping probabilities  $\mathbf{p}=[0.2, 0.6, 0.2]$ . Red cells indicate best performing and blue cells indicate second best models. **Lastly, we note that the results here are not directly comparable to Table 3.1 due to some differences in preprocessing and data.**

heads and pelvic bone region) in our model when compared against all the baselines, including Cross-Stitch. The bone voxel intensities are the most difficult to synthesise from an input MR scan as task uncertainty in the MR to CT mapping at the bone is often highest [104]. Our model was able to disentangle features specific to the bone intensity mapping (Fig. 4.7) without supervision of the pelvic location, which allowed it to learn a more accurate mapping of an intrinsically difficult task.

### 4.6.3 Learned architectures

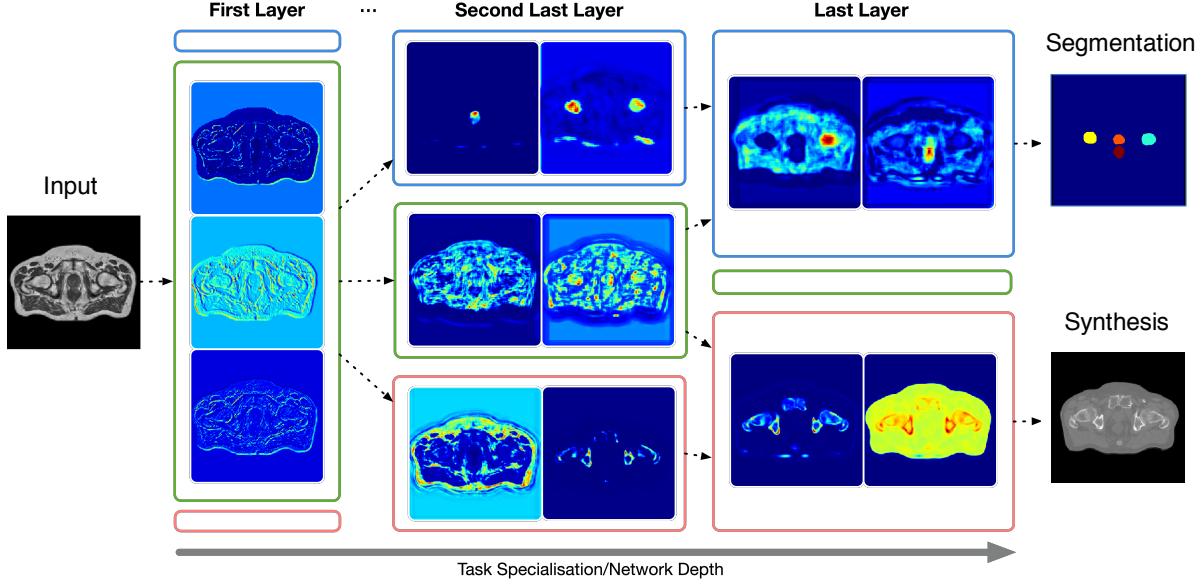


Figure 4.7: Activation maps from example kernels in the learned task-specific and shared filter groups,  $G_1^{(l)}, G_2^{(l)}, G_s^{(l)}$  (enclosed in blue, green and pink funnels) in the first, the second last and the last convolution layers in the SFG-HighResNet model trained on the medical imaging dataset. The results from convolution kernels with low entropy (i.e. high ‘confidence’) of group assignment probabilities  $\mathbf{p}^{(l)}$  are shown for the respective layers.

Analysis of the grouping probabilities of a network embedded with SFG modules permits visualisation of the network connectivity and thus the learned MTL architecture. To analyse the group allocation of kernels at each layer, we computed the sum of class-wise probabilities per layer. Learned groupings for both SFG-VGG11 network trained on UTKFace and the SFG-HighResNet network trained on prostate scans are presented in Fig. 4.8. These figures illustrate increasing task specialisation in the kernels with network depth.

Such behaviour can be also visually confirmed by looking at the activation maps of specialist and generalist kernels in the respective layers. To classify each kernel according to the group (task 1, task 2 or shared), we selected the group with the respective maximum assignment probability. The corresponding activation maps for various input images in the medical imaging dataset are then visualised.

We first analysed the activation maps generated by kernels with low entropy of  $\mathbf{p}$  (i.e. highly confident group assignment). Fig. 4.7 shows that, at the first layer, all kernels are classified as shared ( $\mathbf{p} = [0, 1, 0]$ ), and account for low-order features such as edges or contrast of the images. On the other hand, at deeper layers, higher-order representations are learned, which describe various salient features specific to the tasks such as organs for segmentation, and bones for CT-synthesis. Note that the bones are generally the most difficult region to synthesise CT intensities from an input MR scan [104]. More examples are given in Fig. 4.9.

In addition, we also looked at activation maps from kernels with high entropy of  $\mathbf{p}$  (i.e. highly uncertain group assignment) in Fig. 4.10. In contrast to Fig. 4.9, the learned features do not appear to capture any meaningful structures for both synthesis and segmentation tasks. Of particular note is the dead kernel in the top row of the figure; displaying that a high uncertainty in group allocation correlates with non-informative features.

Notably, the learned connectivity of both models shows striking similarities to hard-parameter sharing architectures commonly used in MTL. Generally, there is a set of shared layers, which aim to learn a feature set common to both tasks. Task-specific branches then learn a mapping from this feature space for task-specific predictions. Our models are able to automatically learn this structure whilst allowing asymmetric allocation of task-specific kernels with no priors on the network structure.

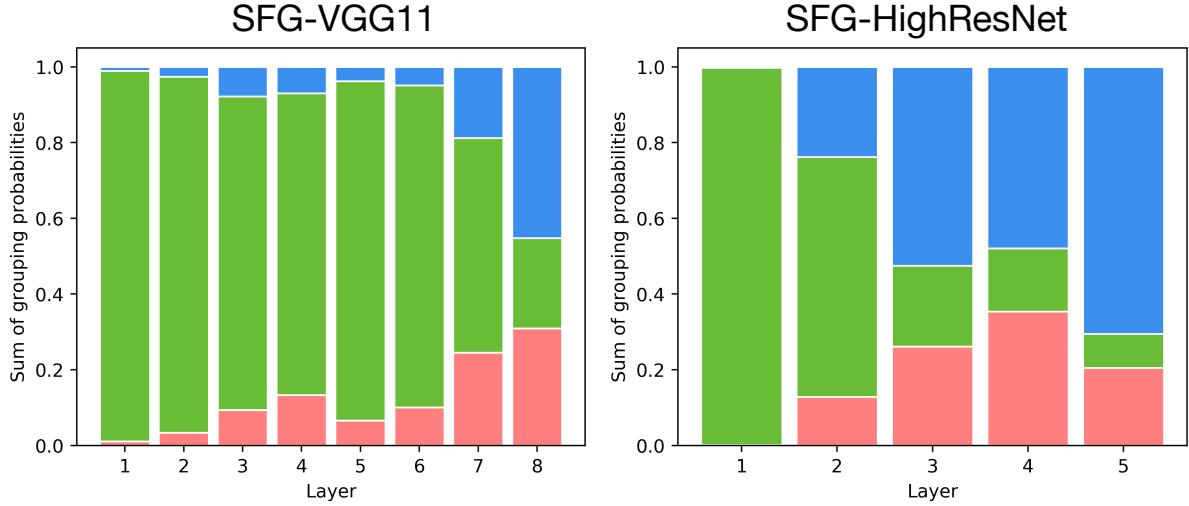


Figure 4.8: Learned kernel grouping in a) SFG-VGG11 network on UTKFace and b) SFG-HighResNet on medical scans. The proportions of task-1, shared and task-2 filter groups are shown in blue, green and pink. Within SFG-VGG11, task-1 age regression and task-2 is gender classification. For SFG-HighResNet, task-1 is CT synthesis and task-2 is organ segmentation.

#### 4.6.4 Learned grouping probability plots

Fig. 4.11 and Fig. 4.12 illustrate the density plots of the learned grouping probabilities  $\mathbf{p}$  for each trained network. We also plot the training trajectories of grouping probabilities  $\mathbf{p}$  of all kernels in each layer. These are colour coded by iteration number—blue for low and yellow for high iteration number. This shows that some grouping probabilities are quickly learned in comparison to others.

Fig. 4.11 and Fig. 4.12 show that most kernels are in the shared group at earlier layers of the network where mostly low-order generic features are learned (as illustrated in Fig. 4.9, SFG layer 1). They converge quickly to the shared vertex of the 2-simplex as evidenced by the colour of the trajectory plots. As the network depth increases, task-specialisation in the kernels increases (see Fig. 4.9, SFG layer  $\geq 4$ ). This is illustrated by high density clusters at task-specific vertices and by the trajectory plots.

#### 4.6.5 Effect of $\mathbf{p}$ initialisation

Fig. 4.13 shows the layer-wise proportion of the learned kernel groups on the UTKFace dataset for four different initialization schemes of grouping probabilities  $\mathbf{p}$ : (i) “dominantly shared”, with  $\mathbf{p} = [0.2, 0.6, 0.2]$ , (ii) “dominantly task-specific”, with  $\mathbf{p} = [0.45, 0.1, 0.45]$ , (iii) “random”, where  $\mathbf{p}$  is drawn from  $\text{Dirichlet}(1, 1, 1)$ , (iv) “start with MT-constant mask”, where an equal number of kernels in each layer are set to probabilities of  $\mathbf{p} = [1, 0, 0], [0, 1, 0]$  and  $[0, 0, 1]$ . In all cases, the same set of hyperparameters, including the annealing rate of the temperature term in GSM approximation and the coefficient of the entropy regularizer  $\mathcal{H}(\mathbf{p})$ , were used during training. We observe that the kernel grouping of respective layers in (i), (ii) and (iii) all converge to a very similar configuration observed in Sec. 4.6.3, highlighting the robustness of our method to different initialisations of  $\mathbf{p}$ . In case (iv), the learning of  $\mathbf{p}$  were much slower than the remaining cases, due to weaker gradients, and we speculate that a higher entropy regularizer is necessary to facilitate its convergence.

#### 4.6.6 Learned filter groups on duplicate tasks

We analysed the dynamics of a network with SFG modules when trained with two duplicates of the same CT regression task (instead of two distinct tasks). Fig. 4.14 visualises the learned grouping and trajectories of the grouping probabilities during training. In the first 3 SFG layers (layers 1, 6 and 11 of the network), all the kernels are grouped as shared. In the penultimate SFG layer (layer 16), either kernels are grouped as shared or with probability  $\mathbf{p}=[1/2, 0, 1/2]$ , signifying that the kernels can belong

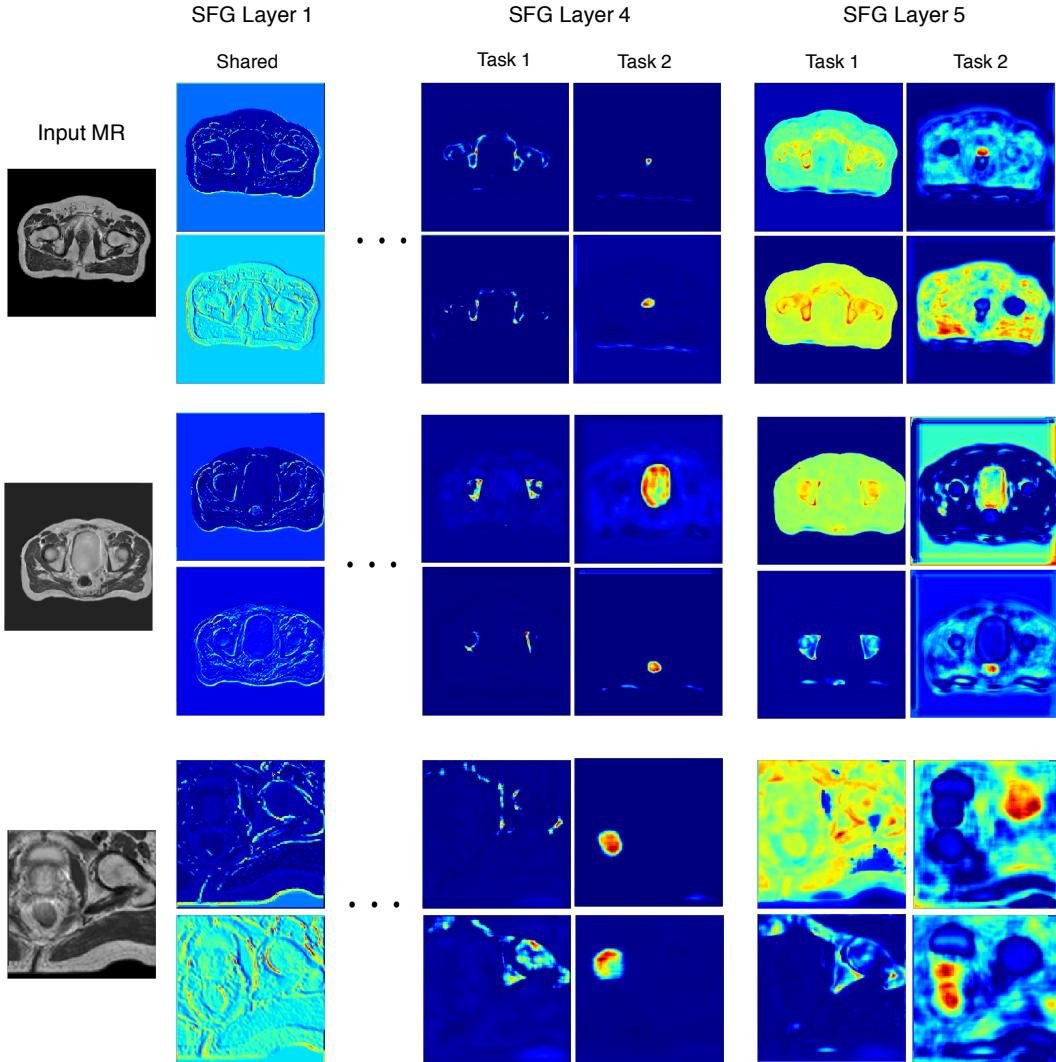


Figure 4.9: Example activations for kernels with low entropy of  $\mathbf{p}$  (i.e. group assignment with high confidence) for three input MR slices in the SFG-HighResNet multi-task network. Columns “Shared”, “Task 1” & “Task 2” display the results from the shared, CT-synthesis and organ-segmentation specific filter groups in respective layers. We illustrate activations stratified by group in layer 1 (SFG layer 1), layer 16 (SFG layer 4) and layer 17 (SFG layer 5).

to either task. The final SFG layer (layer 17) shows that most kernels have probabilities  $\mathbf{p}=[1/3, 1/3, 1/3]$ . Kernels thus have equal probability of being task-specific or shared. This is expected as we are training on duplicate tasks and therefore the kernels are equally likely to be useful across all groups.

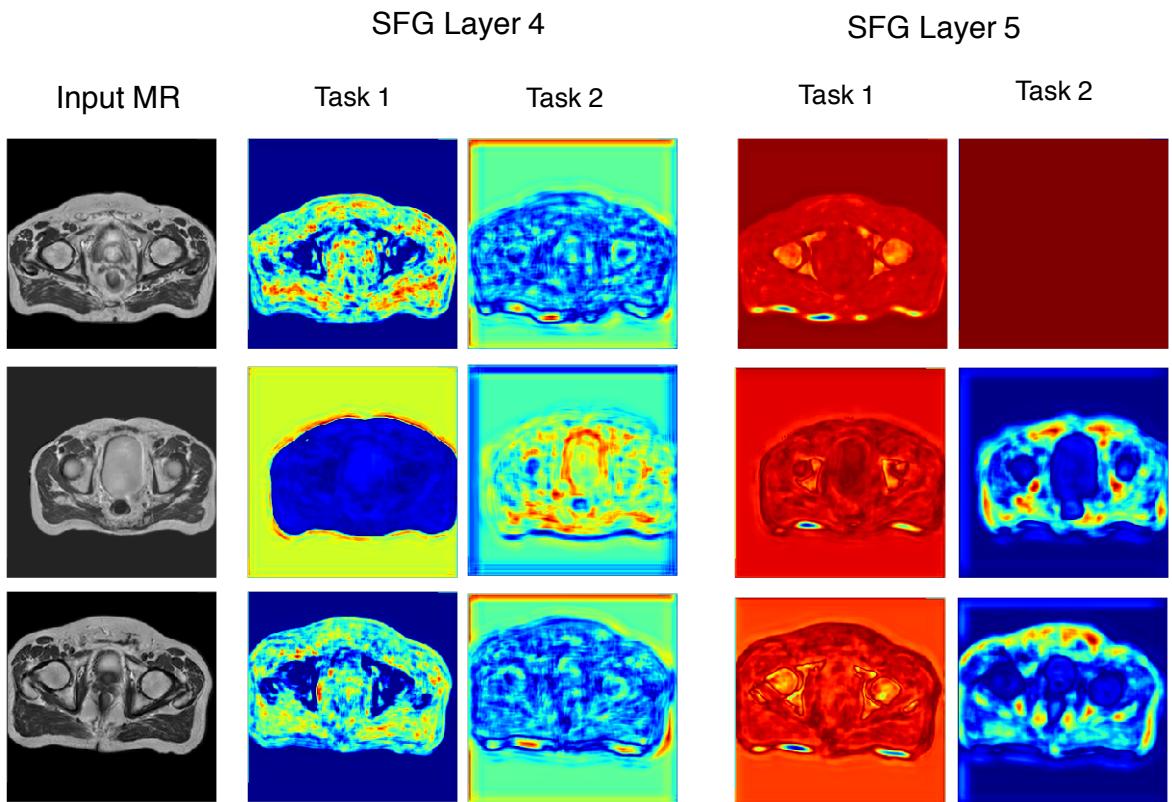


Figure 4.10: Example activations for kernels with high entropy (i.e. group assignment with low confidence) for three input MR slices in the SFG-HighResNet multi-task network. Columns “Shared”, “Task 1” & “Task 2” display the results from the shared, CT-synthesis and organ-segmentation specific filter groups in respective layers. We illustrate activations stratified by group in layer 16 (SFG layer 4) and layer 17 (SFG layer 5).

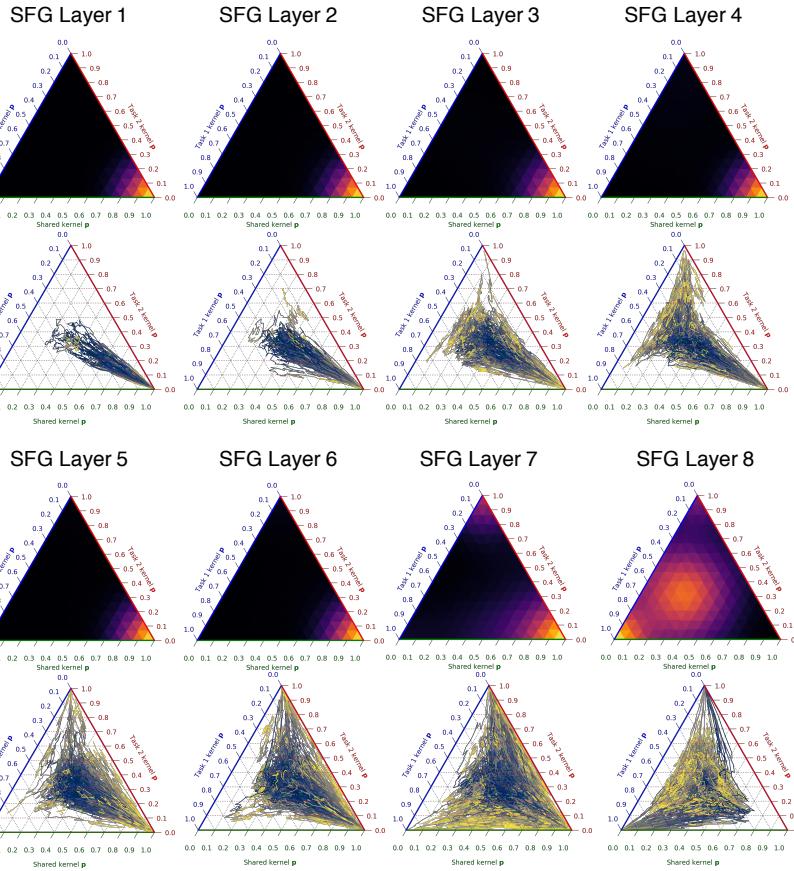


Figure 4.11: Density plots and trajectory plots of the learned grouping probabilities for the SFG-VGG11 architecture. The density plots represent the final learned probabilities per layer for each kernel. The trajectory plots represent how the grouping probabilities are learned during training and thus how the connectivity is determined. Histograms of the grouping probabilities were smoothed with a Gaussian kernel with  $\sigma = 1$ . The densities are mapped to and visualised in the 2-simplex using `python-ternary` [286]. This figure represents that as the network depth increases, task-specialisation in the kernels increases.

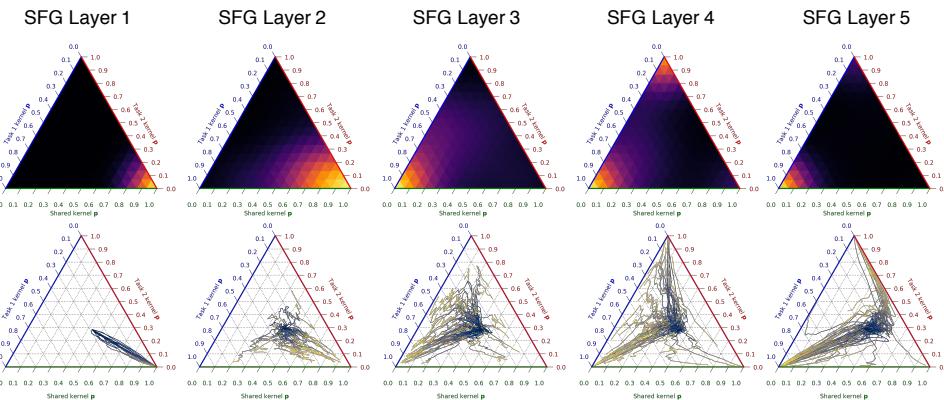


Figure 4.12: Density plots and trajectory plots of the learned grouping probabilities for the SFG-HighResNet architecture. The density plots represent the final learned probabilities per layer for each kernel. The trajectory plots represent how the grouping probabilities are learned during training and thus how the connectivity is determined.

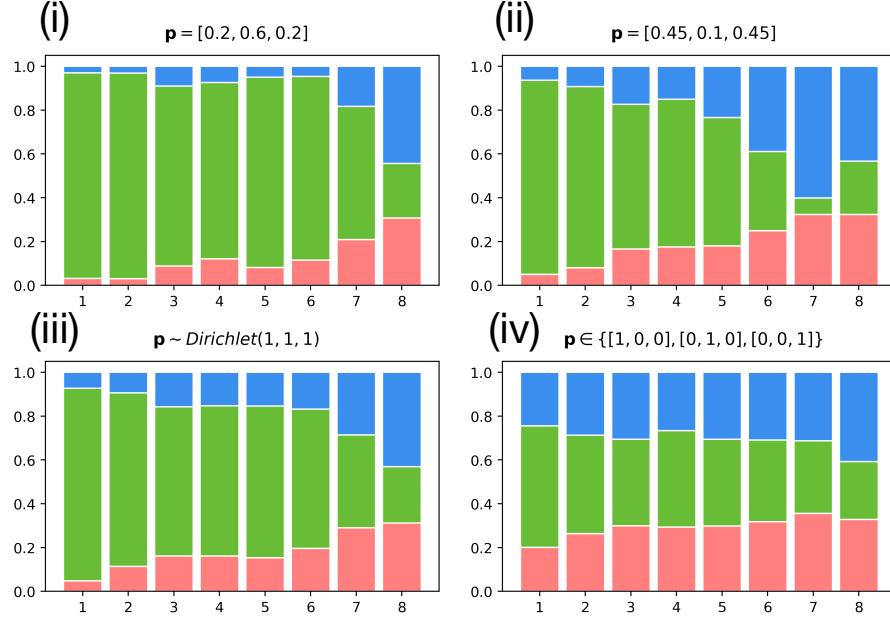


Figure 4.13: Effect of the initial values of grouping probabilities  $\mathbf{p}$  on the learned kernel allocation after convergence for the medical imaging dataset. The proportions of task-1, shared and task-2 filter groups are shown in blue, green and pink.

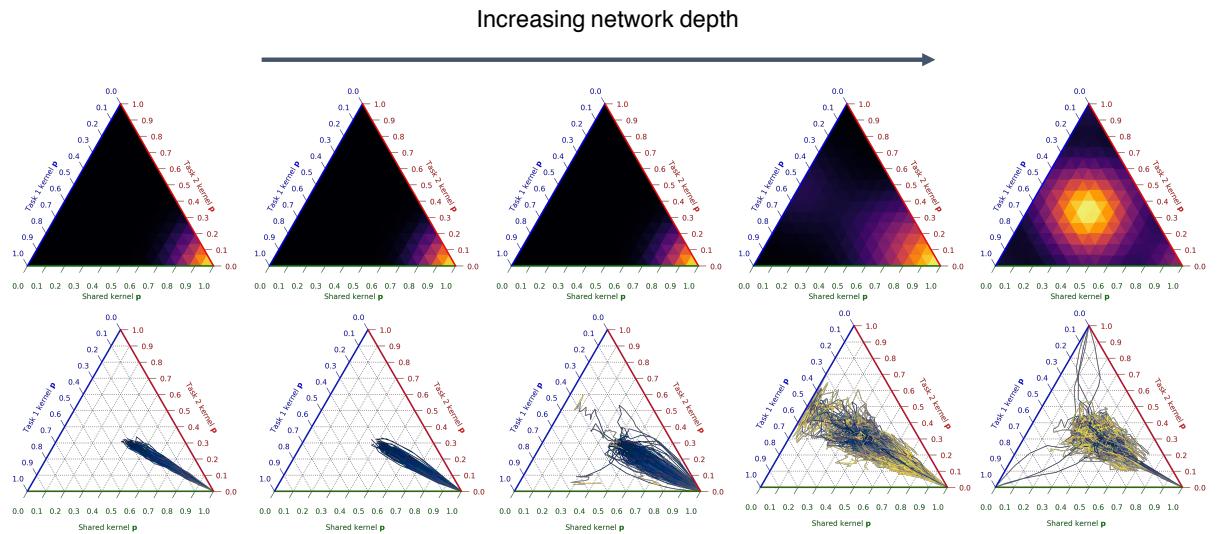


Figure 4.14: Top: density plots for the learned grouping probabilities at each SFG layer in a model where we trained on duplicate tasks i.e. task 1 is CT synthesis and task 2 is also CT synthesis. Bottom: trajectories of the grouping probabilities during training.

## 4.7 Discussion

In this chapter, we have proposed *stochastic filter groups* (SFGs) to disentangle *task-specific* and *generalist* features. SFGs probabilistically defines the grouping of kernels and thus the connectivity of features in a CNNs. We use variational inference to approximate the distribution over connectivity given training data and sample over possible architectures during training. Our method can be considered as a probabilistic form of multi-task architecture learning [287], as the learned posterior embodies the optimal MTL architecture given the data.

Our model learns structure in the representations. The learned shared (generalist) features may be exploited either in a transfer learning or continual learning scenario. As seen in [288], an effective prior learned from multiple tasks can be a powerful tool for learning new, unrelated tasks. Our model consequently offers the possibility to exploit the learned task-specific and generalist features when faced with situations where a third task is needed, which may suffer from unbalanced or limited training data. This is particularly relevant in the medical field, where training data is expensive to acquire as well as laborious. We will investigate this in further work.

Lastly, a network composed of SFG modules can be seen as a superset of numerous MTL architectures. Depending on the data and the analysed problem, SFGs can recover many different architectures such as single task networks, traditional hard-parameter sharing, equivalent allocation across tasks, and asymmetrical grouping (Fig. 4.3). Note, however, that the proposed SFG module only learns connectivity between neighbouring layers. Non-parallel ordering of layers, a crucial concept of MTL models [260, 259], was not investigated. Future work will look to investigate the applicability of SFG modules for learning connections across grouped kernels between non-neighbouring layers.

# Chapter 5

## How to Learn Network Architecture like a Decision Tree

**Abstract:** Deep neural networks and decision trees operate on largely separate paradigms; typically, the former performs representation learning with pre-specified architectures, while the latter is characterised by learning hierarchies over pre-specified features with learned architectures. We unite the two via *adaptive neural trees* (ANTs) that incorporates representation learning into edges, routing functions and leaf nodes of a decision tree, along with a backpropagation-based training algorithm that adaptively grows the architecture from primitive modules (e.g., convolutional layers). We demonstrate that, whilst achieving competitive performance on classification and regression datasets, ANTs benefit from (i) lightweight inference via conditional computation, (ii) hierarchical separation of features useful to the task e.g. learning meaningful class associations, such as separating natural vs. man-made objects, and (iii) a mechanism to adapt the architecture to the size and complexity of the training dataset. This chapter is based on [106].

### 5.1 Introduction

Neural networks (NNs) and decision trees (DTs) are both powerful classes of machine learning models with proven successes in academic and commercial applications. The two approaches, however, typically come with mutually exclusive benefits and limitations, as illustrated in Fig. 5.1.

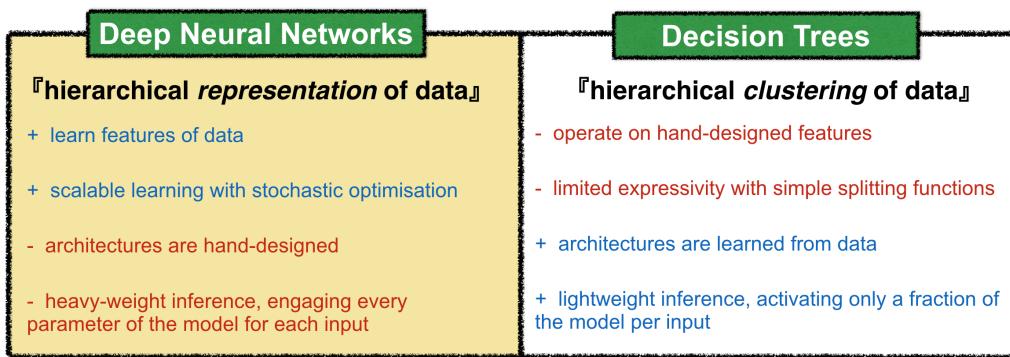


Figure 5.1: Comparison of two machine learning paradigms.

NNs are characterised by learning hierarchical representations of data through the composition of nonlinear transformations [289, 290], which has alleviated the need for feature engineering, in contrast to many other machine learning models. In addition, NNs are trained with stochastic optimisers, such as stochastic gradient descent (SGD), allowing training to scale to large datasets. Consequently, with modern hardware, we can train NNs of many layers on large datasets, solving numerous problems ranging from object detection to speech recognition with unprecedented accuracy [1]. However, their architectures

typically need to be designed by hand and fixed per task or dataset, requiring domain expertise [291]. Inference can also be heavy-weight for large models, as each sample engages every part of the network, i.e., increasing capacity causes a proportional increase in computation [292].

Alternatively, DTs are characterised by learning hierarchical clusters of data [293]. A DT learns how to split the input space, so that in each subset, linear models suffice to explain the data. In contrast to standard NNs, the architectures of DTs are optimised based on training data, and are particularly advantageous in data-scarce scenarios. DTs also enjoy lightweight inference as only a single root-to-leaf path on the tree is used for each input sample. However, successful applications of DTs often require hand-engineered features of data. We can ascribe the limited expressivity of single DTs to the common use of simplistic routing functions, such as splitting on axis-aligned features. The loss function for optimising hard partitioning is non-differentiable, which hinders the use of gradient descent-based optimization and thus complex splitting functions. Current techniques for increasing capacity include ensemble methods such as random forests (RFs) [294] and gradient-boosted trees (GBTs) [295], which are known to achieve state-of-the-art performance in various tasks, including medical applications and financial forecasting [296, 297, 69, 298].

The goal of this work is to combine NNs and DTs to gain the complementary benefits of both approaches. To this end, we propose *adaptive neural trees* (ANTs), which generalise previous work that attempted the similar unification [299, 300, 301, 302, 303, 304, 305] and address their limitations (see Tab. 5.1). ANTs represent routing decisions and root-to-leaf computational paths within the tree structures as NNs, which lets them benefit from hierarchical representation learning, rather than being restricted to partitioning the raw data space. On the other hand, unlike the fully distributed representation of standard NN models, the tree topology of ANTs acts as a strong structural prior that enforces sparse structures by which features are shared and separated in a hierarchical fashion. In addition, we propose a backpropagation-based training algorithm to grow ANTs based on a series of decisions between making the ANT deeper—the central NN paradigm—or partitioning the data—the central DT paradigm (see Fig. 5.2 (Right)). This allows the architectures of ANTs to adapt to the data available. By our design, ANTs inherit the following desirable properties from both DTs and NNs:

- **Representation learning:** as each root-to-leaf path in an ANT is an NN, features can be learned end-to-end with gradient-based optimisation. Combined with the tree structure, an ANT can learn such features which are hierarchically shared and separated. The training algorithm is also amenable to SGD.
- **Architecture learning:** by progressively growing ANTs, the architecture adapts to the availability and complexity of data, embodying Occam’s razor. The growth procedure can be viewed as architecture search with a hard constraint over the model class.
- **Lightweight inference:** at inference time, ANTs perform conditional computation, selecting a single root-to-leaf path on the tree on a per-sample basis, activating only a subset of the parameters of the model.

We empirically validate these benefits for regression and classification through experiments on the SARCOS [306], MNIST [307] and CIFAR-10 [308] datasets. The best performing methods on the SARCOS multivariate regression dataset are all tree-based, with ANTs achieving the lowest mean squared error. On the other hand, along with other forms of neural networks, ANTs far outperform state-of-the-art RF [309] and GBT [310] methods on image classification, with architectures achieving over 99% accuracy on MNIST and over 90% accuracy on CIFAR-10. Our ablations on all three datasets consistently show that the combination of feature learning and data partitioning are required for the best predictive performance of ANTs. In addition, we show that ANTs can learn meaningful hierarchical partitionings of data, e.g., grouping man-made and natural objects (see Fig. 5.3) useful to the end task. ANTs also have reduced time and memory requirements during inference, thanks to such hierarchical structure. In one case, we discover an architecture that achieves over 98% accuracy on MNIST using approximately the same number of parameters as a linear classifier on raw image pixels, showing the benefits of tree-shaped hierarchical sharing and separation of features in enhancing both computational and predictive performance. Finally, we demonstrate the benefits of architecture learning by training ANTs on subsets of CIFAR-10 of varying sizes. The method can construct architectures of adequate size, leading to better generalisation, particularly on small datasets.

Table 5.1: Comparison of tree-structured NNs. The first column denotes if each path on the tree is a NN, and the second column denotes if the routers learn features. The last column shows if the method grows an architecture, or uses a pre-specified one.

Method	Feature learning?		Grown?
	Path	Routers	
SDT [299]	✗	✗	✓
SDT 2 / HME [311]	✗	✓	✗
SDT 3 [300]	✗	✓	✓
SDT 4 [304]	✗	✓	✗
RDT [312]	✗	✓	✗
BT [316]	✗	✓	✓
Conv DT [301]	✗	✓	✗
NDT [302]	✗	✓	✓
NDT 2 [305]	✓	✓	✗
NDF [303]	✓	✓	✗
CNet [315]	✓	✓	✗
<b>ANT (ours)</b>	✓	✓	✓

## 5.2 Related work

Our work is primarily related to research into combining DTs and NNs. Here we explain how ANTs subsume a large body of such prior work as specific cases and address their limitations. We also include additional reviews of work in conditional computation, neural architecture search, and an early form of feature learning with DTs based on cascading.

**Combining Decisions Trees and Neural Networks:** The very first soft decision tree (SDT) introduced in [299] is a specific case where in our terminology the routers are axis-aligned features, the transformers are identity functions, and the routers are static distributions over classes or linear functions. The hierarchical mixture of experts (HMEs) proposed by [311] is a variant of SDTs whose routers are linear classifiers and the tree structure is fixed; [312] recently proposed a more computationally efficient training method that is able to directly optimise hard-partitioning by differentiating through stochastic gradient estimators. More modern SDTs in [302, 301, 304] used multilayer perceptrons (MLPs) or convolutional layers in the routers to learn more complex partitionings of the input space. However, the simplicity of identity transformers used in these methods means that input data is never transformed and thus each path on the tree does not perform representation learning, limiting their performance.

More recent work suggested that integrating non-linear transformations of data into DTs would enhance model performance. The neural decision forest (NDF) [303], which held cutting-edge performance on ImageNet [313] in 2015, is an ensemble of DTs, each of which is also an instance of ANTs where the whole GoogLeNet architecture [314] (except for the last linear layer) is used as the root transformer, prior to learning tree-structured classifiers with linear routers. [305] employed a similar approach with a MLP at the root transformer, and is optimised to minimise a differentiable information gain loss. The conditional network proposed in [315] sparsified CNN architectures by distributing computations on hierarchical structures based on directed acyclic graphs with MLP-based routers, and designed models with the same accuracy with reduced compute cost and number of parameters. However, in all cases, the model architectures are pre-specified and fixed.

In contrast, ANTs satisfy all criteria in Tab. 5.1; they provide a general framework for learning tree-structured models with the capacity of representation learning along each path and within routing functions, and a mechanism for learning its architecture.

Architecture growth is a key facet of DTs [293], and typically performed in a greedy fashion with a termination criteria based on validation set error [299, 300]. Previous works in DT research have made attempts to improve upon this greedy growth strategy. Decision jungles [317] employ a training mechanism to merge partitioned input spaces between different sub-trees, and thus to rectify suboptimal “splits” made due to the locality of optimisation. [316] proposes budding trees, which are grown and

pruned incrementally based on global optimisation of existing nodes. While our training algorithm, for simplicity, grows the architecture by greedily choosing the best option between going “deeper” and “splitting” the input space (see Fig. 5.2), it is certainly amenable to these advances.

**Conditional Computation:** in NNs, computation of each sample engages every parameter of the model. In contrast, DTs route each sample to a single path, only activating a small fraction of the model. [290] advocated for this notion of conditional computation to be integrated into NNs, and this has become a topic of growing interest. Rationales for using conditional computation ranges from attaining better capacity-to-computation ratio [292, 318, 319, 320] to adapting the required computation to the difficulty of the input and task [319, 321, 322, 323, 324, 325]. We view the growth procedure of ANTs as having a similar motivation with the latter—processing raw pixels is suboptimal for computer vision tasks, but we have no reason to believe that the hundreds of convolutional layers in current state-of-the-art architectures [326, 327] are necessary either. Growing ANTs adapts the architecture complexity to the dataset as a whole, with routers determining the computation needed on a per-sample basis.

**Neural Architecture Search:** the ANT growing procedure is related to the progressive growing of NNs [328, 329, 330, 331, 332, 333, 334, 335], or more broadly, the field of neural architecture search [291, 336, 337]. This approach, mainly via greedy layerwise training, has historically been one solution to optimising NNs [328, 329]. However, nowadays it is possible to train NNs in an end-to-end fashion. One area which still uses progressive growing is lifelong learning, in which a model needs to adapt to new tasks while retaining performance on previous ones [330, 333]. In particular, [330] introduced a method that grows a tree-shaped network to accommodate new classes. However, their method never transforms the data before passing it to the children classifiers, and hence never benefit from the parent’s representations.

Whilst we learn the architecture of an ANT in a greedy, layerwise fashion, several other methods search globally. Based on a variety of techniques, including evolutionary algorithms [338, 339], reinforcement learning [291], sequential optimisation [340] and boosting [337], these methods find extremely high-performance yet complex architectures. In our case, we constrain the search space to simple tree-structured NNs, retaining desirable properties of DTs such as data-dependent computation and interpretable structures, while keeping the space and time requirement of architecture search tractable thanks to the locality of our growth procedure.

**Cascaded trees and forests:** another noteworthy strand of work for feature learning with tree-structured models is cascaded forests—stacks of RFs where the outputs of intermediate models are fed into the subsequent ones [73, 341, 309]. It has been shown how a cascade of DTs can be mapped to NNs with sparse connections [342], and more recently [343] extended this argument to RFs. However, the features obtained in this approach are the intermediate outputs of respective component models, which are not optimised for the target task, and cannot be learned end-to-end, thus limiting its representational quality. Recently, [344] introduced a method to jointly train a cascade of gradient boosted trees (GBTs) to improve the limited representation learning ability of such previous work. A variant of target propagation [345] was designed to enable the end-to-end training of cascaded GBTs, each of which is non-differentiable and thus not amenable to back-propagation.

## 5.3 Adaptive Neural Trees

We now formalise the definition of Adaptive Neural Trees (ANTs), which are a form of DTs enhanced with deep, learned representations. We focus on supervised learning, where the aim is to learn the conditional distribution  $p(\mathbf{y}|\mathbf{x})$  from a set of  $N$  labelled samples  $(\mathbf{x}^{(1)}, \mathbf{y}^{(1)}), \dots, (\mathbf{x}^{(N)}, \mathbf{y}^{(N)}) \in \mathcal{X} \times \mathcal{Y}$  as training data.

### 5.3.1 Model Topology and Operations

In short, an ANT is a tree-structured model, characterized by a set of hierarchical partitions of the input space  $\mathcal{X}$ , a series of nonlinear transformations, and separate predictive models in the respective

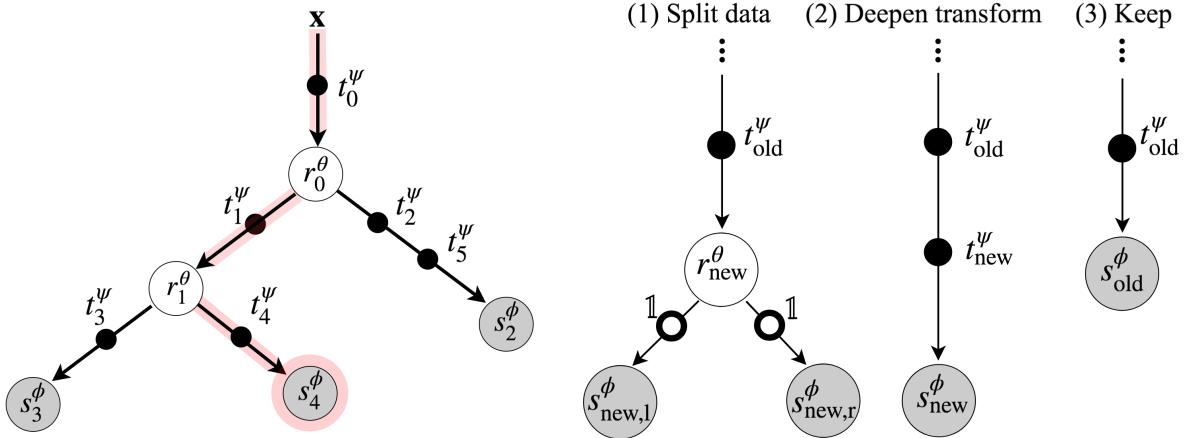


Figure 5.2: (Left). An example ANT. Data is passed through transformers (black circles on edges), routers (white circles on internal nodes), and solvers (gray circles on leaf nodes). The red shaded path shows routing of  $\mathbf{x}$  to reach leaf node 4. Input  $\mathbf{x}$  undergoes a series of selected transformations  $\mathbf{x} \rightarrow \mathbf{x}_0^\psi := t_0^\psi(\mathbf{x}) \rightarrow \mathbf{x}_1^\psi := t_1^\psi(\mathbf{x}_0^\psi) \rightarrow \mathbf{x}_4^\psi := t_4^\psi(\mathbf{x}_1^\psi)$  and the solver module yields the predictive distribution  $p_{4,\psi}^\phi(\mathbf{y}) := s_4^\phi(\mathbf{x}_4^\psi)$ . The probability of selecting this path is given by  $\pi_2^{\psi,\theta}(\mathbf{x}) := r_0^\theta(\mathbf{x}_0^\psi) \cdot (1 - r_1^\theta(\mathbf{x}_1^\psi))$ . (Right). Three growth options at a given node: *split data*, *deepen transform* & *keep*. The small white circles on the edges denote identity transformers.

component regions. More formally, we define an ANT as a pair  $(\mathbb{T}, \mathcal{O})$  where  $\mathbb{T}$  defines the model topology, and  $\mathcal{O}$  denotes the set of operations on it.

Table 5.2: Primitive module specifications for MNIST, CIFAR-10 and SARCOS datasets. “conv5-40” denotes a 2D convolution with 40 kernels of spatial size  $5 \times 5$ . “GAP”, “FC”, “LC” and “LR” stand for global-average-pooling, fully connected layer, linear classifier and linear regressor. “Downsample Freq” denotes the frequency at which  $2 \times 2$  max-pooling is applied.

Model	Router, $\mathcal{R}$	Transformer, $\mathcal{T}$	Solver, $\mathcal{S}$	Downsample Freq.
ANT-SARCOS	$1 \times \text{FC} + \text{Sigmoid}$	$1 \times \text{FC} + \tanh$	LR	0
ANT-MNIST-A	$1 \times \text{conv5-40} + \text{GAP} + 2 \times \text{FC} + \text{Sigmoid}$	$1 \times \text{conv5-40} + \text{ReLU}$	LC	1
ANT-MNIST-B	$1 \times \text{conv3-40} + \text{GAP} + 2 \times \text{FC} + \text{Sigmoid}$	$1 \times \text{conv3-40} + \text{ReLU}$	LC	2
ANT-MNIST-C	$1 \times \text{conv5-5} + \text{GAP} + 2 \times \text{FC} + \text{Sigmoid}$	$1 \times \text{conv5-5} + \text{ReLU}$	LC	2
ANT-CIFAR10-A	$2 \times \text{conv3-128} + \text{GAP} + 1 \times \text{FC} + \text{Sigmoid}$	$2 \times \text{conv3-128} + \text{ReLU}$	LC	1
ANT-CIFAR10-B	$2 \times \text{conv3-96} + \text{GAP} + 1 \times \text{FC} + \text{Sigmoid}$	$2 \times \text{conv3-96} + \text{ReLU}$	LC	1
ANT-CIFAR10-C	$2 \times \text{conv3-72} + \text{GAP} + 1 \times \text{FC} + \text{Sigmoid}$	$2 \times \text{conv3-72} + \text{ReLU}$	GAP + LC	1

We restrict the model topology  $\mathbb{T}$  to be instances of *binary trees*, defined as a set of graphs whose each node is either an internal node or a leaf, and is the child of exactly one parent node, except the root node at the top. We define the topology of a tree as  $\mathbb{T} := \{\mathcal{N}, \mathcal{E}\}$  where  $\mathcal{N}$  is the set of all nodes, and  $\mathcal{E}$  is the set of edges between them. Nodes with no children are leaf nodes,  $\mathcal{N}_{leaf}$ , and all others are internal nodes,  $\mathcal{N}_{int}$ . Every internal node  $j \in \mathcal{N}_{int}$  has exactly two children nodes, represented by  $\text{left}(j)$  and  $\text{right}(j)$ . Unlike standard trees,  $\mathcal{E}$  contains an edge which connects input data  $\mathbf{x}$  with the root node, as shown in Fig.5.2 (Left).

Every node and edge is assigned with operations which acts on the allocated samples of data (Fig.5.2). Starting at the root, each sample gets transformed and traverses the tree according to the set of operations  $\mathcal{O}$ . An ANT is constructed based on three primitive modules of differentiable operations:

1. **Routers,  $\mathcal{R}$ :** each internal node  $j \in \mathcal{N}_{int}$  holds a *router* module,  $r_j^\theta : \mathcal{X}_j \rightarrow [0, 1] \in \mathcal{R}$ , parametrised by  $\theta$ , which sends samples from the incoming edge to either the left or right child. Here  $\mathcal{X}_j$  denotes the representation at node  $j$ . We use *stochastic routing*, where the decision (1 for the left and 0 for the right branch) is sampled from Bernoulli distribution with mean  $r_j^\theta(\mathbf{x}_j)$  for input  $\mathbf{x}_j \in \mathcal{X}_j$ . As an example,  $r_j^\theta$  can be defined as a small CNN.

2. **Transformers**,  $\mathcal{T}$ : every edge  $e \in \mathcal{E}$  of the tree has one or a composition of multiple *transformer* module(s). Each transformer  $t_e^\psi \in \mathcal{T}$  is a nonlinear function, parametrised by  $\psi$ , that transforms samples from the previous module and passes them to the next one. For example,  $t_e^\psi$  can be a single convolutional layer followed by ReLU [346]. Unlike in standard DTs, edges transform data and are allowed to “grow” by adding more operations (Sec. 5.4), learning “deeper” representations as needed.
3. **Solvers**,  $\mathcal{S}$ : each leaf node  $l \in \mathcal{N}_{leaf}$  is assigned to a *solver* module,  $s_l^\phi : \mathcal{X}_l \rightarrow \mathcal{Y} \in \mathcal{S}$ , parametrised by  $\phi$ , which operates on the transformed input data and outputs an estimate for the conditional distribution  $p(\mathbf{y}|\mathbf{x})$ . For classification tasks, we can define, for example,  $s_l^\phi$  as a linear classifier on the feature space  $\mathcal{X}_l$ , which outputs a distribution over classes.

Defining operations on the graph  $\mathbb{T}$  amounts to a specification of the triplet  $\mathbb{O} = (\mathcal{R}, \mathcal{T}, \mathcal{S})$ . For example, given image inputs, we would choose the operations of each module to be from the set of operations commonly used in CNNs (examples are given in Tab. 5.2). In this case, every computational path on the resultant ANT, as well as the set of routers that guide inputs to one of these paths, are given by CNNs. Lastly, many existing tree-structured models [299, 300, 301, 302, 303, 304, 305] are instantiations of ANTs with limitations which we will address with our model (see Sec. 5.2 for a more detailed discussion).

### 5.3.2 Probabilistic Model and Inference

An ANT ( $\mathbb{T}, \mathbb{O}$ ) models the conditional distribution  $p(\mathbf{y}|\mathbf{x})$  as a hierarchical mixture of experts (HMEs) [311], each of which is defined as an NN and is a root-to-leaf path in the tree. Standard HMEs are a special case of ANTs where transformers are the identity function. As a result, the representations within experts are hierarchically shared between similar experts, unlike the independent representations within experts in standard HMEs. In addition, ANTs come with a growth mechanism to determine the number of needed experts and their complexity, as discussed in Sec. 5.4.

Each input to the ANT,  $\mathbf{x}$ , stochastically traverses the tree based on decisions of routers and undergoes a sequence of transformations until it reaches a leaf node where the corresponding solver predicts the label  $\mathbf{y}$ . Suppose we have  $L$  leaf nodes, the full predictive distribution, with parameters  $\Theta = (\boldsymbol{\theta}, \boldsymbol{\psi}, \boldsymbol{\phi})$ , is given by

$$p(\mathbf{y}|\mathbf{x}, \Theta) = \sum_{l=1}^L \underbrace{p(z_l = 1|\mathbf{x}, \boldsymbol{\theta}, \boldsymbol{\psi})}_{\text{Leaf-assignment prob. } \pi_l^{\boldsymbol{\theta}, \boldsymbol{\psi}}} \underbrace{p(\mathbf{y}|\mathbf{x}, z_l = 1, \boldsymbol{\phi}, \boldsymbol{\psi})}_{\text{Leaf-specific prediction. } s_l^{\boldsymbol{\phi}, \boldsymbol{\psi}}} \quad (5.1)$$

where  $\mathbf{z} \in \{0, 1\}^L$  is an  $L$ -dimensional binary latent variable such that  $\sum_{l=1}^L z_l = 1$ , which describes the choice of leaf node (e.g.  $z_l = 1$  means that leaf  $l$  is used). Here  $\boldsymbol{\theta}, \boldsymbol{\psi}, \boldsymbol{\phi}$  summarise the parameters of router, transformer and solver modules in the tree. The mixing coefficient  $\pi_l^{\boldsymbol{\theta}, \boldsymbol{\psi}}(\mathbf{x}) := p(z_l = 1|\mathbf{x}, \boldsymbol{\psi}, \boldsymbol{\theta})$  quantifies the probability that  $\mathbf{x}$  is assigned to leaf  $l$  and is given by a product of decision probabilities over all router modules on the unique path  $\mathcal{P}_l$  from the root to leaf node  $l$ :

$$\pi_l^{\boldsymbol{\psi}, \boldsymbol{\theta}}(\mathbf{x}) = \prod_{r_j^{\boldsymbol{\theta}} \in \mathcal{P}_l} r_j^{\boldsymbol{\theta}}(\mathbf{x}_j^{\boldsymbol{\psi}})^{\mathbb{1}_{l \prec j}} \cdot (1 - r_j^{\boldsymbol{\theta}}(\mathbf{x}_j^{\boldsymbol{\psi}}))^{\mathbb{1}_{l \not\prec j}} \quad (5.2)$$

where  $l \prec j$  is a binary relation and is only true if leaf  $l$  is in the left subtree of internal node  $j$ , and  $\mathbf{x}_j^{\boldsymbol{\psi}}$  is the feature representation of  $\mathbf{x}$  at node  $j$ . Let  $\mathcal{T}_j = \{t_{e_1}^{\boldsymbol{\psi}}, \dots, t_{e_n}^{\boldsymbol{\psi}}\}$  denote the ordered set of the  $n$  transformer modules on the path from the root to node  $j$ , the feature vector  $\mathbf{x}_j^{\boldsymbol{\psi}}$  is given by

$$\mathbf{x}_j^{\boldsymbol{\psi}} := (t_{e_n}^{\boldsymbol{\psi}} \circ \dots \circ t_{e_2}^{\boldsymbol{\psi}} \circ t_{e_1}^{\boldsymbol{\psi}})(\mathbf{x}).$$

On the other hand, the leaf-specific conditional distribution  $p_l^{\boldsymbol{\phi}, \boldsymbol{\psi}}(\mathbf{y}) := p(\mathbf{y}|\mathbf{x}, z_l = 1, \boldsymbol{\phi}, \boldsymbol{\psi})$  in (5.1) yields an estimate for the distribution over target  $\mathbf{y}$  for leaf node  $l$  and is given by its solver’s output  $s_l^{\boldsymbol{\phi}}(\mathbf{x}_{\text{parent}(l)}^{\boldsymbol{\psi}})$ .

We consider two inference schemes based on a trade-off between accuracy and computation, which we refer to as *multi-path* and *single-path* inference. The multi-path inference uses the *full predictive distribution* given in (5.1) as estimate for  $p(\mathbf{y}|\mathbf{x})$ . However, computing this quantity requires averaging the distributions over all the leaves involving computing all operations at all nodes and edges of the tree,

which is expensive for a large ANT. On the other hand, the single-path inference scheme only uses the predictive distribution at the leaf node chosen by greedily traversing the tree in the directions of highest confidence of the routers. This approximation constrains computations to a single path, allowing for more memory- and time-efficient inference.

## 5.4 Optimisation

Training of an ANT proceeds in two stages: 1) *growth phase* during which the model architecture is learned based on *local* optimisation, and 2) *refinement phase* which further tunes the parameters of the model discovered in the first phase based on *global* optimisation. Algorithm 5.1 shows a pseudocode of the training algorithm.

---

### Algorithm 5.1 ANT Optimisation

---

```

Initialise topology  $\mathbb{T}$  and parameters  $\mathbb{O}$   $\triangleright \mathbb{T}$  is set to a root node with one solver and one transformer
Optimise parameters in  $\mathbb{O}$  via gradient descent on NLL  $\triangleright$  Learning root classifier
Set the root node “suboptimal”
while true do  $\triangleright$  Growth of  $\mathbb{T}$  begins
    Freeze all parameters  $\mathbb{O}$ 
    Pick next “suboptimal” leaf node  $l \in \mathcal{N}_{leaf}$  in the breadth-first order
    Add (1) router to  $l$  and train new parameters  $\triangleright$  Split data
    Add (2) transformer to  $l$  and train new parameters  $\triangleright$  Deepen transform
    Add (1) or (2) to  $\mathbb{T}$  if validation error decreases, otherwise set  $l$  to “optimal”
    Add any new modules to  $\mathbb{O}$ 
    if no “suboptimal” leaves remain then
        Break
    end if
end while
Unfreeze and train all parameters in  $\mathbb{O}$   $\triangleright$  Global refinement with fixed  $\mathbb{T}$ 

```

---

### 5.4.1 Loss function: optimising parameters of $\mathbb{O}$

For both phases, we use the negative log-likelihood (NLL) as the common objective function to minimise:

$$-\log p(\mathbf{Y}|\mathbf{X}, \Theta) = -\sum_{n=1}^N \log \left( \sum_{l=1}^L \pi_l^{\theta, \psi}(\mathbf{x}^{(n)}) p_l^{\phi, \psi}(\mathbf{y}^{(n)}) \right)$$

where  $\mathbf{X} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}$ ,  $\mathbf{Y} = \{\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(N)}\}$  denote the training inputs and targets. As all component modules (routers, transformers and solvers) are differentiable with respect to their parameters  $\Theta = (\boldsymbol{\theta}, \boldsymbol{\psi}, \boldsymbol{\phi})$ , we can use gradient-based optimisation. Given an ANT with fixed topology  $\mathbb{T}$ , we use backpropagation [347] for gradient computation and use gradient descent to minimise the NLL for learning the parameters.

### 5.4.2 Growth phase: learning architecture $\mathbb{T}$

We next describe our proposed method for growing the tree  $\mathbb{T}$  to an architecture of adequate complexity for the given training data. Starting from the root, we choose one of the leaf nodes in breadth-first order and incrementally modify the architecture by adding computational modules to it. In particular, we evaluate 3 choices (Fig. 5.2 (Right)) at each leaf node; (1) “split data” extends the current model by splitting the node with an addition of a new router; (2) “deepen transform” increases the depth of the incoming edge by adding a new transformer; (3) “keep” retains the current model. We then locally optimise the parameters of the newly added modules in the architectures of (1) and (2) by minimising NLL via gradient descent, while fixing the parameters of the previous part of the computational graph. Lastly, we select the model with the lowest validation NLL if it improves on the previously observed

lowest NLL, otherwise we execute (3). This process is repeated to all new nodes level-by-level until no more ‘‘split data’’ or ‘‘deepen transform’’ operations pass the validation test.

The rationale for evaluating the two choices is to give the model a freedom to choose the most effective option between ‘‘going deeper’’ or splitting the data space. Splitting a node is equivalent to a soft partitioning of the feature space of incoming data, and gives birth to two new leaf nodes (left and right children solvers). In this case, the added transformer modules on the two branches are identity functions. Deepening an edge on the other hand seeks to learn richer representation via an extra nonlinear transformation, and replaces the old solver with a new one. Local optimisation is efficient in time and space; gradients only need to be computed for the parameters of the new parts of the architecture, reducing computation, while forward activations prior to the new parts do not need to be stored in memory, saving space.

### 5.4.3 Refinement phase: global tuning of $\mathbb{O}$

Once the model topology is determined in the growth phase, we finish by performing global optimisation to refine the parameters of the model, now with a fixed architecture. This time, we perform gradient descent on the NLL with respect to the parameters of all modules in the graph, jointly optimising the hierarchical grouping of data to paths on the tree and the associated expert NNs. The refinement phase can correct suboptimal decisions made during the local optimisation of the growth phase, and empirically improves the generalisation error (see Sec. 5.5.4).

## 5.5 Experiments and Results

We evaluate ANTs using the SARCOS multivariate regression dataset [306], and the MNIST [307] and CIFAR-10 [308] classification datasets. We run ablation studies to show that our different components are vital for the best performance. We then assess the ability of ANTs to automatically learn meaningful hierarchical structures in data. Next, we examine the effects of refinement phase on ANTs, and show that it can automatically prune the tree. Finally, we demonstrate that our proposed training procedure adapts the model size appropriately under varying amounts of labelled data. All of our models are implemented in PyTorch [348] and is available at <https://github.com/rtanno21609/AdaptiveNeuralTrees>.

### 5.5.1 Set-up details

**Data:** we perform our experiments on the SARCOS robot inverse dynamics dataset<sup>1</sup>, the MNIST digit classification task [307] and the CIFAR-10 object recognition task [308]. The SARCOS dataset consists of 44,484 training and 4,449 testing examples, where the goal is to map from the 21-dimensional input space (7 joint positions, 7 joint velocities and 7 joint accelerations) to the corresponding 7 joint torques [306]. No dataset preprocessing or augmentation is used. The MNIST dataset consists of 60,000 training and 10,000 testing examples, all of which are  $28 \times 28$  grayscale images of digits from 0 to 9 (10 classes). The dataset is preprocessed by subtracting the mean, but no data augmentation is used. The CIFAR-10 dataset consists of 50,000 training and 10,000 testing examples, all of which are  $32 \times 32$  coloured natural images drawn from 10 classes. We adopt an augmentation scheme widely used in the literature [349, 350, 351, 326, 327] where images are zero-padded with 4 pixels on each side, randomly cropped and horizontally mirrored. For all three datasets, we hold out 10% of training images as a validation set. The best model is selected based on the validation accuracy over the course of ANT training, spanning both the growth phase and the refinement phase, and its test accuracy is reported.

**Training:** both the growth phase and the refinement phase of ANTs are performed on a single Titan X GPU on all three datasets. For all the experiments in this paper, we employ the following training protocol: (1) optimise parameters using Adam [190] with initial learning rate of  $10^{-3}$  and  $\beta = [0.9, 0.999]$ , with minibatches of size 512; (2) during the growth phase, employ early stopping with a patience of 5, that is, training is stopped after 5 epochs of no progress on the validation set; (3) during the refinement

---

<sup>1</sup><http://www.gaussianprocess.org/gpml/data/>

phase, train for 300 epochs for SARCOS, 100 epochs for MNIST and 200 epochs for CIFAR-10, decreasing the learning rate by a factor of 10 at every multiple of 50.

We observe that the patience level is an important hyperparameter which affects the quality of the growth phase; very low or high patience levels result in new modules underfitting or overfitting locally, thus preventing meaningful further growth and limiting the accuracy of the resultant models. We tuned this hyperparameter using the validation sets, and set the patience level to 5, which produced consistently good performance on SARCOS, MNIST and CIFAR-10 datasets across different specifications of primitive modules. A quantitative evaluation on CIFAR-10 is given in Supp. Sec. ??.

In the SARCOS experiments, all the non-NN-based methods were trained using scikit-learn [352]. Hidden layers in the baseline MLPs are followed by tanh non-linearities and contain 256 units to be consistent with the complexity of transformer modules.

**Primitive modules:** we train ANTs with a range of primitive modules as shown in Tab. 2 in the main text. For simplicity, we define the modules based on three types of NN layers: convolutional, global-average-pooling (GAP) and fully-connected (FC). Solver modules are fixed as linear models e.g. linear classifier and linear regression. Router modules are binary classifiers with a sigmoid output. All convolutional and FC layer are followed by ReLU or tanh non-linearities, except in the last layers of solvers and routers. For image classification experiments, we also apply  $2 \times 2$  max-pooling to feature maps after every  $d$  transformer modules where  $d$  is the downsample frequency. For the SARCOS regression experiment, hidden layers in the routers and transformers contain 256 units. We balance the number of parameters in the router and transformer modules to be of the same order of magnitude to avoid favouring either partitioning the data or learning more expressive features.

### 5.5.2 Model Performance

We compare the performance of ANTs (Tab. 5.2) against a range of DT and NN models (Tab. 5.3), where notably the relative performance of these two classes of models differs between datasets. ANTs inherit from both and achieve the lowest error on SARCOS, and perform favourably on MNIST and CIFAR-10. In general, DT methods without feature learning, such as RFs [294, 309] and GBTs [310], perform poorly on image classification tasks [308]. In comparison with CNNs without shortcut connections [307, 349, 350, 351], different ANTs balance between stronger performance with comparable numbers of trainable parameters, and comparable performance with smaller amount of parameters. At the other end of the spectrum, state-of-the-art NNs [356, 327] contain significantly more parameters.

**Conditional computation:** Tab. 5.3 compares the errors and number of parameters of different ANTs for both multi-path and single-path inference schemes. While reducing the number of parameters (from Params (multi-path) to Params (single-path)) across all ANT models, we observe only a small difference in error (between Error (multi-path) and Error (single-path)), with the largest deviations being 0.06% for classification and 0.158 for regression.

In addition, Tab. 5.4 shows that the single-path inference reduces the floating point operations per second (FLOPS). This means that single-path inference gives an accurate approximation of the multi-path inference, while being more efficient to compute. This close approximation comes from the confident splitting probabilities of routers, being close to 0 or 1 (see blue histograms in Fig. 5.3(b)).

**Ablation study:** we compare the predictive errors of different variants of ANTs in cases where the options for adding transformer or router modules are disabled (see Tab. 5.5). In the first case, the resulting models are equivalent to SDTs [299] or HMEs [311] with locally grown architectures, while the second case is equivalent to standard CNNs, grown adaptively layer by layer. We observe that either ablation consistently leads to higher errors across different module configurations on all three datasets, justifying the combination of feature learning and hierarchical partitioning in ANTs.

**SARCOS multivariate regression:** Tab. 5.3 shows that ANT-SARCOS outperforms all other methods in mean squared error (MSE) with the full set of parameters. With the single-path inference, GBTs performs slightly better than a single ANT while requiring fewer parameters. We note that the top 3 methods are all tree-based, with the third best method being an SDT (with MLP routers). On the other hand, ANT and GBTs outperform the best standard NN model with less than a half of the parameter

Table 5.3: Comparison of performance of different models on SARCOS, MNIST and CIFAR-10. The columns “Error (multi-path)” and “Error (single-path)” indicate the classification (%) or regression (MSE) errors of predictions based on the multi-path and the single-path inference. The columns “Params. (multi-path)” and “Params. (single-path)” respectively show the total number of parameters in the model and the average number of parameters used during single-path inference. “Ensemble Size” indicates the size of ensemble used. An entry of “–” indicates that no value was reported. Methods marked with  $\dagger$  are from our implementations trained in the same experimental setup. \* indicates that the parameters are initialised with a pre-trained CNN.

	Method	Error (multi-path)	Error (single-path)	Params. (multi-path)	Params. (single-path)	Ensemble Size
SARCOS	Linear regression	10.693	N/A	154	N/A	1
	MLP with 2 hidden layers [353]	5.111	N/A	31,804	N/A	1
	Decision tree	3.708	3.708	319,591	25	1
	MLP with 1 hidden layer	2.835	N/A	7,431	N/A	1
	Gradient boosted trees	2.661	2.661	391,324	2,083	$7 \times 30$
	MLP with 5 hidden layers	2.657	N/A	270,599	N/A	1
	Random forest	2.426	2.426	40,436,840	4,791	200
	Random forest	2.394	2.394	141,540,436	16,771	700
	MLP with 3 hidden layers	2.129	N/A	139,015	N/A	1
	SDT (with MLP routers)	2.118	2.246	28,045	10,167	1
MNIST	Gradient boosted trees	1.444	1.444	988,256	6,808	$7 \times 100$
	ANT-SARCOS	1.384	1.542	103,823	61,640	1
	ANT-SARCOS (ensemble)	1.226	1.372	598,280	360,766	8
	Linear classifier	7.91	N/A	7,840	N/A	1
	RDT [312]	5.41	–	–	–	1
CIFAR-10	Random Forests [294]	3.21	3.21	–	–	200
	Compact Multi-Class Boosted Trees [310]	2.88	–	–	–	100
	Alternating Decision Forest [354]	2.71	2.71	–	–	20
	Neural Decision Tree [305]	2.10	–	1,773,130	502,170	1
	ANT-MNIST-C	1.62	1.68	39,670	7,956	1
	MLP with 2 hidden layers [355]	1.40	N/A	1,275,200	N/A	1
	LeNet-5 $\dagger$ [307]	0.82	N/A	431,000	N/A	1
	gcForest [309]	0.74	0.74	–	–	500
	ANT-MNIST-B	0.72	0.73	76,703	50,653	1
	Neural Decision Forest [303]	0.70	–	544,600	463,180	10
CIFAR-10	ANT-MNIST-A	0.64	0.69	100,596	84,935	1
	ANT-MNIST-A (ensemble)	0.29	0.30	850,775	655,449	8
	CapsNet [356]	0.25	–	8.2M	N/A	1
	Compact Multi-Class Boosted Trees [310]	52.31	–	–	–	100
	Random Forests [294]	50.17	50.17	–	–	2000
	gcForest [309]	38.22	38.22	–	–	500
	MaxOut [349]	9.38	N/A	6M	N/A	1
	ANT-CIFAR10-C	9.31	9.34	0.7M	0.5M	1
	ANT-CIFAR10-B	9.15	9.18	0.9M	0.6M	1
	Network in Network [350]	8.81	N/A	1M	N/A	1
CIFAR-10	All-CNN $\dagger$ [351]	8.71	N/A	1.4M	N/A	1
	ANT-CIFAR10-A	8.31	8.32	1.4M	1.0M	1
	ANT-CIFAR10-A (ensemble)	7.71	7.79	8.7M	7.4M	8
	ANT-CIFAR10-A*	6.72	6.74	1.3M	0.8M	1
	ResNet-110 [326]	6.43	N/A	1.7M	N/A	1
	DenseNet-BC (k=24) [327]	3.74	N/A	27.2M	N/A	1

count. This highlights the value of hierarchical clustering for predictive performance and inference speed. Meanwhile, we still reap the benefits of representation learning, as shown by both ANT-SARCOS and the SDT (which is a specific form of ANT with identity transformers) requiring fewer parameters than the best-performing GBT configuration. Finally, we note that deeper NNs (5 vs. 3 hidden layers) can

Table 5.4: Comparison of FLOPs. The results for ResNet110 and DenseNet were retrieved from [357] and [358], respectively. The FLOPs of all other models were computed using TorchStat toolbox available at <https://github.com/Swall0w/torchstat>.

	Model	FLOPS (multi-path)	FLOPS (single-path)
MNIST	Linear Classifier	8K	-
	LeNet-5	231 K	-
	ANT-MNIST-C	99K	83K
	ANT-MNIST-B	346K	331K
	ANT-MNIST-A	382K	380K
CIFAR-10	Net-in-Net	222M	-
	All-CNN	245M	-
	ResNet-110	256M	-
	DenseNet-BC (k=24)	9388M	-
	ANT-CIFAR10-C	66M	61M
	ANT-CIFAR10-B	163M	149M
	ANT-CIFAR10-A	254M	243M

overfit on this small dataset, which makes the adaptive growth procedure of tree-based methods ideal for finding a model that exhibits good generalisation.

**MNIST digit classification:** we observe that ANT-MNIST-A outperforms state-of-the-art GBT [310] and RF [309] methods in accuracy. This performance is attained despite the use of a single tree, while RF methods operate with ensembles of classifiers (the size shown in Tab. 5.2). In particular, the NDF [303] has a pre-specified architecture where LeNet-5 [307] is used as the root transformer module, and 10 trees of fixed depth 5 are built on this base features. On the other hand, ANT-MNIST-A is constructed in a data-driven manner from primitive modules, and displays an improvement over the NDF both in terms of accuracy and number of parameters. In addition, reducing the size of convolution kernels (ANT-MNIST-B) reduces the total number of parameters by 25% and the path-wise average by almost 40% while only increasing the error by < 0.1%.

We also compare against the LeNet-5 CNN [307], comprised of the same types of operations used in our primitive modules (i.e. convolutional, max-pooling and FC layers). For a fair comparison, the network is trained with the same protocol as that of the ANT refinement phase, achieving an error rate of 0.82%. Both ANT-MNIST-A and ANT-MNIST-B attain better accuracy with a smaller number of parameters than LeNet-5. The current state-of-the-art, capsule networks (CapsNets) [356], have more parameters than ANT-MNIST-A by almost two orders of magnitude.<sup>2</sup> By ensembling ANTs, we can reach similar performance (0.29% versus 0.25%) with an order of magnitude less parameters.

Lastly, we highlight the observation that ANT-MNIST-C, with the simplest primitive modules, achieves an error rate of 1.68% with single-path inference, which is significantly better than that of the linear classifier (7.91%), while engaging almost the same number of parameters (7,956 vs. 7,840) on average. To isolate the benefit of convolutions, we took one of the root-to-path CNNs on ANT-MNIST-C and increased the number of kernels to adjust the number of parameters to the same value. We observe a higher error rate of 3.55%, which indicates that while convolutions are beneficial, data partitioning has additional benefits in improving accuracy. This result demonstrates the potential of ANT growth protocol for constructing performant models with lightweight inference. See Sec. G in the supplementary materials for the architecture of ANT-MNIST-C.

**CIFAR-10 object recognition:** we see that ANTs largely outperform the state-of-the-art DT method, gcForest [309], achieving over 90% accuracy, demonstrating the benefit of representation learning in tree-structured models. Secondly, with fewer number of parameters in single-path inference, ANT-CIFAR-A achieves higher accuracy than CNN models without shortcut connections [349, 350, 351] that held the state-of-the-art performance in respective years. With simpler primitive modules we learn more compact models (ANT-MNIST-B and -C) with a marginal compromise in accuracy. In addition,

<sup>2</sup>Notably, CapsNets also feature a routing mechanism, but with a significantly different mechanism and motivation.

Table 5.5: Ablation study on regression (MSE) and classification (%) errors. ‘‘CNN’’ refers to the case where the ANT is grown without routers while ‘‘SDT/HME’’ refers to the case where transformer modules on the edges are disabled.

Model	Error (multi-path)			Error (single-path)		
	ANT (default)	CNN (no $\mathcal{R}$ )	HME (no $\mathcal{T}$ )	ANT (default)	CNN (no $\mathcal{R}$ )	HME (no $\mathcal{T}$ )
SARCOS	1.38	2.51	2.12	1.54	2.51	2.25
MNIST-A	0.64	0.74	3.18	0.69	0.74	4.19
MNIST-B	0.72	0.80	4.63	0.73	0.80	3.62
MNIST-C	1.62	3.71	5.70	1.68	3.71	6.96
CIFAR10-A	8.31	9.29	39.29	8.32	9.29	40.33
CIFAR10-B	9.15	11.08	43.09	9.18	11.08	44.25
CIFAR10-C	9.31	11.61	48.59	9.34	11.61	50.02

initialising the parameters of transformers and routers from a pre-trained single-path CNN further reduced the error rate of ANT-MNIST-A by 20% (see ANT-MNIST-A\* in Tab. 5.3), indicating room for improvement in our proposed optimisation method.

Shortcut connections [328] have recently lead to leaps in performance in deep CNNs [326, 327]. We observe that our best network, ANT-MNIST-A\*, has a comparable error rate and half the parameter count (with single-path inference) to the best-performing residual network, ResNet-110 [326]. Densely connected networks leads to better accuracy, but with an order of magnitude more parameters [327]. We expect shortcut connections to improve ANT performance, and leave integrating them to future work.

**Training times:** Tab. 5.6 summarises the time taken on a single Titan X GPU for the growth phase and refinement phase of various ANTs trained on the CIFAR10 dataset, and compares against the training time of All-CNN [351]. Local optimisation during the growth phase means that the gradient computation is constrained to the newly added component of the graph, allowing us to grow a good candidate model under 3 hours on one GPU.

Table 5.6: Training time comparison. Time and number of epochs taken for the growth and refinement phase are shown. along with the time required to train the baseline, All-CNN [351].

Model	Growth		Fine-tune	
	Time	Epochs	Time	Epochs
All-CNN (baseline)	—	—	1.1 (hr)	200
ANT-CIFAR10-A	1.3 (hr)	236	1.5 (hr)	200
ANT-CIFAR10-B	0.8 (hr)	313	0.9 (hr)	200
ANT-CIFAR10-C	0.7 (hr)	285	0.8 (hr)	200

### 5.5.3 Interpretability

The growth procedure of ANTs is capable of discovering hierarchical structures in the data that are useful to the end task. Without any regularization imposed on routers, the learned hierarchies often display strong specialisation of paths to certain classes or categories of data on both the MNIST and CIFAR-10 datasets—Fig. 5.4 visualises all the ANT architectures discovered on the MNIST (i-iii) and CIFAR-10 (iv-vi) datasets. Fig. 5.3 (a) displays an example with particularly ‘‘human-interpretable’’ partitions e.g. man-made versus natural objects, and road vehicles versus other types of vehicles. It should, however, be noted that human intuitions on relevant hierarchical structures do not necessarily equate to optimal representations, particularly as datasets may not necessarily have an underlying hierarchical structure, e.g., MNIST. Rather, what needs to be highlighted is the ability of ANTs to learn when to share or separate the representation of data to optimise end-task performance, which gives rise to automatically discovering such hierarchies.

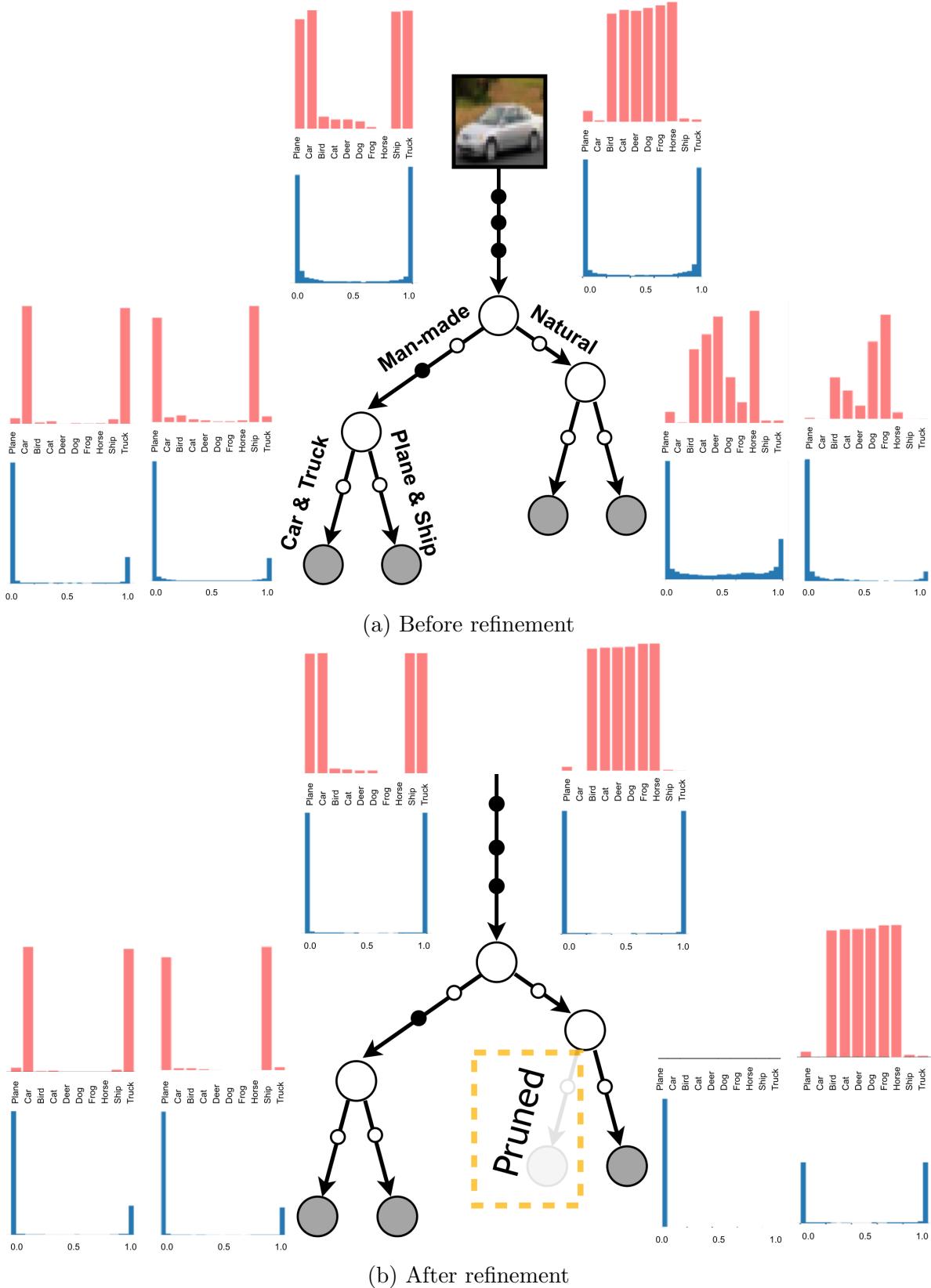


Figure 5.3: Visualisation of class distributions (red) and path probabilities (blue) computed over the whole test set at respective nodes of an example ANT (a) before and (b) after the refinement phase. (a) shows that the model captures an interpretable hierarchy, grouping semantically similar images on the same branches. (b) shows that the refinement phase polarises path probabilities, pruning a branch.

To further attest that the model learns a meaningful routing strategy, we also present the test accuracy of the predictions from the leaf node with the smallest reaching probability in Tab. 5.7. We observe that using the least likely “expert” leads to a substantial drop in classification accuracy to close to that of random guess or even worse for large trees (ANT-MNIST-C and ANT-CIFAR10-C). This demonstrates that features in ANTs become specialised to the subsets of the partitioned input space at lower levels in the tree hierarchy.

Table 5.7: Comparison of classification performance between the default single-path inference scheme and the prediction based on the least likely expert. between the

Module Spec.	Error % (Selected path)	Error % (Least likely path)
ANT-MNIST-A	0.69	86.18
ANT-MNIST-B	0.73	81.98
ANT-MNIST-C	1.68	98.84
ANT-CIFAR10-A	8.32	74.28
ANT-CIFAR10-B	9.18	89.74
ANT-CIFAR10-C	9.34	97.52

In addition, most learned trees are unbalanced. This property of adaptive computation is plausible since certain types of images may be easier to classify than others, as seen in prior work [324]. This property is reflected by traditional DT algorithms, but not “neural” tree-structured models with pre-specified architectures [301, 304, 303, 315].

### 5.5.4 Effect of global refinement

We observe that global refinement phase improves the generalisation error. Fig. 5.5 (Right) shows the generalisation error of various ANT models on CIFAR-10, with vertical dotted lines indicating the epoch when the models enter the refinement phase. As we switch from optimising parts of the ANT in isolation to optimising all parameters, we shift the optimisation landscape, resulting in an initial drop in performance. However, they all consistently converge to higher test accuracy than the best value attained during the growth phase. This provides evidence that refinement phase remedies suboptimal decisions made during the locally-optimised growth phase. In many cases, we observed that global optimisation polarises the decision probability of routers, which occasionally leads to the effective “pruning” of some branches. For example, in the case of the tree shown in Fig. 5.3(b), we observe that the decision probability of routers are more concentrated near 0 or 1 after global refinement, and as a result, the empirical probability of visiting one of the leaf nodes, calculated over the validation set, reduces to 0.09%—meaning that the corresponding branch could be pruned without a negligible change in the network’s accuracy. The resultant model attains lower generalisation error, showing the pruning has resolved a suboptimal partitioning of data.

### 5.5.5 Adaptive model complexity

Overparametrised models, trained without regularization, are vulnerable to overfitting on small datasets. Here we assess the ability of our proposed ANT training method to adapt the model complexity to varying amounts of labelled data. We run classification experiments on CIFAR-10 and train three variants of ANTs, All-CNN [351] and linear classifier on subsets of the dataset of sizes 50, 250, 500, 2.5k, 5k, 25k and 45k (the full training set). Here we choose All-CNN as the baseline as it has similar number of parameters when trained on the full dataset and is the closest in terms of constituent operations (convolutional, GAP and FC layers). Fig. 5.5 (Left) shows the corresponding test performances. The best model is picked based on the performance on the same validation set of 5k examples as before. As the dataset gets smaller, the margin between the test accuracy of the ANT models and All-CNN/linear classifier increases (up to 13%). Fig. 5.5 (Middle) shows the model size of discovered ANTs as the dataset size varies. For different settings of primitive modules, the number of parameters generally increases as a function of the dataset size. All-CNN has a fixed number of parameters, consistently larger than the

discovered ANTs, and suffers from overfitting, particularly on small datasets. The linear classifier, on the other hand, underfits to the data. Our method constructs models of adequate complexity, leading to better generalisation. This shows the value of our tree-building algorithm over using models of fixed-size structures.

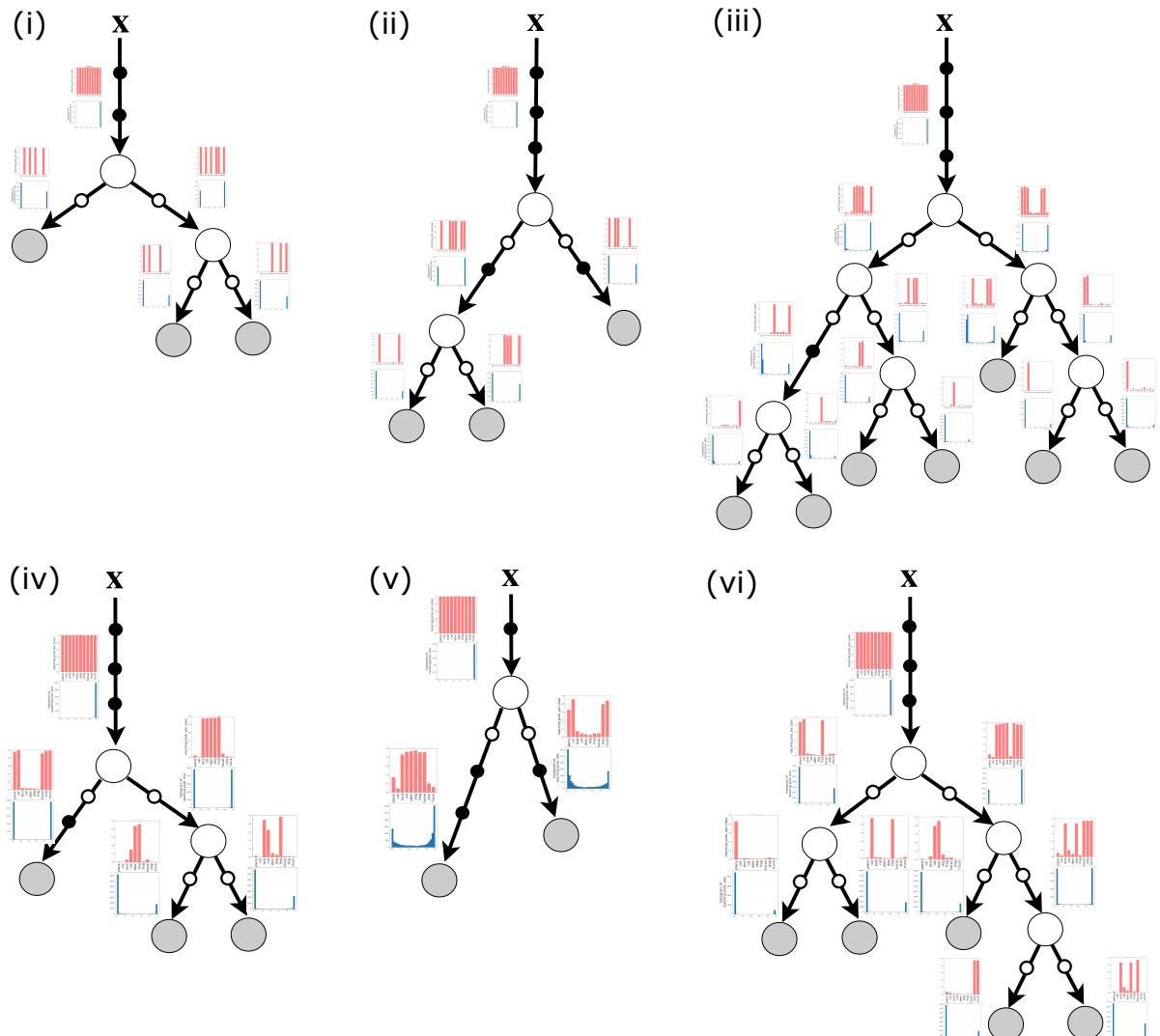


Figure 5.4: Illustration of discovered ANT architectures. (i) ANT-MNIST-A, (ii) ANT-MNIST-B, (iii) ANT-MNIST-C, (iv) ANT-CIFAR10-A, (v) ANT-CIFAR10-B, (vi) ANT-CIFAR10-C. Histograms in red and blue show the class distributions and path probabilities at respective nodes. Small black circles on the edges represent transformers, circles in white at the internal nodes represent routers, and circles in gray are solvers. The small white circles on the edges denote specific cases where transformers are identity functions.

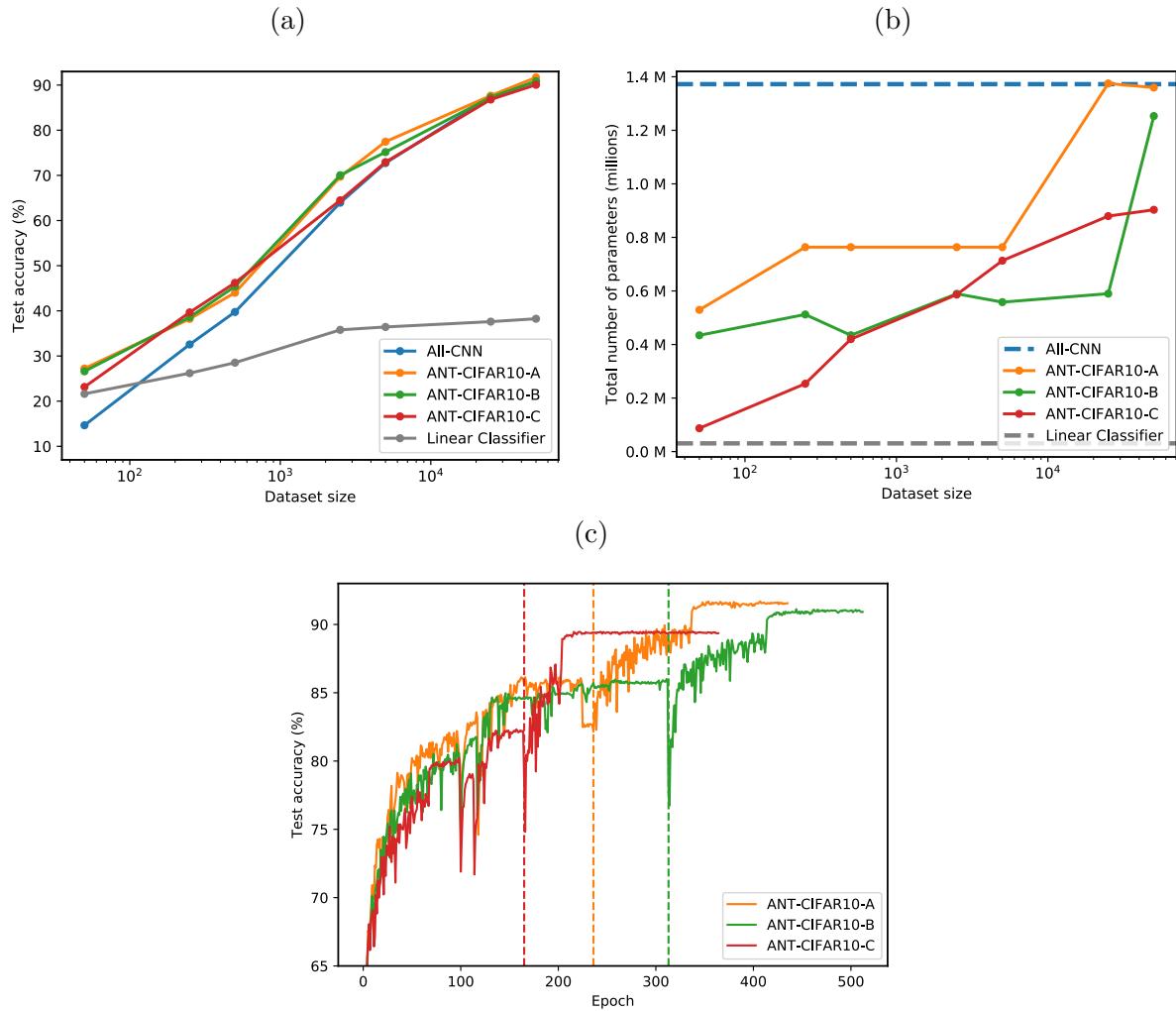


Figure 5.5: (a) Test accuracy on CIFAR-10 of ANTs for varying amounts of training data. (b) The complexity of the grown ANTs increases with dataset size. (c) Refinement improves generalisation; the dotted lines show where the refinement phase starts.

## 5.6 Discussion and Conclusion

In this chapter, we introduced Adaptive Neural Trees (ANTs), a simple way to marry the architecture learning, conditional computation and hierarchical clustering of decision trees (DTs) with the hierarchical representation learning and gradient descent optimization of deep neural networks (DNNs). Our proposed training algorithm optimises both the parameters and architectures of ANTs through progressive growth, tuning them to the size and complexity of the training dataset. Together, these properties make ANTs a generalisation of previous work attempting to unite NNs and DTs. Finally, we validated the claimed benefits of ANTs for regression (SARCOS dataset) and classification (MNIST & CIFAR10 datasets), whilst still achieving high performance.

Future work will aim to scale up ANTs to medical imaging applications which typically involve larger, higher dimensional datasets. This requires training potentially wider and deeper ANTs, which necessitates a more effective and efficient optimisation algorithm. Firstly, the current growth procedure is greedy and suboptimal. The sub-optimality of the local decisions only gets worse as the tree gets more complex. In the future, we will look into different ways to alleviate this issue such as (1) more global optimisation during growth phase, (2) merging of branches [317] or (3) training in entangled settings as already done in decision tree research [73]. Secondly, for large trees, the global refinement phase would be computationally expensive since the associated computation involves every part of the tree structure. One potential solution to this is to employ a Monte-Carlo approximation of the training by stochastically traversing the tree, while using gradient estimators for learning the parameters of the routers such as REINFORCE and straight-through (ST) estimators [359] or more modern methods such as Gumbel ST estimator [274], REBAR [360] and RELAX [361]. We have also limited ourselves to relatively simple NN components in the design of the ANT’s primitive modules. As future work, we aim to extend ANTs to use more recent advances in deep learning, particularly residual learning [326] and dense connections [327], as these have enabled the successful optimisation of more performant NNs.

Another important challenge is to explore the value of the learned hierarchical structures of ANTs in terms of “interpretability” (or algorithmic transparency). I believe that the tree-shaped hierarchy provides a new means to understand its internal decision making process—for instance, given an image where the ANT fail to classify, you could compare against a large number of correctly predicted examples, and potentially localise a point of routing failure where the ANT makes the wrong decision. On the other hand, such localisation of failures is more difficult in conventional CNNs with a single fully distributed representation. In addition, such hierarchical and decomposed interpretations are complementary to existing visualisation techniques such as saliency based attribution methods [362]. For example, providing saliency maps of the series of routing decisions may potentially provide a layer of transparency into the decisions made in the raw feature space.

# Chapter 6

## Modelling Human Uncertainty (I): Classification

The predictive performance of supervised learning algorithms depends on the quality of labels. In a typical label collection process, multiple annotators provide subjective noisy estimates of the “truth” under the influence of their varying skill-levels and biases. Blinely treating these noisy labels as the ground truth limits the accuracy of learning algorithms in the presence of strong disagreement. This problem is critical for applications in domains such as medical imaging where both the annotation cost and inter-observer variability are high. In this work, we present a method for simultaneously learning the individual annotator model and the underlying true label distribution, using only noisy observations. Each annotator is modeled by a confusion matrix that is jointly estimated along with the classifier predictions. We propose to add a regularization term to the loss function that encourages convergence to the true annotator confusion matrix. We provide a theoretical argument as to how the regularization is essential to our approach both for the case of single annotator and multiple annotators. Despite the simplicity of the idea, experiments on image classification tasks with both simulated and real labels show that our method either outperforms or performs on par with the state-of-the-art methods and is capable of estimating the skills of annotators even with a single label available per image. This chapter is based on the publication [107].

### 6.1 Introduction

In many practical applications, supervised learning algorithms are trained on noisy labels obtained from multiple annotators of varying skill levels and biases. When there is a substantial amount of disagreement in the labels, conventional training algorithms that treat such labels as the “truth” lead to models with limited predictive performance. To mitigate such variation, practitioners typically abide by the principle of “wisdom of crowds” [363] and aggregate labels by computing the majority vote. However, this approach has limited efficacy in applications where the number of annotations is modest or the tasks are ambiguous. In particular, many vision applications in medical image analysis [2] require annotations from clinical experts, which incur high costs and commonly suffer from high inter-reader variability [364, 365, 366, 60] — e.g., the average variability in the range 74-85% has been reported for glioblastoma segmentation [367]. While medical imaging data is now extremely abundant due to over two decades of digitisation, the world still remains relatively short of access to clean data with well-curated labels, that is amenable to machine learning [368], necessitating an intelligent method to learn robustly from noisy annotations.

In theory, even in the presence of large label noise, if the exact process by which each annotator generates the labels was known, we could correct the annotations accordingly and thus train our model on a cleaner set of data. Furthermore, this additional knowledge of the annotators’ skills can be utilized to decide on which examples to be labeled by which annotators [369, 370, 371]. Therefore, methods that can accurately model the label noise of annotators are useful for improving not only the accuracy of the trained model, but also the quality of labels in the future.

Previous work, therefore, proposed various methods for jointly estimating the skills of the annotators and the ground truth (GT) labels. We categorize these methods into two groups: (1) *two-stage* approach and (2) *simultaneous* approach. Methods in the first category perform label aggregation and training of

a supervised learning model in two separate steps. The noisy labels  $\tilde{\mathbf{Y}}$  are first aggregated by building a probabilistic model of annotators. The observable variables are the noisy labels  $\tilde{\mathbf{Y}}$ , and the latent variables/parameters to be estimated are the annotator skills and GT labels  $\mathbf{Y}$ . Then, a machine learning model is trained on the pairs of aggregated labels  $\mathbf{Y}$  and input examples  $\mathbf{X}$  (e.g. images) to perform the task of interest. The initial attempt was made in [372] in the early 1970s and more recently, numerous lines of research [373, 60, 374, 375, 376] proposed extensions of this work e.g. by estimating the difficulty of each example. However, in all these cases, information about the raw inputs  $\mathbf{X}$  is completely neglected in the generative model of noisy labels used in the aggregation step, and this highly limits the quality of estimated true labels in practice.

The *simultaneous* approaches [377, 378, 379, 380] address this issue by integrating the prediction of the supervised learning model (i.e. distribution  $p(\mathbf{Y}|\mathbf{X})$ ) into the probabilistic model of noisy labels, and have been shown to improve the predictive performance. These methods employ variants of the expectation-maximization (EM) algorithm during training, and require a reasonable number of labels for each example. However, in most real world applications, it is practically prohibitive to collect a large number of labels per example, and this requirement limits their applications. A notable exception is the Model Bootstrapped EM (MBEM) algorithm presented in [381] that is capable of learning even with little label redundancy.

In this chapter, we propose a more effective alternative to these EM-based approaches for jointly modeling the annotator skills and GT label distribution. Our method separates the annotation noise from true labels by (1) ensuring high fidelity with the data by minimizing the cross entropy loss and (2) encouraging the estimated annotators to be maximally unreliable by minimizing the trace of the estimated confusion matrices. Our method is also simpler to implement, only requiring an addition of a regularization term to the cross-entropy loss. Furthermore, we provide a theoretical result that such regularization is capable of recovering the annotation noise as long as the average confusion matrix (CM) over annotators is diagonally dominant (i.e., every diagonal element is larger than any other off-diagonal element in the corresponding row).

Experiments on image classification tasks with both simulated and real noisy labels demonstrate that our method, despite being much simpler, leads to better or comparable performance with MBEM [381] and generalized EM [377, 184], and is capable of recovering CMs even when there is only one label available per example. We simulated a diverse range of annotator types on MNIST and CIFAR10 data sets while we used a ultrasound dataset for cardiac view classification to test the efficacy in a real-world application. We also show importance of modeling individual annotators by comparing against various modern noise-robust methods [382, 383, 384, 385], when the inter-annotator variability is high.

### 6.1.1 Other Related Works.

More broadly, our work is related to methods for robust learning in the presence of label noise. There is a large body of literature that do not explicitly model individual annotators unlike our method.

The effects of label noise are well studied in common classifiers such as SVMs and logistic regression, and robust variants have been proposed [386, 387, 388]. More recently, various attempts have been made to train deep neural networks under label noise. Reed et al. [382] developed a robust loss to model “prediction consistency”, which was later extended by [389]. In [390] and [383], label noise was parametrized in the form of a transition matrix and incorporated into neural networks for binary and multi-way classification. A more effective alternative for estimating such transition matrix was proposed in [391], and a method for capturing image dependency of label noise was shown in [392]. We will later compare our model to several of these methods to test the value of modelling individual annotators in gaining robustness to label noise.

Multiple lines of work have shown that a small portion of clean labels improves robustness. [393] proposed to learn from clean labels to correct the labels of noisy examples. [394] proposed a method for learning to weigh examples during each training iteration by using the validation loss on clean labels as the meta-objective. [395] employs a similar approach, but trains a separate network that proposes weighting. However, curating a set of clean labels of sufficient size is expensive for many applications, and this work focuses on the scenario of learning from purely noisy labels.

## 6.2 Methods

We assume that a set of images  $\{\mathbf{x}_i\}_{i=1}^N$  are assigned with noisy labels  $\{\tilde{y}_i^{(r)}\}_{i=1,\dots,N}^{r=1,\dots,R}$  from multiple annotators where  $\tilde{y}_i^{(r)}$  denotes the label from annotator  $r$  given to example  $\mathbf{x}_i$ , but no ground truth (GT) labels  $\{y_i\}_{i=1,\dots,N}$  are available. In this work, we present a new procedure for multiclass classification problem that can simultaneously estimate the annotator noise and GT label distribution  $p(y|\mathbf{x})$  from such noisy set of data  $\mathcal{D} = \{\mathbf{x}_i, \tilde{y}_i^{(1)}, \dots, \tilde{y}_i^{(R)}\}_{i=1,\dots,N}$ . The method only requires adding a regularization term, that is the average accuracy of all annotator models, to the cross-entropy loss function. Intuitively, the method biases our models of each annotator to be as inaccurate as possible while having the model still explain the data. We will show that this is capable of decoupling the annotation noise from the true label distribution, as long as the average labels of the real annotators are “sufficiently” correct (which we formalize in Sec. 6.2.3). For simplicity, we first describe the method in the *dense label* scenario in which each image has labels from all annotators, and then extend to scenarios with *missing* labels where only a subset of annotators label each image. As we shall see later, the method works even when each image is only labelled by a single annotator.

### 6.2.1 Noisy Observation Model

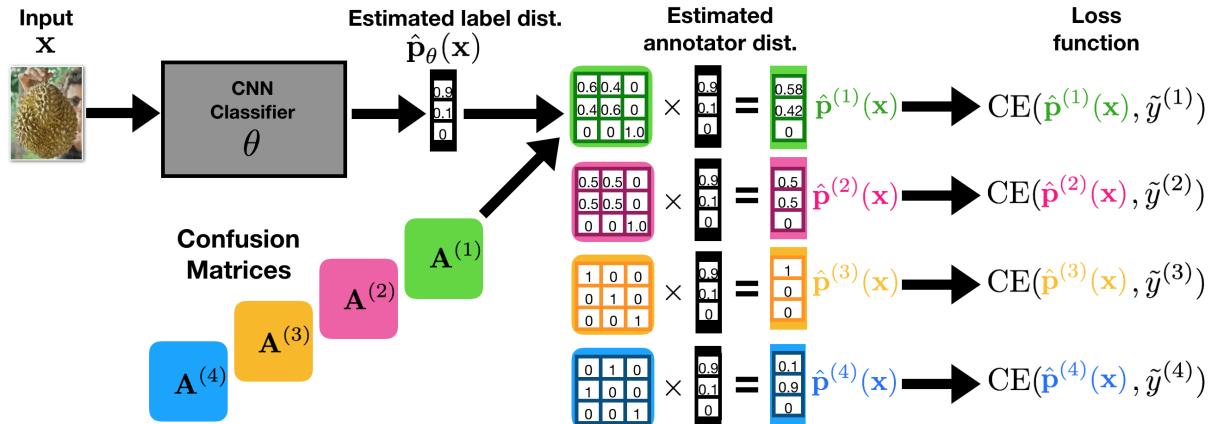


Figure 6.1: General schematic of the model (eq. 6.2) in the presence of 4 annotators. Given input image  $\mathbf{x}$ , the classifier parametrised by  $\theta$  generates an estimate of the ground truth class probabilities,  $\mathbf{p}_\theta(\mathbf{x})$ . Then, the class probabilities of respective annotators  $\mathbf{p}^{(r)}(\mathbf{x}) := \mathbf{A}^{(r)}\mathbf{p}_\theta(\mathbf{x})$  for  $r \in \{1, 2, 3, 4\}$  are computed. The model parameters  $\{\theta, \mathbf{A}^{(1)}, \mathbf{A}^{(2)}, \mathbf{A}^{(3)}, \mathbf{A}^{(4)}\}$  are optimized to minimize the sum of four cross-entropy losses between each estimated annotator distribution  $\mathbf{p}^{(r)}(\mathbf{x})$  and the noisy labels  $\tilde{y}^{(r)}$  observed from each annotator. The probability that each annotator provides accurate labels can be estimated by taking the average diagonal elements of the associated confusion matrix (CM), which we refer to as the “skill level” of the annotator.

We first describe our probabilistic model of the observed noisy labels from multiple annotators. In particular, we make two key assumptions: (1) annotators are statistically independent, (2) annotation noise is independent of the input image. By assumption (1), the probability of observing noisy labels  $\{\tilde{y}^{(1)}, \dots, \tilde{y}^{(R)}\}$  on image  $\mathbf{x}$  can be written as:

$$p(\tilde{y}^{(1)}, \dots, \tilde{y}^{(R)} | \mathbf{x}) = \prod_{r=1}^R \int_{y \in \mathcal{Y}} p(\tilde{y}^{(r)} | y, \mathbf{x}) \cdot p(y | \mathbf{x}) dy \quad (6.1)$$

where  $p(y | \mathbf{x})$  denotes the true label distribution of the image, and  $p(\tilde{y}^{(r)} | y, \mathbf{x})$  describes the noise model by which annotator  $r$  corrupts the ground truth label  $y$ . For classification problems, the label  $y$  takes a discrete value in  $\mathcal{Y} = \{1, \dots, L\}$ . From assumption (2), the probability that annotator  $r$  corrupts the GT label  $y = i$  to  $\tilde{y}^{(r)} = j$  is independent of the image  $\mathbf{x}$  i.e.  $p(\tilde{y}^{(r)} = j | y = i, \mathbf{x}) = p(\tilde{y}^{(r)} = j | y = i) =: a_{ji}^{(r)}$ .

Here we refer to the associated  $L \times L$  transition matrix  $\mathbf{A}^{(r)} = (a_{ji}^{(r)})$  as the *confusion matrix* (CM) of annotator  $r$ . The joint probability over the noisy labels is simplified to:

$$p(\tilde{y}^{(1)}, \dots, \tilde{y}^{(R)} | \mathbf{x}) = \prod_{r=1}^R \sum_{y=1}^L a_{\tilde{y}^{(r)}, y}^{(r)} \cdot p(y | \mathbf{x}) \quad (6.2)$$

Fig. 6.1 provides a schematic of our overall architecture, which models the different constituents in the above joint probability distribution. In particular, the model consists of two components: the *base classifier* which estimates the ground truth class probability vector  $\hat{\mathbf{p}}_\theta(\mathbf{x})$  whose  $i^{\text{th}}$  element approximates  $p(y = i | \mathbf{x})$ , and the set of the CM estimators  $\{\hat{\mathbf{A}}^{(r)}\}_{r=1}^R$  which approximate  $\{\mathbf{A}^{(r)}\}_{r=1}^R$ . Each product  $\hat{\mathbf{p}}^{(r)}(\mathbf{x}) := \hat{\mathbf{A}}^{(r)} \hat{\mathbf{p}}_\theta(\mathbf{x})$  represents the estimated class probability vector of the corresponding annotator. At inference time, we use the most confident class in  $\hat{\mathbf{p}}_\theta(\mathbf{x})$  as the final classification output. Next, we describe our optimization algorithm for jointly learning the parameters of the base classifier,  $\theta$  and the CMs,  $\{\hat{\mathbf{A}}^{(r)}\}_{r=1}^R$ .

## 6.2.2 Joint Estimation of Confusion and True labels

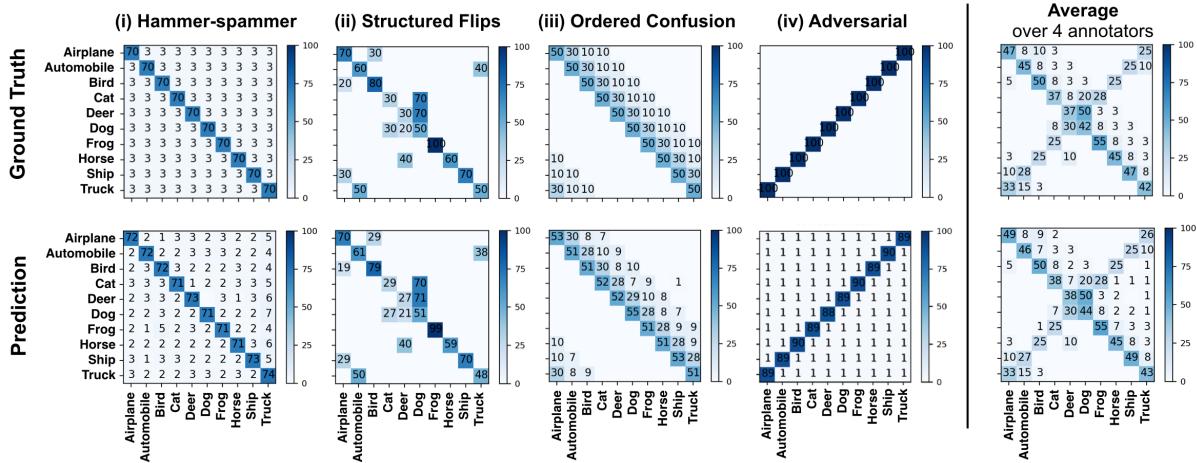


Figure 6.2: A diverse set of 4 simulated annotators on CIFAR-10. The top row shows the ground truths while the bottom row shows the estimation from our method, trained with only one label per image.

Given training inputs  $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^N$  and noisy labels  $\tilde{\mathbf{Y}}^{(r)} = \{\tilde{y}_i^{(r)}\}_{i=1}^N$  for  $r = 1, \dots, R$ , we optimize the parameters  $\{\theta, \hat{\mathbf{A}}^{(r)}\}$  by minimizing the negative log-likelihood (NLL),  $-\log p(\tilde{\mathbf{Y}}^{(1)}, \dots, \tilde{\mathbf{Y}}^{(R)} | \mathbf{X})$ . From eq. 6.2, this optimization objective equates to the sum of cross-entropy losses between the observed labels and the estimated annotator label distributions:

$$-\log p(\tilde{\mathbf{Y}}^{(1)}, \dots, \tilde{\mathbf{Y}}^{(R)} | \mathbf{X}) = \sum_{i=1}^N \sum_{r=1}^R \text{CE}(\mathbf{A}^{(r)} \hat{\mathbf{p}}_\theta(\mathbf{x}_i), \tilde{y}_i^{(r)}). \quad (6.3)$$

Minimizing above encourages each annotator-specific prediction  $\hat{\mathbf{p}}^{(r)}(\mathbf{x}) := \hat{\mathbf{A}}^{(r)} \hat{\mathbf{p}}_\theta(\mathbf{x})$  to be as close as possible to the noisy label distribution of the corresponding annotator  $\mathbf{p}^{(r)}(\mathbf{x})$ . However, this loss function alone is not capable of separating the annotation noise from the true label distribution; there are infinite combinations of  $\{\hat{\mathbf{A}}^{(r)}\}_{r=1}^R$  and classification model  $\hat{\mathbf{p}}_\theta$  such that  $\hat{\mathbf{p}}^{(r)}$  perfectly matches the annotator's label distribution  $\mathbf{p}^{(r)}$  for any input  $\mathbf{x}$ .

To formalize this problem, we denote the CM of the estimated true label distribution<sup>1</sup>  $\hat{\mathbf{p}}_\theta$  by  $\mathbf{P}$ . The CM of the estimated annotator's label distribution  $\hat{\mathbf{p}}^{(r)}$  is then given by the product  $\hat{\mathbf{A}}^{(r)} \mathbf{P}$ . Minimizing the cross-entropy loss (eq. 6.3) encourages  $\hat{\mathbf{A}}^{(r)} \mathbf{P}$  to converge to the true CM of the corresponding

<sup>1</sup> $\mathbf{P}_{ji} = \int_{\mathbf{x} \in \mathcal{X}} p(\text{argmax}_k [\hat{\mathbf{p}}_\theta(\mathbf{x})]_k = j | y = i) p(\mathbf{x}) d\mathbf{x}$

annotator  $\mathbf{A}^{(r)}$  i.e.  $\hat{\mathbf{A}}^{(r)} \mathbf{P} \rightarrow \mathbf{A}^{(r)}$ . However, there are infinitely many solutions pairs  $(\hat{\mathbf{A}}^{(r)}, \mathbf{P})$  that satisfy the equality  $\hat{\mathbf{A}}^{(r)} \mathbf{P} = \mathbf{A}^{(r)}$ . This means that we need to regularize the optimization to encourage convergence to the desired solutions i.e.  $\hat{\mathbf{A}}^{(r)} \rightarrow \mathbf{A}^{(r)}$  and  $\mathbf{P} \rightarrow \mathbf{I}$ .

To combat this problem, we propose to add the trace of the estimated CMs to the loss in eq. 6.3. Extending to the “missing labels” regime in which only a subset of annotators label each example, we derive the combined loss:

$$\sum_{i=1}^N \sum_{r=1}^R \mathbb{1}(\tilde{y}_i^{(r)} \in \mathcal{S}(\mathbf{x}_i)) \cdot \text{CE}(\hat{\mathbf{A}}^{(r)} \hat{\mathbf{p}}_\theta(\mathbf{x}_i), \tilde{y}_i^{(r)}) + \lambda \sum_{r=1}^R \text{tr}(\hat{\mathbf{A}}^{(r)}) \quad (6.4)$$

where  $\mathcal{S}(\mathbf{x})$  denotes the set of all labels available for image  $\mathbf{x}$ , and  $\text{tr}(\mathbf{A})$  denotes the trace of matrix  $\mathbf{A}$ . We simply perform gradient descent on this loss to learn  $\{\theta, \hat{\mathbf{A}}^{(1)}, \dots, \hat{\mathbf{A}}^{(R)}\}$ .

Numerous previous work have considered the same observation model, but proposed various optimization schemes. The original work [377, 184] employed the generalized EM algorithm to estimate  $\{\theta, \hat{\mathbf{A}}^{(1)}, \dots, \hat{\mathbf{A}}^{(R)}\}$ , and more recent work [379, 380] employed variants of hard-EM to optimize the same model. Khetan et al.,[381] proposed a method called model-bootstrapped EM (MBEM) in which the predictions of the base neural network classifier are used in the M-step update of CMs to learn from singly labelled data, which was not viable with the prior work. However, in all of the above EM-based methods, each M-step for the parameters of NN is not available in closed form and thus performed via gradient descent. This means that every M-step requires a training of the CNN classifier, rendering each iteration of EM expensive. A naive solution to this is to perform only few iterations of gradient descent in each E-step, however, this could limit the performance if sufficient convergence is not achieved. Our approach directly maximizes the likelihood with the trace regularizer and does not suffer from these issues. In Sec. 4, we show empirically this approach leads to an improvement both in terms of accuracy and convergence rate over the previous methods on noisy labels with high inter-annotator variability.

We also provide pseudo-codes of our method (Algorithm 6.1), generalized EM [377] (Algorithm 6.2) and model-bootstrapped EM [381] (Algorithm 6.3) to clarify the differences between different methods for jointly learning the true label distribution and confusion matrices of annotators in eq. (6.2) in the main text. Given the training set  $\mathcal{D} = \{\mathbf{x}_n, \tilde{y}_n^{(1)}, \dots, \tilde{y}_n^{(R)}\}_{n=1}^N$ , each example may not be labelled by all the annotators. In such cases, for ease of notation, we assign pseudo class  $\tilde{y}_n^{(r)} = -1$  to fill the missing labels. The comparison between these three algorithms illustrates the implementational simplicity of our method, despite the comparable or superior performance demonstrated on all three datasets.

---

### Algorithm 6.1 Our method

---

**Inputs:**  $\mathcal{D} = \{\mathbf{x}_n, \tilde{y}_n^{(1)}, \dots, \tilde{y}_n^{(R)}\}_{n=1}^N, \lambda$  : scale of trace regularizer

**Initialize the confusion matrices  $\{\hat{\mathbf{A}}^{(r)}\}_{r=1}^R$  to identity matrices**

**Initialize the parameters of the base classifier  $\theta$**

**Learn  $\theta$  and  $\{\hat{\mathbf{A}}^{(r)}\}_{r=1}^R$  by performing minibatch SGD on the combined loss:**

$$\theta, \{\hat{\mathbf{A}}^{(r)}\}_{r=1}^R \leftarrow \underset{\theta, \{\hat{\mathbf{A}}^{(r)}\}}{\text{argmin}} \left[ \sum_{i=1}^N \sum_{r=1}^R \mathbb{1}(\tilde{y}_i^{(r)} \neq -1) \cdot \text{CE}(\hat{\mathbf{A}}^{(r)} \hat{\mathbf{p}}_\theta(\mathbf{x}_i), \tilde{y}_i^{(r)}) + \lambda \sum_{r=1}^R \text{tr}(\hat{\mathbf{A}}^{(r)}) \right]$$

**Return:**  $\hat{\mathbf{p}}_\theta$  and  $\{\hat{\mathbf{A}}^{(r)}\}_{r=1}^R$

---

**Algorithm 6.2** Generalized EM [377]

**Inputs:**  $\mathcal{D} = \{\mathbf{x}_n, \tilde{y}_n^{(1)}, \dots, \tilde{y}_n^{(R)}\}_{n=1}^N$ ,  $T$ : # EM steps,  $G$ : # SGD in each M-step  
**Initialize posterior distribution by the mean labels:** for  $j = 1, \dots, L, n = 1, \dots, N$

$$q_{nj}^{(0)} := p(y_n = j | \mathbf{x}_n, \{\tilde{y}_n^{(r)}\}_r, \theta^{(0)}) \leftarrow R^{-1} \sum_{r=1}^R \mathbb{1}(\tilde{y}_n^{(r)} = j)$$

**Initialize the parameters of the base classifier  $\theta$**

**Repeat  $T$  times:**

**M-step for  $\theta$ .** Learn the base classifier  $\hat{\mathbf{p}}_\theta$  by performing minibatch SGD for  $G$  iterations

$$\theta^{(t+1)} \leftarrow \operatorname{argmin}_\theta \left[ - \sum_{n=1}^N \sum_{l=1}^L q_{nj}^{(t)} \cdot \log p(y_n = l | \mathbf{x}_n, \theta) \right]$$

**M-step for  $\{\hat{\mathbf{A}}^{(r)}\}_{r=1}^R$ .** Estimate the confusion matrices

$$\hat{a}_{ji}^{(r),t+1} \leftarrow \frac{\sum_{n=1}^N \mathbb{1}(\tilde{y}_n^{(r)} \neq -1) \cdot \mathbb{1}(\tilde{y}_n^{(r)} = i) \cdot q_{nj}^{(t)}}{\sum_{n=1}^N \mathbb{1}(\tilde{y}_n^{(r)} \neq -1) \cdot q_{nj}^{(t)}}$$

**E-step.** Estimate the posterior label distribution

$$q_{nj}^{(t+1)} \leftarrow \frac{p(y_n = j | \mathbf{x}_n, \theta^{(t+1)}) \cdot \prod_{r=1}^R (\hat{a}_{j\tilde{y}_n^{(r)}}^{(r),t+1})^{\mathbb{1}(\tilde{y}_n^{(r)} \neq -1)}}{\sum_{l=1}^L p(y_n = l | \mathbf{x}_n, \theta^{(t+1)}) \cdot \prod_{r=1}^R (\hat{a}_{l\tilde{y}_n^{(r)}}^{(r),t+1})^{\mathbb{1}(\tilde{y}_n^{(r)} \neq -1)}}$$

**Return:**  $\hat{\mathbf{p}}_{\theta^{(T)}}$  and  $\{\hat{\mathbf{A}}^{(r),T}\}_{r=1}^R$

**Algorithm 6.3** Model-Bootstrapped EM [381]

**Inputs:**  $\mathcal{D} = \{\mathbf{x}_n, \tilde{y}_n^{(1)}, \dots, \tilde{y}_n^{(R)}\}_{n=1}^N$ ,  $T$ : # EM steps,  $G$ : # SGD in each M-step  
**Initialize posterior distribution by the mean labels:** for  $j = 1, \dots, L, n = 1, \dots, N$

$$q_{nj}^{(0)} := p(y_n = j | \mathbf{x}_n, \{\tilde{y}_n^{(r)}\}_r, \theta^{(0)}) \leftarrow R^{-1} \sum_{r=1}^R \mathbb{1}(\tilde{y}_n^{(r)} = j)$$

**Initialize the parameters of the base classifier  $\theta$**

**Repeat  $T$  times:**

**M-step for  $\theta$ .** Learn the base classifier  $\hat{\mathbf{p}}_\theta$  by performing minibatch SGD for  $G$  iterations

$$\theta^{(t+1)} \leftarrow \operatorname{argmin}_\theta \left[ - \sum_{n=1}^N \sum_{l=1}^L q_{nj}^{(t)} \cdot \log p(y_n = l | \mathbf{x}_n, \theta) \right]$$

**Predict on training examples.** for  $n = 1, \dots, N$ :

$$c_n \leftarrow \operatorname{argmax}_{l \in \{1, \dots, L\}} p(y_n = l | \mathbf{x}_n, \theta^{(t+1)})$$

**M-step for  $\{\hat{\mathbf{A}}^{(r)}\}_{r=1}^R$ .** Estimate the CMs. For  $i, j = 1, \dots, L$  and  $r = 1, \dots, R$ :

$$\hat{a}_{ji}^{(r),t+1} \leftarrow \frac{\sum_{n=1}^N \mathbb{1}(\tilde{y}_n^{(r)} \neq -1) \cdot \mathbb{1}(\tilde{y}_n^{(r)} = i) \cdot \mathbb{1}(c_n = j)}{\sum_{n=1}^N \mathbb{1}(\tilde{y}_n^{(r)} \neq -1) \cdot \mathbb{1}(c_n = j)}$$

**Update prior label distribution.** for  $l = 1, \dots, L$ :

$$p_l \leftarrow N^{-1} \sum_{n=1}^N \mathbb{1}(c_n = l)$$

**E-step.** Estimate the posterior label distribution

$$q_{nj}^{(t+1)} \leftarrow \frac{p_j \cdot \prod_{r=1}^R (\hat{a}_{j\tilde{y}_n^{(r)}}^{(r),t+1})^{\mathbb{1}(\tilde{y}_n^{(r)} \neq -1)}}{\sum_{l=1}^L p_l \cdot \prod_{r=1}^R (\hat{a}_{l\tilde{y}_n^{(r)}}^{(r),t+1})^{\mathbb{1}(\tilde{y}_n^{(r)} \neq -1)}}$$

**Return:**  $\hat{\mathbf{p}}_{\theta^{(T)}}$  and  $\{\hat{\mathbf{A}}^{(r),T}\}_{r=1}^R$

### 6.2.3 Motivation for Trace Regularization

Here we intend to motivate the addition of the trace regularizer in eq. 6.4. In the last section, we saw that minimizing cross-entropy loss alone encourages  $\hat{\mathbf{A}}^{(r)} \mathbf{P} \rightarrow \mathbf{A}^{(r)}$ . Therefore, if we could devise a regularizer which, when minimized, uniquely ensures the convergence  $\hat{\mathbf{A}}^{(r)} \rightarrow \mathbf{A}^{(r)}$ , then this would make  $\mathbf{P}$  tend to the identity matrix, implying that the base model fully captures the true label distribution i.e.  $\text{argmax}_k [\hat{\mathbf{p}}(\mathbf{x})_\theta]_k = y \forall \mathbf{x}$ . We describe below the trace regularizer is indeed a such regularizer when both  $\hat{\mathbf{A}}^{(r)}$  and  $\mathbf{A}^{(r)}$  satisfy some conditions. We first show this result assuming that there is a single annotator, and then extend to the scenario with multiple annotators.

**Lemma 6.1** (Single Annotator). *Let  $\mathbf{P}$  be the CM of the estimated true labels  $\hat{\mathbf{p}}_\theta$  and  $\hat{\mathbf{A}}$  be the estimated CM of the annotator. If the model matches the noisy label distribution of the annotator i.e.  $\hat{\mathbf{A}}\mathbf{P} = \mathbf{A}$ , and both  $\hat{\mathbf{A}}$  and  $\mathbf{A}$  are diagonally dominant ( $a_{ii} > a_{ij}$ ,  $\hat{a}_{ii} > \hat{a}_{ij}$ ) for all  $i \neq j$ , then  $\hat{\mathbf{A}}$  with the minimal trace uniquely coincides with the true confusion matrix  $\mathbf{A}$ .*

*Proof.* We show that each diagonal element in the true CM  $\mathbf{A}$  forms a lower bound to the corresponding element in its estimation.

$$a_{ii} = \sum_j \hat{a}_{ij} p_{ji} \leq \sum_j \hat{a}_{ii} p_{ji} = \hat{a}_{ii} (\sum_j p_{ji}) = \hat{a}_{ii} \quad (6.5)$$

for all  $i \in \{1, \dots, L\}$ . It therefore follows that  $\text{tr}(\mathbf{A}) \leq \text{tr}(\hat{\mathbf{A}})$ . We now show that the equality  $\hat{\mathbf{A}} = \mathbf{A}$  is uniquely achieved when the trace is the smallest i.e.  $\text{tr}(\mathbf{A}) = \text{tr}(\hat{\mathbf{A}}) \Rightarrow \mathbf{A} = \hat{\mathbf{A}}$ . From (6.5), if the trace of  $\mathbf{A}$  and  $\hat{\mathbf{A}}$  are the same, we see that their diagonal elements also match i.e.  $a_{ii} = \hat{a}_{ii} \forall i \in \{1, \dots, L\}$ . Now, the non-negativity of all elements in CMs  $\mathbf{P}$  and  $\hat{\mathbf{A}}$ , and the equality  $a_{ii} = \sum_j \hat{a}_{ij} p_{ji}$  imply that  $p_{ji} = \mathbb{1}[i = j]$  i.e.  $\mathbf{P}$  is the identity matrix.  $\square$

We note that the above result was also mentioned in [383] in a more general context of label noise modelling (that neglects annotator information). Here we further augment their proof by showing the uniqueness of solutions (i.e.  $\text{tr}(\mathbf{A}) = \text{tr}(\hat{\mathbf{A}}) \Rightarrow \mathbf{A} = \hat{\mathbf{A}}$ ). We also note that the trace regularization was, in fact, never used in practice in [383] — for implementation reason, the Frobenius norm was used in all their experiments. We now extend this to the multiple annotator regime. We will show later that minimizing the mean trace of all annotators indeed enhances the estimation quality of both CM and true label distributions, particularly in the presence of high annotator disagreement.

**Theorem 6.1** (Multiple Annotators). *Let  $\hat{\mathbf{A}}^{(r)}$  be the estimated CM of annotator  $r$ . If  $\hat{\mathbf{A}}^{(r)} \mathbf{P} = \mathbf{A}^{(r)}$  for  $r = 1, \dots, R$ , and the average true and estimated CMs  $\mathbf{A}^* := R^{-1} \sum_{r=1}^R \mathbf{A}^{(r)}$  and  $\hat{\mathbf{A}}^* := R^{-1} \sum_{r=1}^R \hat{\mathbf{A}}^{(r)}$  are diagonally dominant (**every diagonal element is larger than any other off-diagonal element in the corresponding row**), then  $\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(R)} = \text{argmin}_{\hat{\mathbf{A}}^{(1)}, \dots, \hat{\mathbf{A}}^{(R)}} [\text{tr}(\hat{\mathbf{A}}^*)]$  and such solutions are unique. In other words, when the trace of the mean CM is minimized, the estimation of respective annotator's CMs match the true values.*

*Proof.* As the average CMs  $\mathbf{A}^*$  and  $\hat{\mathbf{A}}^*$  are diagonally dominant and we have  $\mathbf{A}^* = \hat{\mathbf{A}}^* \mathbf{P}$ , Lemma 6.1 yields that  $\text{tr}(\mathbf{A}^*) \leq \text{tr}(\hat{\mathbf{A}}^*)$  with equality if and only if  $\mathbf{A}^* = \hat{\mathbf{A}}^*$ . Therefore, when the trace of the average CM of annotators is minimized i.e.  $\text{tr}(\hat{\mathbf{A}}^*) = \text{tr}(\mathbf{A}^*)$ , the estimated CM of the true label distribution  $\mathbf{P}$  reduces to identity, giving  $\hat{\mathbf{A}}^{(r)} = \mathbf{A}^{(r)}$  for all  $r \in \{1, \dots, R\}$ .  $\square$

The above result shows that if each estimated annotator's distribution  $\hat{\mathbf{A}}^{(r)} \hat{\mathbf{p}}_\theta(\mathbf{x})$  is very close to the true noisy distribution  $\mathbf{p}^{(r)}(\mathbf{x})$  (which is encouraged by minimizing the cross-entropy loss), and on average for each class  $c$ , the number of correctly labelled examples exceeds the number of examples of every other class  $c'$  that are mislabelled as  $c$  (the mean CM is diagonally dominant), then minimizing its trace will drive the estimates of CMs towards the true values. To encourage  $\{\hat{\mathbf{A}}^{(1)}, \dots, \hat{\mathbf{A}}^{(R)}\}$  to be also diagonally dominant, we initialize them with identity matrices. Intuitively, the combination of the trace term and cross-entropy separates the true distribution from the annotation noise by finding the maximal amount of confusion which can explain the noisy observations well.

## 6.3 Experiments and Results

We now aim to verify the proposed method on various image recognition tasks. Particularly, we demonstrate (1) advantage of our simpler optimization scheme compared to EM-based approaches (Sec. 6.3.3), (2) importance of modeling multiple annotators (Sec. 6.3.4) and (3) the applicability of the model in a challenging real world application (Sec. 6.3.5). We address the first two questions by testing the proposed method on MNIST and CIFAR-10 datasets with a diverse set of simulated annotators. To answer the final question, we evaluate our approach on the task of cardiac view classification using ultrasound images where the labels are noisy and sparse, and are acquired from multiple annotators of varying levels of expertise.

### 6.3.1 Datasets, training and architectures

**Datasets.** We evaluated our method on three classification datasets: MNIST digit classification dataset [307]; CIFAR-10 object recognition dataset [308]; the cardiac view classification (CVC) dataset from a handheld ultra-sound probe. The MNIST dataset consists of 60,000 training and 10,000 testing examples, all of which are  $28 \times 28$  grayscale images of digits from 0 to 9. The CIFAR-10 dataset consists of 50,000 training and 10,000 testing examples, all of which are  $32 \times 32$  coloured natural images drawn from 10 classes.

The CVC data set contains 26,000 training and 20,000 test examples, which are grayscale images of size  $96 \times 96$  from 6 different cardiac views (see Figure 6.9 (a)). Each image is labelled by a subset of 8 annotators (6 sonographers and 2 non-experts). A committee of sonographers (with varying levels of experience) were tasked with providing the cardiac view labels to a large volume of ultra-sound images, and each example is only labelled by a subset of them. To acquire ground truth in this setting, we chose those samples where the three most experienced sonographers agreed on a given label. The resulting data set consists of noisy labels provided by the remaining less experienced 6 sonographers for a total of 240,000 training images and 22,000 validation images. In addition, we also acquired labels from two non-expert users and included in the training data.

**Training.** For all experiments, we employ the same training scheme unless otherwise stated. We optimize parameters using Adam [190] with initial learning rate of  $10^{-3}$  and  $\beta = [0.9, 0.999]$ , with minibatches of size 50 and train for 200 epochs. For our method, we set the scale of the trace regularization to  $\lambda = 0.01$ . For the training of the EM-based approaches (Model-Bootstrapped EM [381] and generalized EM [377]), we train the base classifier for 200 epochs in total over the course of the EM steps, following the same protocol above. For CIFAR-10, we performed two iterations of EM algorithm ( $T = 2$ ) and 100 epochs worth of gradient descent steps during each E-step to update the parameters of the base classifier ( $G = 100$  epochs), following the original implementation in [381]. For the experiments on the CVC data set, we run more rounds of EM with  $T = 10$  and  $G = 20$  epochs. In all cases, we hold out 10% of training images as a validation set and best model is selected based on the validation accuracy over the course of training. No data augmentation is performed during training in all three data sets. We note that for CIFAR-10, we, in addition, decreased the learning rate by a factor of 10 at every multiple of 50 in a similar fashion to the schedule used in [396, 326, 351].

**Architectures.** For MNIST, the base classifier was defined as a CNN architecture comprised of 4 convolution layers, each with  $3 \times 3$  kernels follower by ReLU. The number of kernerls in respective layers are  $\{32, 32, 64, 64\}$ . After the first two convolution layers, we perform  $2 \times 2$  max-pooling, and after the last one, we further down-sample the features with Global Average Pooling (GAP) prior to the final fully connected layer. For the CVC dataset, we employed the same architecture, but with increased number of kernels i.e.  $\{128, 128, 128, 128\}$ . For CIFAR-10, we used a 50-layer ResNet [326].

### 6.3.2 Set-Up

We focus on a regime in which models have only access to noisy labels from multiple annotators. For MNIST and CIFAR-10 data sets, we simulate noisy labels from a range of annotators with different skill levels and biases.

**MNIST Experiments.** We consider two different models of annotator types: (i) *pairwise-flipper*: each annotator is correct with probability  $p$  or flips the label of each class to another label (the flipping target is chosen uniformly at random for each class), (ii) *hammer-spammer*: each annotator is always correct with probability  $p$  or otherwise chooses labels uniformly at random [381]. Example confusion matrices for both cases are shown in Fig. 6.3.

For each annotator type and skill level  $p$ , we create a group of 5 annotators by generating confusion matrices (CMs) from the associated distribution. More specifically, each CM is generated by perturbing the mean skill level  $p$  by injecting a small Gaussian noise  $\epsilon \sim \text{Normal}(0, 0.01)$  and choosing the flipping target class randomly in the case of a *pairwise-flipper*. Given the GT labels, we generate noisy labels as defined by the CM per annotator. These noisy labels are used during training.

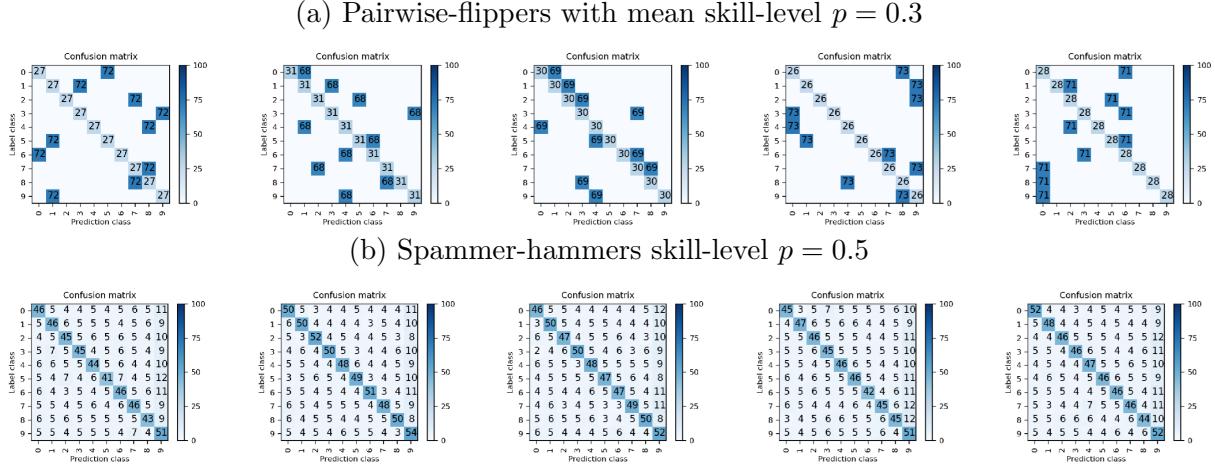


Figure 6.3: Examples of annotator groups. The value of diagonal entries are fixed constant for each annotator and is drawn from  $\text{Normal}(p, 10^{-2})$ .

**CIFAR-10 Experiments.** We consider a diverse group of 4 annotators with different patterns of CMs as shown in Fig. 6.2: (i) is a “hammer-spammer” as defined above, (ii) tends to mix up semantically similar categories of images e.g. cats and dogs, and automobiles and trucks, (iii) is likely to confuse “neighbouring” classes and (iv) is an adversarial annotator who has a wrong association of class names to object categories. On average, labels generated by these annotators are correct only 45% of the time.

In synthetic experiments, we assume that equal number of labels are generated by each annotator on average. We also note that all models are trained on noisy labels and do not have access to the ground truth. Unless otherwise stated, we hold out 10% of training images as a validation set, on which the best performing model is selected. We also perform no data augmentation during training. Full details of training and model architectures are provided in the supplementary materials of [107]. In Sec. 6.3.3 and Sec. 6.3.4 below, we compare our model against two separate sets of baselines to address different questions.

### 6.3.3 Comparing with EM-based Approaches

This section examines the ability of our method in learning the CMs of annotators and the GT label distribution on MNIST and CIFAR-10. In particular, we compare against two prior methods: (1) generalized EM [184], the first method for end-to-end training of the CM model in the presence of multiple annotators, and (2) Model Bootstrapped EM (MBEM) [381], the present state-of-the-art method. We analyze the performance in two cases, one in which all labels from 5 annotators are available for each image (“dense labels”), and another where only one randomly selected annotator labels each example (“1 label per image”). We quantify the error of CM estimation by the average Frobenius norm between each CM and its estimate over the annotators, and this metric is normalized to be in the range  $[0, 1]$  by dividing by the number of classes  $L$  i.e.  $R^{-1}L^{-1} \sum_r \sum_{i,j} \|a_{ij}^{(r)} - \hat{a}_{ij}^{(r)}\|^2$ .

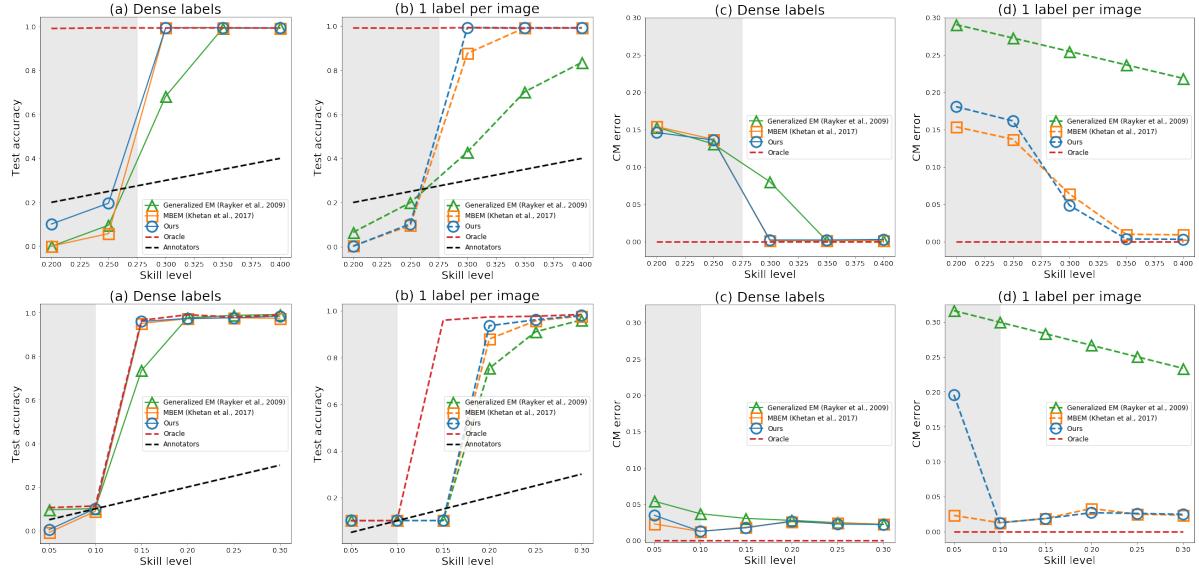


Figure 6.4: Comparison between our method, generalized EM, MBEM trained on noisy labels on MNIST from “pairwise flippers” (Top row) and “hammer-spammers” (Bottom row) for a range of mean skill level  $p$ . (a), (b) show classification accuracy in two cases, one where all annotators label each example and the other where only one label is available per example. (c), (d) quantify the CM recovery error as the annotator-wise average of the normalized Frobenius norm between each ground truth CM and its estimate. The shaded areas represent the cases where the average CM over the annotators are not diagonally dominant.

**Performance Comparison.** Fig. 6.4 compares the classification accuracy and the error of CM estimation on MNIST for a range of mean skill-levels  $p$  where labels are generated by a group of 5 “pairwise-flippers” and “hammer-spammers”. The “oracle” model is the idealistic scenario where CMs of the annotators are a priori known to the model while “annotators” indicate the average labeling accuracy of each annotator group.

The top row in Fig. 6.4 shows a strong correlation between the classification accuracy and the error of CM estimation. We observe our model displays consistently better or comparable performance in terms of both classification accuracy and estimation of CMs with dense labels (Fig. 6.4(a) and (c)). When each example receives only one label from one of the annotators, we observe the same trend as long as the mean CMs are diagonally dominant (Fig. 6.4(b,d)). We also observe that when the diagonal dominance holds, all three methods perform better than the annotators. On the other hand, when the diagonal dominance does not hold (see the grey regions), all models undergo a steep drop in classification accuracy due to the inability to estimate CMs accurately as reflected in Fig. 6.4(c,d), which is consistent with Theorem 7.1. Fig. 6.5 also visualizes the average of the estimated CMs at this break point. We also note that with only one label per image, the generalized EM algorithm [377, 184] is not capable of recovering CMs at all and predict identity matrices (Fig. 6.5), which renders the model equivalent to a vanilla classifier directly trained on noisy labels. A similar set of results are also obtained in the “spammer-hammer” case as shown in the second row of Fig. 6.4.

On CIFAR-10 dataset, Tab. 6.1 shows that our method outperforms MBEM and the generalized EM in terms of both classification accuracy and CM estimation by a large margin. In addition, the standard deviations of these metrics are generally smaller for our method than for the baselines. Fig. 6.2 illustrates that our method can estimate CMs of the 4 very different annotators even when each image receives only one label.

**Ablation of Trace Regularization.** Interestingly, Tab. 6.1 shows that even removing the trace norm can achieve reasonably high classification accuracy and low CM estimation error. We believe that this is due to the natural robustness of CNN classifiers to label noise; an accurate estimation of the label distribution  $\hat{\mathbf{p}}_\theta$  regularizes the estimation of confusion matrices  $\{\hat{\mathbf{A}}^{(r)}\}_{r=1}^R$ . Nevertheless, adding the trace norm still improves the performance.

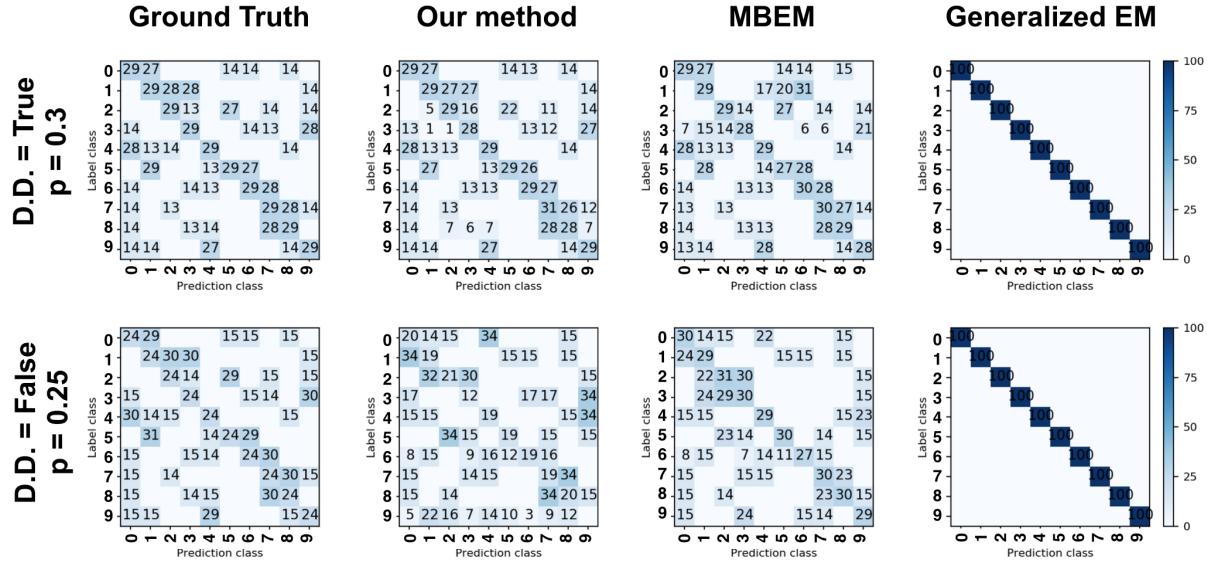


Figure 6.5: Visualization of the mean CM estimates when the diagonal dominance (D.D.) holds (mean skill level,  $p = 0.3$ ) and does not hold ( $p = 0.25$ ). In all cases, only one label is provided per image. The numbers are rounded to nearest integers. Here the respective models are trained on the noisy labels from 5 “pairwise flippers”. Note that when each image receives only 1 label, the generalised EM [377] completely fails to recover the CM due to the failure of M-step for updating the confusion matrices (see Algorithm 6.2).

We observe on MNIST that such improvement is much pronounced in the presence of larger noise. We compare our method on MNIST against the case where the trace norm regularization is removed (results on CIFAR-10 and CVC datasets are given in the main text). Fig. 6.6 shows that adding the trace norm generally improves the performance in terms of both classification accuracy and CM estimation error, and such improvement is pronounced in the presence of larger noise i.e. lower skill levels of annotators. We also observe that when the noise level is low, our model still attains very high accuracy even without trace norm regularization. This can be explained by the natural robustness of the CNN classifier; if the amount of label noise is sufficiently small, the base classifier is still capable of learning the true label distribution well. This, in turn, allows the model to separate annotation noise from true label distribution, improving the quality of CM estimation and thus the overall performance. However, in the presence of large label noise, having trace-norm regularization shows evident benefits.

Lastly, we note that our method, without the trace norm, performs almost as poorly as the naive classifier on the real-world medical imaging dataset (Fig. 6.9 (d)).

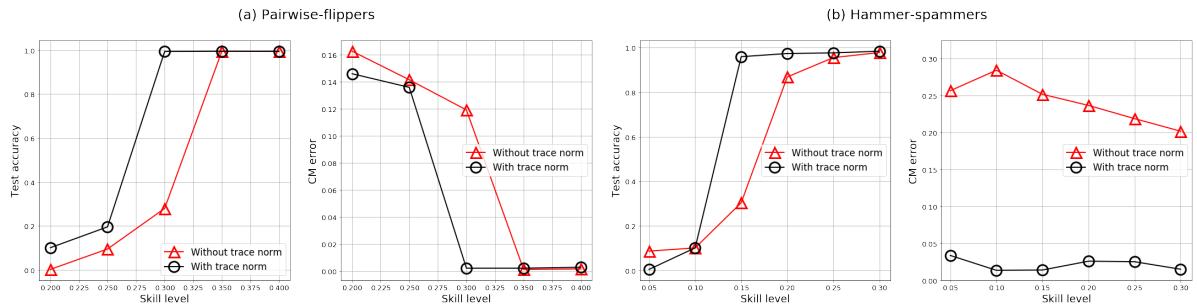


Figure 6.6: Comparison between our method with and without trace norm on MNIST. Results for two annotator groups, consisting of “hammer-spammers” and “pairwise-flippers” are shown for a range of mean skill level  $p$ .

(a) Dense labels			(b) 1 label per image		
Method	Accuracy	CM error	Method	Accuracy	CM error
Our method	<b>81.23 ± 0.21</b>	<b>0.72 ± 0.01</b>	Our method	<b>77.65 ± 0.31</b>	<b>1.22 ± 0.01</b>
Our method (no trace norm)	80.29 ± 0.65	1.37 ± 0.12	Our method (no trace norm)	76.31 ± 0.49	1.46 ± 0.27
MBEM [381]	73.33 ± 0.46	2.53 ± 0.24	MBEM [381]	55.97 ± 1.23	4.58 ± 0.64
generalized EM [377]	70.49 ± 0.23	6.13 ± 0.28	generalized EM [377]	53.38 ± 0.71	4.47 ± 0.64
Single CM [383]	68.82 ± 2.27	-	Single CM [383]	59.91 ± 0.98	-
Weighted Doctor Net [385]	60.11 ± 1.80	-	Weighted Doctor Net [385]	57.98 ± 0.14	-
Soft-bootstrap [382]	54.73 ± 1.33	-	Soft-bootstrap [382]	42.91 ± 1.08	-
Vanilla CNN [382]	52.33 ± 0.31	-	Vanilla CNN [382]	36.04 ± 1.04	-

Table 6.1: Mean classification accuracy and CM estimation errors ( $\times 10^{-2}$ ) on CIFAR-10 with dense labels. Average annotator accuracy is 45%. Standard deviations are computed based on 3 runs with varied weight initialization.

**Sensitivity to Hyper-parameters.** We next study the robustness of our method against the generalized EM and MBEM to the specification of hyper-parameters. We used the group of five pairwise-flippers with the mean skill level  $p = 0.35$  to generate noisy labels on MNIST data set. For our model, we compare the effects of the scaling  $\lambda$  of the trace-norm in eq. 6.4 on the trajectory of classification accuracy on the validation set and the quality of CM estimation. For the baselines, we experiment by varying the number of EM steps (denoted by  $T$ ) and the number of stochastic gradient descent for each E-step (denoted by  $G$ ) while fixing the total number of training iterations at 100,000. We observed our model presents robustness to different values of  $\lambda$  as long as the trace-norm loss is not larger than the cross-entropy loss (where the estimated CMs will start to diffuse too much), and Fig. 6.7 shows the stability of the validation curves for  $\lambda \in \{0.1, 0.01, 0.001\}$ . Both the MBEM and generalized EM show evident dependence on the values of  $T$  and  $G$  and by and large display slower convergence than our method. We also observe that if too few gradient descents are performed ( $G = 1000$ ) during each E-step, the model converges to a lower accuracy in both classification and CM estimation.

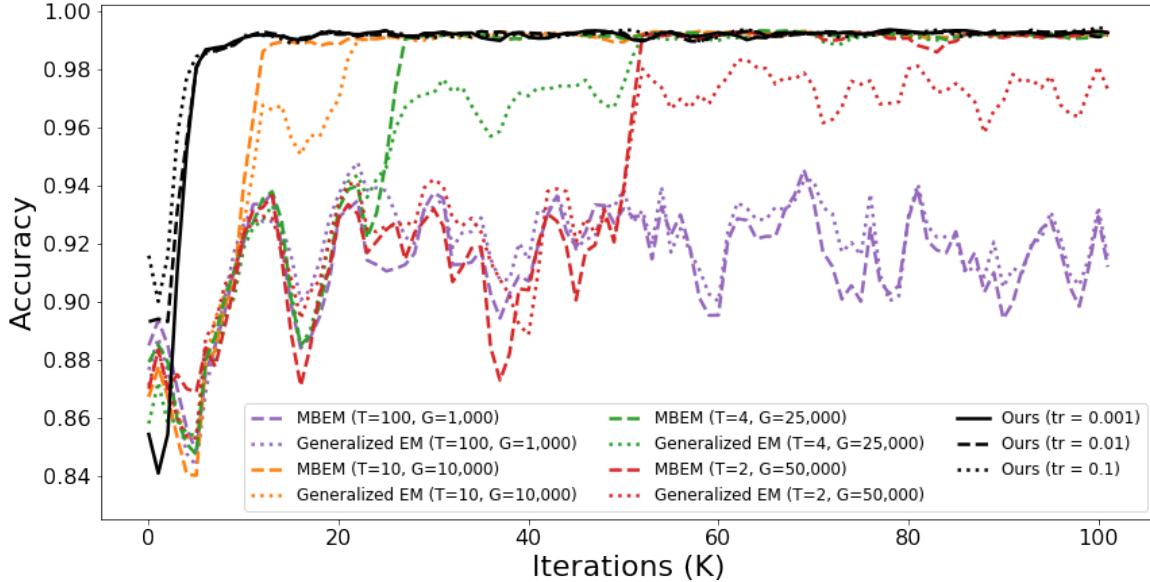


Figure 6.7: Curves of validation accuracy during training of our method, generalized EM and MBEM for a range of hyper-parameters. For our method, the scaling of the trace regularizer is varied in  $[0.001, 0.01, 0.1]$ , while, for EM and MBEM, we vary the number of EM steps ( $T$ ), and the number of gradient descent steps per E-step ( $G$ ) while fixing the total number of training iterations at 100,000.

### 6.3.4 Value of Modelling Individual Annotators

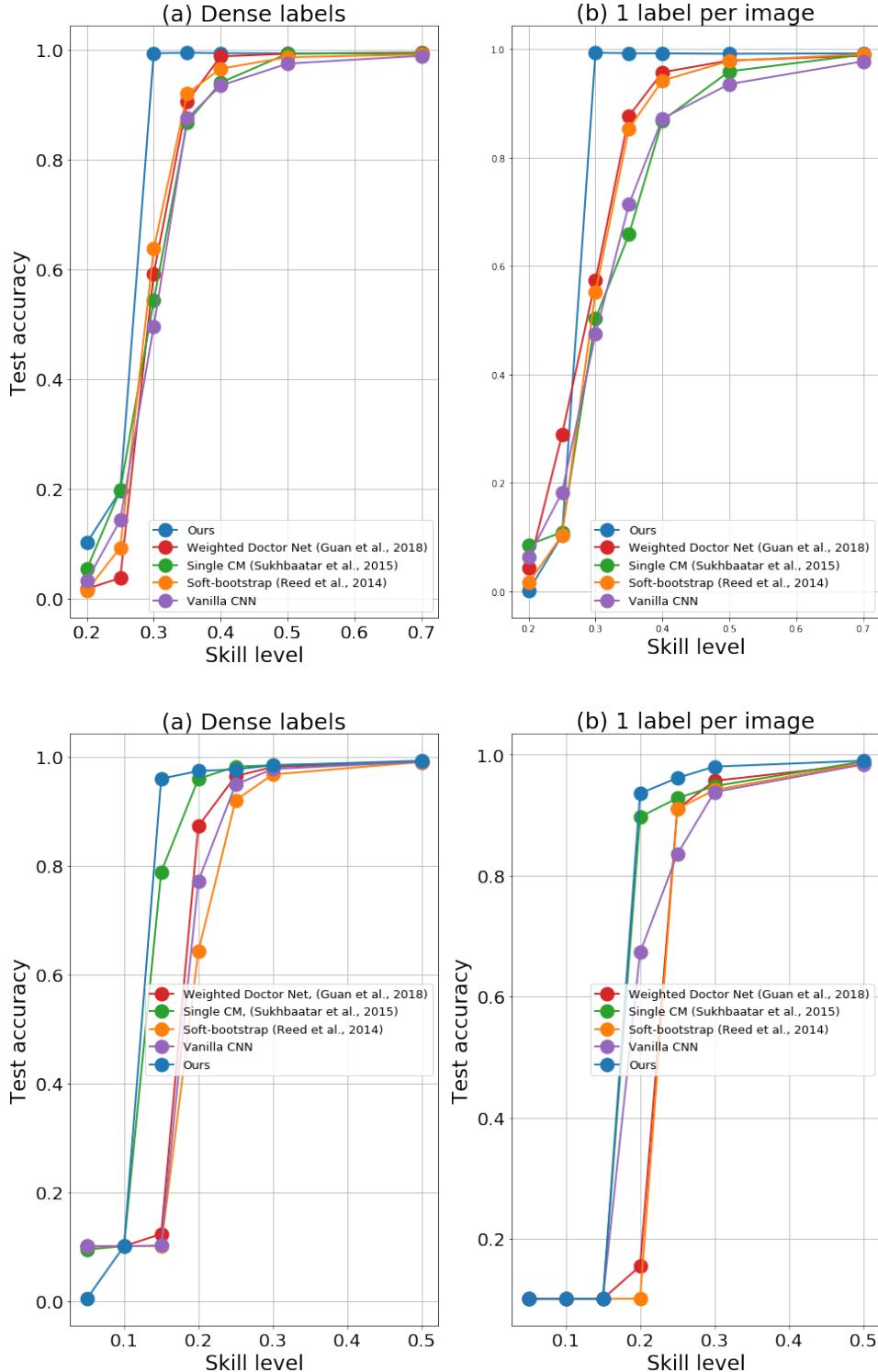


Figure 6.8: Classification accuracy on MNIST of different noise-robust models as a function of the mean annotator skill level  $p$  in two cases. (Top row): For each mean skill-level  $p$ , a group of 5 “pairwise flippers” is formed and used to generate labels. (a). each example receives labels from all the annotators. (b). each example is labelled by only 1 randomly selected annotator. (Bottom row): the equivalent results for a group of 5 “hammer-spammers” with varying mean skill-levels are presented.

Now, we compare the performance of our method against the prior work that aim to improve robustness to noisy labels without explicitly modelling the individual annotators. The first baseline is the vanilla classifier trained on the majority vote labels. We also compare against the noise robust approaches proposed in [382] and [383]. Reed *et al.* [382] adds to the cross-entropy loss a label consistency term based on the negative entropy of the softmax outputs, and we used the default hyper-parameter  $\beta = 0.95$  for comparison. Sukhbaatar *et al.* [383] explicitly accounts for the label noise with a single CM, but does not model individual annotators. We add the trace-norm of the same scaling used in our method ( $\lambda = 0.01$ ) to the loss function for training. We also include Weighted Doctor Net architecture (WDN) [385] in the comparison, a recent method that models the annotators individually and then learns averaging weights for combining them. It should be noted that this model considers a different observation model of the labels and does not explicitly model the true label distribution. When we have access to multiple labels per example, with the exception of WDN, we aggregated the labels by computing the majority vote and trained all models. This is because we observed a consistent improvement on validation accuracy (thus poses a tougher challenge against our method) and this would be a more realistic utilization of such data set. For both MNIST and CIFAR-10 experiments, we test on the same set of simulated labels as used in Sec. 6.3.3.

Fig. 6.8 shows better or comparable classification accuracy than all the baselines when the diagonal dominance of the mean CM holds. In particular, our methods show significant improvement when the mean skill level of the annotators are relatively low (e.g.  $p = 0.3$  and  $0.35$ ). The results are pronounced in the case with only one label available per image for which the baseline methods undergo a steep drop in accuracy (see Fig. 6.8(b)). Similarly on CIFAR-10 data set, Tab. 6.1 shows that our method improves the classification accuracy upon the baselines. Such improvement is pronounced in the case of sparse labels. On the other hand, a vanilla CNN with only L2 weight decay overfits to the training data very quickly in the presence of such high noise.

### 6.3.5 Experiments on Cardiac View Classification

Lastly, we illustrate the results of our approach for a real data set with sparse and noisy labels from the medical domain. This data set consists of images of the cardiac region in different views, acquired using a hand-held ultrasound probe. The task is to classify a given ultrasound image into one of six different view classes (see Fig. 6.9(a)). The process of obtaining a cardiac view label is crucial for guiding the user to the correct locations of measurements, and affects the quality of the downstream cardiac tasks.

We estimated the skill-level of each annotator by computing the average value of the diagonal elements in the corresponding learned CM, and Fig. 6.9(b) shows that the group of experts can be separated from the two non-experts with varying levels of experinces (one is less competent than the other). Fig. 6.9(c) shows that confusion between  $A3C$  and  $A5C$ , even common among experts, can be detected (see the result for ‘Expert 1’) while clearly capturing the patterns of mistakes for the non-experts. In addition, Fig. 6.9(d) shows that our model outperforms MBEM [381] again in classification accuracy and the quality of CM estimation. Lastly, the higher classification accuracy of our model with respect to the other baseline models illustrates again that modelling individual annotators improves robustness to label noise.

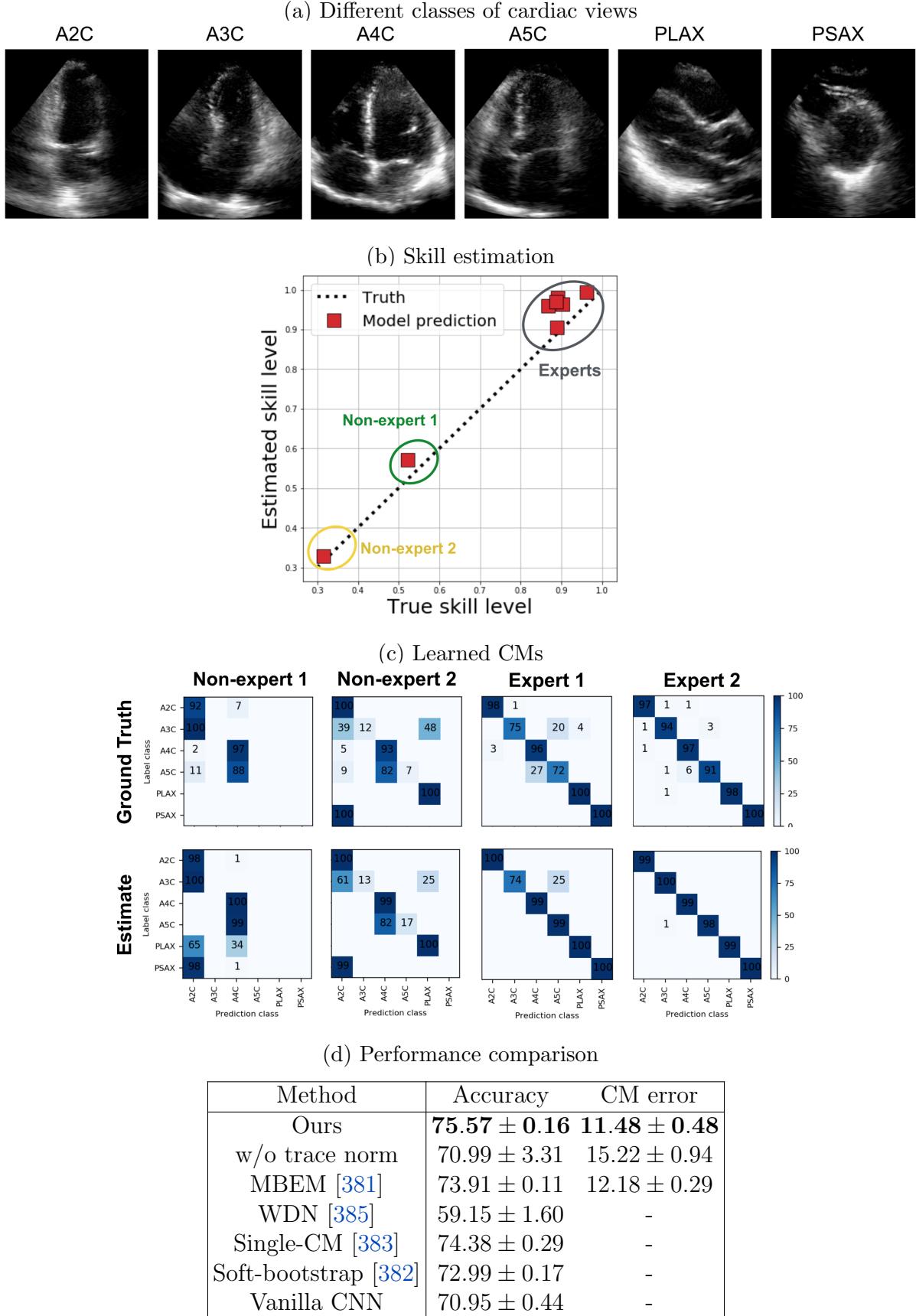


Figure 6.9: Results on the cardiac view classification dataset: (a) illustrates examples of different cardiac view images. (b) plots the estimated skill level of each annotator (average of the diagonal elements of its estimated CM) against the ground truth (c) compares the estimated CMs of the two least skilled and two most skilled annotators according the GT labels (d) summarizes the classification accuracy and error of CM estimation for different methods.

## 6.4 Discussion and Conclusion

We introduced a new theoretically grounded algorithm for simultaneously recovering the label noise of multiple annotators and the ground truth label distribution. Our method enjoys implementation simplicity, requiring only adding a regularization term to the loss function. Experiments on both synthetic and real data sets have shown superior performance over the common EM-based methods in terms of both classification accuracy and the quality of confusion matrix estimation. Comparison against the other modern noise-robust methods demonstrates that the modelling individual annotators improves robustness to label noise. Furthermore, the method is capable of estimating annotation noise even when there is a single label per image.

There are many exciting avenues of future research. Our work was primarily motivated by medical imaging applications for which the number of classes are mostly limited to below 10. However, future work shall consider imposing structures on the confusion matrices to broaden up the applicability to massively multi-class scenarios e.g. introducing taxonomy based sparsity [380] and low-rank approximation.

We also assumed that there is only one ground truth for each input; this no longer holds true when the input images are truly ambiguous—recent advances in modelling multi-modality of label distributions [397, 182] potentially facilitate relaxation of such assumption. Another limiting assumption is the image independence of the annotator’s label noise. The majority of disagreement between annotators arise in the difficult cases. Integrating such input dependence of label noise [378, 398] is also a valuable next step.

Another simplification to note is that the annotators are assumed statistically independent. In many data collection processes, such assumption is not realistic (e.g., communication between experts may be encouraged to improve the quality of labels). Recently, Yuan et al., [399] showed that modelling annotator correlations, improves accuracy in high-noise crowdsourcing applications. Our model would certainly benefit from such approach but efficient training and inference schemes need to be considered.

Lastly, we are also interested in incorporating our approach into active learning schemes, where the current knowledge of the confusion matrices may be used to decide, not only which examples to annotate, but which annotator to route them to. Similarly, such information about annotators’ characteristics could be potentially used for further education of less experienced annotators. We believe such applications of the human uncertainty modelling have utility in designing a more efficient data collection pipeline.

# Chapter 7

## Modelling Human Uncertainty (II): Semantic Segmentation

**Abstract:** In this chapter, we extend the method introduced in Chapter 6 to the more challenging task of semantic segmentation where every pixel in a given input image is classified. To this end, we model the annotation mistakes of each annotator as a pixel-wise confusion matrix that is a function of the input image. Analogous to the classification approach, the separation between the true labels and the annotation mistakes is achieved by encouraging the estimated annotators to be maximally unreliable while achieving high fidelity with the noisy training data. We first define a toy segmentation dataset based on the MNIST dataset and empirically study the behaviours of the proposed algorithm. We then demonstrate the utility of the method on three public medical imaging segmentation datasets with simulated (when necessary) and real diverse annotations: 1) MSLSC (multiple-sclerosis lesions); 2) BraTS (brain tumours); 3) LIDC-IDRI (lung abnormalities). In all cases, our method outperforms competing methods and relevant baselines particularly in cases where the number of annotations is small and the amount of disagreement is large. The experiments also show strong ability to capture the complex spatial characteristics of annotators' mistakes. Our code is available at [https://github.com/moucheng2017/Learn\\_Noisy\\_Labels\\_Medical\\_Images](https://github.com/moucheng2017/Learn_Noisy_Labels_Medical_Images). This chapter is based on a joint work [108] with Le Zhang at UCL where I primarily contributed to the method development, theoretical results, implementation and experiment design.

### 7.1 Introduction

Segmentation of anatomical structures in medical images is known to suffer from high inter-reader variability [366, 364, 365, 367, 62], influencing the performance of downstream supervised machine learning models. This problem is particularly prominent in the medical domain where the labelled data is commonly scarce due to the high cost of annotations. For instance, accurate identification of multiple sclerosis (MS) lesions in MRIs is difficult even for experienced experts due to variability in lesion location, size, shape and anatomical variability across patients [400]. Another example [367] reports the average inter-reader variability in the range 74-85% for glioblastoma (a type of brain tumour) segmentation. Further aggravated by differences in biases and levels of expertise, segmentation annotations of structures in medical images suffer from high annotation variations [401]. In consequence, despite the present abundance of medical imaging data thanks to over two decades of digitisation, the world still remains relatively short of access to data with curated labels [368], that is amenable to machine learning, necessitating intelligent methods to learn robustly from such noisy annotations.

To mitigate inter-reader variations, different pre-processing techniques are commonly used to curate segmentation annotations by fusing labels from different experts. The most basic yet popular approach is based on the majority vote where the most representative opinion of the experts is treated as the ground truth (GT). A smarter version that accounts for similarity of classes has proven effective in aggregation of brain tumour segmentation labels [367]. A key limitation of such approaches, however, is that all experts are assumed to be equally reliable. Warfield *et al.*[60] proposed a label fusion method, called STAPLE that explicitly models the reliability of individual experts and uses that information to “weigh” their opinions in the label aggregation step. After consistent demonstration of its superiority over the

standard majority-vote pre-processing in multiple applications, STAPLE has become the go-to label fusion method in the creation of public medical image segmentation datasets e.g., ISLES [402], MSSeg [403], Gleason’19 [404] datasets. Asman *et al.* later extended this approach in [245] by accounting for voxel-wise consensus to address the issue of under-estimation of annotators’ reliability. In [405], another extension was proposed in order to model the reliability of annotators across different pixels in images. More recently, within the context of multi-atlas segmentation problems [406] where image registration is used to warp segments from labeled images (“atlases”) onto a new scan, STAPLE has been enhanced in multiple ways to encode the information of the underlying images into the label aggregation process. A notable example is STEP proposed in Cardoso *et al.* [407] who designed a strategy to further incorporate the local morphological similarity between atlases and target images, and different extensions of this approach such as [408, 409] have since been considered. However, these previous label fusion approaches have a common drawback—they critically lack a mechanism to integrate information across different training images. This fundamentally limits the remit of applications to cases where each image comes with a reasonable number of annotations from multiple experts, which can be prohibitively expensive in practice. Moreover, relatively simplistic functions are used to model the relationship between observed noisy annotations, true labels and reliability of experts, which may fail to capture complex characteristics of human annotators.

In this work, we introduce the first instance of an end-to-end supervised segmentation method that jointly estimates, from noisy labels alone, the reliability of multiple human annotators and true segmentation labels. The proposed architecture (Fig. 7.1) consists of two coupled CNNs where one estimates the true segmentation probabilities and the other models the characteristics of individual annotators (e.g., tendency to over-segmentation, mix-up between different classes, etc) by estimating the pixel-wise confusion matrices (CMs) on a per image basis. Unlike STAPLE [60] and its variants, our method models, and disentangles with deep neural networks, the complex mappings from the input images to the annotator behaviours and to the true segmentation label. Furthermore, the parameters of the CNNs are “global variables” that are optimised across different image samples; this enables the model to disentangle robustly the annotators’ mistakes and the true labels based on correlations between similar image samples, even when the number of available annotations is small per image (e.g., a single annotation per image). In contrast, this would not be possible with STAPLE [60] and its variants [405, 407] where the annotators’ parameters are estimated on every target image separately.

For evaluation, we first simulate a diverse range of annotator types on the MNIST dataset by performing morphometric operations with Morpho-MNIST framework [410]. Then we demonstrate the potential in several real-world medical imaging datasets, namely (i) MS lesion segmentation dataset (MSLSC) from the ISBI 2015 challenge [411], (ii) Brain tumour segmentation dataset (BraTS) [367] and (iii) Lung nodule segmentation dataset (LIDC-IDRI) [412]. Experiments on all datasets demonstrate that our method consistently leads to better segmentation performance compared to widely adopted label-fusion methods and other relevant baselines, especially when the number of available labels for each image is low and the degree of annotator disagreement is high.

## 7.2 Related Works

The majority of algorithmic innovations in the space of *label aggregation for segmentation* have uniquely originated from the medical imaging community, partly due to the prominence of the inter-reader variability problem in the field, and the wide-reaching values of reliable segmentation methods [405]. The aforementioned methods based on the STAPLE-framework such as [60, 245, 405, 407, 413, 408, 408, 409, 61] are based on generative models of human behaviours, where the latent variables of interest are the unobserved true labels and the “reliability” of the respective annotators. Our method can be viewed as an instance of translation of the STAPLE-framework to the supervised learning paradigm. As such, our method produces a model that can segment test images without needing to acquire labels from annotators or atlases unlike STAPLE and its local variants. Another key difference is that our method is jointly trained on many different subjects while the STAPLE-variants are only fitted on a per-subject basis. This means that our method is able to learn from correlations between different subjects, which previous works have not attempted—for example, our method uniquely can estimate the reliability and true labels even when there is only one label available per input image as shown later.

Our work also relates to a recent strand of methods that aim to generate a set of diverse and plausible segmentation proposals on a given image. Notably, probabilistic U-net [182] and its recent variants, PHiSeg [414] have shown that the aforementioned inter-reader variations in segmentation labels can be

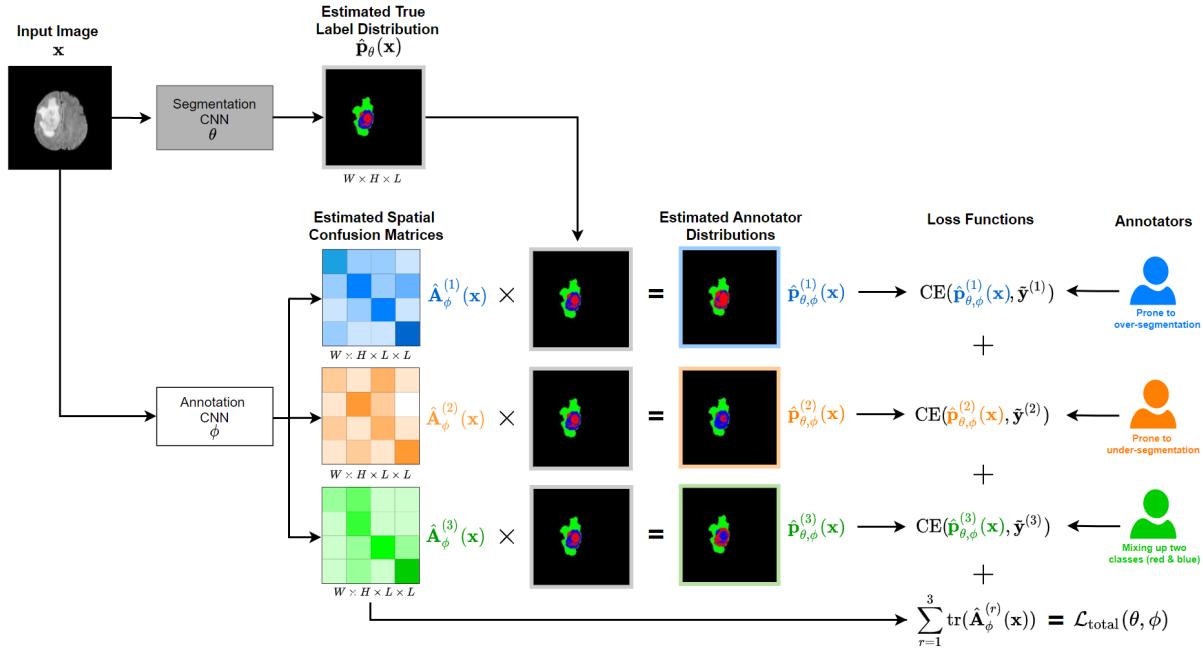


Figure 7.1: An architecture schematic in the presence of 3 annotators of varying characteristics (over-segmentation, under-segmentation and confusing between two classes, red and blue). The model consists of two parts: (1) *segmentation network* parametrised by  $\theta$  that generates an estimate of the unobserved true segmentation probabilities,  $\mathbf{p}_\theta(\mathbf{x})$ ; (2) *annotator network*, parametrised by  $\phi$ , that estimates the pixelwise confusion matrices (CMs),  $\{\mathbf{A}_\phi^{(r)}(\mathbf{x})\}_{r=1}^3$  of the annotators for the given input image  $\mathbf{x}$ . During training, the estimated annotator distributions  $\hat{\mathbf{p}}_{\theta,\phi}^{(r)}(\mathbf{x}) := \mathbf{A}_\phi^{(r)}(\mathbf{x}) \cdot \mathbf{p}_\theta(\mathbf{x})$  are computed, and the parameters  $\{\theta, \phi\}$  are learned by minimizing the sum of their cross-entropy losses with respect to the acquired noisy segmentation labels  $\tilde{\mathbf{y}}^{(r)}$ , and the trace of the estimated CMs. At test time, the output of the segmentation network,  $\hat{\mathbf{p}}_\theta(\mathbf{x})$  is used to yield the prediction.

modelled with sophisticated forms of probabilistic CNNs. Such approaches, however, fundamentally differ from ours in that variable annotations from many experts in the training data are assumed to be all realistic instances of the true segmentation; we assume, on the other hand, that there is a single, unknown, true segmentation map of the underlying anatomy, and each individual annotator produces a noisy approximation to it with variations that reflect their individual characteristics. The latter assumption may be reasonable in the context of segmentation problems since there exists only one true boundary of the physical objects captured in an image while multiple hypothesis can arise from ambiguities in human interpretations.

We also note that, in standard classification problems, a plethora of different works have shown the utility of modelling the labeling process of human annotators in restoring the true label distribution [184, 381, 107]. Such approaches can be categorized into two groups: (1) *two-stage* approach [372, 373, 374, 375, 376], and (2) *simultaneous* approach [377, 378, 379, 380, 381, 107, 415]. In the first category, the noisy labels are first curated through a probabilistic model of annotators, and subsequently, a supervised machine-learning model is trained on the curated labels. The initial attempt [372] was made in the early 1970s, and numerous advances such as [373, 374, 375, 376] since built upon this work e.g. by estimating sample difficulty and human biases. In contrast, models in the second category aim to curate labels and learn a supervised model jointly in an end-to-end fashion [377, 378, 379, 380, 381, 107] so that the two components inform each other. Although the evidence still remains limited to the simple classification task, these *simultaneous* approaches have shown promising improvements over the methods in the first category in terms of the predictive performance of the supervised model and the sample efficiency (i.e., fewer labels are required per input). However, to date very little attention has been paid to the same problem in more complicated, structured prediction tasks where the outputs are high dimensional. In this work, we propose the first *simultaneous* approach to addressing such a problem for image segmentation, while drawing inspirations from the STAPLE framework [60] which would fall into the *two-stage* approach category.

## 7.3 Method

### 7.3.1 Problem Set-up

In this work, we consider the problem of learning a supervised segmentation model from noisy labels acquired from multiple human annotators. Specifically, we consider a scenario where set of images  $\{\mathbf{x}_n \in \mathbb{R}^{W \times H \times C}\}_{n=1}^N$  (with  $W, H, C$  denoting the width, height and channels of the image) are assigned with noisy segmentation labels  $\{\tilde{\mathbf{y}}_n^{(r)} \in \mathcal{Y}^{W \times H}\}_{n=1, \dots, N}^{r \in S(\mathbf{x}_n)}$  from multiple annotators where  $\tilde{\mathbf{y}}_n^{(r)}$  denotes the label from annotator  $r \in \{1, \dots, R\}$  and  $S(\mathbf{x}_n)$  denotes the set of all annotators who labelled image  $\mathbf{x}_n$  and  $\mathcal{Y} = [1, 2, \dots, L]$  denotes the set of classes.

Here we assume that every image  $\mathbf{x}$  annotated by at least one person i.e.,  $|S(\mathbf{x})| \geq 1$ , and no GT labels  $\{\mathbf{y}_n \in \mathcal{Y}^{W \times H}\}_{n=1, \dots, N}$  are available. The problem of interest here is to *learn the unobserved true segmentation distribution  $p(\mathbf{y} | \mathbf{x})$  from such noisy labelled dataset  $\mathcal{D} = \{\mathbf{x}_n, \tilde{\mathbf{y}}_n^{(r)}\}_{n=1, \dots, N}^{r \in S(\mathbf{x}_n)}$*  i.e., the combination of images, noisy annotations and experts' identities for labels (which label was obtained from whom).

We also emphasise that *the goal at inference time is to segment a given unlabelled test image* but not to fuse multiple available labels as is typically done in multi-atlas segmentation approaches [406].

### 7.3.2 Probabilistic Model and Proposed Architecture

Here we describe the probabilistic model of the observed noisy labels from multiple annotators. We make two key assumptions: (1) annotators are statistically independent, (2) annotations over different pixels are independent given the input image. Under these assumptions, the probability of observing noisy labels  $\{\tilde{\mathbf{y}}^{(r)}\}_{r \in S(\mathbf{x})}$  on  $\mathbf{x}$  factorises as:

$$p(\{\tilde{\mathbf{y}}^{(r)}\}_{r \in S(\mathbf{x})} | \mathbf{x}) = \prod_{r \in S(\mathbf{x})} p(\tilde{\mathbf{y}}^{(r)} | \mathbf{x}) = \prod_{r \in S(\mathbf{x})} \prod_{w \in \{1, \dots, W\} \atop h \in \{1, \dots, H\}} p(\tilde{y}_{wh}^{(r)} | \mathbf{x}) \quad (7.1)$$

where  $\tilde{y}_{wh}^{(r)} \in [1, \dots, L]$  denotes the  $(w, h)^{\text{th}}$  elements of  $\tilde{\mathbf{y}}^{(r)} \in \mathcal{Y}^{W \times H}$ . Now we rewrite the probability of observing each noisy label on each pixel  $(w, h)$  as:

$$p(\tilde{y}_{wh}^{(r)} | \mathbf{x}) = \sum_{y_{wh}=1}^L p(\tilde{y}_{wh}^{(r)} | y_{wh}, \mathbf{x}) \cdot p(y_{wh} | \mathbf{x}) \quad (7.2)$$

where  $p(y_{wh} | \mathbf{x})$  denotes the GT label distribution over the  $(w, h)^{\text{th}}$  pixel in the image  $\mathbf{x}$ , and  $p(\tilde{y}_{wh}^{(r)} | y_{wh}, \mathbf{x})$  describes the noisy labelling process by which annotator  $r$  corrupts the true segmentation label. In particular, we refer to the  $L \times L$  matrix whose each  $(i, j)^{\text{th}}$  element is defined by the second term  $\mathbf{a}^{(r)}(\mathbf{x}, w, h)_{ij} := p(\tilde{y}_{wh}^{(r)} = i | y_{wh} = j, \mathbf{x})$  as the CM of annotator  $r$  at pixel  $(w, h)$  in image  $\mathbf{x}$ .

We introduce a CNN-based architecture which models the different constituents in the above joint probability distribution  $p(\{\tilde{\mathbf{y}}^{(r)}\}_{r \in S(\mathbf{x})} | \mathbf{x})$  as illustrated in Fig. 7.1. The model consists of two components: (1) *Segmentation Network*, parametrised by  $\theta$ , which estimates the GT segmentation probability map,  $\hat{\mathbf{p}}_\theta(\mathbf{x}) \in \mathbb{R}^{W \times H \times L}$  whose each  $(w, h, i)^{\text{th}}$  element approximates  $p(y_{wh} = i | \mathbf{x})$ ; (2) *Annotator Network*, parametrised by  $\phi$ , that generate estimates of the pixel-wise CMs of respective annotators as a function of the input image,  $\{\hat{\mathbf{A}}_\phi^{(r)}(\mathbf{x}) \in [0, 1]^{W \times H \times L \times L}\}_{r=1}^R$  whose each  $(w, h, i, j)^{\text{th}}$  element approximates  $p(\tilde{y}_{wh}^{(r)} = i | y_{wh} = j, \mathbf{x})$ . Each product  $\hat{\mathbf{p}}_{\theta, \phi}^{(r)}(\mathbf{x}) := \hat{\mathbf{A}}_\phi^{(r)}(\mathbf{x}) \cdot \hat{\mathbf{p}}_\theta(\mathbf{x})$  represents the estimated segmentation probability map of the corresponding annotator. Note that here “ $\cdot$ ” denotes the element-wise matrix multiplications in the spatial dimensions  $W, H$ . At inference time, we use the output of the segmentation network  $\hat{\mathbf{p}}_\theta(\mathbf{x})$  to segment test images.

We note that each spatial CM,  $\hat{\mathbf{A}}_\phi^{(r)}(\mathbf{x})$  contains  $WHL^2$  variables, and calculating the corresponding annotator's prediction  $\hat{\mathbf{p}}_{\theta, \phi}^{(r)}(\mathbf{x})$  requires  $WH(2L - 1)L$  floating-point operations, potentially incurring a large time/space cost when the number of classes is large.

### 7.3.3 Learning Confusion Matrices and True Segmentation

Next, we describe how we jointly optimise the parameters of segmentation network,  $\theta$  and the parameters of annotator network,  $\phi$ . In short, we minimise the negative log-likelihood of the probabilistic model

plus a regularisation term via stochastic gradient descent. A detailed description is provided below.

Given training input  $\mathbf{X} = \{\mathbf{x}_n\}_{n=1}^N$  and noisy labels  $\tilde{\mathbf{Y}}^{(r)} = \{\tilde{\mathbf{y}}_n^{(r)} : r \in S(\mathbf{x}_n)\}_{n=1}^N$  for  $r = 1, \dots, R$ , we optimise the parameters  $\{\theta, \phi\}$  by minimizing the negative log-likelihood (NLL),  $-\log p(\tilde{\mathbf{Y}}^{(1)}, \dots, \tilde{\mathbf{Y}}^{(R)} | \mathbf{X})$ . From eqs. (7.1) and (7.2), this optimization objective equates to the sum of cross-entropy losses between the observed noisy segmentations and the estimated annotator label distributions:

$$-\log p(\tilde{\mathbf{Y}}^{(1)}, \dots, \tilde{\mathbf{Y}}^{(R)} | \mathbf{X}) = \sum_{n=1}^N \sum_{r=1}^R \mathbb{1}(r \in S(\mathbf{x}_n)) \cdot \text{CE}(\hat{\mathbf{A}}_\phi^{(r)}(\mathbf{x}_n) \cdot \hat{\mathbf{p}}_\theta(\mathbf{x}_n), \tilde{\mathbf{y}}_n^{(r)}) \quad (7.3)$$

Minimizing the above encourages each annotator-specific predictions  $\hat{\mathbf{p}}_{\theta, \phi}^{(r)}(\mathbf{x})$  to be as close as possible to the true noisy label distribution of the annotator  $\mathbf{p}^{(r)}(\mathbf{x})$ . However, this loss function alone is not capable of separating the annotation noise from the true label distribution; there are many combinations of pairs  $\hat{\mathbf{A}}_\phi^{(r)}(\mathbf{x})$  and segmentation model  $\hat{\mathbf{p}}_\theta(\mathbf{x})$  such that  $\hat{\mathbf{p}}_{\theta, \phi}^{(r)}(\mathbf{x})$  perfectly matches the true annotator's distribution  $\mathbf{p}^{(r)}(\mathbf{x})$  for any input  $\mathbf{x}$  (e.g., permutations of rows in the CMs). To combat this problem, inspired by Tanno *et al.*[107], which addressed an analogous issue for the classification task, we add the trace of the estimated CMs to the loss function in Eq. (7.3) as a regularisation term (see Sec 7.3.4). We thus optimize the combined loss:

$$\mathcal{L}_{\text{total}}(\theta, \phi) := \sum_{n=1}^N \sum_{r=1}^R \mathbb{1}(r \in S(\mathbf{x}_n)) \cdot [\text{CE}(\hat{\mathbf{A}}_\phi^{(r)}(\mathbf{x}_n) \cdot \hat{\mathbf{p}}_\theta(\mathbf{x}_n), \tilde{\mathbf{y}}_n^{(r)}) + \lambda \cdot \text{tr}(\hat{\mathbf{A}}_\phi^{(r)}(\mathbf{x}_n))] \quad (7.4)$$

where  $S(\mathbf{x})$  denotes the set of all labels available for image  $\mathbf{x}$ , and  $\text{tr}(\mathbf{A})$  denotes the trace of matrix  $\mathbf{A}$ . The mean trace represents the average probability that a randomly selected annotator provides an accurate label. Intuitively, minimising the trace encourages the estimated annotators to be maximally unreliable while minimising the cross entropy ensures fidelity with observed noisy annotators. We minimise this combined loss via stochastic gradient descent to learn both  $\{\theta, \phi\}$ .

### 7.3.4 Justification for the Trace Norm

Here we provide a further justification for using the trace regularisation. Tanno *et al.*[107] showed that if the average CM of annotators is *diagonally dominant*, and the cross-entropy term in the loss function is zero, minimising the trace of the estimated CMs uniquely recovers the true CMs. However, their results concern properties of the average CMs of both the annotators and the classifier over the data population, rather than individual data samples. We show a similar but slightly weaker result in the sample-specific regime, which is more relevant as we estimate CMs of respective annotators on every input image.

First, let us set up the notations. For brevity, for a given input image  $\mathbf{x} \in \mathbb{R}^{W \times H \times C}$ , we denote the ground-truth CM of annotator  $r$  at  $(i, j)^{\text{th}}$  pixel and its estimate by  $\mathbf{A}^{(r)} := [\mathbf{A}^{(r)}(\mathbf{x})_{ij}]$  and  $\hat{\mathbf{A}}^{(r)} := [\hat{\mathbf{A}}^{(r)}(\mathbf{x})_{ij}] \in [0, 1]^{L \times L}$ , respectively. We also define the mean CM  $\mathbf{A}^* := \sum_{r=1}^R \pi_r \mathbf{A}^{(r)}$  and its estimate  $\hat{\mathbf{A}}^* := \sum_{r=1}^R \pi_r \hat{\mathbf{A}}^{(r)}$  where  $\pi_r \in [0, 1]$  is the probability that the annotator  $r$  labels image  $\mathbf{x}$ . Lastly, as we stated earlier, we assume there is a single GT segmentation label per image — thus the true  $L$ -dimensional probability vector at pixel  $(i, j)$  takes the form of a one-hot vector i.e.,  $\mathbf{p}(\mathbf{x}) = \mathbf{e}_k$  for, say, class  $k \in [1, \dots, L]$ . Then, the followings result motivates the use of the trace regularisation:

**Theorem 7.1.** *If the annotator's segmentation probabilities are perfectly modelled by the model for the given image i.e.,  $\hat{\mathbf{A}}^{(r)} \hat{\mathbf{p}}_\theta(\mathbf{x}) = \mathbf{A}^{(r)} \mathbf{p}(\mathbf{x}) \forall r = 1, \dots, R$ , and the average true confusion matrix  $\mathbf{A}^*$  at a given pixel and its estimate  $\hat{\mathbf{A}}^*$  satisfy that  $a_{kk}^* > a_{kj}^*$  for  $j \neq k$  and  $\hat{a}_{ii}^* > \hat{a}_{ij}^*$  for all  $i, j$  such that  $j \neq i$ , then  $\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(R)} = \text{argmin}_{\hat{\mathbf{A}}^{(1)}, \dots, \hat{\mathbf{A}}^{(R)}} [\text{tr}(\hat{\mathbf{A}}^*)]$  and such solutions are **unique** in the  $k^{\text{th}}$  column where  $k$  is the correct pixel class.*

The corresponding proof is provided in the original paper [108]. The above result shows that if each estimated annotator's distribution  $\hat{\mathbf{A}}^{(r)} \hat{\mathbf{p}}_\theta(\mathbf{x})$  is very close to the true noisy distribution  $\mathbf{p}^{(r)}(\mathbf{x})$  (which is encouraged by minimizing the cross-entropy loss), and for a given pixel, the average CM has the  $k^{\text{th}}$

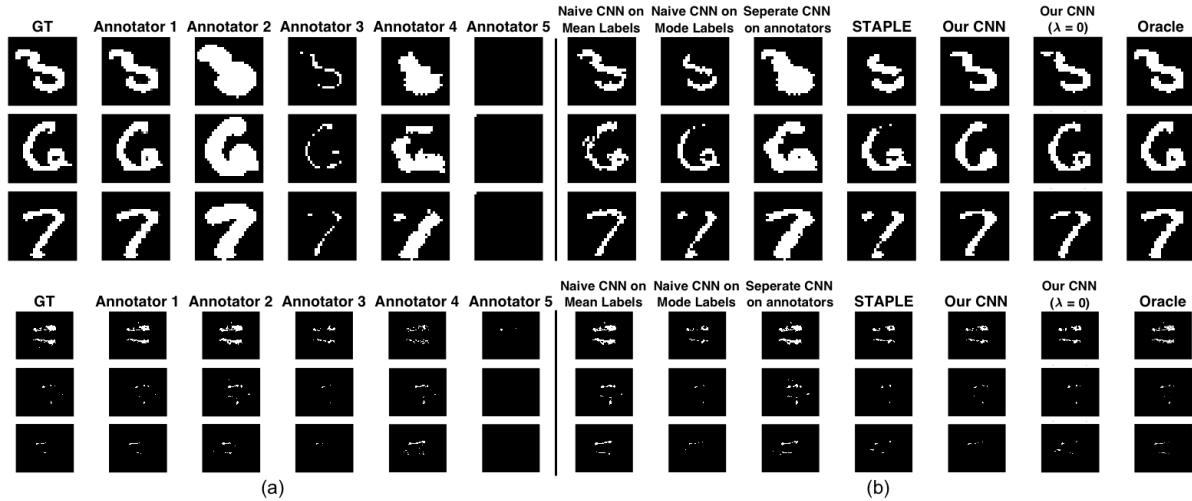


Figure 7.2: Visualisation of segmentation labels on two datasets: (a) ground-truth (GT) and segmentation labels from simulated annotators (Annotators 1 - 5); (b) the predictions from the supervised models.

diagonal entry larger than any other entries in the same row<sup>1</sup>, then minimizing its trace will drive the estimates of the  $k^{\text{th}}$  ('correct class') columns in the respective annotator's CMs to match the true values. Although this result is weaker than what was shown in [107] for the population setting rather than the individual samples, the single-ground-truth assumption means that the remaining values of the CMs are uniformly equal to  $1/L$ , and thus it suffices to recover the column of the correct class.

To encourage  $\{\hat{\mathbf{A}}^{(1)}, \dots, \hat{\mathbf{A}}^{(R)}\}$  to be also diagonally dominant, we initialize them with identity matrices by training the *annotation network* to maximise the trace for sufficient iterations as a warm-up period. Intuitively, the combination of the trace term and cross-entropy separates the true distribution from the annotation noise by finding the maximal amount of confusion which explains the noisy observations well.

## 7.4 Experiments

**Datasets** We evaluate our method on a variety of datasets including both synthetic and real-world scenarios: 1) for MNIST segmentation and ISBI2015 MS lesion segmentation challenge dataset [416], we apply morphological operations to simulate different types of annotation noise in binary segmentation tasks; 2) for the BraTS 2019 brain tumour segmentation dataset [367], we perform a similar simulation but in the multi-class setting; 3) we also consider the LIDC-IDRI dataset which contains real noisy annotations acquired from 4 different clinical experts.

More specifically, we simulate a group of 5 annotators of disparate characteristics by performing morphological transformations (e.g., thinning, thickening, fractures, etc) on the ground-truth (GT) segmentation labels, using the Morpho-MNIST software [410]. Fig. 7.2 shows the examples of the respective annotators; the first annotator provides faithful segmentation labels (“good-segmentation”) that match the GTs, the second tends to over-segment (“over-segmentation”), the third tends to under-segment (“under-segmentation”), the fourth is prone to the combination of small fractures and over-segmentation (“wrong-segmentation”) and the fifth always annotates everything as the background (“blank-segmentation”). To create synthetic noisy labels in multi-class scenario, we first choose a target class, create a binary mask and then apply morphological operations in a similar manner to simulate 4 annotators of different characteristics, namely, over-segmentation, under-segmentation, wrong segmentation (mixing up between the target and another class) and good segmentation (see Fig. 7.6 for examples). We create noisy training data by deriving labels from the simulated annotators. Lastly, we also experiment with varying the levels of morphological operations on MNIST and MS lesion datasets, to test the robustness of our methods

<sup>1</sup>For the standard “majority vote” label to capture the correct true labels, one requires the  $k^{\text{th}}$  diagonal element in the average CM to be larger than the sum of the remaining elements in the same row, which is a more strict condition.

to varying degrees of annotation noise.

**Baselines** Our experiments are based on the assumption that no ground-truth (GT) labels are not known a priori, hence, we compare our method against multiple label fusion methods. In particular, we consider four label fusion baselines: a) mean of all of the noisy labels; b) mode labels by taking the “majority vote”; c) label fusion via the original STAPLE method [60]; d) Spatial STAPLE, a more recent extension of c) that accounts for spatial variations in CMs. After curating the noisy annotations via the above methods, we train the segmentation network and report the results. For c) and d), we used the toolkit<sup>2</sup>. To get an upper-bound performance, we also include the *oracle* model that is directly trained on the ground-truth annotations. To test the value of the proposed image-dependent spatial CMs, we also include “Global CM” model where a single CM is learned per annotator but fixed across pixels and images (analogous to *et al.*[377, 381, 107], but in segmentation task). Lastly, we also compare against a recent method called Probabilistic U-net as another baseline, which has been shown to capture inter-reader variations accurately.

**Metrics** For evaluation metrics, we use: 1) root-MSE between estimated CMs and real CMs; 2) Dice coefficient (DICE) between estimated segmentation and true segmentation; 3) The generalized energy distance proposed in [182] to measure the quality of the estimated annotator’s labels.

**Implementation** Our method as well as the above baselines are implemented in Pytorch. Our network is based on a 4 down-sampling stages 2D U-net [417], the channel numbers for each encoders are 32, 64, 128, 256, we also replaced the batch normalisation layers with instance normalisation. Our segmentation network and annotator network share the same parameters apart from the last layer in the decoder of U-net, essentially, the overall architecture is implemented as an U-net with multiple output last layers: one for prediction of true segmentation; others for predictions of noisy segmentation respectively. For segmentation network, the output of the last layer has  $c$  channels where  $c$  is the number of classes. On the other hand, for annotator network, by default, the output of the last layer has  $L \times L$  number of channels for estimating confusion matrices at each spatial location; when low-rank approximation is used, the output of the last layer has  $2 \times L \times l$  number of channels. The Probabilistic U-net implementation is adopted from <https://github.com/stefanknegt/Probabilistic-Unet-Pytorch>, for fair comparison, we adjusted the number of the channels and the depth of the U-net backbone in Probabilistic U-net to match with our networks. All of the models were trained on a NVIDIA RTX 208 for at least 3 times with different random initialisations to compute the mean performance and its standard deviation (run 3 times of the experiments with the same initialization). The Adam [190] optimiser was used in all experiments with the default hyper-parameter settings. We also provide all of the hyper-parameters of the experiments for each data set in Table 7.1. We also kept the training details the same between the baselines and our method.

Data set	Learning Rate	Epoch	Batch Size	Augmentation	weight for regularisation ( $\lambda$ )
MNIST	1e-4	60	2	Random flip	0.7
MS	1e-4	55	2	Random flip	0.7
BraTS	1e-4	60	8	Random flip	1.5
LIDC	1e-4	75	4	Random flip	0.9

Table 7.1: Hyper-parameters used for respective datasets.

### 7.4.1 MNIST and MS lesion segmentation datasets

MNIST dataset consists of 60,000 training and 10,000 testing examples, all of which are  $28 \times 28$  grayscale images of digits from 0 to 9, and we derive the segmentation labels by thresholding the intensity values at 0.5. The MS dataset is publicly available and comprises 21 3D scans from 5 subjects. All scans are split into 10 for training and 11 for testing. We hold out 20% of training images as a validation set for both datasets. On both datasets, our proposed model achieves a higher dice similarity coefficient

<sup>2</sup><https://www.nitrc.org/projects/masi-fusion/>

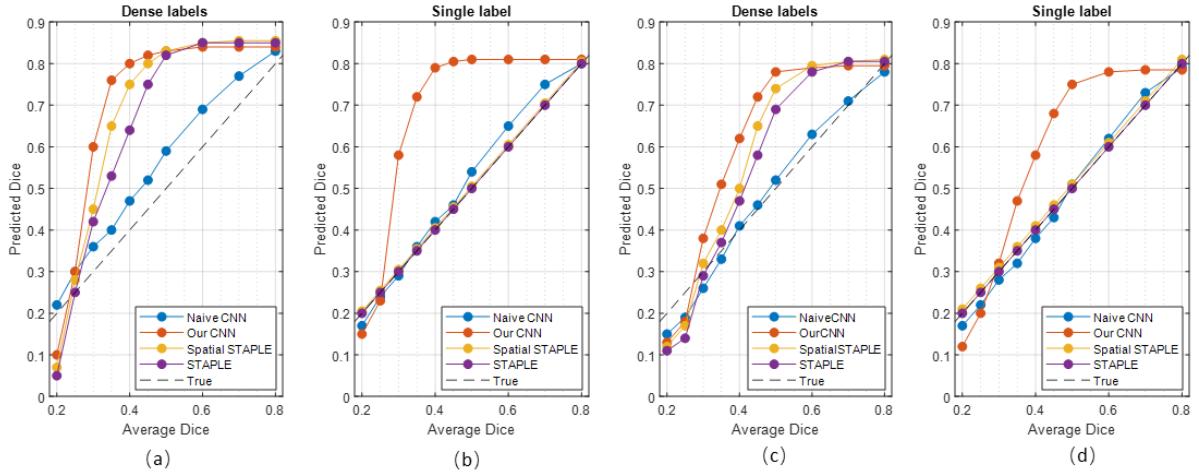


Figure 7.3: Segmentation accuracy of different models on MNIST (a, b) and MS (c, d) dataset for a range of annotation noise (measured in averaged Dice with respect to GT).

than STAPLE on the dense label case and, even more prominently, on the single label (i.e., randomly choose 1 label per image, aka, “one label per image”) case (shown in Tables. 7.2&7.3 and Fig. 7.2). In addition, our model outperforms STAPLE without or with trace norm, in terms of CM estimation, specifically, we could achieve an increase at 6.3%. Additionally, we include the performance on a set of different regularisation coefficients—Fig. 7.4 shows that the presented method is quite robust to this hyper-parameter. Fig. 7.3 compares the segmentation accuracy on MNIST and MS lesion for a range of average dice where labels are generated by a group of 5 simulated annotators. Fig. 7.5 illustrates our model can capture the patterns of mistakes for each annotator. We also notice that our model is consistently more accurate than the global CM model, indicating the value of image-dependent pixel-wise CMs.

Models	MNIST	MNIST	MSLesion	MSLesion
	DICE (%)	CM estimation	DICE (%)	CM estimation
Naive CNN on mean labels	$38.36 \pm 0.41$	n/a	$46.55 \pm 0.53$	n/a
Naive CNN on mode labels	$62.89 \pm 0.63$	n/a	$47.82 \pm 0.76$	n/a
Probabilistic U-net [182]	$65.12 \pm 0.83$	n/a	$46.15 \pm 0.59$	n/a
Separate CNNs on annotators	$70.44 \pm 0.65$	n/a	$46.84 \pm 1.24$	n/a
STAPLE [60]	$78.03 \pm 0.29$	$0.1241 \pm 0.0011$	$55.05 \pm 0.53$	$0.1502 \pm 0.0026$
Spatial STAPLE [405]	$78.96 \pm 0.22$	$0.1195 \pm 0.0013$	$58.37 \pm 0.47$	$0.1483 \pm 0.0031$
Ours with Global CMs	$79.21 \pm 0.41$	$0.1132 \pm 0.0028$	$61.58 \pm 0.59$	$0.1449 \pm 0.0051$
Ours without Trace	$79.63 \pm 0.53$	$0.1125 \pm 0.0037$	$65.77 \pm 0.62$	$0.1342 \pm 0.0053$
Ours	$82.92 \pm 0.19$	$0.0893 \pm 0.0009$	$67.55 \pm 0.31$	$0.0811 \pm 0.0024$
Oracle (Ours but with known CMs)	$83.29 \pm 0.11$	$0.0238 \pm 0.0005$	$78.86 \pm 0.14$	$0.0415 \pm 0.0017$

Table 7.2: Comparison of segmentation accuracy (DICE) and quality of confusion matrix (CM) estimation (MSE) for different methods with dense labels (mean  $\pm$  standard deviation).

Models	MNIST DICE (%)	MNIST CM estimation	MSLesion DICE (%)	MSLesion CM estimation
Naive CNN	$32.79 \pm 1.13$	n/a	$27.41 \pm 1.45$	n/a
STAPLE [60]	$54.07 \pm 0.68$	$0.2617 \pm 0.0064$	$35.74 \pm 0.84$	$0.2833 \pm 0.0081$
Spatial STAPLE [405]	$56.73 \pm 0.53$	$0.2384 \pm 0.0061$	$38.21 \pm 0.71$	$0.2591 \pm 0.0074$
Ours with Global CMs	$59.01 \pm 0.65$	$0.1953 \pm 0.0041$	$40.32 \pm 0.68$	$0.1974 \pm 0.0063$
Ours without Trace	$74.48 \pm 0.37$	$0.1538 \pm 0.0029$	$54.76 \pm 0.66$	$0.1745 \pm 0.0044$
Ours	$76.48 \pm 0.25$	$0.1329 \pm 0.0012$	$56.43 \pm 0.47$	$0.1542 \pm 0.0023$

Table 7.3: Comparison of segmentation accuracy (DICE) and error of CM estimation (MSE) for different methods with one label per image (mean  $\pm$  standard deviation). We note that ‘Naive CNN’ is a baseline trained by simply minimising the cross-entropy between the predictions and the noisy labels.

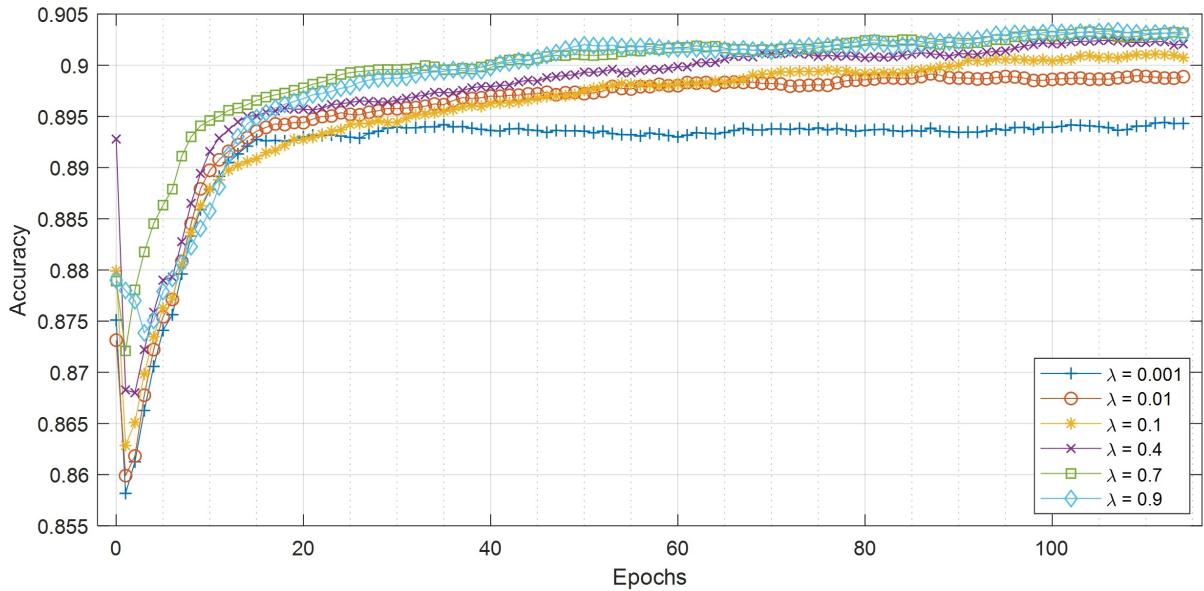


Figure 7.4: Curves of validation accuracy during training of our model on MNIST for a range of hyperparameters. For our method, the scaling of trace regularizer is varied in  $[0.001, 0.01, 0.1, 0.4, 0.7, 0.9]$ .

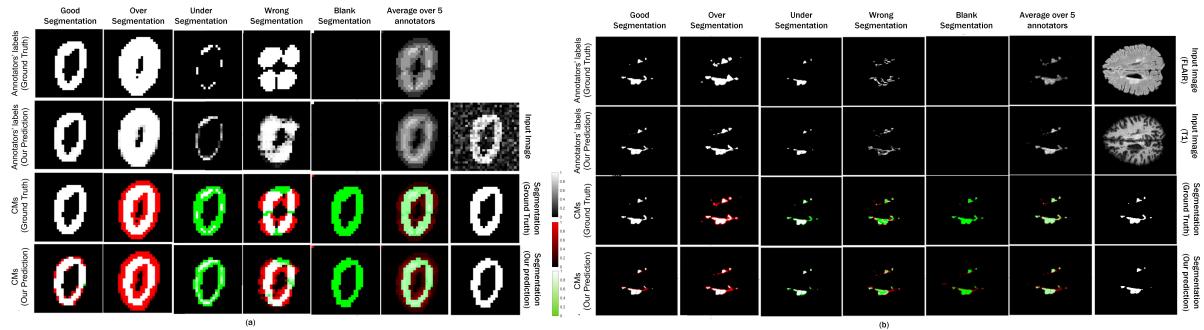


Figure 7.5: Visualisation of the estimated true labels and the estimated pixel-wise confusion matrices on MNIST/MS datasets along with their targets (best viewed in colour). White is the true positive, green is the false negative, red is the false positive and black is the true negative.

### 7.4.2 BraTS Dataset and LIDC-IDRI Dataset

We also evaluate our model on a multi-class segmentation task, using all of the 259 high grade glioma (HGG) cases in training data from 2019 multi-modal Brain Tumour Segmentation Challenge (BraTS). We extract each slice as 2D images and split them at case-wise to have, 1600 images for training, 300 for validation and 500 for testing. Pre-processing includes: concatenation of all of available modalities; centre cropping to 192 x 192; normalisation for each case at each modality.

Models	BraTS DICE (%)	BraTS CM estimation	LIDC-IDRI DICE (%)	LIDC-IDRI CM estimation
Naive CNN on mean labels	29.42 ± 0.58	n/a	56.72 ± 0.61	n/a
Naive CNN on mode labels	34.12 ± 0.45	n/a	58.64 ± 0.47	n/a
Probabilistic U-net [182]	40.53 ± 0.75	n/a	61.26 ± 0.69	n/a
STAPLE [60]	46.73 ± 0.17	0.2147 ± 0.0103	69.34 ± 0.58	0.0832 ± 0.0043
Spatial STAPLE [405]	47.31 ± 0.21	0.1871 ± 0.0094	70.92 ± 0.18	0.0746 ± 0.0057
Ours with Global CMs	47.33 ± 0.28	0.1673 ± 0.1021	70.94 ± 0.19	0.1386 ± 0.0052
Ours without Trace	49.03 ± 0.34	0.1569 ± 0.0072	71.25 ± 0.12	0.0482 ± 0.0038
Ours	<b>53.47 ± 0.24</b>	<b>0.1185 ± 0.0056</b>	<b>74.12 ± 0.19</b>	<b>0.0451 ± 0.0025</b>
Oracle (Ours but with known CMs)	67.13 ± 0.14	0.0843 ± 0.0029	79.41 ± 0.17	0.0381 ± 0.0021

Table 7.4: Comparison of segmentation accuracy and error of CM estimation for different methods trained with **dense labels** (mean ± standard deviation). The best results are shown in bold. Note that we count out the Oracle from the model ranking as it forms a theoretical upper-bound on the performance where true labels are known on the training data.

Models	BraTS DICE (%)	BraTS CM estimation	LIDC-IDRI DICE (%)	LIDC-IDRI CM estimation
Naive CNN on mean & mode labels	36.12 ± 0.93	n/a	48.36 ± 0.79	n/a
STAPLE [60]	38.74 ± 0.85	0.2956 ± 0.1047	57.32 ± 0.87	0.1715 ± 0.0134
Spatial STAPLE [405]	41.59 ± 0.74	0.2543 ± 0.0867	62.35 ± 0.64	0.1419 ± 0.0207
Ours with Global CMs	41.76 ± 0.71	0.2419 ± 0.0829	63.25 ± 0.66	0.1382 ± 0.0175
Ours without Trace	43.74 ± 0.49	0.1825 ± 0.0724	66.95 ± 0.51	0.0921 ± 0.0167
Ours	<b>46.21 ± 0.28</b>	<b>0.1576 ± 0.0487</b>	<b>68.12 ± 0.48</b>	<b>0.0587 ± 0.0098</b>

Table 7.5: Comparison of segmentation accuracy and error of CM estimation for different methods trained with only **one label available per image** (mean ± standard deviation). The best results are shown in bold.

Lastly, we use the LIDC-IDRI dataset to evaluate the method in the scenario where multiple labels are actually acquired from different clinical experts. The dataset contains 1018 lung CT scans from 1010 lung patients with manual lesion segmentations from four experts. For each scan, 4 radiologists provided annotation masks for lesions that they independently detected and considered to be abnormal. For our experiments, we used the same method in [182] to pre-process all scans. We split the dataset at case-wise into a training (722 patients), validation (144 patients) and testing (144 patients). We then resampled the CT scans to  $1mm \times 1mm$  in-plane resolution. We also centre cropped 2D images ( $180 \times 180$  pixels) around lesion positions, in order to focus on the annotated lesions. The lesion positions are those where at least one of the experts segmented a lesion. We hold 5000 images in the training set, 1000 images in the validation set and 1000 images in the test set. Since the dataset does not provide a single curated ground-truth for each image, we created a “gold standard” by aggregating the labels via spatial STAPLE [405], a recent variant of the STAPLE framework employed in the creation of public medical image segmentation datasets e.g., ISLES [402], MSSeg [403], Gleason’19 [404] datasets. We further note that, as before, we assume labels are only available to the model during training, but not at test time, thus label aggregation methods cannot be applied on the test examples.

On both BraTS and LIDC-IDRI datasets, our proposed model achieves a higher dice similarity coefficient than STAPLE and Spatial STAPLE on both of the dense labels and single label scenarios (shown in Table. 7.4 and Table. 7.5). In addition, our model (with trace) outperforms STAPLE in terms of CM estimation by a large margin at 14.4% on BraTS. In Fig. 7.6, we visualize the segmentation results on BraTS and the corresponding annotators’ predictions. Fig. 7.8 presents three examples of the segmentation results and the corresponding four annotator contours, as well as the consensus. As shown in both figures, our model successfully predicts both the segmentation of lesions and the variations of each annotator in different cases. We also measure the inter-reader consensus levels by computing the IoU of

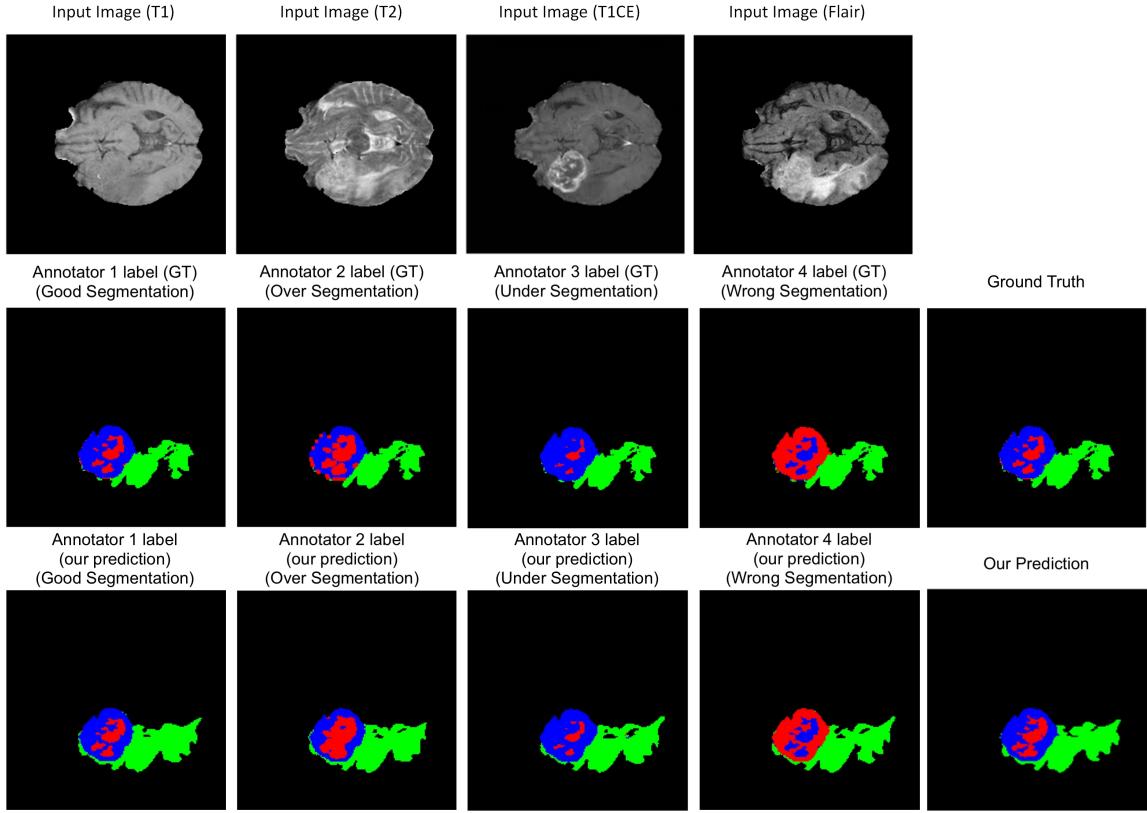


Figure 7.6: Visualisation of the estimated true segmentation on the BraTS dataset and the estimated annotations of the respective annotators (best viewed in colour). The tumour core (red) is the target class on which annotation mistakes are simulated.

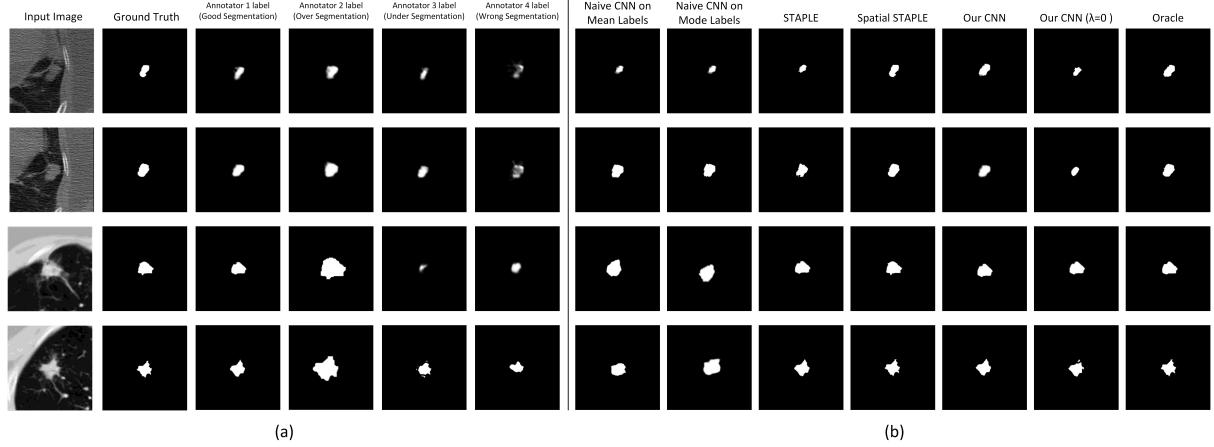


Figure 7.7: Visualisation of segmentation labels on LIDC-IDRI dataset: (a) “Ground Truth” (GT) and simulated annotator’s segmentations (Annotator 1 - 5); (b) the predictions from the supervised models. We note that the GT labels are created by aggregating multiple labels via spatial STAPLE [405].)

multiple annotations, and compare the segmentation performance in 3 subgroups of different consensus levels (low, medium and high). Fig. 7.9 illustrates that our method attains consistent improvement over the baselines in all cases, indicating its ability to segment more robustly even the hard examples where the experts in reality have disagreed to a large extent.

Additionally, as shown in Table 7.6, our model consistently outperforms Probabilistic U-Net on generalized energy distance across the four test different datasets, indicating our method can better capture the inter-annotator variations than the baseline Probabilistic U-Net. This result shows that the information about which labels are acquired from whom is useful in modelling the variability in the observed segmentation labels.

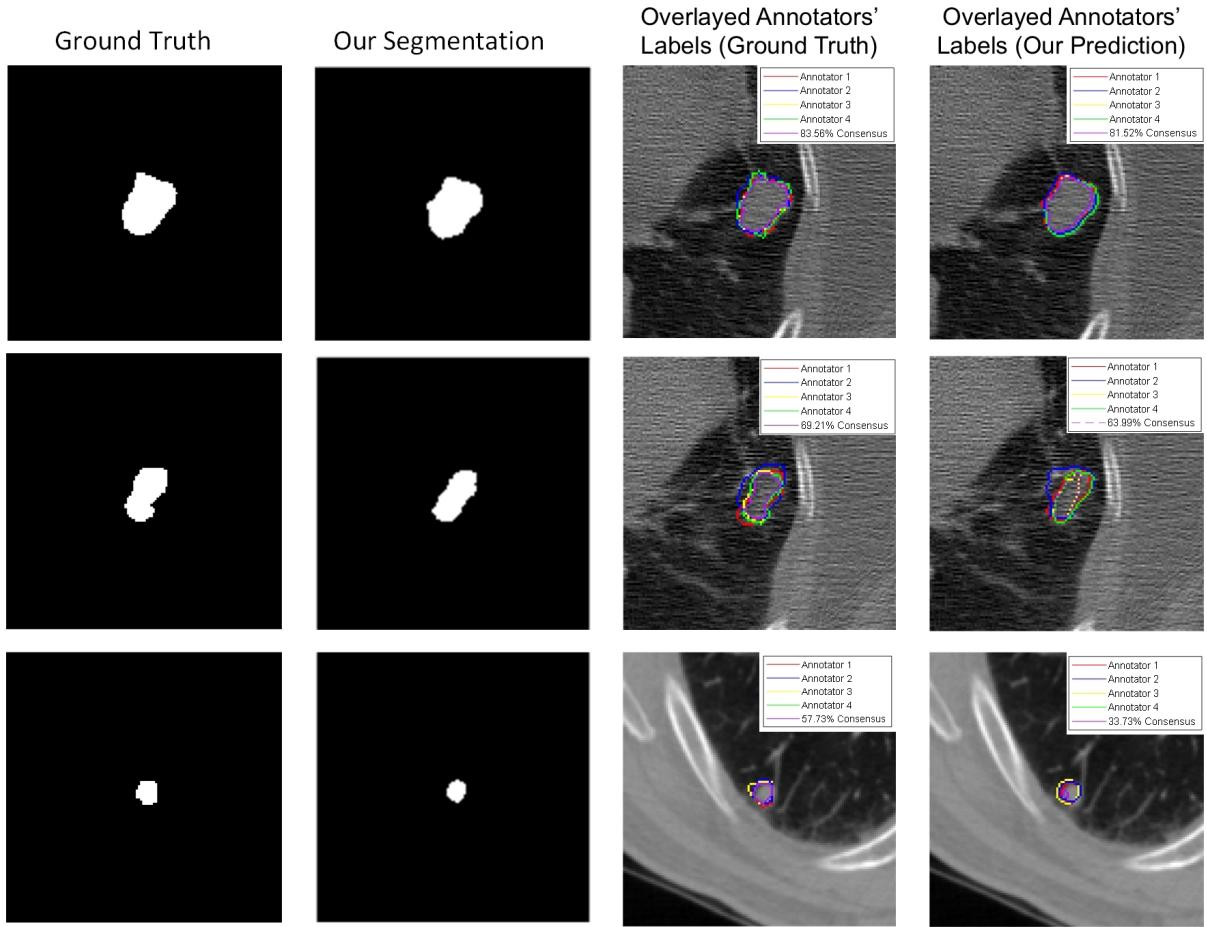


Figure 7.8: Segmentation results on LIDC-IDRI dataset and the visualization of each annotator contours and the consensus. The bottom row shows an interesting example in which annotator 4 (green) misses the abnormality completely, which is also predicted by our model.

Models	MNIST	MS	BraTS	LIDC-IDRI
Probabilistic U-net [182]	$1.46 \pm 0.04$	$1.91 \pm 0.03$	$3.23 \pm 0.07$	$1.97 \pm 0.03$
Ours	<b><math>1.24 \pm 0.02</math></b>	<b><math>1.67 \pm 0.03</math></b>	<b><math>3.14 \pm 0.05</math></b>	<b><math>1.87 \pm 0.04</math></b>

Table 7.6: Comparison of Generalised Energy Distance (GED) on different datasets (mean  $\pm$  standard deviation). The distance metric used here is the DICE score.

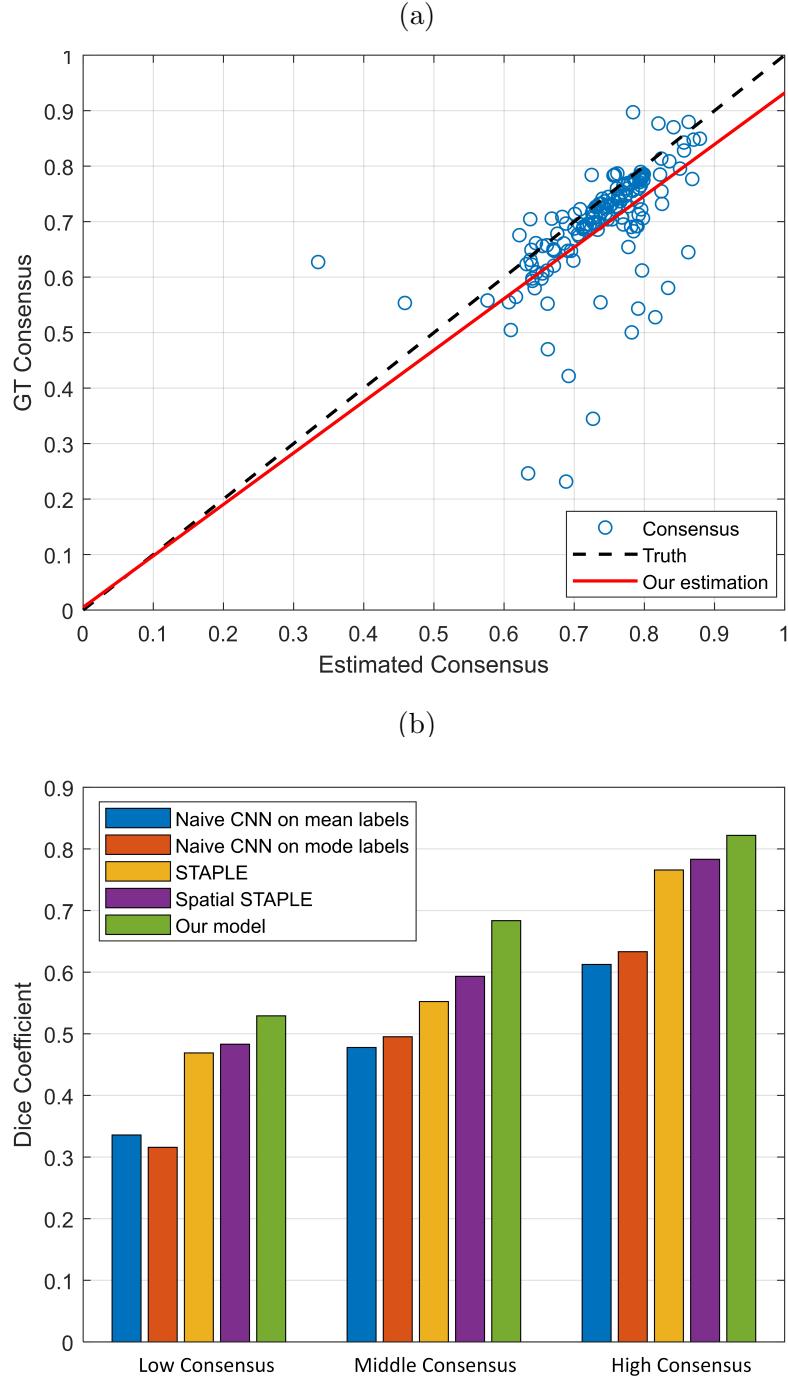


Figure 7.9: (a) The consensus level amongst the estimated annotators is plotted against the ground truth on LIDC-IDRI dataset. Here the inter-reader consensus level is measured as the IoU of available annotations. The strong positive linear correlation shows that the variation in the inter-reader variability on different input examples (e.g., some examples are more ambiguous than others) is captured well. We do note, however, that the inter-reader variation seems more under-estimated for “easy” (i.e., higher consensus) samples. (b) Segmentation performance on 3 different subgroups of the LIDC-IDRI dataset with varying levels of inter-reader agreement: 1) Low-consensus (30% to 65% IoU); 2) Middle consensus (65% to 75% IoU); 3) High consensus (75% to 90% IoU). Our method shows *consistent* improvement over the baselines and the competing methods in all groups, showing its enhanced ability to segment challenging examples (i.e., low-consensus cases).

## 7.5 Discussion and Conclusion

We introduced the first supervised segmentation method for jointly estimating the spatial characteristics of labelling errors from multiple human annotators and the ground-truth label distribution. We demonstrated this method on real-world datasets with both synthetic and real-world annotations. Our method is capable of estimating individual annotators and thereby improving robustness against label noise. Experiments have shown our model achieves considerable improvement over the traditional label fusion approaches including averaging, the majority vote and the widely used STAPLE framework and its recent extensions, in terms of both segmentation accuracy and the quality of confusion matrix (CM) estimation.

In the future, we plan to accommodate meta-information of annotators (e.g., number of years of experience), and non-image data (e.g., genetics) that may influence the pattern of the underlying segmentation label such as lesion appearance, in our framework. We are also interested in assessing the utility of our approach in downstream applications. Of particular interest is the design of active data collection schemes where the segmentation model is used to select which samples to annotate (“active learning”), and the annotator models are used to decide which experts to label them (“active labelling”) [378]. Another exciting avenue of applications is education of inexperienced annotators; the estimated spatial characteristics of segmentation mistakes provide further insights into their annotation behaviours, which they may benefit from in improving their annotation quality.

# Chapter 8

## Conclusions and Future Work

This thesis has explored the utility of modelling and reasoning with different types of uncertainty in deep learning models in high-dimensional, challenging medical imaging applications where safety is critical. The paucity of methodological literature upon embarking on this journey in 2015 motivated the development of new algorithms presented in the thesis, many of which are attained by translating ideas from the established paradigm of probabilistic machine learning to the developing world of deep learning. While most of the research questions I have attempted to answer derive from specific applications, the resultant solutions are general and their applicability extends to other problems in medical imaging and beyond.

**Predictive Uncertainty and its Constituents** Chapter 2 has investigated the importance of quantifying predictive uncertainty and its constituents (namely, aleatoric and parameter uncertainty) in the context of MRI super-resolution application. In particular we have demonstrated that rather simple attempts at modelling uncertainty (with the mean-field Gaussian likelihood and variational approximation of the weight posterior) confer tangible benefits such as 1) performance improvement: e.g., the generalisation to out-of-distribution data and robustness to input noise and outliers; 2) reliability assessment of prediction: e.g., certification of performance based on uncertainty-thresholding and detection of unfamiliar structures. Moreover, assuming sufficient flexibility of the model structure, we also introduce a way to decompose the predictive uncertainty into its orthogonal sources i.e. aleatoric and parameter uncertainty. Preliminary results indicate the potential utility of such decoupling in providing quantitative “explanations” into the model’s under-performance (e.g., if one observes high parameter and low aleatoric uncertainty around a certain image feature, then acquiring more training data of similar cases may be recommended). However, much work still is needed to design systematic means to elicit actionable insights from such uncertainty decomposition.

For example, a recent work from Antoran *et al.*[418] has explored a way to use a generative model to explain the obtained uncertainty estimates. In a nutshell, the method generates a “counterfactual” explanation by synthesizing a version of the given input image with a minimal difference that causes the target uncertainty estimate to dwindle considerably. Scaling up this method to a much more high-dimensional medical imaging datasets would be likely enabled by advances in the generative models of such datasets such as [419]. More broadly, despite the numerous attempts made to explain the predictions of deep learning models [362], little attention has been paid to explaining uncertainty estimates. This is rather strange given that the cases in which which model is uncertain may benefit more from further explanations than the cases where the correct predictions are obvious (assuming that the model is well-calibrated). Thinking about whether or not adapting the existing “interpretability” methods to the analysis of uncertainty metrics poses any new challenges would be worthwhile in my view.

Chapter 3 has shown that the same methods of uncertainty modelling could be naturally extended to the multi-task learning setting. Such adaptation leads to a mechanism which automatically determines, in a spatially adaptive fashion, the relative weighting between the task losses, which is a key determinant in the efficacy of multi-task learning. We demonstrate similar benefits in performance improvement and reliability assessment in another challenging multi-task image translation problem where the synthetic CT image and the segmentation of organs at risk are simultaneously predicted from the input MR image for the downstream radio-therapy planning.

**Structural Uncertainty** Motivated by the observation that multi-task learning with deep learning hinges on the design of feature sharing between tasks within the architecture, Chapter 4 develops a Bayesian approach to learning the connectivity structures in a given neural network. More specifically, the method circumvents the need of handcrafting an architecture by introducing a mechanism that learns probabilities to assign convolution kernels in each layer to “specialist” or “generalist” groups, which are specific to or shared across different tasks. Variational inference is employed to learn the posterior distribution over the possible grouping of kernels and network parameters, leading to the state-of-the-art performance in the radio-therapy planning application that is also addressed in Chapter 3. In all tested benchmarks, learning the grouping of kernels seems to enhance the overall performance, corroborating an important connection between the sparse structure and the representation quality in multi-task learning.

In contrast, Chapter 5 moves away from the probabilistic treatment of architecture learning, and instead draws inspirations from deterministic algorithms used to grow the structure of decision trees. The proposed method, Adaptive Neural Trees (ANTs) progressively grows the structure of a neural network architecture from simple building blocks, and adapt to the given availability of data and the complexity of the task. Results on classification benchmarks show the growth procedure construct architectures of adequate size, leading to better generalisation, particularly in cases with limiting training data where overfitting is prominent with complex models. In a sense, the progressive architecture growth embodies the principles of Occam’s razor [420] which states preferences over simpler models when equal fidelity with data is observed. Another noteworthy aspect is the implications of the ANT’s tree topology as a structural prior. Due to this strong structural constraint, a trained ANT can be construed as a hierarchical form of deep ensemble models [178] where the routing functions determine which transformations should be shared and specific to a given input data. Our results empirically demonstrate such structured sparsity achieves a better trade-off between accuracy and computation than standard densely distributed deep learning models, which indicates room for designing better models that may be achieved by inducing appropriate sparse structures in over-parametrised models. Lastly, I would like to note that such insight is also consistent with the recent observation such as the lottery ticket hypothesis [421] and other works on model pruning [422, 423] that have shown one can prune away 70-80% of the connections in a neural network without adversely affecting performance.

**Human Uncertainty** In many medical imaging applications, acquiring reliable annotations is often challenging due to the high cost and rarity of highly experienced clinical experts. As a result, many medical imaging datasets in practice are contaminated with systematic annotation noise (a form of “measurement error”) which are complex functions of biases (e.g., personality, hospital policies) and competence levels (e.g., amount of experiences) of individual human experts. However, annotation noise is commonly seen as a form of (irreducible) aleatoric uncertainty, and unwillingly accepted as a part of the training data without any additional measure.

Contrary to such popular stance, Chapter 6 develops a method that explicitly models the annotation noise (e.g., biases and skill levels) of human experts, and thereby better infer the (possibly unobserved) true label distribution, from noisy labels alone, by “inverting” such noise models. A well-grounded and practical optimisation method is designed to disentangle the annotation noise and the true labels, leading to a considerable improvement in accuracy in a real-world classification task where the annotations are very noisy and sparse. Chapter 7 extends this idea to the more challenging task of semantic segmentation where every pixel in the input image is classified. The modified noise model combined with an effective optimisation algorithm improve substantially the robustness of the model to annotation noise in three established benchmarks in comparison with the blind approach based on the majority vote and more established label curation methods. These empirical evidence in both classification and segmentation problems suggest that if one has a good understanding of the noise contaminating in the data, incorporating them explicitly into the model helps and even allow ones to learn a model that is more accurate than the available data.

Future work remains, however, to extend this research both in theory and applications. A key simplifying assumption we have made is that there is a single, unknown, true segmentation map of the underlying anatomy, and each individual annotator produces a noisy approximation to it with variations that reflect their individual characteristics. This is in stark contrast with many recent advances (e.g., Probabilistic U-net [182] and PHiSeg [414]) that assume variable annotations from experts are all realistic instances of the true segmentation. One could argue that single-truth assumption may be sensible in the context of segmentation problems since there exists only one true boundary of the physical objects captured in an image while multiple hypothesis can arise from ambiguities in human interpretations. However, we believe that the reality lies somewhere between i.e., some variations are indeed intrinsic

while some are specific to human imperfections. Separation of the two could be potentially addressed by using some prior knowledge about the individual annotators (e.g., meta-information such as the years of experiences, etc) [377] or using a small portion of dataset with curated annotations as a reference set which can be assumed to come from the true label distribution.

Another exciting avenue of research is the application of the annotation models in downstream tasks. Of particular interest is the design of active data collection schemes where the segmentation model is used to select which samples to annotate (“active learning”), and the annotator models are used to decide which experts to label them (“active labelling”—e.g., extending Yan *et al.*[424] from simple classification task to segmentation remains future work. Another exciting application is education of inexperienced annotators; the estimated spatial characteristics of segmentation mistakes provide further insights into their annotation behaviours, and as a result, potential help them improve the quality of their annotations in the next data acquisition.



# Bibliography

- [1] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521(7553):436, 2015.
- [2] Geert Litjens, Thijs Kooi, Babak Ehteshami Bejnordi, Arnaud Arindra Adiyoso Setio, Francesco Ciompi, Mohsen Ghafoorian, Jeroen AWM van der Laak, Bram van Ginneken, and Clara I Sánchez. A survey on deep learning in medical image analysis. *Medical image analysis*, 42:60–88, 2017.
- [3] Dinggang Shen, Guorong Wu, and Heung-Il Suk. Deep learning in medical image analysis. *Annual review of biomedical engineering*, 19:221–248, 2017.
- [4] Andre Esteva, Brett Kuprel, Roberto A Novoa, Justin Ko, Susan M Swetter, Helen M Blau, and Sebastian Thrun. Dermatologist-level classification of skin cancer with deep neural networks. *Nature*, 542(7639):115, 2017.
- [5] Varun Gulshan, Lily Peng, Marc Coram, Martin C Stumpe, Derek Wu, Arunachalam Narayanaswamy, Subhashini Venugopalan, Kasumi Widner, Tom Madams, Jorge Cuadros, et al. Development and validation of a deep learning algorithm for detection of diabetic retinopathy in retinal fundus photographs. *Jama*, 316(22):2402–2410, 2016.
- [6] Pranav Rajpurkar, Jeremy Irvin, Kaylie Zhu, Brandon Yang, Hershel Mehta, Tony Duan, Daisy Ding, Aarti Bagul, Curtis Langlotz, Katie Shpanskaya, et al. Chexnet: Radiologist-level pneumonia detection on chest x-rays with deep learning. *arXiv preprint arXiv:1711.05225*, 2017.
- [7] Nan Wu, Jason Phang, Jungkyu Park, Yiqiu Shen, Zhe Huang, Masha Zorin, Stanislaw Jastrzebski, Thibault Fevry, Joe Katsnelson, Eric Kim, et al. Deep neural networks improve radiologists’ performance in breast cancer screening. 2019.
- [8] Zoubin Ghahramani. Probabilistic machine learning and artificial intelligence. *Nature*, 521(7553):452, 2015.
- [9] Stephen C Hora. Aleatory and epistemic uncertainty in probability elicitation with an example from hazardous waste management. *Reliability Engineering & System Safety*, 54(2-3):217–223, 1996.
- [10] Armen Der Kiureghian and Ove Ditlevsen. Aleatory or epistemic? does it matter? *Structural Safety*, 31(2):105–112, 2009.
- [11] David Roxbee Cox. *Principles of statistical inference*. Cambridge university press, 2006.
- [12] Francisco J Samaniego. *A comparison of the Bayesian and frequentist approaches to estimation*. Springer Science & Business Media, 2010.
- [13] Timothy John Sullivan. *Introduction to uncertainty quantification*, volume 63. Springer, 2015.
- [14] Robert S Ledley and Lee B Lusted. Reasoning foundations of medical diagnosis. *Science*, 130(3366):9–21, 1959.
- [15] Edward H Shortliffe, Bruce G Buchanan, and Edward A Feigenbaum. Knowledge engineering for medical decision making: A review of computer-based clinical decision aids. *Proceedings of the IEEE*, 67(9):1207–1224, 1979.
- [16] Casimir A Kulikowski. Artificial intelligence methods and systems for medical consultation. *IEEE Transactions on pattern analysis and Machine Intelligence*, (5):464–476, 1980.
- [17] Richard O Duda and Edward H Shortliffe. Expert systems research. *Science*, 220(4594):261–268, 1983.

- [18] Klaus-Peter Adlassnig, Gernot Kolarz, Werner Scheithauer, Harald Effenberger, and Georg Grabner. Cadiag: Approaches to computer-assisted medical diagnosis. *Computers in biology and medicine*, 15(5):315–335, 1985.
- [19] Klaus-Peter Adlassnig. Fuzzy set theory in medical diagnosis. *IEEE Transactions on Systems, Man, and Cybernetics*, 16(2):260–265, 1986.
- [20] Eta S Berner and Mark L Gruber. Overconfidence as a cause of diagnostic error in medicine. *The American journal of medicine*, 121(5):S2–S23, 2008.
- [21] Friedrich Steimann, Klaus P Adlassnig, et al. Fuzzy medical diagnosis. *Handbook of fuzzy computation*, page G13, 1998.
- [22] Robert Ivor John and Peter R Innocent. Modeling uncertainty in clinical diagnosis using fuzzy logic. *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)*, 35(6):1340–1350, 2005.
- [23] Ewa Straszecka. Combining uncertainty and imprecision in models of medical diagnosis. *Information Sciences*, 176(20):3026–3059, 2006.
- [24] PK Anooj. Clinical decision support system: Risk level prediction of heart disease using weighted fuzzy rules. *Journal of King Saud University-Computer and Information Sciences*, 24(1):27–40, 2012.
- [25] Markos G Tsipouras, Themis P Exarchos, Dimitrios I Fotiadis, Anna P Kotsia, Konstantinos V Vakalis, Katerina K Naka, and Lampros K Michalis. Automated diagnosis of coronary artery disease based on data mining and fuzzy modeling. *IEEE Transactions on Information Technology in Biomedicine*, 12(4):447–458, 2008.
- [26] Judea Pearl. *Probabilistic reasoning in intelligent systems: networks of plausible inference*. Elsevier, 1984.
- [27] Daniel Nikovski. Constructing bayesian networks for medical diagnosis from incomplete and partially correct statistics. *IEEE Transactions on Knowledge & Data Engineering*, (4):509–516, 2000.
- [28] Francisco Javier Diez, José Mira, E Iturralde, and S Zubillaga. Diaval, a bayesian expert system for echocardiography. *Artificial Intelligence in Medicine*, 10(1):59–73, 1997.
- [29] Subramani Mani, Suzanne McDermott, and Marco Valtorta. Mentor: a bayesian model for prediction of mental retardation in newborns. *Research in Developmental Disabilities*, 18(5):303–318, 1997.
- [30] Charles E Kahn Jr, Linda M Roberts, Katherine A Shaffer, and Peter Haddawy. Construction of a bayesian network for mammographic diagnosis of breast cancer. *Computers in biology and medicine*, 27(1):19–29, 1997.
- [31] Paul Sajda, Clay Spence, and Lucas Parra. A multi-scale probabilistic network model for detection, synthesis and compression in mammographic image analysis. *Medical image analysis*, 7(2):187–204, 2003.
- [32] Elizabeth S Burnside, Daniel L Rubin, Ross D Shachter, Rita E Sohlich, and Edward A Sickles. A probabilistic expert system that provides automated mammographic–histologic correlation: initial experience. *American Journal of Roentgenology*, 182(2):481–488, 2004.
- [33] Flávio Luiz Seixas, Bianca Zadrozny, Jerson Laks, Aura Conci, and Débora Christina Muchaluat Saade. A bayesian network decision model for supporting the diagnosis of dementia, alzheimer? s disease and mild cognitive impairment. *Computers in biology and medicine*, 51:140–158, 2014.
- [34] Kunio Doi, ML Giger, RM Nishikawa, KR Hoffmann, H MacMahon, RA Schmidt, and K-G Chua. Digital radiography: A useful clinical tool for computer-aided diagnosis by quantitative analysis of radiographic images. *Acta Radiologica*, 34(5):426–439, 1993.
- [35] Paul Sajda. Machine learning for detection and diagnosis of disease. *Annu. Rev. Biomed. Eng.*, 8:537–565, 2006.
- [36] JB Antoine Maintz and Max A Viergever. A survey of medical image registration. *Medical image analysis*, 2(1):1–36, 1998.

- [37] Ben Glocker, Aristeidis Sotiras, Nikos Komodakis, and Nikos Paragios. Deformable medical image registration: setting the state of the art with discrete methods. *Annual review of biomedical engineering*, 13:219–244, 2011.
- [38] Aristeidis Sotiras, Christos Davatzikos, and Nikos Paragios. Deformable medical image registration: A survey. *IEEE transactions on medical imaging*, 32(7):1153, 2013.
- [39] Koen Van Leemput. Encoding probabilistic brain atlases using bayesian inference. *IEEE Transactions on Medical Imaging*, 28(6):822–837, 2008.
- [40] Petter Risholm, Steve Pieper, Egil Samset, and William M Wells. Summarizing and visualizing uncertainty in non-rigid registration. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 554–561. Springer, 2010.
- [41] Dana Cobzas and Abhishek Sen. Random walks for deformable image registration. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 557–565. Springer, 2011.
- [42] Tayebeh Lotfi, Lisa Tang, Shawn Andrews, and Ghassan Hamarneh. Improving probabilistic image registration via reinforcement learning and uncertainty evaluation. In *International Workshop on Machine Learning in Medical Imaging*, pages 187–194. Springer, 2013.
- [43] Ivor JA Simpson, Julia A Schnabel, Adrian R Groves, Jesper LR Andersson, and Mark W Woolrich. Probabilistic inference of regularisation in non-rigid registration. *NeuroImage*, 59(3):2438–2451, 2012.
- [44] Petter Risholm, Firdaus Janoos, Isaiah Norton, Alex J Golby, and William M Wells III. Bayesian characterization of uncertainty in intra-subject non-rigid registration. *Medical image analysis*, 17(5):538–555, 2013.
- [45] Karteek Popuri, Dana Cobzas, and Martin Jägersand. A variational formulation for discrete registration. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 187–194. Springer, 2013.
- [46] Miaomiao Zhang, Nikhil Singh, and P Thomas Fletcher. Bayesian estimation of regularization and atlas building in diffeomorphic image registration. In *International conference on information processing in medical imaging*, pages 37–48. Springer, 2013.
- [47] Demian Wassermann, Matthew Toews, Marc Niethammer, and William Wells. Probabilistic diffeomorphic registration: Representing uncertainty. In *International Workshop on Biomedical Image Registration*, pages 72–82. Springer, 2014.
- [48] Ivor JA Simpson, Manuel Jorge Cardoso, Marc Modat, David M Cash, Mark W Woolrich, Jesper LR Andersson, Julia A Schnabel, Sébastien Ourselin, Alzheimer’s Disease Neuroimaging Initiative, et al. Probabilistic non-linear registration with spatially adaptive regularisation. *Medical image analysis*, 26(1):203–216, 2015.
- [49] Mattias P Heinrich, Ivor JA Simpson, Bartłomiej W Papież, Michael Brady, and Julia A Schnabel. Deformable image registration by combining uncertainty estimates from supervoxel belief propagation. *Medical image analysis*, 27:57–71, 2016.
- [50] Loic Le Folgoc, Hervé Delingette, Antonio Criminisi, and Nicholas Ayache. Quantifying registration uncertainty with sparse bayesian modelling. *IEEE transactions on medical imaging*, 36(2):607–617, 2016.
- [51] Petter Risholm, James Balter, and William M Wells. Estimation of delivered dose in radiotherapy: the influence of registration uncertainty. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 548–555. Springer, 2011.
- [52] Ivor JA Simpson, Mark W Woolrich, and Julia A Schnabel. Probabilistic segmentation propagation from uncertainty in registration. In *Medical Image Analysis and Understanding (MIUA)*, 2011.
- [53] Ivor JA Simpson, Mark W Woolrich, Jesper L R Andersson, Adrian R Groves, and Julia A Schnabel. Ensemble learning incorporating uncertain registration. *IEEE transactions on medical imaging*, 32(4):748–756, 2012.

- [54] Jie Luo, Sarah F. Frisken, Karteek Popuri, Dana Cobzas, Frank Preiswerk, Matthew Toews, Miao-miao Zhang, Hongyi Ding, Polina Golland, Alexandra J. Golby, Masashi Sugiyama, and William M. Wells III. On the ambiguity of registration uncertainty. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2018.
- [55] David T Gering, Arya Nabavi, Ron Kikinis, Noby Hata, Lauren J O'Donnell, W Eric L Grimson, Ferenc A Jolesz, Peter M Black, and William M Wells III. An integrated visualization system for surgical planning and guidance using image fusion and an open mr. *Journal of Magnetic Resonance Imaging: An Official Journal of the International Society for Magnetic Resonance in Medicine*, 13(6):967–975, 2001.
- [56] Gloria P Mazzara, Robert P Velthuizen, James L Pearlman, Harvey M Greenberg, and Henry Wagner. Brain tumor target volume determination for radiation treatment planning through automated mri segmentation. *International Journal of Radiation Oncology\* Biology\* Physics*, 59(1):300–312, 2004.
- [57] Bruce Fischl, David H Salat, Evelina Busa, Marilyn Albert, Megan Dieterich, Christian Haselgrove, Andre Van Der Kouwe, Ron Killiany, David Kennedy, Shuna Klaveness, et al. Whole brain segmentation: automated labeling of neuroanatomical structures in the human brain. *Neuron*, 33(3):341–355, 2002.
- [58] Marcel Prastawa, John H Gilmore, Weili Lin, and Guido Gerig. Automatic segmentation of mr images of the developing newborn brain. *Medical image analysis*, 9(5):457–466, 2005.
- [59] Alex P Zijdenbos, Reza Forghani, Alan C Evans, et al. Automatic” pipeline” analysis of 3-d mri data for clinical trials: application to multiple sclerosis. *IEEE Trans. Med. Imaging*, 21(10):1280–1291, 2002.
- [60] Simon K Warfield, Kelly H Zou, and William M Wells. Simultaneous truth and performance level estimation (staple): an algorithm for the validation of image segmentation. *IEEE transactions on medical imaging*, 23(7):903–921, 2004.
- [61] Leo Joskowicz, D Cohen, N Caplan, and Jacob Sosna. Automatic segmentation variability estimation with segmentation priors. *Medical image analysis*, 50:54–64, 2018.
- [62] Leo Joskowicz, D Cohen, N Caplan, and J Sosna. Inter-observer variability of manual contour delineation of structures in ct. *European radiology*, 29(3):1391–1399, 2019.
- [63] Rolf A Heckemann, Joseph V Hajnal, Paul Aljabar, Daniel Rueckert, and Alexander Hammers. Automatic anatomical brain mri segmentation combining label propagation and decision fusion. *NeuroImage*, 33(1):115–126, 2006.
- [64] Zhuowen Tu, Katherine L Narr, Piotr Dollár, Ivo Dinov, Paul M Thompson, and Arthur W Toga. Brain anatomical structure segmentation by hybrid discriminative/generative models. *IEEE transactions on medical imaging*, 27(4):495–508, 2008.
- [65] Juan Eugenio Iglesias, Ender Konukoglu, Albert Montillo, Zhuowen Tu, and Antonio Criminisi. Combining generative and discriminative models for semantic segmentation of ct scans via active learning. In *Biennial International Conference on Information Processing in Medical Imaging*, pages 25–36. Springer, 2011.
- [66] Antonio Criminisi, Krishna Juluru, and Sayan Pathak. A discriminative-generative model for detecting intravenous contrast in ct images. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 49–57. Springer, 2011.
- [67] Bjoern H Menze, Koen Van Leemput, Danial Lashkari, Tammy Riklin-Raviv, Ezequiel Geremia, Esther Alberts, Philipp Gruber, Susanne Wegener, Marc-André Weber, Gabor Székely, et al. A generative probabilistic model and discriminative extensions for brain lesion segmentation?with application to tumor and stroke. *IEEE transactions on medical imaging*, 35(4):933–946, 2015.
- [68] Michael Wels, Gustavo Carneiro, Alexander Aplas, Martin Huber, Joachim Hornegger, and Dorin Comaniciu. A discriminative model-constrained graph cuts approach to fully automated pediatric brain tumor segmentation in 3-d mri. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 67–75. Springer, 2008.

- [69] Loic Le Folgoc, Aditya V Nori, Siddharth Ancha, and Antonio Criminisi. Lifted auto-context forests for brain tumour segmentation. In *International Workshop on Brainlesion: Glioma, Multiple Sclerosis, Stroke and Traumatic Brain Injuries*, pages 171–183. Springer, 2016.
- [70] Mithun Prasad and Arcot Sowmya. Multi-level classification of emphysema in hrct lung images using delegated classifiers. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 59–66. Springer, 2008.
- [71] Ezequiel Geremia, Olivier Clatz, Bjoern H Menze, Ender Konukoglu, Antonio Criminisi, and Nicholas Ayache. Spatial decision forests for ms lesion segmentation in multi-channel magnetic resonance images. *NeuroImage*, 57(2):378–390, 2011.
- [72] Antonio Criminisi, Jamie Shotton, and Stefano Bucciarelli. Decision forests with long-range spatial context for organ localization in ct volumes. In *Medical Image Computing and Computer-Assisted Intervention (MICCAI)*, pages 69–80, 2009.
- [73] Albert Montillo, Jamie Shotton, John Winn, Juan Eugenio Iglesias, Dimitri Metaxas, and Antonio Criminisi. Entangled decision forests and their application for semantic segmentation of ct images. In *Biennial International Conference on Information Processing in Medical Imaging*, pages 184–196. Springer, 2011.
- [74] Jonathan H Morra, Zhiwen Tu, Liana G Apostolova, Amity E Green, Arthur W Toga, and Paul M Thompson. Automatic subcortical segmentation using a contextual model. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 194–201. Springer, 2008.
- [75] Li Wang, Yaozong Gao, Feng Shi, Gang Li, John H Gilmore, Weili Lin, and Dinggang Shen. Links: Learning-based multi-source integration framework for segmentation of infant brain images. *NeuroImage*, 108:160–172, 2015.
- [76] Tianhu Lei and Wilfred Sewchand. Statistical approach to x-ray ct imaging and its applications in image analysis. ii. a new stochastic model-based image segmentation technique for x-ray ct image. *IEEE Transactions on Medical Imaging*, 11(1):62–69, 1992.
- [77] Zhengrong Liang, James R MacFall, and Donald P Harrington. Parameter estimation and tissue segmentation from multispectral mr images. *IEEE transactions on Medical Imaging*, 13(3):441–449, 1994.
- [78] William M Wells, W Eric L Grimson, Ron Kikinis, and Ferenc A Jolesz. Adaptive segmentation of mri data. *IEEE transactions on medical imaging*, 15(4):429–442, 1996.
- [79] Koen Van Leemput, Frederik Maes, Dirk Vandermeulen, and Paul Suetens. Automated model-based tissue classification of mr images of the brain. *IEEE transactions on medical imaging*, 18(10):897–908, 1999.
- [80] Yongyue Zhang, Michael Brady, and Stephen Smith. Segmentation of brain mr images through a hidden markov random field model and the expectation-maximization algorithm. *IEEE transactions on medical imaging*, 20(1):45–57, 2001.
- [81] Paul P Wyatt and J Alison Noble. Map mrf joint segmentation and registration of medical images. *Medical Image Analysis*, 7(4):539–552, 2003.
- [82] John Ashburner and Karl J Friston. Unified segmentation. *Neuroimage*, 26(3):839–851, 2005.
- [83] Kilian M Pohl, John Fisher, W Eric L Grimson, Ron Kikinis, and William M Wells. A bayesian model for joint segmentation and registration. *NeuroImage*, 31(1):228–239, 2006.
- [84] Koen Van Leemput, Akram Bakkour, Thomas Benner, Graham Wiggins, Lawrence L Wald, Jean Augustinack, Bradford C Dickerson, Polina Golland, and Bruce Fischl. Automated segmentation of hippocampal subfields from ultra-high resolution in vivo mri. *Hippocampus*, 19(6):549–557, 2009.
- [85] Mert R Sabuncu, BT Thomas Yeo, Koen Van Leemput, Bruce Fischl, and Polina Golland. A generative model for image segmentation based on label fusion. *IEEE transactions on medical imaging*, 29(10):1714–1729, 2010.
- [86] M Jorge Cardoso, Matthew J Clarkson, Gerard R Ridgway, Marc Modat, Nick C Fox, Sébastien Ourselin, Alzheimer’s Disease Neuroimaging Initiative, et al. Load: a locally adaptive cortical segmentation algorithm. *NeuroImage*, 56(3):1386–1397, 2011.

- [87] Mark William Woolrich and Timothy E Behrens. Variational bayes inference of spatial mixture models for segmentation. *IEEE Transactions on Medical Imaging*, 25(10):1380–1391, 2006.
- [88] GuangJian Tian, Yong Xia, Yanning Zhang, and Dagan Feng. Hybrid genetic and variational expectation-maximization algorithm for gaussian-mixture-model-based brain mr image segmentation. *IEEE transactions on information technology in biomedicine*, 15(3):373–380, 2011.
- [89] Claudia Blaiotta, M Jorge Cardoso, and John Ashburner. Variational inference for medical image segmentation. *Computer Vision and Image Understanding*, 151:14–28, 2016.
- [90] Juan Eugenio Iglesias, Mert Rory Sabuncu, and Koen Van Leemput. Incorporating parameter uncertainty in bayesian segmentation models: Application to hippocampal subfield volumetry. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 50–57. Springer, 2012.
- [91] Juan Eugenio Iglesias, Mert Rory Sabuncu, Koen Van Leemput, Alzheimer’s Disease Neuroimaging Initiative, et al. Improved inference in bayesian segmentation using monte carlo sampling: Application to hippocampal subfield volumetry. *Medical image analysis*, 17(7):766–778, 2013.
- [92] Juan Eugenio Iglesias, Ender Konukoglu, Darko Zikic, Ben Glocker, Koen Van Leemput, and Bruce Fischl. Is synthesizing mri contrast useful for inter-modality analysis? In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 631–638. Springer, 2013.
- [93] Dong Hye Ye, Darko Zikic, Ben Glocker, Antonio Criminisi, and Ender Konukoglu. Modality propagation: coherent synthesis of subject-specific scans with data-driven regularization. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 606–613. Springer, 2013.
- [94] François Rousseau. Brain hallucination. In *ECCV 2008*, pages 497–508. Springer, 2008.
- [95] Yu Zhang, Guorong Wu, Pew-Thian Yap, Qianjin Feng, Jun Lian, Wufan Chen, and Dinggang Shen. Hierarchical patch-based sparse representation? a new approach for resolution enhancement of 4d-ct lung data. *IEEE transactions on medical imaging*, 31(11):1993–2005, 2012.
- [96] Wolfgang Wein, Shelby Brunke, Ali Khamene, Matthew R Callstrom, and Nassir Navab. Automatic CT-ultrasound registration for diagnostic imaging and image-guided intervention. *Medical image analysis*, 12(5):577–585, 2008.
- [97] Ninon Burgos, M Jorge Cardoso, Kris Thielemans, Marc Modat, Stefano Pedemonte, John Dickson, Anna Barnes, Rebekah Ahmed, Colin J Mahoney, Jonathan M Schott, et al. Attenuation correction synthesis for hybrid pet-mr scanners: application to brain studies. *IEEE transactions on medical imaging*, 33(12):2332–2341, 2014.
- [98] M Jorge Cardoso, Carole H Sudre, Marc Modat, and Sébastien Ourselin. Template-based multi-modal joint generative model of brain data. In *International conference on information processing in medical imaging*, pages 17–29. Springer, 2015.
- [99] Nicolas Cordier, Hervé Delingette, Matthieu Lê, and Nicholas Ayache. Extended modality propagation: Image synthesis of pathological cases. *IEEE transactions on medical imaging*, 35(12):2598–2608, 2016.
- [100] Ryutaro Tanno, Aurobrata Ghosh, Francesco Grussu, Enrico Kaden, Antonio Criminisi, and Daniel C Alexander. Bayesian image quality transfer. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 265–273. Springer, 2016.
- [101] Andre Esteva, Alexandre Robicquet, Bharath Ramsundar, Volodymyr Kuleshov, Mark DePristo, Katherine Chou, Claire Cui, Greg Corrado, Sebastian Thrun, and Jeff Dean. A guide to deep learning in healthcare. *Nature medicine*, 25(1):24, 2019.
- [102] Ryutaro Tanno, Daniel E Worrall, Aurobrata Ghosh, Enrico Kaden, Stamatios N Sotiroopoulos, Antonio Criminisi, and Daniel C Alexander. Bayesian image quality transfer with cnns: Exploring uncertainty in dmri super-resolution. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 611–619. Springer, 2017.
- [103] Ryutaro Tanno, Daniel Worrall, Enrico Kaden, Aurobrata Ghosh, Francesco Grussu, Alberto Bizzi, Stamatios N Sotiroopoulos, Antonio Criminisi, and Daniel C Alexander. Uncertainty quantification in deep learning for safer neuroimage enhancement. 2019.

- [104] Felix Bragman, Ryu Tanno, Zach Eaton-Rosen, Wenqi Li, David Hawkes, Sebastien Ourselin, Daniel Alexander, Jamie McClelland, and M. Jorge Cardoso. Uncertainty in multitask learning: joint representations for probabilistic mr-only radiotherapy planning. In *Medical Image Computing and Computer-Assisted Interventions (MICCAI)*, 2018.
- [105] Felix JS Bragman, Ryutaro Tanno, Sebastien Ourselin, Daniel C Alexander, and M Jorge Cardoso. Stochastic filter groups for multi-task cnns: Learning specialist and generalist convolution kernels. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2019.
- [106] Ryutaro Tanno, Kai Arulkumaran, Daniel Alexander, Antonio Criminisi, and Aditya Nori. Adaptive neural trees. In *Proceedings of the 36th International Conference on Machine Learning (ICML)*, 2019.
- [107] Ryutaro Tanno, Ardavan Saeedi, Swami Sankaranarayanan, Daniel C Alexander, and Nathan Silberman. Learning from noisy labels by regularized estimation of annotator confusion. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [108] Le Zhang, Ryutaro Tanno, Mou-Cheng Xu, Joseph Jacob, Olga Ciccarelli, Frederik Barkhof, and Daniel C Alexander. Disentangling human error from the ground truth in segmentation of medical images. *arXiv preprint arXiv:2007.15963*, 2020.
- [109] Konstantinos Kamnitsas, Christian Ledig, Virginia FJ Newcombe, Joanna P Simpson, Andrew D Kane, David K Menon, Daniel Rueckert, and Ben Glocker. Efficient multi-scale 3d cnn with fully connected crf for accurate brain lesion segmentation. *Medical image analysis*, 36:61–78, 2017.
- [110] Holger R Roth, Le Lu, Ari Seff, Kevin M Cherry, Joanne Hoffman, Shijun Wang, Jiamin Liu, Evrim Turkbey, and Ronald M Summers. A new 2.5 d representation for lymph node detection using random sets of deep convolutional neural network observations. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 520–527. Springer, 2014.
- [111] Teresa Araújo, Guilherme Aresta, Eduardo Castro, José Rouco, Paulo Aguiar, Catarina Eloy, António Polónia, and Aurélio Campilho. Classification of breast cancer histology images using convolutional neural networks. *PloS one*, 12(6):e0177544, 2017.
- [112] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1125–1134, 2017.
- [113] Ozan Oktay, Wenjia Bai, Matthew Lee, Ricardo Guerrero, Konstantinos Kamnitsas, Jose Caballero, Antonio de Marvao, Stuart Cook, Declan O'Regan, and Daniel Rueckert. Multi-input cardiac image super-resolution using convolutional neural networks. In *MICCAI*. Springer, 2016.
- [114] Yuhua Chen, Feng Shi, Anthony G Christodoulou, Yibin Xie, Zhengwei Zhou, and Debiao Li. Efficient and accurate mri super-resolution using a generative adversarial network and 3d multi-level densely connected network. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 91–99. Springer, 2018.
- [115] Daniele Ravì, Agnieszka Barbara Szczotka, Stephen P Pereira, and Tom Vercauteren. Adversarial training with cycle consistency for unsupervised super-resolution in endomicroscopy. *Medical image analysis*, 53:123–131, 2019.
- [116] Dong Nie, Xiaohuan Cao, Yaozong Gao, Li Wang, and Dinggang Shen. Estimating ct image from mri data using 3d fully convolutional networks. In *Deep Learning and Data Labeling for Medical Applications*, pages 170–178. Springer, 2016.
- [117] Eunhee Kang, Junhong Min, and Jong Chul Ye. A deep convolutional neural network using directional wavelets for low-dose x-ray ct reconstruction. *Medical physics*, 44(10), 2017.
- [118] A Benou, R Veksler, A Friedman, and T Riklin Raviv. Ensemble of expert deep neural networks for spatio-temporal denoising of contrast-enhanced mri sequences. *Medical image analysis*, 42:145–159, 2017.
- [119] Hu Chen, Yi Zhang, Weihua Zhang, Peixi Liao, Ke Li, Jiliu Zhou, and Ge Wang. Low-dose ct via convolutional neural network. *Biomedical optics express*, 8(2):679–694, 2017.

- [120] Suheyla Cetin Karayumak, Marek Kubicki, and Yogesh Rathi. Harmonizing diffusion mri data across magnetic field strengths. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 116–124. Springer, 2018.
- [121] Chantal MW Tax, Francesco Grussu, Enrico Kaden, Lipeng Ning, Umesh Rudrapatna, John Evans, Samuel St-Jean, Alexander Leemans, Simon Koppers, Dorit Merhof, et al. Cross-scanner and cross-protocol diffusion mri data harmonisation: A benchmark database and evaluation of algorithms. *NeuroImage*, 2019.
- [122] Jian Sun, Huibin Li, Zongben Xu, et al. Deep admm-net for compressive sensing mri. In *Advances in neural information processing systems*, pages 10–18, 2016.
- [123] Kyong Hwan Jin, Michael T McCann, Emmanuel Froustey, and Michael Unser. Deep convolutional neural network for inverse problems in imaging. *IEEE Transactions on Image Processing*, 26(9):4509–4522, 2017.
- [124] Kerstin Hammernik, Teresa Klatzer, Erich Kobler, Michael P Recht, Daniel K Sodickson, Thomas Pock, and Florian Knoll. Learning a variational network for reconstruction of accelerated mri data. *Magnetic resonance in medicine*, 79(6):3055–3071, 2018.
- [125] Jo Schlemper, Jose Caballero, Joseph V Hajnal, Anthony N Price, and Daniel Rueckert. A deep cascade of convolutional neural networks for dynamic mr image reconstruction. *IEEE transactions on medical imaging*, 37(2):491–503, 2018.
- [126] Bo Zhu, Jeremiah Z Liu, Stephen F Cauley, Bruce R Rosen, and Matthew S Rosen. Image reconstruction by domain-transform manifold learning. *Nature*, 555(7697):487, 2018.
- [127] Guang Yang, Simiao Yu, Hao Dong, Greg Slabaugh, Pier Luigi Dragotti, Xujiong Ye, Fangde Liu, Simon Arridge, Jennifer Keegan, Yike Guo, et al. Dagan: Deep de-aliasing generative adversarial networks for fast compressed sensing mri reconstruction. *IEEE transactions on medical imaging*, 37(6):1310–1321, 2018.
- [128] Yeo Hun Yoon, Shujaat Khan, Jaeyoung Huh, and Jong Chul Ye. Efficient b-mode ultrasound image reconstruction from sub-sampled rf data using deep learning. *IEEE transactions on medical imaging*, 38(2):325–336, 2019.
- [129] Hessam Sokooti, Bob de Vos, Floris Berendsen, Boudeijn PF Lelieveldt, Ivana Išgum, and Marius Staring. Nonrigid image registration using multi-scale 3d convolutional neural networks. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 232–239. Springer, 2017.
- [130] Guha Balakrishnan, Amy Zhao, Mert R Sabuncu, John Guttag, and Adrian V Dalca. An unsupervised learning model for deformable medical image registration. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9252–9260, 2018.
- [131] Lingyun Wu, Jie-Zhi Cheng, Shengli Li, Baiying Lei, Tianfu Wang, and Dong Ni. Fuiqa: Fetal ultrasound image quality assessment with deep convolutional networks. *IEEE transactions on cybernetics*, 47(5):1336–1349, 2017.
- [132] Steven J Esses, Xiaoguang Lu, Tiejun Zhao, Krishna Shanbhogue, Bari Dane, Mary Bruno, and Hersh Chandarana. Automated image quality evaluation of t2-weighted liver mri utilizing deep learning architecture. *Journal of Magnetic Resonance Imaging*, 47(3):723–728, 2018.
- [133] Joseph Paul Cohen, Margaux Luck, and Sina Honari. Distribution matching losses can hallucinate features in medical image translation. *arXiv preprint arXiv:1805.08841*, 2018.
- [134] Edmon Begoli, Tanmoy Bhattacharya, and Dimitri Kusnezov. The need for uncertainty quantification in machine-assisted medical decision making. *Nature Machine Intelligence*, 1(1):20, 2019.
- [135] David A Nix and Andreas S Weigend. Estimating the mean and variance of the target probability distribution. In *IEEE WCCI*, volume 1, pages 55–60. IEEE, 1994.
- [136] Diederik P Kingma, Tim Salimans, and Max Welling. Variational dropout and the local reparameterization trick. In *NIPS*, pages 2575–2583, 2015.
- [137] Daniel C Alexander and et al. Image quality transfer via random forest regression: applications in diffusion MRI. In *MICCAI 2014*, pages 225–232. Springer, 2014.

- [138] Daniel C Alexander, Darko Zikic, Aurobrata Ghosh, Ryutaro Tanno, Viktor Wottschel, Jiaying Zhang, Enrico Kaden, Tim B Dyrby, Stamatios N Sotropoulos, Hui Zhang, et al. Image quality transfer and applications in diffusion mri. *Neuroimage*, 152:283–298, 2017.
- [139] Stefano B Blumberg, Ryutaro Tanno, Iasonas Kokkinos, and Daniel C Alexander. Deeper image quality transfer: Training low-memory neural networks for 3d images. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 118–125. Springer, 2018.
- [140] Stamatios N Sotropoulos and et al. Advances in diffusion MRI acquisition and processing in the human connectome project. *Neuroimage*, 80:125–143, 2013.
- [141] Michael P Harms, Leah H Somerville, Beau M Ances, Jesper Andersson, Deanna M Barch, Matteo Bastiani, Susan Y Bookheimer, Timothy B Brown, Randy L Buckner, Gregory C Burgess, et al. Extending the human connectome project across ages: Imaging protocols for the lifespan development and aging projects. *NeuroImage*, 183:972–984, 2018.
- [142] Xiao Yang, Roland Kwitt, and Marc Niethammer. Fast predictive image registration. In *Deep Learning and Data Labeling for Medical Applications*, pages 48–57. Springer, 2016.
- [143] Khosro Bahrami, Feng Shi, Islem Rekik, and Dinggang Shen. Convolutional neural network for reconstruction of 7t-like images from 3t mri using appearance and anatomical features. In *MICCAI DLDLM workshop*, pages 39–47. Springer, 2016.
- [144] Ozan Oktay, Enzo Ferrante, Konstantinos Kamnitsas, Matthias Heinrich, Wenjia Bai, Jose Caballero, Stuart A Cook, Antonio de Marvao, Timothy Dawes, Declan P O'Regan, et al. Anatomically constrained neural networks (acnn): application to cardiac image enhancement and segmentation. *IEEE transactions on medical imaging*, 37(2):384–395, 2018.
- [145] Ian Goodfellow, Yoshua Bengio, Aaron Courville, and Yoshua Bengio. *Deep learning*, volume 1. MIT Press, 2016.
- [146] Wenzhe Shi, Jose Caballero, Ferenc Huszár, Johannes Totz, Andrew P Aitken, Rob Bishop, Daniel Rueckert, and Zehan Wang. Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network. In *CVPR*, pages 1874–1883, 2016.
- [147] Augustus Odena, Vincent Dumoulin, and Chris Olah. Deconvolution and checkerboard artifacts. *Distill*, 2016.
- [148] Hong Wang, Dennis M Levi, and Stanley A Klein. Intrinsic uncertainty and integration efficiency in bisection acuity. *Vision research*, 36(5):717–739, 1996.
- [149] David Draper. Assessment and propagation of model uncertainty. *Journal of the Royal Statistical Society: Series B (Methodological)*, 57(1):45–70, 1995.
- [150] Alex Kendall and Yarin Gal. What uncertainties do we need in bayesian deep learning for computer vision? In *Advances in Neural Information Processing Systems*, pages 5580–5590, 2017.
- [151] C Radhakrishna Rao. Estimation of heteroscedastic variances in linear models. *Journal of the American Statistical Association*, 65(329):161–172, 1970.
- [152] Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Insights and applications. In *Deep Learning Workshop, ICML*, 2015.
- [153] David M Blei, Alp Kucukelbir, and Jon D McAuliffe. Variational inference: A review for statisticians. *Journal of the American Statistical Association*, 112(518):859–877, 2017.
- [154] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research*, 15(1):1929–1958, 2014.
- [155] Daniel E Worrall, Clare M Wilson, and Gabriel J Brostow. Automated retinopathy of prematurity case detection with convolutional neural networks. In *MICCAI DLDLM Workshop*, pages 68–76. Springer, 2016.
- [156] Christian Leibig, Vaneeda Allken, Murat Seçkin Ayhan, Philipp Berens, and Siegfried Wahl. Leveraging uncertainty information from deep neural networks for disease detection. *Scientific reports*, 7(1):17816, 2017.

- [157] Abhijit Guha Roy, Sailesh Conjeti, Nassir Navab, Christian Wachinger, Alzheimer’s Disease Neuroimaging Initiative, et al. Bayesian quicknat: Model uncertainty in deep whole-brain segmentation for structure-wise quality control. *NeuroImage*, 2019.
- [158] Zach Eaton-Rosen, Felix Bragman, Sotirios Bisdas, Sébastien Ourselin, and M Jorge Cardoso. Towards safe deep learning: accurately quantifying biomarker uncertainty in neural network predictions. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 691–699. Springer, 2018.
- [159] Tanya Nair, Doina Precup, Douglas L Arnold, and Tal Arbel. Exploring uncertainty measures in deep networks for multiple sclerosis lesion detection and segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 655–663. Springer, 2018.
- [160] Felix JS Bragman, Ryutaro Tanno, Zach Eaton-Rosen, Wenqi Li, David J Hawkes, Sébastien Ourselin, Daniel C Alexander, Jamie R McClelland, and M Jorge Cardoso. Uncertainty in multitask learning: joint representations for probabilistic mr-only radiotherapy planning. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 2018.
- [161] Radford M Neal. Bayesian learning via stochastic dynamics. In *Advances in neural information processing systems*, pages 475–482, 1993.
- [162] Max Welling and Yee W Teh. Bayesian learning via stochastic gradient langevin dynamics. In *Proceedings of the 28th International Conference on Machine Learning (ICML-11)*, pages 681–688, 2011.
- [163] Tianqi Chen, Emily Fox, and Carlos Guestrin. Stochastic gradient hamiltonian monte carlo. In *International Conference on Machine Learning*, pages 1683–1691, 2014.
- [164] Yi-An Ma, Tianqi Chen, and Emily Fox. A complete recipe for stochastic gradient mcmc. In *Advances in Neural Information Processing Systems*, pages 2917–2925, 2015.
- [165] Dmitry Molchanov, Arsenii Ashukha, and Dmitry Vetrov. Variational dropout sparsifies deep neural networks. *arXiv preprint arXiv:1701.05369*, 2017.
- [166] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. In *International Conference on Learning Representations*, 2014.
- [167] Neil A Weiss. *A course in probability*. Addison-Wesley, 2006.
- [168] Clive G Bowsher and Peter S Swain. Identifying sources of variation and the flow of information in biochemical networks. *Proceedings of the National Academy of Sciences*, 2012.
- [169] Can Zhao, Aaron Carass, Blake E Dewey, Jonghye Woo, Jiwon Oh, Peter A Calabresi, Daniel S Reich, Pascal Sati, Dzung L Pham, and Jerry L Prince. A deep learning based anti-aliasing self super-resolution algorithm for mri. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 100–108. Springer, 2018.
- [170] Dwarikanath Mahapatra, Behzad Bozorgtabar, Sajini Hewavitharanage, and Rahil Garnavi. Image super resolution using generative adversarial networks and local saliency maps for retinal image analysis. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 382–390. Springer, 2017.
- [171] Haichao Yu, Ding Liu, Honghui Shi, Hanchao Yu, Zhangyang Wang, Xinchao Wang, Brent Cross, Matthew Bramler, and Thomas S Huang. Computed tomography super-resolution using convolutional neural networks. In *Image Processing (ICIP), 2017 IEEE International Conference on*, pages 3944–3948. IEEE, 2017.
- [172] Dong Nie, Roger Trullo, Jun Lian, Li Wang, Caroline Petitjean, Su Ruan, Qian Wang, and Ding-gang Shen. Medical image synthesis with deep convolutional adversarial networks. *IEEE Transactions on Biomedical Engineering*, 2018.
- [173] Jelmer M Wolterink, Anna M Dinkla, Mark HF Savenije, Peter R Seevinck, Cornelis AT van den Berg, and Ivana Išgum. Deep mr to ct synthesis using unpaired data. In *International Workshop on Simulation and Synthesis in Medical Imaging*, pages 14–23. Springer, 2017.
- [174] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. *arXiv preprint*, 2017.

- [175] Stefano B Blumberg, Marco Palombo, Can Son Khoo, Chantal Tax, Ryutaro Tanno, and Daniel C Alexander. Multi-stage prediction networks for data harmonization. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2019.
- [176] Hu Shi, Daniel Worrall, Bas Veeling, Henkjan Huisman, and Max Welling. Supervised uncertainty quantification for segmentation with multiple annotations. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 2019.
- [177] Adrian V Dalca, Guha Balakrishnan, John Guttag, and Mert R Sabuncu. Unsupervised learning for fast probabilistic diffeomorphic registration. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 2018.
- [178] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. In *Advances in Neural Information Processing Systems*, pages 6402–6413, 2017.
- [179] Jo Schlemper, Guang Yang, Pedro Ferreira, Andrew Scott, Laura-Ann McGill, Zohya Khalique, Margarita Gorodezky, Malte Roehl, Jennifer Keegan, Dudley Pennell, et al. Stochastic deep compressive sensing for the reconstruction of diffusion tensor cardiac mri. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 2018.
- [180] Murat Seckin Ayhan and Philipp Berens. Test-time data augmentation for estimation of heteroscedastic aleatoric uncertainty in deep neural networks. 2018.
- [181] Maithra Raghu, Katy Blumer, Rory Sayres, Ziad Obermeyer, Sendhil Mullainathan, and Jon M. Kleinberg. Direct uncertainty prediction for medical second opinions. 2018.
- [182] Simon Kohl, Bernardino Romera-Paredes, Clemens Meyer, Jeffrey De Fauw, Joseph R Ledsam, Klaus Maier-Hein, SM Ali Eslami, Danilo Jimenez Rezende, and Olaf Ronneberger. A probabilistic u-net for segmentation of ambiguous images. In *Advances in Neural Information Processing Systems*, pages 6965–6975, 2018.
- [183] Christian F. Baumgartner, Kerem C. Tezcan, Krishna Chaitanya, Andreas M. Hötker, Urs J. Muehlematter, Khoschy Schawkat, Anton S. Becker, Olivio Donati, and Ender Konukoglu. PhiSeg: Capturing uncertainty in medical image segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 2019.
- [184] Vikas C Raykar, Shipeng Yu, Linda H Zhao, Gerardo Hermosillo Valadez, Charles Florin, Luca Bognoni, and Linda Moy. Learning from crowds. *Journal of Machine Learning Research*, 11(Apr):1297–1322, 2010.
- [185] David C Van Essen, Stephen M Smith, Deanna M Barch, Timothy EJ Behrens, Essa Yacoub, Kamil Ugurbil, Wu-Minn HCP Consortium, et al. The wu-minn human connectome project: an overview. *Neuroimage*, 80:62–79, 2013.
- [186] Matthew F Glasser, Stamatios N Sotiroopoulos, J Anthony Wilson, Timothy S Coalson, Bruce Fischl, Jesper L Andersson, Junqian Xu, Saad Jbabdi, Matthew Webster, Jonathan R Polimeni, et al. The minimal preprocessing pipelines for the human connectome project. *Neuroimage*, 80:105–124, 2013.
- [187] Matteo Figini, Marco Riva, Mark Graham, Gian Marco Castelli, Bethania Fernandes, Marco Grimaldi, Giuseppe Baselli, Federico Pessina, Lorenzo Bello, Hui Zhang, et al. Prediction of isocitrate dehydrogenase genotype in brain gliomas with mri: single-shell versus multishell diffusion models. *Radiology*, 289(3):788–796, 2018.
- [188] Peter J Basser, James Mattiello, and Denis LeBihan. Mr diffusion tensor spectroscopy and imaging. *Biophysical journal*, 66(1):259–267, 1994.
- [189] Evren Özarslan, Cheng Guan Koay, Timothy M Shepherd, Michal E Komlosh, M Okan İrfanoğlu, Carlo Pierpaoli, and Peter J Basser. Mean apparent propagator (map) mri: a novel diffusion imaging method for mapping tissue microstructure. *NeuroImage*, 78:16–32, 2013.
- [190] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *CoRR*, 2014.
- [191] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pages 448–456, 2015.

- [192] Yarin Gal, Jiri Hron, and Alex Kendall. Concrete dropout. In *Advances in Neural Information Processing Systems*, pages 3581–3590, 2017.
- [193] Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. Tensorflow: a system for large-scale machine learning. In *OSDI*, volume 16, pages 265–283, 2016.
- [194] Zhou Wang, Alan Conrad Bovik, Hamid Rahim Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE Trans Image Proc*, 13(4), 2004.
- [195] Christopher M Bishop. Mixture density networks. Technical report, Citeseer, 1994.
- [196] Abner Guzman-Rivera, Dhruv Batra, and Pushmeet Kohli. Multiple choice learning: Learning to produce multiple structured outputs. In *Advances in Neural Information Processing Systems*, pages 1799–1807, 2012.
- [197] Diane Bouchacourt, Pawan K Mudigonda, and Sebastian Nowozin. Disco nets: Dissimilarity coefficients networks. In *Advances in Neural Information Processing Systems*, pages 352–360, 2016.
- [198] Hsin-Ying Lee, Hung-Yu Tseng, Jia-Bin Huang, Maneesh Singh, and Ming-Hsuan Yang. Diverse image-to-image translation via disentangled representations. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 35–51, 2018.
- [199] Xun Huang, Ming-Yu Liu, Serge Belongie, and Jan Kautz. Multimodal unsupervised image-to-image translation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 172–189, 2018.
- [200] Danilo Jimenez Rezende and Shakir Mohamed. Variational inference with normalizing flows. *arXiv preprint arXiv:1505.05770*, 2015.
- [201] George Papamakarios, Theo Pavlakou, and Iain Murray. Masked autoregressive flow for density estimation. In *Advances in Neural Information Processing Systems*, pages 2338–2347, 2017.
- [202] Augustus Odena, Christopher Olah, and Jonathon Shlens. Conditional image synthesis with auxiliary classifier gans. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 2642–2651. JMLR. org, 2017.
- [203] Siddhartha Chandra and Iasonas Kokkinos. Fast, exact and multi-scale inference for semantic image segmentation with deep gaussian crfs. In *European Conference on Computer Vision*, pages 402–418. Springer, 2016.
- [204] Christos Louizos and Max Welling. Structured and efficient variational deep learning with matrix gaussian posteriors. In *International Conference on Machine Learning*, pages 1708–1716, 2016.
- [205] Changyong Oh, Kamil Adamczewski, and Mijung Park. Radial and directional posteriors for bayesian neural networks. *arXiv preprint arXiv:1902.02603*, 2019.
- [206] David Krueger, Chin-Wei Huang, Riashat Islam, Ryan Turner, Alexandre Lacoste, and Aaron Courville. Bayesian hypernetworks. *arXiv preprint arXiv:1710.04759*, 2017.
- [207] Ruqi Zhang, Chunyuan Li, Jianyi Zhang, Changyou Chen, and Andrew Gordon Wilson. Cyclical stochastic gradient mcmc for bayesian deep learning. *arXiv preprint arXiv:1902.03932*, 2019.
- [208] Nick Pawlowski, Andrew Brock, Matthew CH Lee, Martin Rajchl, and Ben Glocker. Implicit weight uncertainty in neural networks. *arXiv preprint arXiv:1711.01297*, 2017.
- [209] Christos Louizos and Max Welling. Multiplicative normalizing flows for variational bayesian neural networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 2218–2227. JMLR. org, 2017.
- [210] Cheng Zhang, Judith Bütepage, Hedvig Kjellström, and Stephan Mandt. Advances in variational inference. *IEEE transactions on pattern analysis and machine intelligence*, 41(8):2008–2026, 2018.
- [211] Burr Settles. Active learning literature survey. Technical report, University of Wisconsin-Madison Department of Computer Sciences, 2009.

- [212] Yarin Gal, Riashat Islam, and Zoubin Ghahramani. Deep bayesian active learning with image data. In *Proceedings of the 34th International Conference on Machine Learning- Volume 70*, pages 1183–1192. JMLR. org, 2017.
- [213] Marc Gorriz, Axel Carlier, Emmanuel Faure, and Xavier Giro-i Nieto. Cost-effective active learning for melanoma segmentation. *arXiv preprint arXiv:1711.09168*, 2017.
- [214] Sinno Jialin Pan and Qiang Yang. A survey on transfer learning. *IEEE Transactions on knowledge and data engineering*, 22(10):1345–1359, 2010.
- [215] Konstantinos Kamnitsas, Christian Baumgartner, Christian Ledig, Virginia Newcombe, Joanna Simpson, Andrew Kane, David Menon, Aditya Nori, Antonio Criminisi, Daniel Rueckert, et al. Unsupervised domain adaptation in brain lesion segmentation with adversarial networks. In *International conference on information processing in medical imaging*, pages 597–609. Springer, 2017.
- [216] Neerav Karani, Krishna Chaitanya, Christian Baumgartner, and Ender Konukoglu. A lifelong learning approach to brain mr segmentation across scanners and protocols. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 476–484. Springer, 2018.
- [217] Chaitanya Baweja, Ben Glocker, and Konstantinos Kamnitsas. Towards continual learning in medical imaging. *arXiv preprint arXiv:1811.02496*, 2018.
- [218] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114(13):3521–3526, 2017.
- [219] Cuong V Nguyen, Yingzhen Li, Thang D Bui, and Richard E Turner. Variational continual learning. In *International Conference on Learning Representations*, 2018.
- [220] Chelsea Finn, Kelvin Xu, and Sergey Levine. Probabilistic model-agnostic meta-learning. In *Advances in Neural Information Processing Systems*, pages 9516–9527, 2018.
- [221] Jaesik Yoon, Taesup Kim, Ousmane Dia, Sungwoong Kim, Yoshua Bengio, and Sungjin Ahn. Bayesian model-agnostic meta-learning. In *Advances in Neural Information Processing Systems*, pages 7332–7342, 2018.
- [222] Dan Ma, Vikas Gulani, Nicole Seiberlich, Kecheng Liu, Jeffrey L Sunshine, Jeffrey L Duerk, and Mark A Griswold. Magnetic resonance fingerprinting. *Nature*, 495(7440):187–192, 2013.
- [223] Ouri Cohen, Bo Zhu, and Matthew S Rosen. Mr fingerprinting deep reconstruction network (drone). *Magnetic resonance in medicine*, 80(3):885–894, 2018.
- [224] Yoseob Han, Jaejun Yoo, Hak Hee Kim, Hee Jung Shin, Kyunghyun Sung, and Jong Chul Ye. Deep learning with domain adaptation for accelerated projection-reconstruction mr. *Magnetic resonance in medicine*, 80(3):1189–1205, 2018.
- [225] Amod Jog, Aaron Carass, Snehashis Roy, Dzung L Pham, and Jerry L Prince. MR image synthesis by contrast learning on neighborhood ensembles. *Medical image analysis*, 24(1):63–76, 2015.
- [226] Ninon Burgos, M Jorge Cardoso, Filipa Guerreiro, Catarina Veiga, Marc Modat, Jamie McClelland, Antje-Christin Knopf, Shonit Punwani, David Atkinson, Simon R Arridge, et al. Robust CT synthesis for radiotherapy planning: Application to the head and neck region. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 476–484. Springer, 2015.
- [227] Hengameh Mirzaalian, Lipeng Ning, Peter Savadjiev, Ofer Pasternak, Sylvain Bouix, O Michailovich, G Grant, CE Marx, Rajendra A Morey, LA Flashman, et al. Inter-site and inter-scanner diffusion mri data harmonization. *NeuroImage*, 135:311–323, 2016.
- [228] Rasmus L Christiansen, Henrik R Jensen, and Carsten Brink. Magnetic resonance only workflow and validation of dose calculations for radiotherapy of prostate cancer. *Acta Oncologica*, 56(6):787–791, 2017.
- [229] Neelam Tyagi, Sandra Fontenla, Michael Zelefsky, Marcia Chong-Ton, Kyle Ostergren, Niral Shah, Lizette Warner, Mo Kadbi, Jim Mechakos, and Margie Hunt. Clinical workflow for mr-only simulation and planning in prostate. *Radiation Oncology*, 12(1):119, 2017.

- [230] Mikko Tenhunen, Juha Korhonen, Mika Kapanen, Tiina Seppälä, Lauri Koivula, Juhani Collan, Kauko Saarilahti, and Harri Visapää. MRI-only based radiation therapy of prostate cancer: workflow and early clinical experience. *Acta Oncologica*, 57(7):902–907, 2018.
- [231] Joakim Jonsson, Tufve Nyholm, and Karin Söderkvist. The rationale for mr-only treatment planning for external radiotherapy. *Clinical and Translational Radiation Oncology*, 2019.
- [232] Jens M Edmund and Tufve Nyholm. A review of substitute ct generation for mri-only radiation therapy. *Radiation Oncology*, 12(1):28, 2017.
- [233] Emily Johnstone, Jonathan J Wyatt, Ann M Henry, Susan C Short, David Sebag-Montefiore, Louise Murray, Charles G Kelly, Hazel M McCallum, and Richard Speight. Systematic review of synthetic computed tomography generation methodologies for use in magnetic resonance imaging-only radiation therapy. *International Journal of Radiation Oncology\* Biology\* Physics*, 100(1):199–217, 2018.
- [234] Rich Caruana. Multitask learning. *Machine learning*, 28(1):41–75, 1997.
- [235] CC Parker, A Damyanovich, T Haycocks, M Haider, A Bayley, and CN Catton. Magnetic resonance imaging in the radiation treatment planning of localized prostate cancer using intra-prostatic fiducial markers for computed tomography co-registration. *Radiotherapy and oncology*, 66(2):217–224, 2003.
- [236] Michael Milosevic, Sachi Voruganti, Ralph Blend, Hamideh Alasti, Padraig Warde, Michael McLean, Pamela Catton, Charles Catton, and Mary Gospodarowicz. Magnetic resonance imaging (mri) for localization of the prostatic apex: comparison to computed tomography (ct) and urethrography. *Radiotherapy and oncology*, 47(3):277–284, 1998.
- [237] International Atomic Energy Agency. *Investigation of an accidental exposure of radiotherapy patients in Panama*. International Atomic Energy Agency, 2001.
- [238] G Wack, F Lalande, and MD Seligman. Summary of asn report no. 2006 enstr 019-igas n rm 2007-015p on the epinal radiotherapy accident. *French Nuclear Safety Authority*, 2007.
- [239] N. Burgos et al. Iterative framework for the joint segmentation and ct synthesis of mr images: application to mri-only radiotherapy treatment planning. *Phys. Med. Biol.*, 62, 2017.
- [240] D. Nie et al. Medical image synthesis with context-aware generative adversarial networks. *arXiv:1612.05362*.
- [241] J. Wolterink et al. Deep mr to ct synthesis using unpaired data. In *SASHIMI*, pages 14–22, 2017.
- [242] P. Moeskops et al. Deep learning for multi-task medical image segmentation in multiple modalities. In *MICCAI*, pages 478–486, 2016.
- [243] Ryutaro Tanno, Antonios Makropoulos, Salim Arslan, Ozan Oktay, Sven Mischkewitz, Fouad Al-Noor, Jonas Oppenheimer, Ramin Mandegaran, Bernhard Kainz, and Mattias P Heinrich. Autodvt: Joint real-time classification for vein compressibility analysis in deep vein thrombosis ultrasound diagnostics. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 905–912. Springer, 2018.
- [244] Alex Kendall, Yarin Gal, and Roberto Cipolla. Multi-task learning using uncertainty to weigh losses for scene geometry and semantics. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [245] Andrew J Asman and Bennett A Landman. Robust statistical label fusion through consensus level, labeler accuracy, and truth estimation (collate). *IEEE transactions on medical imaging*, 30(10):1779–1794, 2011.
- [246] Rich Caruana. Multitask learning: A knowledge-based source of inductive bias. In *ICML*, 1993.
- [247] W. Li et al. On the compactness, efficiency, and representation of 3d convolutional networks: brain parcellation as a pretext task. In *IPMI*, pages 348–360, 2017.
- [248] Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, pages 1050–1059, 2016.

- [249] Fisher Yu and Vladlen Koltun. Multi-scale context aggregation by dilated convolutions. *arXiv preprint arXiv:1511.07122*, 2015.
- [250] Gordon Sands, Enda F Fallon, and Wilhelm J van der Putten. The utilisation of probabilistic risk assessment in radiation oncology. *Procedia Manufacturing*, 3:250–257, 2015.
- [251] David Tilly, ?sa Holm, Erik Grusell, and Anders Ahnesjö. Probabilistic optimization of dose coverage in radiotherapy. *Physics and Imaging in Radiation Oncology*, 10:1–6, 2019.
- [252] Pierre Sermanet, David Eigen, Xiang Zhang, Michaël Mathieu, Rob Fergus, and Yann LeCun. Overfeat: Integrated recognition, localization and detection using convolutional networks. In *ICLR*, 2014.
- [253] David Eigen and Rob Fergus. Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture. In *Proceedings of the IEEE international conference on computer vision*, pages 2650–2658, 2015.
- [254] Ishan Misra, Abhinav Shrivastava, Abhinav Gupta, and Martial Hebert. Cross-stitch Networks for Multi-task Learning. In *CVPR*, 2016.
- [255] Iasonas Kokkinos. Ubernet: Training a universal convolutional neural network for low-, mid-, and high-level vision using diverse datasets and limited memory. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6129–6138, 2017.
- [256] Rajeev Ranjan, Vishal M Patel, and Rama Chellappa. Hyperface: A deep multi-task learning framework for face detection, landmark localization, pose estimation, and gender recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 41(1):121–135, 2019.
- [257] Hakan Bilen and Andrea Vedaldi. Integrated perception with recurrent multi-task neural networks. In *Advances in neural information processing systems*, pages 235–243, 2016.
- [258] Sihong Chen, Dong Ni, Jing Qin, Baiying Lei, Tianfu Wang, and Jie-Zhi Cheng. Bridging computational features toward multiple semantic features with multi-task regression: A study of ct pulmonary nodules. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 53–60. Springer, 2016.
- [259] Sebastian Ruder, Joachim Bingel, Isabelle Augenstein, and Anders Søgaard. Latent multi-task architecture learning. 2019.
- [260] Elliot Meyerson and Risto Miikkulainen. Beyond shared hierarchies: Deep multitask learning through soft layer ordering. In *International Conference on Learning Representations*, 2018.
- [261] Junshi Huang, Rogerio S Feris, Qiang Chen, and Shuicheng Yan. Cross-domain image retrieval with a dual attribute-aware ranking network. In *Proceedings of the IEEE international conference on computer vision*, pages 1062–1070, 2015.
- [262] Brendan Jou and Shih-Fu Chang. Deep cross residual learning for multitask visual recognition. In *Proceedings of the 24th ACM international conference on Multimedia*, pages 998–1007. ACM, 2016.
- [263] Amir R. Zamir, Alexander Sax, William B. Shen, Leonidas J. Guibas, Jitendra Malik, and Silvio Savarese. Taskonomy: Disentangling task transfer learning. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2018.
- [264] Zhifei Zhang, Yang Song, and Hairong Qi. Age progression/regression by conditional adversarial autoencoder. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2017.
- [265] Mingsheng Long and Jianmin Wang. Learning multiple tasks with deep relationship networks. In *Advances in Neural Information Processing Systems*, 2017.
- [266] Ya Xue, Xuejun Liao, Lawrence Carin, and Balaji Krishnapuram. Multi-task learning for classification with dirichlet process priors. *Journal of Machine Learning Research*, 8(Jan):35–63, 2007.
- [267] Laurent Jacob, Jean philippe Vert, and Francis R. Bach. Clustered multi-task learning: A convex formulation. In *Advances in Neural Information Processing Systems 21*, 2009.

- [268] Zhuoliang Kang, Kristen Grauman, and Fei Sha. Learning with whom to share in multi-task feature learning. In *Proceedings of the 28th International Conference on International Conference on Machine Learning*, ICML'11, pages 521–528, USA, 2011. Omnipress.
- [269] Yongxi Lu, Abhishek Kumar, Shuangfei Zhai, Yu Cheng, Tara Javidi, and Rogério Schmidt Feris. Fully-adaptive feature sharing in multi-task networks with applications in person attribute classification. In *CVPR*, volume 1, page 6, 2017.
- [270] Youssef A Mejjati, Darren Cosker, and Kwang In Kim. Multi-task learning by maximizing statistical dependence. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3465–3473, 2018.
- [271] Yani Ioannou, Duncan Robertson, Roberto Cipolla, Antonio Criminisi, et al. Deep roots: Improving cnn efficiency with hierarchical filter groups. 2017.
- [272] Yarin Gal. Uncertainty in deep learning. *University of Cambridge*, 2016.
- [273] Chris J Maddison, Andriy Mnih, and Yee Whye Teh. The concrete distribution: A continuous relaxation of discrete random variables. *arXiv preprint arXiv:1611.00712*, 2016.
- [274] Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. *arXiv preprint arXiv:1611.01144*, 2016.
- [275] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition, 2014.
- [276] Wenqi Li, Guotai Wang, Lucas Fidon, Sébastien Ourselin, M. Jorge Cardoso, and Tom Vercauteren. On the compactness, efficiency, and representation of 3d convolutional networks: Brain parcellation as a pretext task. 2017.
- [277] Kerstin Kläser, Paweł Markiewicz, Marta Ranzini, Wenqi Li, Marc Modat, Brian F Hutton, David Atkinson, Kris Thielemans, M Jorge Cardoso, and Sébastien Ourselin. Deep boosted regression for mr to ct synthesis, 2018.
- [278] Zach Eaton-Rosen, Felix Bragman, Sotirios Bisdas, Sébastien Ourselin, and M. Jorge Cardoso. Towards safe deep learning: accurately quantifying biomarker uncertainty in neural network predictions, 2018.
- [279] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *3rd International Conference for Learning Representations*, 2015.
- [280] Nicholas J Tustison, Brian B Avants, Philip A Cook, Yuanjie Zheng, Alexander Egan, Paul A Yushkevich, and James C Gee. N4itk: Improved n3 bias correction. *IEEE Transactions on Medical Imaging*, 29(6):1310–1320, 2010.
- [281] Fabian Isensee, Jens Petersen, Andre Klein, David Zimmerer, Paul F. Jaeger, Simon Kohl, Jakob Wasserthal, Gregor Koehler, Tobias Norajitra, Sebastian Wirkert, and Klaus H. Maier-Hein. nnunet: Self-adapting framework for u-net-based medical image segmentation, 2018.
- [282] L.G. Nyul, J.K. Udupa, and Xuan Zhang. New variants of a method of MRI scale standardization. *IEEE Transactions on Medical Imaging*, 19(2):143–150, 2000.
- [283] Eli Gibson, Wenqi Li, Carole Sudre, Lucas Fidon, Dzhoshkun I. Shakir, Guotai Wang, Zach Eaton-Rosen, Robert Gray, Tom Doel, Yipeng Hu, Tom Whyntie, Parashkev Nachev, Marc Modat, Dean C. Barratt, Sébastien Ourselin, M. Jorge Cardoso, and Tom Vercauteren. NiftyNet: a deep-learning platform for medical imaging. *Computer Methods and Programs in Biomedicine*, 158:113–122, 2018.
- [284] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification, 2015.
- [285] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks, 2016.
- [286] Marc Harper. python-ternary: Ternary plots in python. In *10.5281/zenodo.34938*, 2015.
- [287] Jason Liang, Elliot Meyerson, and Risto Miikkulainen. Evolutionary architecture search for deep multitask networks. In *Proceedings of the Genetic and Evolutionary Computation Conference*, pages 466–473. ACM, 2018.

- [288] Alexandre Lacoste, Boris Oreshkin, Wonchang Chung, Thomas Boquet, Negar Rostamzadeh, and David Krueger. Uncertainty in multitask transfer learning. In *arXiv:1806.07528*, 2018.
- [289] Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In *European conference on computer vision*, pages 818–833. Springer, 2014.
- [290] Yoshua Bengio. Deep learning of representations: Looking forward. In *International Conference on Statistical Language and Speech Processing*, pages 1–37. Springer, 2013.
- [291] Barret Zoph and Quoc V Le. Neural architecture search with reinforcement learning. *ICLR*, 2017.
- [292] Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradients through stochastic neurons for conditional computation. *CoRR*, 2013.
- [293] Antonio Criminisi and Jamie Shotton. *Decision forests for computer vision and medical image analysis*. Springer, 2013.
- [294] Leo Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001.
- [295] Jerome H Friedman. Greedy function approximation: a gradient boosting machine. *Annals of statistics*, pages 1189–1232, 2001.
- [296] Vlad Sandulescu and Mihai Chiru. Predicting the future relevance of research institutions—the winning solution of the kdd cup 2016. *CoRR*, 2016.
- [297] Kaggle.com. Two sigma financial modeling challenge, 2017.
- [298] Maksims Volkovs, Guang Wei Yu, and Tomi Poutanen. Content-based neighbor models for cold start in recommender systems. In *Proceedings of the Recommender Systems Challenge 2017*, page 7. ACM, 2017.
- [299] Alberto Suárez and James F Lutsko. Globally optimal fuzzy decision trees for classification and regression. *IEEE Transactions. PAMI*, 21(12):1297–1311, 1999.
- [300] Ozan İrsoy, Olcay Taner Yıldız, and Ethem Alpaydın. Soft decision trees. In *ICPR*, pages 1819–1822. IEEE, 2012.
- [301] Dmitry Laptev and Joachim M Buhmann. Convolutional decision trees for feature learning and segmentation. In *German Conference on Pattern Recognition*, pages 95–106. Springer, 2014.
- [302] Samuel Rota Bulo and Peter Kortschieder. Neural decision forests for semantic image labelling. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 81–88, 2014.
- [303] Peter Kortschieder, Madalina Fiterau, Antonio Criminisi, and Samuel Rota Bulo. Deep neural decision forests. In *ICCV*, pages 1467–1475, 2015.
- [304] Nicholas Frosst and Geoffrey E Hinton. Distilling a neural network into a soft decision tree. *CoRR*, 2017.
- [305] Han Xiao. Ndt: Neual decision tree towards fully functioned neural graph. *arXiv preprint arXiv:1712.05934*, 2017.
- [306] Sethu Vijayakumar and Stefan Schaal. Locally weighted projection regression: An  $O(n)$  algorithm for incremental real time learning in high dimensional space. In *ICML*, volume 1, pages 288–293, 2000.
- [307] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [308] Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. 2009.
- [309] Zhi-Hua Zhou and Ji Feng. Deep forest: Towards an alternative to deep neural networks. In *IJCAI*, 2017.
- [310] Natalia Ponomareva, Thomas Colthurst, Gilbert Hendry, Salem Haykal, and Soroush Radpour. Compact multi-class boosted trees. In *International Conference on Big Data*, pages 47–56, 2017.
- [311] Michael I Jordan and Robert A Jacobs. Hierarchical mixtures of experts and the em algorithm. *Neural computation*, 1994.

- [312] Aurélia Léon and Ludovic Denoyer. Policy-gradient methods for decision trees. In *ESANN*, 2015.
- [313] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009.
- [314] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, Andrew Rabinovich, et al. Going deeper with convolutions. In *CVPR*, 2015.
- [315] Yani Ioannou, Duncan Robertson, Darko Zikic, Peter Kotschieder, Jamie Shotton, Matthew Brown, and Antonio Criminisi. Decision forests, convolutional networks and the models in-between. *CoRR*, 2016.
- [316] Ozan İrsoy, Olcay Taner Yıldız, and Ethem Alpaydin. Budding trees. In *ICPR*, pages 3582–3587. IEEE, 2014.
- [317] Jamie Shotton, Toby Sharp, Pushmeet Kohli, Sebastian Nowozin, John Winn, and Antonio Criminisi. Decision jungles: Compact and rich models for classification. In *Advances in Neural Information Processing Systems*, pages 234–242, 2013.
- [318] Andrew Davis and Itamar Arel. Low-rank approximations for conditional feedforward computation in deep neural networks. *CoRR*, 2013.
- [319] Emmanuel Bengio, Pierre-Luc Bacon, Joelle Pineau, and Doina Precup. Conditional computation in neural networks for faster models. *CoRR*, 2015.
- [320] Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc V. Le, Geoffrey E. Hinton, and Jeff Dean. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. *ICLR*, abs/1701.06538, 2017.
- [321] Amjad Almahairi, Nicolas Ballas, Tim Cooijmans, Yin Zheng, Hugo Larochelle, and Aaron Courville. Dynamic capacity networks. In *ICML*, 2016.
- [322] Surat Teerapittayanon, Bradley McDanel, and H. T. Kung. Branchynet: Fast inference via early exiting from deep neural networks. In *ICPR*, 2016.
- [323] Alex Graves. Adaptive computation time for recurrent neural networks. *CoRR*, abs/1603.08983, 2016.
- [324] Michael Figurnov, Maxwell D. Collins, Yukun Zhu, Li Zhang, Jonathan Huang, Dmitry P. Vetrov, and Ruslan Salakhutdinov. Spatially adaptive computation time for residual networks. *CVPR*, pages 1790–1799, 2017.
- [325] Andreas Veit and Serge Belongie. Convolutional networks with adaptive computation graphs. *CoRR*, 2017.
- [326] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [327] Gao Huang, Zhuang Liu, Kilian Q Weinberger, and Laurens van der Maaten. Densely connected convolutional networks. In *CVPR*, 2017.
- [328] Scott E Fahlman and Christian Lebiere. The cascade-correlation learning architecture. In *Advances in neural information processing systems*, pages 524–532, 1990.
- [329] Geoffrey E Hinton, Simon Osindero, and Yee-Whye Teh. A fast learning algorithm for deep belief nets. *Neural computation*, 18(7):1527–1554, 2006.
- [330] Tianjun Xiao, Jiaxing Zhang, Kuiyuan Yang, Yuxin Peng, and Zheng Zhang. Error-driven incremental learning in deep convolutional neural network for large-scale image classification. In *ACM Multimedia*, 2014.
- [331] Tianqi Chen, Ian Goodfellow, and Jonathon Shlens. Net2net: Accelerating learning via knowledge transfer. In *ICLR*, 2016.
- [332] Rupesh Kumar Srivastava, Klaus Greff, and Jürgen Schmidhuber. Highway networks. *arXiv preprint arXiv:1505.00387*, 2015.

- [333] Jeongtae Lee, Jaehong Yun, Sungju Hwang, and Eunho Yang. Lifelong learning with dynamically expandable networks. *CoRR*, 2017.
- [334] Han Cai, Tianyao Chen, Weinan Zhang, Yong Yu, and Jun Wang. Efficient architecture search by network transformation. In *AAAI*, 2018.
- [335] Ozan İrsoy and Ethem Alpaydin. Continuously constructive deep neural networks. *arXiv preprint arXiv:1804.02491*, 2018.
- [336] Andrew Brock, Theodore Lim, James M Ritchie, and Nick Weston. Smash: one-shot model architecture search through hypernetworks. *CoRR*, 2017.
- [337] Corinna Cortes, Xavier Gonzalvo, Vitaly Kuznetsov, Mehryar Mohri, and Scott Yang. Adanet: Adaptive structural learning of artificial neural networks. In *ICML*, pages 874–883, 2017.
- [338] Kenneth O Stanley and Risto Miikkulainen. Evolving neural networks through augmenting topologies. *Evolutionary computation*, 2002.
- [339] Esteban Real, Sherry Moore, Andrew Selle, Saurabh Saxena, Yutaka Leon Suematsu, Quoc Le, and Alex Kurakin. Large-scale evolution of image classifiers. *CoRR*, 2017.
- [340] Chenxi Liu, Barret Zoph, Jonathon Shlens, Wei Hua, Li-Jia Li, Li Fei-Fei, Alan Yuille, Jonathan Huang, and Kevin Murphy. Progressive neural architecture search. *CoRR*, 2017.
- [341] Peter Kontschieder, Pushmeet Kohli, Jamie Shotton, and Antonio Criminisi. Geof: Geodesic forests for learning coupled predictors. In *CVPR*, 2013.
- [342] Ishwar Krishnan Sethi. Entropy nets: From decision trees to neural networks. *Proceedings of the IEEE*, pages 1605–1613, 1990.
- [343] David L Richmond, Dagmar Kainmueller, Michael Yang, Eugene W Myers, and Carsten Rother. Mapping stacked decision forests to deep and sparse convolutional neural networks for semantic segmentation. *CoRR*, 2015.
- [344] Ji Feng, Yang Yu, and Zhi-Hua Zhou. Multi-layered gradient boosting decision trees. In *Advances in Neural Information Processing Systems*, pages 3551–3561, 2018.
- [345] Dong-Hyun Lee, Saizheng Zhang, Asja Fischer, and Yoshua Bengio. Difference target propagation. In *Joint european conference on machine learning and knowledge discovery in databases*, pages 498–515. Springer, 2015.
- [346] Vinod Nair and Geoffrey E Hinton. Rectified linear units improve restricted boltzmann machines. In *ICML*, pages 807–814, 2010.
- [347] David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning representations by back-propagating errors. *Nature*, 323(6088):533, 1986.
- [348] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017.
- [349] Ian J Goodfellow, David Warde-Farley, Mehdi Mirza, Aaron Courville, and Yoshua Bengio. Maxout networks. In *ICML*, 2013.
- [350] Min Lin, Qiang Chen, and Shuicheng Yan. Network in network. In *ICLR*, 2014.
- [351] Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, and Martin Riedmiller. Striving for simplicity: The all convolutional net. *CoRR*, 2015.
- [352] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. Scikit-learn: Machine learning in python. *JMLR*, 12(Oct):2825–2830, 2011.
- [353] Han Zhao, Otilia Stretcu, Renato Negrinho, Alex Smola, and Geoff Gordon. Efficient multi-task feature and relationship learning. *arXiv preprint arXiv:1702.04423*, 2017.
- [354] Samuel Schulter, Paul Wohlhart, Christian Leistner, Amir Saffari, Peter M Roth, and Horst Bischof. Alternating decision forests. In *CVPR, 2013*, 2013.

- [355] Patrice Y Simard, David Steinkraus, John C Platt, et al. Best practices for convolutional neural networks applied to visual document analysis. In *ICDAR*, volume 3, pages 958–962, 2003.
- [356] Sara Sabour, Nicholas Frosst, and Geoffrey E Hinton. Dynamic routing between capsules. In *Advances in Neural Information Processing Systems*, pages 3856–3866, 2017.
- [357] Jiaqi Guan, Yang Liu, Qiang Liu, and Jian Peng. Energy-efficient amortized inference with cascaded deep classifiers. *arXiv preprint arXiv:1710.03368*, 2017.
- [358] Gao Huang, Shichen Liu, Laurens Van der Maaten, and Kilian Q Weinberger. Condensenet: An efficient densenet using learned group convolutions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2752–2761, 2018.
- [359] Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and new perspectives. *IEEE Trans. Pattern Anal. Mach. Intell.*, 35(8):1798–1828, 2013.
- [360] George Tucker, Andriy Mnih, Chris J Maddison, John Lawson, and Jascha Sohl-Dickstein. Rebar: Low-variance, unbiased gradient estimates for discrete latent variable models. In *Advances in Neural Information Processing Systems*, pages 2627–2636, 2017.
- [361] Will Grathwohl, Dami Choi, Yuhuai Wu, Geoffrey Roeder, and David Duvenaud. Backpropagation through the void: Optimizing control variates for black-box gradient estimation. *arXiv preprint arXiv:1711.00123*, 2017.
- [362] Riccardo Guidotti, Anna Monreale, Salvatore Ruggieri, Franco Turini, Fosca Giannotti, and Dino Pedreschi. A survey of methods for explaining black box models. *ACM computing surveys (CSUR)*, 51(5):1–42, 2018.
- [363] James Surowiecki. *The wisdom of crowds*. Anchor, 2005.
- [364] Takeyuki Watadani, Fumikazu Sakai, Takeshi Johkoh, Satoshi Noma, Masanori Akira, Kiminori Fujimoto, Alexander A Bankier, Kyung Soo Lee, Nestor L Müller, Jae-Woo Song, et al. Interobserver variability in the ct assessment of honeycombing in the lungs. *Radiology*, 266(3):936–944, 2013.
- [365] Andrew B Rosenkrantz, Ruth P Lim, Mershad Haghghi, Molly B Somberg, James S Babb, and Samir S Taneja. Comparison of interreader reproducibility of the prostate imaging reporting and data system and likert scales for evaluation of multiparametric prostate mri. *American Journal of Roentgenology*, 201(4):W612–W618, 2013.
- [366] Elizabeth Lazarus, Martha B Mainiero, Barbara Schepps, Susan L Koelliker, and Linda S Livingston. Bi-rads lexicon for us and mammography: interobserver variability and positive predictive value. *Radiology*, 239(2):385–391, 2006.
- [367] Bjoern H Menze, Andras Jakab, Stefan Bauer, Jayashree Kalpathy-Cramer, Keyvan Farahani, Justin Kirby, Yuliya Burren, Nicole Porz, Johannes Slotboom, Roland Wiest, et al. The multimodal brain tumor image segmentation benchmark (brats). *IEEE transactions on medical imaging*, 34(10):1993–2024, 2014.
- [368] Hugh Harvey and Ben Glocker. A standardised approach for preparing imaging data for machine learning tasks in radiology. In *Artificial Intelligence in Medical Imaging*, pages 61–72. Springer, 2019.
- [369] Peter Welinder and Pietro Perona. Online crowdsourcing: rating annotators and obtaining cost-effective labels. In *Computer Vision and Pattern Recognition Workshops (CVPRW), 2010 IEEE Computer Society Conference on*, pages 25–32. IEEE, 2010.
- [370] Chengjiang Long, Gang Hua, and Ashish Kapoor. Active visual recognition with expertise estimation in crowdsourcing. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3000–3007, 2013.
- [371] Chengjiang Long and Gang Hua. Multi-class multi-annotator active learning with robust gaussian process for visual recognition. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2839–2847, 2015.
- [372] Alexander Philip Dawid and Allan M Skene. Maximum likelihood estimation of observer error-rates using the em algorithm. *Applied statistics*, pages 20–28, 1979.

- [373] Padhraic Smyth, Usama M Fayyad, Michael C Burl, Pietro Perona, and Pierre Baldi. Inferring ground truth from subjective labelling of venus images. In *Advances in neural information processing systems*, pages 1085–1092, 1995.
- [374] Jacob Whitehill, Ting-fan Wu, Jacob Bergsma, Javier R Movellan, and Paul L Ruvolo. Whose vote should count more: Optimal integration of labels from labelers of unknown expertise. In *Advances in neural information processing systems*, pages 2035–2043, 2009.
- [375] Peter Welinder, Steve Branson, Pietro Perona, and Serge J Belongie. The multidimensional wisdom of crowds. In *Advances in neural information processing systems*, pages 2424–2432, 2010.
- [376] Filipe Rodrigues, Francisco Pereira, and Bernardete Ribeiro. Learning from multiple annotators: distinguishing good from random labelers. *Pattern Recognition Letters*, 34(12):1428–1436, 2013.
- [377] Vikas C Raykar, Shipeng Yu, Linda H Zhao, Anna Jerebko, Charles Florin, Gerardo Hermosillo Valadez, Luca Bogoni, and Linda Moy. Supervised learning from multiple experts: whom to trust when everyone lies a bit. In *Proceedings of the 26th Annual international conference on machine learning*, pages 889–896. ACM, 2009.
- [378] Yan Yan, Rómer Rosales, Glenn Fung, Mark Schmidt, Gerardo Hermosillo, Luca Bogoni, Linda Moy, and Jennifer Dy. Modeling annotator expertise: Learning when everybody knows a bit of something. In *AISTATS*, pages 932–939, 2010.
- [379] Steve Branson, Grant Van Horn, and Pietro Perona. Lean crowdsourcing: Combining humans and machines in an online system. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7474–7483, 2017.
- [380] Grant Van Horn, Steve Branson, Scott Loarie, Serge Belongie, Cornell Tech, and Pietro Perona. Lean multiclass crowdsourcing. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2714–2723, 2018.
- [381] Ashish Khetan, Zachary C Lipton, and Anima Anandkumar. Learning from noisy singly-labeled data. In *International Conference on Learning Representations*, 2018.
- [382] Scott Reed, Honglak Lee, Dragomir Anguelov, Christian Szegedy, Dumitru Erhan, and Andrew Rabinovich. Training deep neural networks on noisy labels with bootstrapping. *arXiv preprint arXiv:1412.6596*, 2014.
- [383] Sainbayar Sukhbaatar, Joan Bruna, Manohar Paluri, Lubomir Bourdev, and Rob Fergus. Training convolutional networks with noisy labels. *arXiv preprint arXiv:1406.2080*, 2014.
- [384] Jacob Goldberger and Ehud Ben-Reuven. Training deep neural-networks using a noise adaptation layer. In *International Conference on Learning Representations*, 2017.
- [385] Melody Y Guan, Varun Gulshan, Andrew M Dai, and Geoffrey E Hinton. Who said what: Modeling individual labelers improves classification. *AAAI*, 2018.
- [386] Benoît Frénay and Michel Verleysen. Classification in the presence of label noise: a survey. *IEEE transactions on neural networks and learning systems*, 25(5):845–869, 2014.
- [387] Nagarajan Natarajan, Inderjit S Dhillon, Pradeep K Ravikumar, and Ambuj Tewari. Learning with noisy labels. In *Advances in neural information processing systems*, pages 1196–1204, 2013.
- [388] Jakramate Bootkrajang and Ata Kabán. Label-noise robust logistic regression and its applications. In *Joint European conference on machine learning and knowledge discovery in databases*, pages 143–158. Springer, 2012.
- [389] Daiki Tanaka, Daiki Ikami, Toshihiko Yamasaki, and Kiyoharu Aizawa. Joint optimization framework for learning with noisy labels. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2018.
- [390] Volodymyr Mnih and Geoffrey E Hinton. Learning to label aerial images from noisy data. In *International conference on machine learning*, pages 567–574, 2012.
- [391] Giorgio Patrini, Alessandro Rozza, Aditya Krishna Menon, Richard Nock, and Lizhen Qu. Making deep neural networks robust to label noise: A loss correction approach. In *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.(CVPR)*, pages 2233–2241, 2017.

- [392] Jacob Goldberger and Ehud Ben-Reuven. Training deep neural-networks using a noise adaptation layer. In *ICLR*, 2017.
- [393] Andreas Veit, Neil Alldrin, Gal Chechik, Ivan Krasin, Abhinav Gupta, and Serge J Belongie. Learning from noisy large-scale datasets with minimal supervision. In *CVPR*, pages 6575–6583, 2017.
- [394] Mengye Ren, Wenyuan Zeng, Bin Yang, and Raquel Urtasun. Learning to reweight examples for robust deep learning. In *International Conference on Machine Learning*, 2018.
- [395] Lu Jiang, Zhengyuan Zhou, Thomas Leung, Li-Jia Li, and Li Fei-Fei. Mentornet: Learning data-driven curriculum for very deep neural networks on corrupted labels. In *International Conference on Machine Learning*, pages 2309–2318, 2018.
- [396] Furong Huang, Jordan T. Ash, John Langford, and Robert E. Schapire. Learning deep resnet blocks sequentially using boosting theory. *CoRR*, 2017.
- [397] Ardavan Saeedi, Matthew D Hoffman, Stephen J DiVerdi, Asma Ghandeharioun, Matthew J Johnson, and Ryan P Adams. Multimodal prediction and personalization of photo edits with deep generative models. *arXiv preprint arXiv:1704.04997*, 2017.
- [398] Tong Xiao, Tian Xia, Yi Yang, Chang Huang, and Xiaogang Wang. Learning from massive noisy labeled data for image classification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2691–2699, 2015.
- [399] Yuan Li, Benjamin Rubinstein, and Trevor Cohn. Exploiting worker correlation for label aggregation in crowdsourcing. In *International Conference on Machine Learning*, pages 3886–3895, 2019.
- [400] Huahong Zhang, Alessandra M Valcarcel, Rohit Bakshi, Renxin Chu, Francesca Bagnato, Russell T Shinohara, Kilian Hett, and Ipek Oguz. Multiple sclerosis lesion segmentation with tiramisu and 2.5 d stacked slices. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 338–346. Springer, 2019.
- [401] Eytan Kats, Jacob Goldberger, and Hayit Greenspan. A soft staple algorithm combined with anatomical knowledge. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 510–517. Springer, 2019.
- [402] Stefan Winzeck, Arsany Hakim, Richard McKinley, Jos AADSR Pinto, Victor Alves, Carlos Silva, Maxim Pisov, Egor Krivov, Mikhail Belyaev, Miguel Monteiro, et al. Isles 2016 and 2017-benchmarking ischemic stroke lesion outcome prediction based on multispectral mri. *Frontiers in neurology*, 9:679, 2018.
- [403] Olivier Commowick, Audrey Istace, Michael Kain, Baptiste Laurent, Florent Leray, Mathieu Simon, Sorina Camarasu Pop, Pascal Girard, Roxana Ameli, Jean-Christophe Ferré, et al. Objective evaluation of multiple sclerosis lesion segmentation using a data management and processing infrastructure. *Scientific reports*, 8(1):1–17, 2018.
- [404] Gleason 2019 challenge. <https://gleason2019.grand-challenge.org/Home/>. Accessed: 2020-02-30.
- [405] Andrew J Asman and Bennett A Landman. Formulating spatially varying performance in the statistical fusion framework. *IEEE transactions on medical imaging*, 31(6):1326–1336, 2012.
- [406] Juan Eugenio Iglesias, Mert Rory Sabuncu, and Koen Van Leemput. A unified framework for cross-modality multi-atlas segmentation of brain mri. *Medical image analysis*, 17(8):1181–1191, 2013.
- [407] M Jorge Cardoso, Kelvin Leung, Marc Modat, Shiva Keihaninejad, David Cash, Josephine Barnes, Nick C Fox, Sebastien Ourselin, Alzheimer’s Disease Neuroimaging Initiative, et al. Steps: Similarity and truth estimation for propagated segmentations and its application to hippocampal segmentation and brain parcelation. *Medical image analysis*, 17(6):671–684, 2013.
- [408] Andrew J Asman and Bennett A Landman. Non-local statistical label fusion for multi-atlas segmentation. *Medical image analysis*, 17(2):194–208, 2013.

- [409] Alireza Akhondi-Asl, Lennox Hoyte, Mark E Lockhart, and Simon K Warfield. A logarithmic opinion pool based staple algorithm for the fusion of segmentations with associated reliability weights. *IEEE transactions on medical imaging*, 33(10):1997–2009, 2014.
- [410] Daniel C. Castro, Jeremy Tan, Bernhard Kainz, Ender Konukoglu, and Ben Glocker. Morpho-MNIST: Quantitative assessment and diagnostics for representation learning. *Journal of Machine Learning Research*, 20, 2019.
- [411] Martin Styner, Joohwi Lee, Brian Chin, M Chin, Olivier Commowick, H Tran, S Markovic-Plese, V Jewells, and S Warfield. 3d segmentation in the clinic: A grand challenge ii: Ms lesion segmentation. *Midas Journal*, 2008:1–6, 2008.
- [412] Samuel G Armato III, Geoffrey McLennan, Luc Bidaut, Michael F McNitt-Gray, Charles R Meyer, Anthony P Reeves, Binsheng Zhao, Denise R Aberle, Claudia I Henschke, Eric A Hoffman, et al. 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(2):915–931, 2011.
- [413] Neil I Weisenfeld and Simon K Warfield. Learning likelihoods for labeling (l3): a general multi-classifier segmentation algorithm. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 322–329. Springer, 2011.
- [414] Christian F Baumgartner, Kerem C Tezcan, Krishna Chaitanya, Andreas M Hötker, Urs J Muehlematter, Khoshy Schawkat, Anton S Becker, Olivio Donati, and Ender Konukoglu. Phiseg: Capturing uncertainty in medical image segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 119–127. Springer, 2019.
- [415] Carole H Sudre, Beatriz Gomez Anson, Silvia Ingala, Chris D Lane, Daniel Jimenez, Lukas Haider, Thomas Varsavsky, Ryutaro Tanno, Lorna Smith, Sébastien Ourselin, et al. Let’s agree to disagree: Learning highly debatable multirater labelling. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 665–673. Springer, 2019.
- [416] Andrew Jesson and Tal Arbel. Hierarchical mrf and random forest segmentation of ms lesions and healthy tissues in brain mri. *Proceedings of the 2015 Longitudinal Multiple Sclerosis Lesion Segmentation Challenge*, pages 1–2, 2015.
- [417] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015.
- [418] Javier Antorán, Umang Bhatt, Tameem Adel, Adrian Weller, and José Miguel Hernández-Lobato. Getting a clue: A method for explaining uncertainty estimates. *arXiv preprint arXiv:2006.06848*, 2020.
- [419] Petru-Daniel Tudosi, Thomas Varsavsky, Richard Shaw, Mark Graham, Parashkev Nachev, Sébastien Ourselin, Carole H Sudre, and M Jorge Cardoso. Neuromorphologicaly-preserving volumetric data encoding using vq-vae. *arXiv preprint arXiv:2002.05692*, 2020.
- [420] Carl Edward Rasmussen and Zoubin Ghahramani. Occam’s razor. In *Advances in neural information processing systems*, pages 294–300, 2001.
- [421] Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. *arXiv preprint arXiv:1803.03635*, 2018.
- [422] Pavlo Molchanov, Stephen Tyree, Tero Karras, Timo Aila, and Jan Kautz. Pruning convolutional neural networks for resource efficient inference. *arXiv preprint arXiv:1611.06440*, 2016.
- [423] Michael Zhu and Suyog Gupta. To prune, or not to prune: exploring the efficacy of pruning for model compression. *arXiv preprint arXiv:1710.01878*, 2017.
- [424] Yan Yan, Romer Rosales, Glenn Fung, and Jennifer G Dy. Active learning from crowds. In *ICML*, 2011.