



Fast MRI registration leveraging k-space properties

*Schnelle MRT-Registrierung mittels k-Raum Eigenschaften*

**Masterarbeit**

verfasst am

**Institut für Medizinische Informatik**

im Rahmen des Studiengangs

**Medizinische Informatik**

der Universität zu Lübeck

vorgelegt von

**Jan Meyer**

ausgegeben und betreut von

**Prof. Dr. Mattias Heinrich**

mit Unterstützung von

**Ziad Al-Haj Hemidi, Eytan Kats**

Lübeck, den 1. Januar 2025

**Eidesstattliche Erklärung**

*Ich erkläre hiermit an Eides statt, dass ich diese Arbeit selbständig verfasst und keine anderen als die angegebenen Quellen und Hilfsmittel benutzt habe.*

---

Jan Meyer

## Zusammenfassung

Irgendwas über MRT und Registration

## Abstract

Something about MRI and registration

## Acknowledgements

Danke an Mattias, Ziad und Eytan für die gute Betreuung.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Contributions of the Thesis	1
1.2	Related Work	1
1.3	Structure of the Thesis	2
<b>2</b>	<b>Basics</b>	<b>3</b>
2.1	Magnetic Resonance Imaging	3
2.2	Image Registration	11
2.3	Deep Learning Architectures	13
2.4	Deep Learning for Image Registration	14
2.5	Motion-Corrected Reconstruction	17
<b>3</b>	<b>Data</b>	<b>19</b>
3.1	OASIS Dataset	19
3.2	CMRxRecon Dataset	19
3.3	ACDC Dataset	22
<b>4</b>	<b>Network Architectures</b>	<b>24</b>
4.1	Fourier-Net	24
4.2	Fourier Net+	31
<b>5</b>	<b>Experiments</b>	<b>34</b>
5.1	Parameter Studies on the CMRxRecon Dataset	34
5.2	Parameter Tests on the ACDC Dataset	35
5.3	Data Consistency Experiments and Network Changes	38
5.4	Integration into a Motion-Corrected Reconstruction Pipeline	39
<b>6</b>	<b>Results and Discussion</b>	<b>40</b>
6.1	Parameter Studies on the CMRxRecon Dataset	40
6.2	Further Tests on the ACDC Dataset	42
<b>7</b>	<b>Conclusion</b>	<b>56</b>



# 1

## Introduction

MRI is a commonly used medical imaging technique that is non-invasive, radiation-free and has great contrast for soft tissue. However, speed is its main weakness as full-body scan can take up to 30 minutes, which is a burden on patients as well as hindering efficiency. Thus, MRI acquisition is usually accelerated by subsampling the k-space in which the data is recorded. This can, however, lead to problematic image artifacts that hinder many processing algorithms. Another problem of the slow data acquisition are motion artifacts, which are especially common for organs like lung and heart. While breathing can be controlled, cardiac motion is involuntarily and should not be stopped. These challenges have been traditionally addressed with computationally intensive algorithms that usually iteratively solve an optimization problem. New approaches, however, are based around the deep neural networks, which have seen a rise in popularity in the recent years in many fields including image processing. In this thesis, the possibility of using an unsupervised deep learning approach to align subsampled MRI data as well as its potential usage in a motion reconstruction pipeline will be investigated.

### 1.1 Contributions of the Thesis

We managed to implement an unsupervised approach to correct motion in subsampled MRI images and much more.

### 1.2 Related Work

There are many papers on image registration in general, however in the context of medical image registration with deep learning their number is reduced due to the specialized nature of the subject at hand. Yet, there are a couple of papers that give a good overview of the topic. They give a brief overview of registration methods, basics of deep learning with already existing networks for image registration as well as covering potential applications and challenges [1, 2, 3, 4, 5]. A lot of different approaches for medical image registration are based on *VoxelMorph* [6], such as both *Fourier-Net* [7] and its successor *Fourier-Net+* [8].

### 1.3 Structure of the Thesis

In chapter 2, the general basis of the thesis containing descriptions of MRI in section 2.1, Image registration in section 2.2 and deep learning architectures as well as their use in image registration in sections 2.3 and 2.4. A special emphasize is put on motion reconstruction in MRI as described in section ??.

The data used in this thesis is described in Chapter 3 with details concerning the *OASIS* and *CMRxRecon* datasets in sections 3.1 and 3.2 respectively.

In chapter 4, the networks used in this thesis, namely *Fourier-Net* and *Fourier-Net+*, are described in section 4.1 and 4.2.

The experiments conducted in the making of this thesis are explained in chapter 5 with results and discussion following in chapter 6.

Lastly, the thesis is briefly summarized in chapter 7.

# 2

## Basics

In this chapter the basics of the thesis are explained.

### 2.1 Magnetic Resonance Imaging

Magnetic Resonance Imaging (MRI) is a non-invasive, radiation-free, tomographic imaging technology based on measurements of a magnetic field. An MRI machine comprises four main components as seen in Figure 2.1. The first component is a strong magnet powerful enough to generate a static magnetic field  $B_0$  that is required to induce nuclear proton polarization. The second is a radio frequency (RF) system which generates an alternating magnetic field  $B_1$  at the resonant frequency  $f$  and detects the MR signal that is returned from the patient. The third component is the set of gradient systems (oriented orthogonally in X, Y and Z directions) that generates linear magnetic field variations, which are then superimposed upon  $B_0$  and are used to spatially encode the MR signal. In clinical MR scanners, the three gradient sets and whole-body RF coils are typically concentrically positioned inside the bore of the magnet. The fourth component is a computer providing the user interface, generating images to be displayed and interpreted on the console [9].

#### Physical Foundation

Individual nuclei precess around the field  $B_0$  at a resonance frequency known as the Larmor frequency  $\omega_0$  of the net magnetization vector, which generally occurs in the RF range of the electromagnetic spectrum and is related to the external magnetic field as:

$$\omega_0 = \gamma \cdot B_0, \quad (2.1)$$

with  $\gamma$  being the gyromagnetic ratio, which is a fixed value depending on the nuclei [10]. When nuclei are placed in the presence of a strong static magnetic field such as  $B_0$  the nuclei split into two energy states, either aligned parallel to the magnetic field  $B_0$  (called a spin-up state) or aligned anti-parallel to  $B_0$  (called a spin-down state). The spin-up state has a slightly lower energy level as compared to the spin-down state and is therefore preferred. This slight difference in the spin states (0.001%) results in an overall net mag-

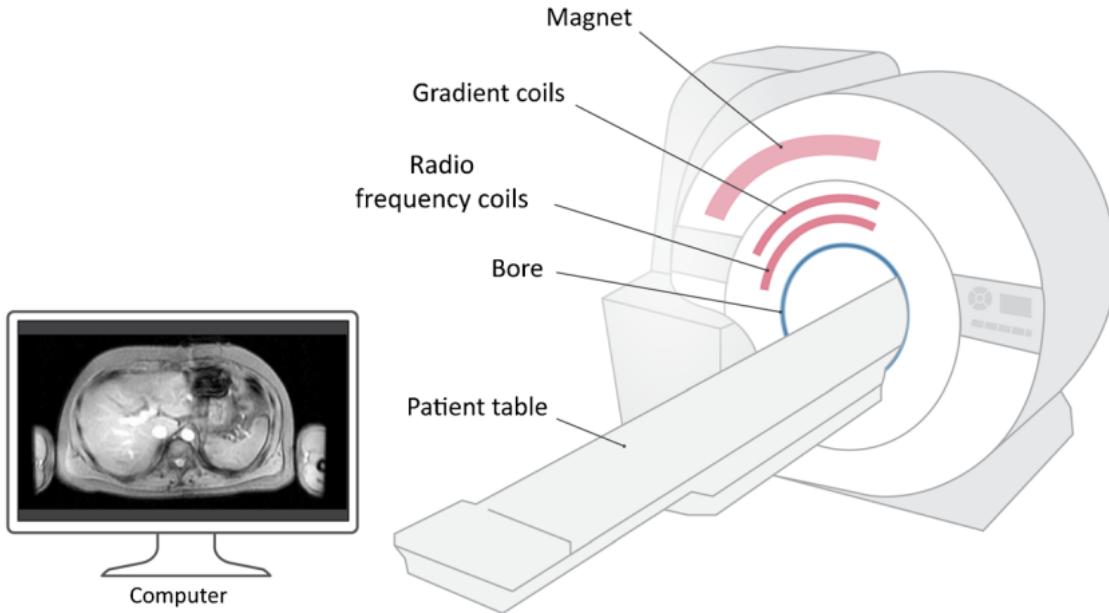


Figure 2.1: Schematic of an MRI scanner taken from [9].

netization  $M$  aligned in the same direction as  $B_0$ . To create an MR signal, the spins are excited out of their resting equilibrium, i.e. tipping  $M$  away from  $B_0$ . To detect the signal from the hydrogen nuclei in the tissues an additional external field  $B_1$  is introduced at the resonant Larmor frequency  $\omega_0$  that can affect magnetization vector, causing it to rotate into a plane orthogonal to its original orientation. The rotated vector continues to precess around the  $B_0$ . The precession of the magnetization vector in the transverse plane can be detected by an RF coil tuned to the resonant frequency  $\omega_0$ . RF coils can be operated in a receive-only mode, in which case the inherent body coil is used as a transmitter; or the RF coils can be both transmit and receive. The purpose of the RF transmit coil is to create a time-varying  $B_1$  field at right angles to  $B_0$  that could be linearly or circularly polarized. The closer the receiving coil is to the source of the MR signal, the better the signal-to-noise ratio (SNR). Receiver coils generally comprise arrays of smaller individual coils or elements; however, each individual coil element has a limited depth penetration. Multiple arrays of coils, termed phased-array coils, can be used together to achieve a higher coverage. The multiple coils are electronically decoupled from one another so that they do not appear as just a single large coil. The images from individual coils are independently reconstructed and then grouped together to create the final image. An RF pulse of amplitude  $B_1$ , called excitation pulse, is applied for a certain time duration to tip the magnetization at an angle away from the  $B_0$  field. The precessing transverse magnetization induces a voltage in the receiver RF coil; this induced voltage is known as the free induction decay. After the pulse, the magnetization returns to thermal equilibrium by processes known as MR relaxation. To fully encode the spatial information within the field of view (FOV), pulse sequences must be iterated numerous times. The time between successive iterations of a pulse sequence is known as the repetition time (TR). The time between the application of the initial RF pulse and the middle of the detected echo is known as the echo time (TE).

Due to this the overall time to acquire an MR image is quite long [9]. Once the RF pulse is turned off,  $M$  continues to precess as it returns to its thermal equilibrium state. During this time, two types of relaxation occur: T1 (longitudinal or spin lattice) and T2 (transverse or spin/spin). One attribute that makes T1 and T2 so valuable for determining the signal in MRI is their sensitivity to the presence and type of tissue. It is this tissue dependence property that gives MR its excellent soft-tissue contrast. T1 relaxation describes the recovery of the longitudinal magnetization back to thermal equilibrium following a perturbation by an RF pulse. The longitudinal component regrows along the Z direction with a time constant T1. In other words, after the RF pulse is turned off, the protons that were disturbed give their energy to the surrounding environment and return back to their original equilibrium state, realigning with  $B_0$ . Hence T1 relaxation is also called longitudinal relaxation. Furthermore, because T1 relaxation involves the loss of energy that was put into the spin system by the RF pulse, it is also referred to as spin-lattice relaxation, the lattice consisting of surrounding macromolecules. This loss of energy is stimulated by the fluctuating magnetic fields associated with the dipole–dipole interactions of neighboring magnetic moments. T1 relaxation can only occur when these magnetic field fluctuations occur at the resonant frequency  $\omega_0$ . The rate at which the spin magnetization  $M_z$  recovers to  $M_0$  at time  $t$  is called T1 relaxation time. It can be expressed as follows:

$$M_z = M_0 \cdot \left(1 - e^{-\frac{t}{T1}}\right). \quad (2.2)$$

A preferred method of measuring T1 relaxation is the Look-Locker method, which is based on the principle that one does not need to wait for the net magnetization vector to equilibrate in order to measure T1. Instead, an RF pulse is used with a small flip, which can be repeatedly applied. The acquisition pulse sequence is designed to generate a train of signals that gradually approaches a steady-state recovery. The recovery curve measured by this technique can be fitted to an exponential curve to provide an effective T1 measurement. This has become one of the most popular T1-mapping method for abdomen and cardiac imaging. The general principle of T1 mapping is to acquire multiple images with different T1 weightings and to fit the signal intensities of the images to the equation of T1 relaxation. T1 relaxation values for tissues can be estimated by fitting the data to the following equation:

$$S = S_0 \cdot \left(1 - A \cdot e^{\frac{T1}{T1}}\right), \quad (2.3)$$

where  $S$  is the signal intensity measured at each inversion time value  $T1$ ,  $S_0$  is the initial signal intensity at time  $t = 0$  and  $A$  is a constant.

T2 relaxation results in the loss of transverse magnetization caused by interactions between the magnetic fields of neighboring hydrogen nuclei. It is not an energy loss process like T1 but is a loss of phase coherence within the spin system. This process, also known as spin-spin relaxation, leads to the destruction of transverse magnetization and causes the magnetic moments of the tissue to dephase. T2 mapping is a method of measuring the T2 value of the tissue. T2 relaxation time can be calculated using a T2 sequence with different echo times (TEs). The most fundamental sequence for T2 mapping is signal measured

with spinecho techniques (multiple sequences with different TE values) [19]. Other 2-D sequences have been used, such as multi-echo spin echo (MSME) and fast spin echo (FSE) [19]. Synthetic MRI is another quantitative method in which a single saturation recovery turbo spin-echo sequence is used to estimate T2 transverse relaxation [20]. T2 relaxation values for tissues can be estimated by fitting the data to the following equation using a mono-exponential decay curve:

$$S(TE) = S_0 \cdot e^{\frac{-TE}{T2}}, \quad (2.4)$$

where  $S(TE)$  is the signal intensity measured at each TE and  $S_0$  is the initial signal intensity at time  $t = 0$ .

### Image Acquisition and K-Space

Going back to the Larmor frequency (Equation 2.1), the position of spins can be determined if they are brought into a magnetic field, which varies linearly across the object as seen in Figure 2.2. Their position can then be derived by frequency analysis of the measured MR signals, which is done by the Fourier transformation (FT). The generalization of this principle to a 2D image is less trivial as the simultaneous application of magnetic field gradients will render an unambiguous assignment only for spins aligned along the direction of the combined gradient. All spins orthogonal to this direction will still see the same field and thus yield identical Larmor frequencies. Thus, even with arbitrarily shaped magnetic field profiles there will always be an infinite number of locations with identical field strength. A direct measurement of a 2D image from an intrinsically 1D signal is therefore not possible [10].

The solution to the problem of encoding a 2D image with an 1D MR signal is using time-variant magnetic fields. If the magnetic field at each point is varied in time in a specific manner, the location of spins can be established as long as the temporal variation is unambiguous for each position. In principle, there are of course an infinite number of ways in which time-varying magnetic fields could be used to encode for each point in the image plane. A very powerful tool to design such imaging strategies is offered by the so-called k-space.

The concept of the k-space is a generalization of the simple relation of a time-variant signal to a spectrum of two or more dimensions. An 2D image is related to a 2D k-space data set by a 2D FT, as seen in Figure 2.3, same as a 1D spectrum is related to the 1D signal by a normal (1D) FT. A 2D FT can be achieved by applying the (1D) FT successively row-by-row and column-by-column. As the FT is an information-preserving operation the k-space data contains exactly the same information as the image data. Thus, in order to get the full image information, the full k-space data has to be measured. The task of forming an image by this approach is then converted to the task of finding a way to measure the necessary corresponding k-space data. An understanding of the physical meaning of the k-space domain can be derived from Figure 2.2. The time-domain signal defines the 1D k-space. Its coordinates are defined by the area under the gradient which the spins have experienced prior to the collection of each k-space data-point. Using the Larmor frequency again from Equation 2.1 the coordinates in 2D k-space will thus be given as  $k_x = \gamma G_x t$  and  $k_y = \gamma G_y t$ .

In order to acquire an MR image, methods need to be developed in order to get from one k-space point to the next. There are only two possibilities as seen in Figure 2.4. The first one is to apply a gradient, in which case the k-space trajectory will be a line defined by the orientation of the gradient (see Figure 2.4a). As long as the gradient is kept constant, the k-space trajectory will be a straight line. This way it is easy to see why simple superposition of two gradients will not produce an image. The resulting k-space trajectory will merely represent a straight line traversing k-space at some skewed angle.

The second method for moving in the k-space is the application of a refocusing pulse, which will lead to a jump of the k-space trajectory around the origin (see Figure 2.4b). It is clear that a spin-echo formation alone will not allow coverage of all of the k-space, since the trajectory would only jump between the two mirror symmetric points. A combination of spin echoes and gradients can, however, produce very efficient k-space trajectories. All of the existing k-space-based imaging techniques are based on some combination of these two basic means of moving through the k-space.

It should also be noted that the signals are measured not continuously, but at discrete time intervals called the dwell-time of data acquisition. This discrete sampling leads to an ambiguous assignment of frequencies above a given threshold which is called the Nyquist frequency. The Nyquist frequency therefore determines the acquisition bandwidth inside which the signal should occur. The definition of the k-space coordinates implies that these are invariant with respect to the actual strength of the gradient used, as long as  $G$  and  $t$  are constant. Mathematically, data acquisition under a strong gradient and in a shorter acquisition time will yield the same k-space data and thus the same image as acquisition under a weaker gradient and a longer acquisition time. A shorter acquisition time and a closer spacing of sampling points are equivalent to a higher-acquisition bandwidth. Since the received noise grows with the square root of the bandwidth, a faster imaging technique will therefore by principle always carry the penalty of a lower signal-to-noise ratio (SNR), even if all other factors influencing high-speed imaging are neglected.

For a conventional spin-echo imaging sequence, which is still widely used, a typical acquisition time of 5 to 10 ms for each phase-encoding step is used. Acquisition of typically 192 to 256 phase-encoding steps therefore requires a total net acquisition time of 1 to 2 s. This is much shorter than the total acquisition time of such a sequence, which is determined by the repetition time defined by the T1-contrast. From the aforementioned general considerations it immediately follows that the SNR of a fast imaging sequence leading to acquisition times of approximately 50 ms, will be lower by at least a factor of 5 to 10, regardless of the actual sequence used. Fast imaging therefore requires improved SNR provided by radio-frequency coil designs. The improvement in SNR driven by the need for fast imaging will also benefit slower techniques, providing images with better spatial resolution and/or higher SNR.

The discrete data sampling has some additional important consequences regarding the sampling density and coverage of data in k-space. Different parts of the k-space encode different features of the image: The center of the k-space hold the lower frequencies which represent the image contrast, whereas the outer parts encode the higher frequencies for sharp structures. Due to the symmetry of the 2D FT, this statement can also be reversed: The image center will be encoded by low-resolution k-space data. Therefore,

sampling of sparsely distributed k-space data will reduce the effective field of view of the final image.

Mathematically, the final image will look exactly the same, irrespective of the way data are sampled in the k-space. In practice, the different approaches to sampling can have a major impact on image quality. This is due to the fact that sampling data has to proceed sequentially. The observed spins evolve not only as a function of the gradients defining the intended k-space trajectory, but are also influenced by other mechanisms unrelated to image encoding. It is noteworthy that the coordinates in k-space are expressed in units of a phase angle across space. Any mechanism affecting the phase of the signal along this trajectory will therefore alter the k-space trajectory. Some such commonly occurring mechanisms are flow and motion effects, magnetic field inhomogeneities, as well as susceptibility and chemical shift effects. The consequence of such phase effects in terms of imaging properties is dependent on the particular sequence and data sampling speed. In addition to phase effects, the signal decay with T2 needs to be taken into account also, especially if the data acquisition time is similar or even longer than the T2 of the observed tissues under construction [10].

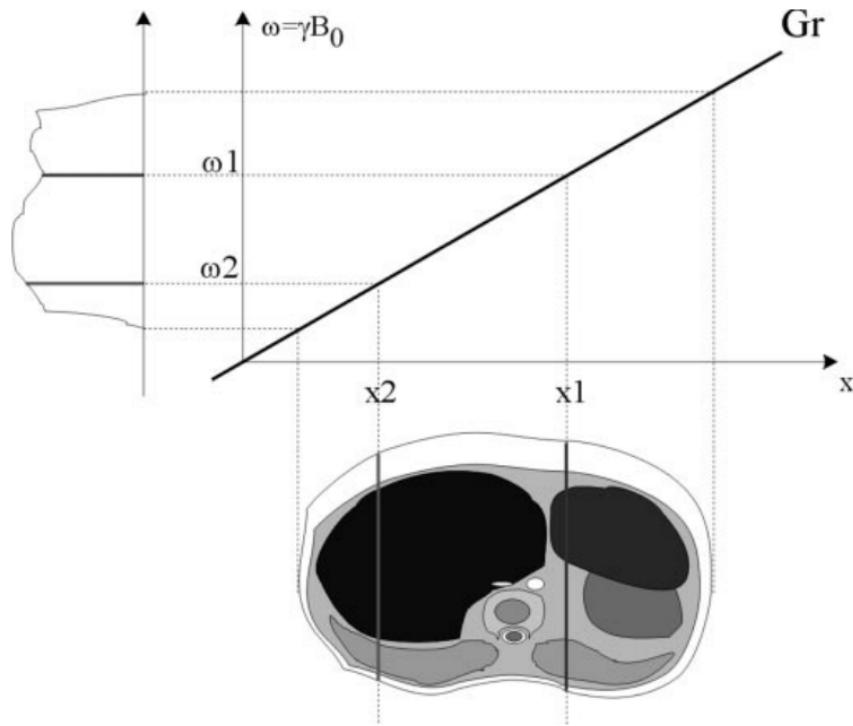


Figure 2.2: Schematic for the spatial encoding using a constant gradient taken from [10]. Gr spind that are located at positions  $x_1$  and  $x_2$  have different Larmor frequencies  $\omega_1$  and  $\omega_2$ . Spins located orthogonally to  $Gr$  have the same x-coordinates and thus the same Larmor frequency

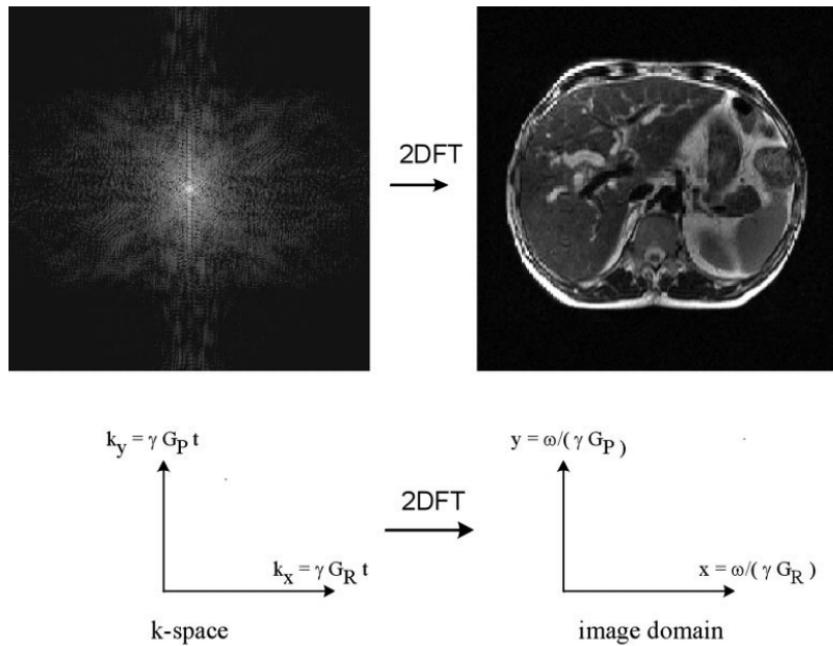


Figure 2.3: Correspondence of the image domain and the k-space via a 2D FT, taken from [10]. The coordinates of each domain are defined by the Larmor frequency according to Equation 2.1.

### Rectilinear K-Space Sampling

Almost all MR imaging sequences used in clinical routine presently are based on rectilinear k-space sampling, i. e., the sampling points in k-space are placed on a rectangular grid. This reflects the ability to use a fast Fourier transform (FFT) for such a sampling strategy. It allows image reconstruction in considerably under 1 s. With special reconstruction processors, the reconstruction speed can be as high as 10 to 25 images per second. The general 2D FT for transformation of a comparable number of non-rectilinear grid points in comparison requires 5 to 50 min per image, using identical computer hardware [10].

Using the two means of moving through the k-space, the basic techniques for rectilinear sampling can be divided into spin-echo-based techniques and gradient-echo-based techniques. In all rectilinear sampling techniques, the gradient which is used during signal acquisition is commonly called the readout gradient, whereas gradients used to bring the k-space trajectory to a certain starting point before data acquisition are called phase-encoding gradients. A gradient echo is formed using a reversal of the readout gradient. In conventional gradient-echo imaging, only the refocusing part of the k-space trajectory is measured, whereas the signal dephasing as well as phase encoding is performed prior to data acquisition (Fig. 8). The experiment is repeated after the recovery time TR, using different offsets of the phase-encoding gradient prior to each acquisition step. Using the k-space, it can be easily seen how gradient-echo imaging can be generalized to a faster technique, using less excitation periods by sampling more than one k-space line after each excitation. Repeating gradient reversal and the application of a phase-encoding

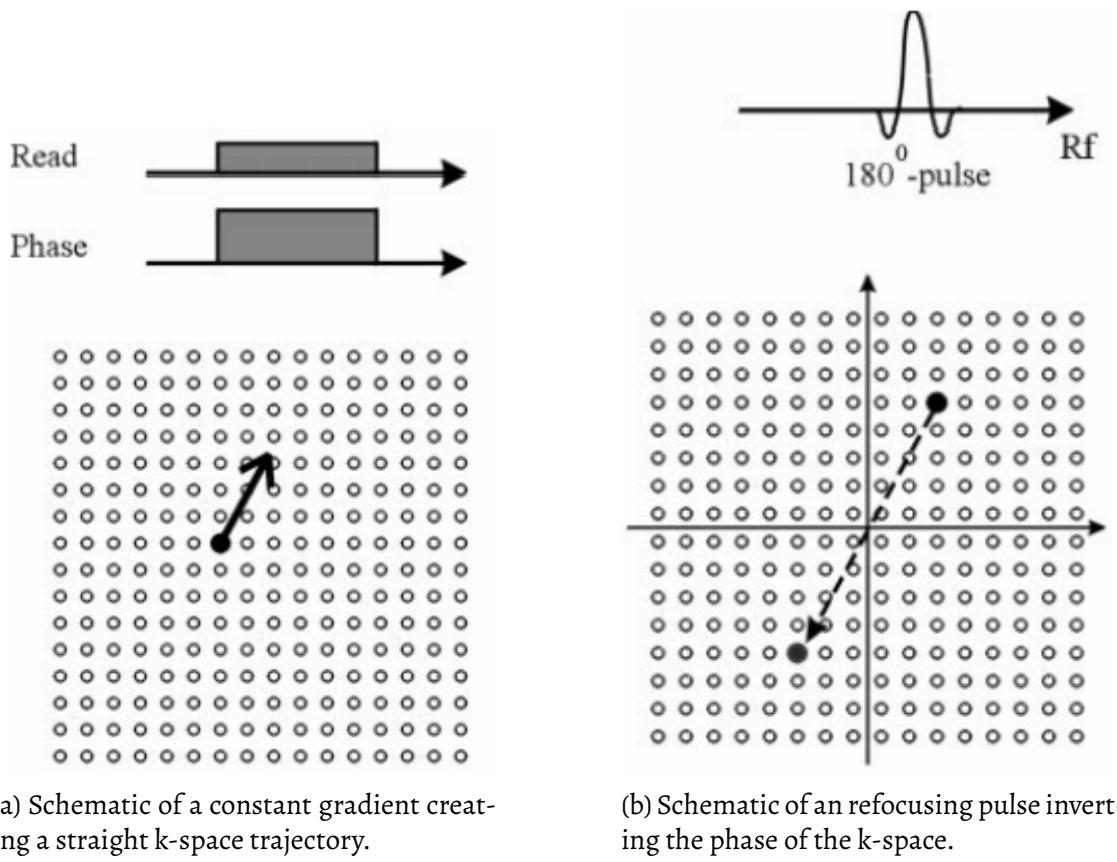


Figure 2.4: Methods of moving in the k-space using a) constant gradients to move along straight lines or b) mirror the k-space at the origin using a refocusing pulse. Taken from [10].

gradient can be used to read out several k-lines per acquisition with a concordant reduction in total imaging time. In the extreme case, all k-space data can be read out following a single excitation. This is called the single-shot, echo-planar imaging (EPI) sequence [11]. Despite the close conceptional similarity between EPI and gradient echo techniques with respect to their k-space trajectories, both types of sequences have vastly different imaging and hardware implications. This difference is caused by the sequential sampling of k-space points and the resulting influences on the signal of the mechanisms mentioned previously. As the signal decays with T2 the signal amplitude is reduced during data acquisition. For conventional gradient-echo imaging, even with moderate gradient strength (and thus longer acquisition times), this effect is inconsequential since the typical acquisition times of 5 to 10 ms are considerably shorter than the typical T2 of tissues. For EPI, however, the total acquisition time for n k-space lines is n times longer, even if the time necessary for gradient reversal and phase encoding is neglected. For the moderate acquisition times noted previously, this will lead to total acquisition times of 1 s or more, which exceeds most T2 values observed in human tissues. Dephasing of spins across the selected slice will lead to T2 of the order of 100 ms at most. For identical acquisition times,

as in our gradient-echo example, the signal will have vanished long before the end of data acquisition. Depending on whether data acquisition has commenced at the center or in the outskirts of k-space, the resulting image will either be severely blurred or reveal merely contours without contrast. This illustrates that EPI absolutely requires very fast data acquisition in order to acquire all of the k-space data within not much more than 100 ms at most; otherwise, image quality will deteriorate significantly. In the given time range, state-of-the-art scanners can acquire 64 to 128 phase-encoding steps. Longer acquisition times allowing higher nominal resolution are of course technically feasible, but the T2-decay of the signal will prohibit realization of the nominal resolution. Since cerebral spinal fluid (CSF) has a long T2, EPI images of the head with long acquisition times can look deceptively sharp. It must be noted, however, that the perceived sharpness applies only to CSF-filled structures, whereas the true resolution of parenchyma with a T2 of approximately 60 ms is considerably worse. With even stronger and faster gradients it is of course feasible to acquire true high-resolution images with short acquisition times. The penalty in SNR for such applications is, however, severe and has prohibited its use for most forms of morphological imaging. The EPI technique has emerged, however, as an ideal vehicle for functional imaging either in BOLD-based fMRI, in first-pass perfusion measurements or in diffusion imaging, but it is not a viable technique for high-resolution morphological imaging [10].

Even for short acquisition times, EPI suffers from severe problems regarding its imaging properties, which are caused by the phase-altering mechanisms mentioned previously. Chemical shift effects and susceptibility both lead to signal phase shifts which increase linearly with time. Since data acquisition along the readout gradient is fast, there is little effect in the readout direction. In the phase-encoding direction, however, the dwell time amounts to the acquisition time of one k-space line. For the acquisition of 64 k-space lines in 64 ms, the field of view in frequency units will be 1000 Hz. For a frequency shift of fat vs water of approximately 200 Hz at 1.5 T, this will mean that the signal of fatty tissue will be shifted by 20% of the image width. The resulting misregistration artifact necessitates the use of fat suppression for practically all EPI applications. Field inhomogeneity and susceptibility will lead to a similar shift. Since the resulting frequency shifts are not as discrete as chemical shift effects, the misregistration will vary continuously leading to highly non-linear image deformations. An additional problem in EPI relates to the fact that the k-space trajectory runs in opposite directions in every other k-space line due to the alternating polarity of the readout gradient. Any asymmetry in the acquired data will therefore lead to some slight perturbation in any other line. The final image will therefore be a superposition of the true image with an image representing the perturbation in every other k-space line. From this discussion of the k-space trajectory, it is easy to see that this perturbation will generate an image which is displaced by half the field of view in the phase-encoding direction. This is the so-called 180° artifact. All the confounding factors discussed herein will produce some such artifact. With appropriate post-processing routines based on reference signals, some of the sources of artifacts can be corrected for. In general, the 180° ghost will become stronger whenever any of the mechanisms mentioned previously become significant and will therefore deteriorate the image quality further. Motion can also cause phase-effects on the signal. With an alternating readout gradient, motion along this direction will be refocused in every

even refocusing period. Consequently, a flow-dependent 180° ghost will occur. Motion along the phase-encoding direction will lead to a velocity-dependent increase in the signal phase and consequently to a velocity-dependent signal displacement. Signal from a round vessel with laminar flow will therefore appear as a boomerang-shaped artifact displaced from the vessel wall. Any vascular diagnosis based on EPI images therefore has to be made with the utmost care. The inherent artifact behavior limits the practical use of EPI to examinations of the head. Abdominal or thoracic EPI and even cardiac EPI have not been developed beyond a feasibility state. They require considerable efforts in region-selective shimming and even post-processing, using actual field maps of the patient such that the time advantage of EPI is lost. For the time being, EPI therefore has to be regarded as a brain-only technique where it serves as a vehicle for the fMRI applications, including diffusion and perfusion imaging.

Conventional gradient-echo imaging suffers from no such effects and can therefore be applied all over the body with no or only mild susceptibility misregistration. Chemical shift effects will, however, influence the contrast of gradient-echo images by the chemical-shift-dependent dephasing of fat and water signals. Signals from fat and water protons located at the same position will periodically be in- and out of phase with each other in a time interval determined by their difference in Larmor frequencies. For 1.5 T the time difference between the two extreme conditions will be approximately 2.5 ms. Images acquired with TE<sub>s</sub> differing by 2.5 ms will look dramatically different due to this effect, although T2-dependent effects for such small TE differences will be negligible. As shown in Fig. 11, the out-of phase image will show a considerable reduction in signal amplitude in tissues containing both fat and water protons such as bone marrow. In addition, an artificial black line will appear at the interface between fatty and water tissues due to the zero crossing of the signal between these antiphase images. When (im)properly used, this artifact can be used to beautify images, giving them a crisp and well-defined appearance. Care must be taken, however, not to mis-read this artifact as an anatomical structure. With respect to the total acquisition time, gradient-echo imaging is somewhat slower than EPI due to the necessity of repetitive application of excitation pulses. Gradient-echo sequences, which are optimized for speed, typically achieve imaging times of 200 to 400 ms, using gradient systems comparable to those of EPI. Due to the strong signal saturation caused by the very short repetition times, the intrinsic SNR is comparably low (but for many applications still satisfactory). The SNR can be optimized by using low-excitation flip angles and/or total signal refocusing (true FISP). In combination with bolus application of contrast agent, gradient-echo sequences allow the acquisition of MR angiograms especially of the thorax, abdomen, and the extremities with unprecedented image quality. This relatively new application has led to whole new field of applications for gradient-echo imaging.

In order to reach a compromise between the more efficient sampling of EPI and the much more benign artifact behavior of gradient echoes, segmented approaches have been developed. In these approaches, more than one excitation is used and several k-lines are sampled per excitation. Depending more on the perspective of their creators than on anything else, such approaches are called multi-echo gradient echo, multi-shot EPI, or segmented EPI. The optimum number of k-lines per excitation is strongly dependent on the particular application; it can range from 3 for MR angiography and some cardiac imaging

applications to 64 or more for high-resolution head imaging.

### Imaging Acceleration

To alleviate the slow image acquisition times of the MR imaging, different MRI acceleration techniques were used. Most of these included a subsampling of the k-space, meaning the omission of certain higher frequencies, from which a images is reconstructed. While there are different subsampling techniques, most involve fully sampling a center region containing the lower frequencies and dropping higher frequencies, which are deemed less important. However, all of these methods create artifacts during image reconstruction leading to different technologies being developed to minimize these effects such as compressed sensing (CS) and parallel imaging (PI).

## 2.2 Image Registration

Image registration is a challenging, yet important task for image processing. It can be described as the process of transforming different image datasets into one coordinate system with matched imaging contents [2]. In the medical field this can be used for clinical applications such as disease diagnosis and monitoring, image-guided treatment delivery, and post-operative assessment. Medical image registration is typically used to pre-process data for tasks like object detection (for e.g. tumor growth monitoring) and segmentation (for e.g. organ atlas creation) where variation in spatial resolution is common between modalities like CT and MRI and patients. Thus the performance of these methods is dependent on the quality of image registration [1].

Medical image registration was often done manually by clinicians, however, registration tasks are often challenging and the quality of manual alignments is dependent on the expertise of the user. These manual registrations are thus not only time consuming, but also hardly reproducible leading to high interobserver-variability. The need for automatic registration is very much apparent, but this task remained hard to solve for a long time, requiring a lot of computational power and time for computer algorithms to solve the problem. While neural networks also require a lot of computational power and time to train, they promise fast execution after training. With the rise of deep learning these network gained popularity and now pose a real alternative to conventional algorithms and manual registration [2]. We will discuss these new approaches in the next section, but first we need to formally define our problem.

In pair-wise image registration two images ( $F$  and  $M$ ) are to be aligned, with  $F$  denoting the fixed and  $M$  the moving image.  $T$  is the desired spatial transformation that aligns the two images. This can be posed as an optimization problem:

$$T' = \arg \max S(F, T(M)), \quad (2.5)$$

with  $T'$  being the best transformation that maximizes the similarity  $S$  between the two images. This process is done iteratively improving estimates for the desired  $T$ , such that the defined similarity in the cost function is maximized [1].

Transformations can be categorized as rigid, affine, and deformable. A rigid transformation consists of rotation and translation; an affine transformation includes translations, rotations, scaling, and sheering; the two kinds of transformations are described as a 2D single matrix. Unlike rigid and affine transformation, deformable transformation is a high-dimension problem that we need to formulate by a 3D matrix for 2D deformable registration i.e., a so-called deformation field. While rigid and affine registration algorithms have already achieved good performance in many applications, deformable registration is still a challenging task due to its intrinsic complexity, particularly when the deformation is large. However, these are also the transformations most likely encountered in clinical practice as it can be utilized to fuse information from different modalities such as MRI and CT [4]. Additionally, deformable image registration can also be utilized for various computer-assisted interventions like biopsy [12] and (MRI-guided) radiotherapy [13, 14].

Intuitively, deformable image registration is an ill-posed problem, making it fundamentally different from other computer vision tasks such as object localization, segmentation or classification. Given two images, deformable image registration aims to find a spatial transformation that warps the moving image to match the fixed image as closely as possible. However, there is no ground-truth available for the desired deformation field and without enforcing any constraints on the properties of the spatial transformation, the resulting cost function is ill-conditioned and highly non-convex. In order to address the latter and ensure tractability, all image registration algorithms regularize the estimated deformation field, based on some prior assumptions on the properties of the underlying unknown deformation [1].

Many methods have been proposed for medical image registration to deal with the complex challenges of this task. Popular conventional registration methods include optical flow [15], demons [16] and many more. However, most of these still lack accuracy and computation speed, which makes newer deep learning approaches all the more interesting [3].

### 2.3 Deep Learning Architectures

Neural networks, despite the theoretical concepts being around for decades, have seen a meteoric rise in popularity over the last few years as constraints on computational power have been alleviated. Especially deep neural networks, which are often summarized under the term deep learning (DL). Recent years have witnessed an almost exponential growth in the development and use of DL algorithms, sustained thus far by rapid improvements in computational hardware (e.g. GPUs). Consequently, clinical applications requiring image classification, segmentation, registration, or object detection/localization, have witnessed significant improvements in algorithmic performance, in terms of accuracy and/or efficiency [1]. The following network architecture are widely used for different tasks including medical image registration.

Some basic stuff about network training, testing and different architectures that are relevant for the later Chapters.

## Convolutional Neural Networks

Convolutional neural networks (CNNs) are a type of deep neural networks with regularized multilayer perceptron, which are mainly used for image processing. CNNs use convolution operations instead of general matrix multiplications in typical neural networks. These convolutional filters make CNNs very suitable for visual signal processing. Because of their excellent feature extraction ability, CNNs are some of the most successful models for image analysis. Different variants of CNN have been proposed and have achieved the-state-of-art performances in various image processing tasks. A typical CNN usually consists of multiple convolutional layers, max pooling layers, batch normalization layers, sometimes dropout layers, a sigmoid or softmax layer. In each convolutional layer, multiple channels of feature maps are extracted by sliding trainable convolutional kernels across the input feature maps. Hierarchical features with high-level abstraction are extracted using multiple convolutional layers. These feature maps usually go through multiple fully connected layer before reaching the final decision layer. Max pooling layers are often used to reduce the image sizes and to promote spatial invariance of the network. Batch normalization is used to reduce internal covariate shift among the training samples. Weight regularization and dropout layers are used to alleviate data overfitting [3]. The loss function is often defined as the difference between the predicted and the target output. CNNs are usually trained by minimizing the loss via gradient back propagation using optimization methods like Adam [17].

## U-Net

The U-Net [18] architecture is an extension of the typical CNNs structure typically used for image segmentation, however it can also be used for image registration tasks. It adopts symmetrical contractive and expansive paths with skip connections between them. The encoding blocks on the left extract important features from the image using convolution layers and max pooling, which are then stored in the latent space in the middle. From there it is reconstructed using upsampling and convolutions in the decoding blocks on the right. Additionally, skip connections are used to improve the spatial resolution of the segmentation. This architecture allows effective feature learning from a small number of training datasets [3].

## Autoencoders

An autoencoder (AE) is a type of CNNs that learns to reconstruct an image from its input without supervision. AEs usually consists of an encoder which extracts the input features, which are stored a low-dimensional latent state space, similar to a U-Net, and a decoder which restore the original input from the latent space. To prevent an AE from learning an identity function, regularized autoencoders were invented, which can be used for e.g. denoising AEs. Variational AEs (VAEs) are generative models that learn latent representation using a variational approach, which constrains the variability of the outputs. VAEs can been used for anomaly detection and image generation [3].

## Generative Adversarial Networks

Generative adversarial networks (GANs) consist of two competing networks, a generator and a discriminator. The generator is trained to generate artificial data that approximate a target data distribution from a low-dimensional latent space similar to an AE. The discriminator is trained to distinguish the artificial data from actual data. The discriminator encourages the generator to predict realistic data by penalizing unrealistic predictions via learning. Therefore, the discriminative loss could be considered as a dynamic network-based loss term. The generator and discriminator both are getting better during training to reach Nash equilibrium. In medical imaging, GANs have been used to perform image synthesis for inter- or intra-modality, such as MRI to synthetic CT and vice versa. In medical image registration, GANs are usually used to either provide additional regularization or translate multi-modal registration to uni-modal registration [3].

## 2.4 Deep Learning for Image Registration

Recently, there has been a surge in the use of deep learning based approaches for medical image registration. Their success is largely due to their ability to perform fast inference, and the flexibility to leverage auxiliary information such as anatomical masks as part of the training process. The most effective methods, such as *VoxelMorph* [6], typically employ a U-Net style architecture to estimate dense spatial deformation fields. These methods require only one forward pass during inference, making them orders of magnitude faster than traditional iterative methods. Following the success of *VoxelMorph*, numerous deep neural networks have been proposed for various registration tasks [8]. Other approaches also utilize CNNs, AEs and GANs. Typical strategies are discussed in more detail in the following sections.

### Supervised Registration

Supervised registration describes training a network with a ground truth displacement field that is either real (created by hand) or synthetic (generated via traditional iterative registration algorithms). Thus the loss can easily be calculated as the difference in the displacement fields of the network prediction and the ground truth (see Figure 2.5 for a visual overview). These methods have achieved notable results with real displacement fields as supervision. However, this approach is very limited by the size and the diversity of the dataset. As the displacement fields are often calculated by conventional algorithms their effectiveness might be limited for difficult problems with which the traditional algorithms struggle. Fully supervised methods are widely studied and have notable results, but the generation of real or synthetic displacement fields is hard, and these displacements fields might be different from the real ground truth, which can impact the accuracy and efficiency of these kinds of methods [4]. Notable approaches include *BIRNet* [19] and *LAPNet* [20], however the latter works in the k-space domain.

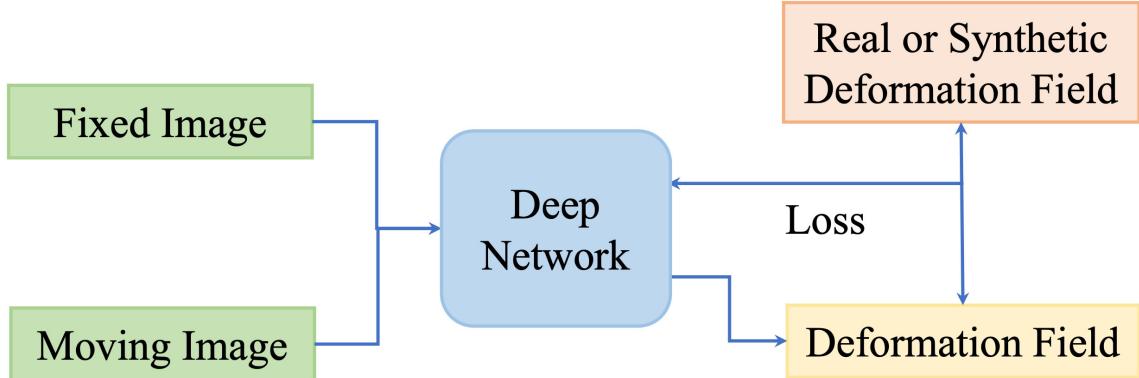


Figure 2.5: Example graph illustrating the training process of a supervised network, taken from [4].

### Unsupervised Registration

As the preparation of the ground truth displacement field for supervised methods is inconvenient, limitations in generalizing results in different domains and various registration tasks are inevitable. Thus, unsupervised registration has a more convenient training process with paired images as inputs, but without a ground truth. Generally, unsupervised learning consists of similarity-based (see Figure 2.6) and GAN-based methods (see Figure 2.7), where the loss function computes the similarity between the aligned images and the smoothness of the displacement field, rather than the difference to a ground truth [4]. Well known example are *IC-Net* [21], *VoxelMorph* [6], *TransMorph* [22] and *SYM-Net* [23].

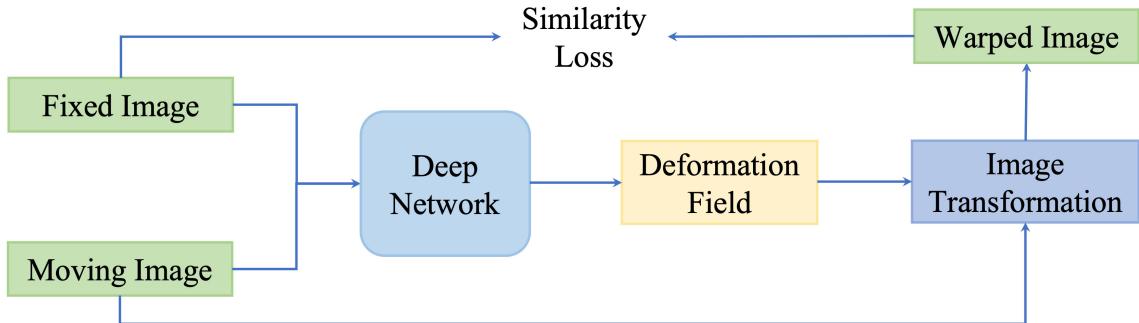


Figure 2.6: Example graph illustrating the training process of an unsupervised network with only a image similarity loss, taken from [4].

### Evaluation Metrics

Different metrics can be used for evaluation of registration performance including similarity metrics and other stuff. The mean squared error (MSE) is calculated between a

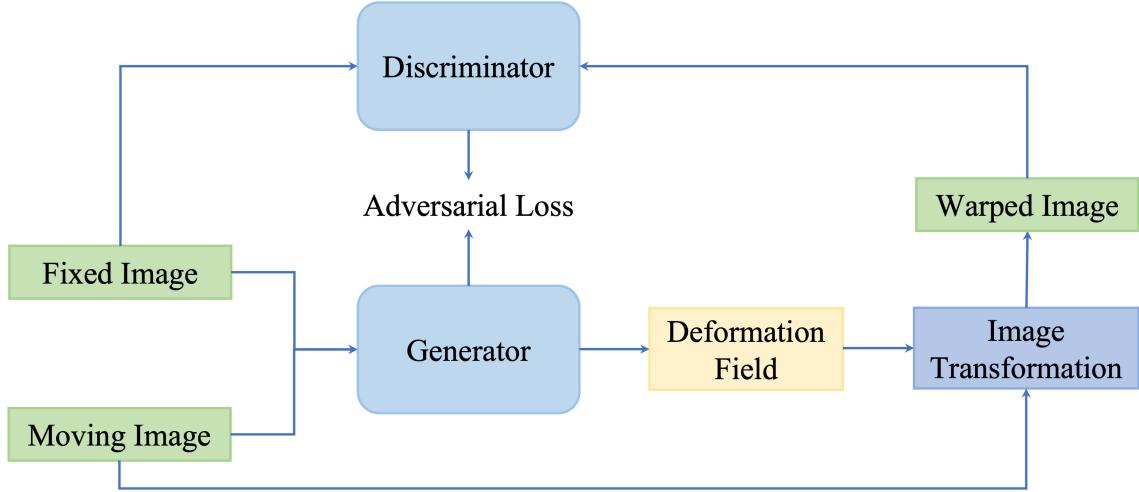


Figure 2.7: Example graph illustrating the training process of a GAN, taken from [4].

fixed (ground truth) image and the warped image giving a pixel-wise comparison:

$$\text{MSE} = \frac{1}{N} \sum_{i=1}^N |F(x,y) - W(x,y)|^2, \quad (2.6)$$

where  $F$  is the fixed image and  $W$  is the warped image with  $N$  being the number of pixels in the images. The lower the MSE the higher the similarity with 0 being a perfect match (the images are exactly the same).

Another common image metric is the structural similarity index measure (SSIM), which operates between 1 (complete similarity) and 0 (no similarity) and tries to estimate the general similarity instead of a pixel-wise comparison making it more robust against contrast changes compared to e.g. MSE.

$$\text{SSIM} = \frac{(2\mu_F\mu_W + c_1) \cdot (2\sigma_{FW} + c_2)}{(\mu_F^2 + \mu_W^2 + c_1) \cdot (\sigma_F^2 + \sigma_W^2 + c_2)}, \quad (2.7)$$

where  $\mu_F, \mu_W$  are the mean values of the images  $F$  and  $W$ ;  $\sigma_F^2, \sigma_W^2$  the variances of  $F$  and  $W$  as well as  $\sigma_{FW}$  the covariance for  $F$  and  $W$ , with  $c_1, c_2$  being constants derived from the dynamic range of the images.

All of these metrics can also be used as loss functions for network training.

For comparison of segmentations the Dice score is a commonly used metric to estimate the similarity of two segmentations. A score of 1 indicates a complete overlap/match and a score of 0 indicates no overlapping of the segmentations. The Dice score is calculated as follows:

$$\text{Dice} = \frac{2|M_F \cap M_W|}{|M_F| + |M_W|}, \quad (2.8)$$

with  $M_F$  and  $M_W$  being segmentations corresponding to  $F$  and  $W$  with  $M_F$  being a manual segmentation which is warped to obtain  $M_W$  [24].

Aside from image similarity measures and the evaluation of image segmentations, the displacement field itself can also be evaluated. This is usually based on the assumption that the displacement should be smooth as, for example, folding the image could result in physically unrealistic anatomic structures, indicating errors. The Jacobian determinant of the deformation  $|J_\phi|$  must be positive everywhere to avoid folding and thus the percentage of non-positive Jacobian determinant of the deformation  $\% |J_\phi| \leq 0$  can be used to evaluate the quality of the generated displacement field [5].

## 2.5 Motion-Corrected Reconstruction

Patient motion during acquisition is one of the major impediments of high-quality MRI scans. This is especially true for thoracic and abdominal imaging, as organs move during breathing. This motion can induce several consequences on MR signal formation. Intraview and interview motion have to be distinguished between: motion is intraview when occurring during individual MR experiments (between RF excitation and echo formation), whereas motion is interview when occurring between individual MR experiments. Whenever the period of motion is slow compared to the period of MR acquisition defined by the repetition time  $TR$ , the assumption can be made that motion is interview. This is often a reasonable assumption when considering pseudo-periodic motion induced by respiration, and also possibly by cardiac contraction, which are the two common sources of motion in cardiac and abdominal imaging (typically, the adult respiratory period is about 4–5 s, and  $TR \approx 10\text{ms}$  for fast imaging). Interview motion results in spatial encoding inconsistencies, and thus in image deterioration which can take complex forms (blurring/ghosting artifacts) as acquisition is performed in a Fourier space. Several strategies can be employed in order to handle patient motion better. Patient cooperation is the most commonly used method. However, breath-holds cannot last much longer than 20s and physiological drifts cannot be completely avoided. This leads to a limitation on the time-period of signal recording and thus, signal-to-noise ratio (SNR). Moreover, the position of organs in successive breath-holds may not be reproducible. Synchronization techniques are well-established and systematically used in clinical protocols, but they require a high-level of motion reproducibility. This is often a limiting factor considering heart rate variability (whether in free breathing or during a breathhold), and respiratory variability in terms of amplitude and frequency. This is why alternative techniques have been put forward with the aim of inverting the process of spatial encoding of moving structures that underlies artifacts. While it is possible to correct for motion prospectively, by modulating the magnetic field gradients and RF fields in order to cancel the effect of motion in the Bloch equations, the method is limited to correcting of, at best, affine motion, due to magnetic field gradient systems being linear. Motion can also be compensated for in reconstruction, however, many methods are restricted to rigid or affine motion due either to computational issues or to the difficulty of modeling complex displacement fields [25]. The latter however seems to become less of a problem as image registration improves with e.g. the help of deep learning.

Motion-resolved data acquisition for these applications are usually accelerated by Parallel Imaging or Compressed Sensing techniques yielding sub-Nyquist sampled (subsampled)

pled) k-space data. In order to reconstruct aliasing-free images, these methods rely on reconstruction schemes that, for example incorporate sparsity or low-rank constraints to solve the ill-posed problem. Fixed sparsity assumptions in Compressed Sensing are often too restrictive and incapable of fully modeling spatio-temporal dynamics. Careful fine-tuning between regularization and data consistency is required and especially in highly subsampled cases residual aliasing may remain in the image (under-regularization) or staircasing and blurring artifacts can occur (over-regularization) which affect the image registration. After reconstruction, non-rigid motion fields can be estimated in image space from reconstructed images by solving a registration problem. A particular interest and challenge lies in the derivation of reliable motion fields which capture the spatio-temporal non-rigid deformations, such as respiratory or cardiac movement. Instead of performing these two steps sequentially, motion-compensated image reconstruction schemes like *GRICS* [25] integrate both motion field estimation and motion correction into the reconstruction process. These methods require reliable motion-resolved images from which the motion fields can be estimated. Motion field estimation can be controlled or supported by external motion surrogate signals, initial motion field estimates, from motion-aliased images or low-frequency image contents. Moreover, spatio-temporal redundancies can be exploited to achieve an aliasing-free image. While these methods have been proven to be more robust against registration errors, they can require a significantly increased computational demand and/or limit imaging acceleration [26].

There are also approaches for learning new subsampling strategies in a data-driven manner (pruning unimportant k-space frequencies) [27] as well as deep learning based radial [28] and non-Cartesian [29] subsampling for MRI acceleration.

# 3

## Data

In the following chapter, the datasets used in this thesis are presented and potential pre-processing steps, as well as uses, discussed. The first dataset called *OASIS*, containing MRI brain scans, was used in the original publications for training and evaluation of both *Fourier-Net* and *Fourier-Net+*. It was thus used to get familiar with these networks and their respective work-flows as well as reproducing some results from the paper to ensure that the code works as intended.

For our major experiment the *CMRxRecon* dataset was used. The cardiac MRI k-space data was used for extending the networks to frame-to-frame registration in order for them to be integrated into a motion reconstruction pipeline. It already contains subsampled data, but does not provide segmentations for multi-coil data.

Thus a third dataset, the *ACDC* cardiac dataset, was used. While it contains no k-space data segmentations are given for end-systolic and end-diastolic frames, which can thus be used for the evaluation of registration performance. Most parameter studies and optimizations of the networks were done with this dataset.

### 3.1 OASIS Dataset

The *OASIS-1* dataset [30] contains T1-weighted MRI brain scans from 454 subjects and was mainly used for training and testing the original *Fourier-Net/Fourier-Net+*. The brain scans were further pre-processed by [31] for the *Learn2Reg-Challenge* [32]. This enables subject-to-subject brain registration, as all MRI scans were bias-corrected, skull-stripped, aligned, and cropped to the size of  $160 \times 192 \times 224$ . The images were stored in the *NIFTI* format, which makes it really easy to use. Examples can be seen in Figures 3.1b and 3.1a with slices from the center of the x-, y- and z-axis.

### 3.2 CMRxRecon Dataset

The *CMRxRecon* dataset [33] from the *CMRxRecon2023* challenge specializes in Cardiac magnetic resonance imaging (CMR). The dataset includes fully sampled and subsampled multi-coil k-space data, as well as auto-calibration lines. This includes imaging of differ-

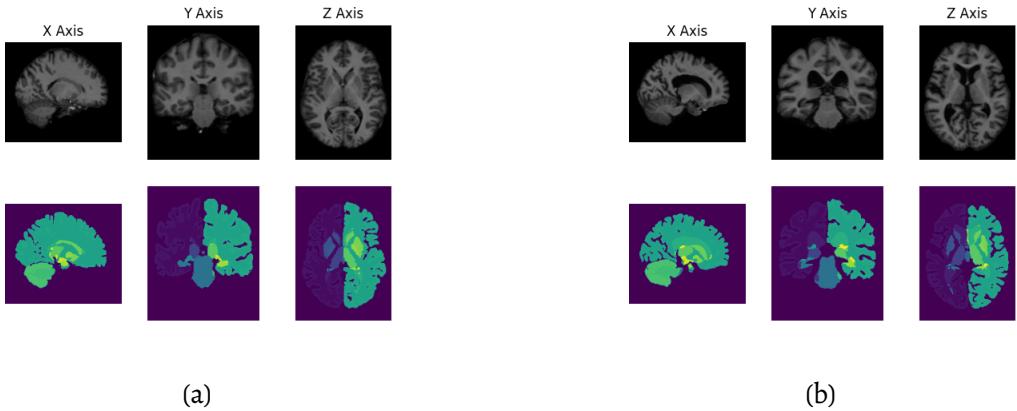


Figure 3.1: Example images (upper row) from the *OASIS* dataset with corresponding labels (bottom row).

ent anatomical views like long-axis (2-chamber, 3-chamber, and 4-chamber) and short-axis (SAX). There is a total of 120 training data, 60 validation data, and 120 test data from healthy volunteers. One of the goals of the challenge is the reconstruction from subsampling and motion correction for the heart movement.

### Image Reconstruction

As the data was recorded in the k-space (see Figure 3.3) and stored as *.mat* files we first need to reconstruct the images to use them for training and testing of *Fourier-Net/Fourier-Net+*. For this, multi-coil reconstruction algorithms like *SENSE* [34] can be used. The sensitivity of the different coils can be seen in Figure 3.2, where each coil focuses on a specific area of the image. These images are then stitched together in reconstruction to produce a combined image with great overall contrast. The dataset, as mentioned before, contains fully sampled as well as subsampled k-space data (see Figure 3.3a and Figure 3.3b), with the latter being typically used to accelerate the MRI acquisition process. This is done by not using all of the available k-space data, but rather masking the signal to achieve the subsampling. The center region of the k-space is always fully sampled, but the outer regions are subsampled depending on the sampling strategy [10]. The dataset contains subsampled k-space data for 4x and 8x acceleration with the latter of course leading to more distortions in the reconstructed image as the subsampling can induce image reconstruction artifacts. This can be seen in Figure 3.3d, where the image reconstructed from 4x accelerated k-space data seems blurred when compared to the fully sampled one in Figure 3.3c. For all experiments the short-axis view data was used. As image sizes between patients can vary slightly, interpolation was used to standardize the images to a size of  $246 \times 512$  to avoid further problems with e.g. loss calculation. Additionally, min-max-normalization is used to standardize the data range for all images to  $[0, 1]$  as this was also varying. The reconstructed images were stored for every image slice for every patient. Thus, all frames for a specific image slice were in a single folder to enable easy access for later data load-in for frame-to-frame registration.

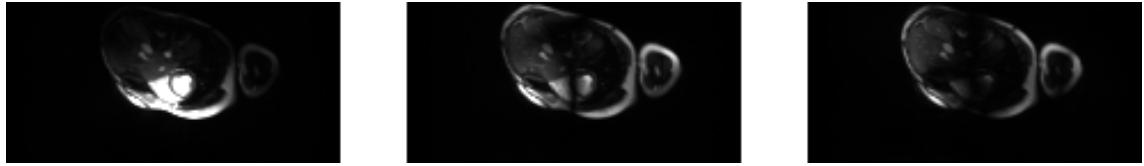
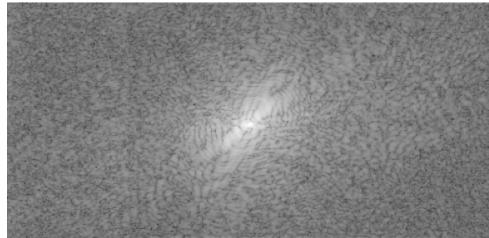
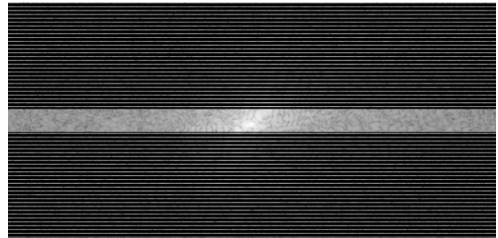


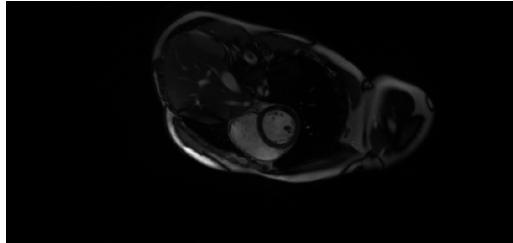
Figure 3.2: MRI images for three different coils from the *CMRxRecon* dataset [33].



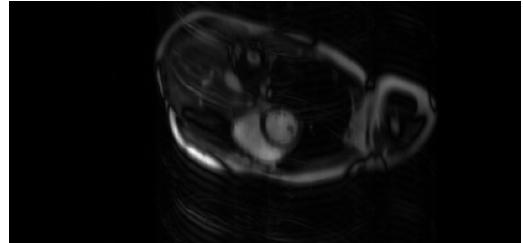
(a) Example for a fully sampled k-space with the low frequencies in the middle.



(b) Example of a subsampled k-space for 4x acceleration.



(c) Corresponding image reconstructed from the fully sampled k-space.



(d) Corresponding image reconstructed from the subsampled k-space.

Figure 3.3: Fully sampled and subsampled k-space data from the *CMRxRecon* dataset [33] with corresponding images.

## Pre-Processing

As seen before, a large part of the image is just black background that contains no usable information. The interesting cardiac region in the image, however, is much smaller and thus image cropping can be used to extract it. The fairly large image size can thus be cut down without losing important image information, which also helps with memory efficiency. One problem however remains: image alignment. Whereas in the *OASIS* dataset the brain scan were already aligned, the cardiac data in the *CMRxRecon* dataset is not. Thus, one cannot simply center-crop all images as the cardiac region might be located in very different parts of the image. To resolve this issue, adaptive cropping with automatic detection of the cardiac region in the images is needed.

To begin, some assumption about the data needs to be made: First, the cardiac region is at the same position for all slices and frames from one patient. Second, most of the movement happens in the cardiac region between frames. From the latter we can conclude that we need to examine the frames, thus the difference between consecutive frames is

calculated and averaged for each slice. According to the first assumption this process is repeated for all slices and their results summed to create a sort of heatmap as seen in Figure 3.4c with a clearly higher intensity in the cardiac region compared to the background. To further isolate this region a simple thresholding is performed followed by an opening operation, which removes any smaller fragments still left in the image. The cardiac region is now clearly separated from the rest of the image as seen in Figure 3.4d. The only step left is to calculate the approximate center of this region and create a crop around it. To drastically reduce the image size a cropping factor of one third was used, reducing the image size from  $246 \times 512$  to  $82 \times 170$ . The result of this drastic crop can be seen in Figure 3.4b when compared to the original image in Figure 3.4a. Note, however, that the cropped image in Figure 3.4 appears larger than it actually is, when compared to the original, to allow for a more detailed look. Thus both images are not presented to scale, but in order to show the most image detail.

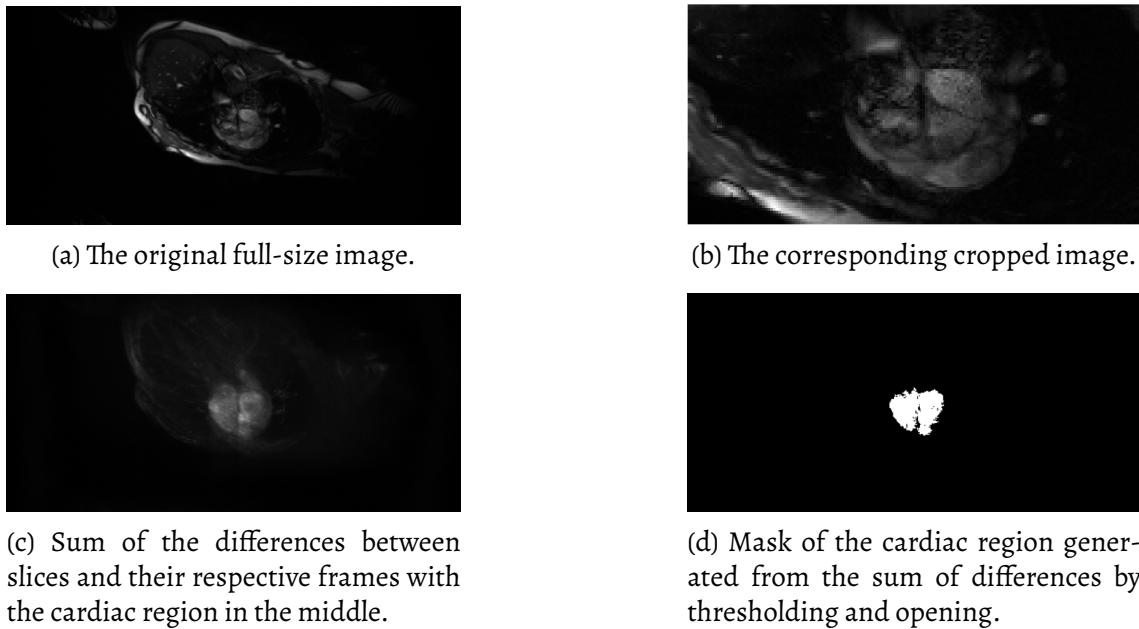


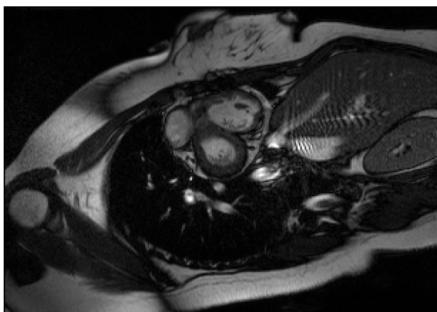
Figure 3.4: Cropping of an example image by locating the movement between frames for each slice. Note that the cropped image appears larger than it actually is compared to the original.

### 3.3 ACDC Dataset

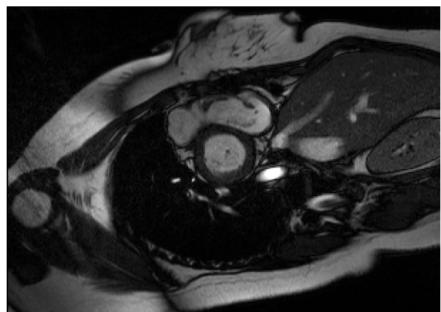
The ACDC dataset [35] from the *Automated Cardiac Diagnosis Challenge* (ACDC) during the MICCAI 2017 conference. It contains cardiac cine-MRI short-axis data from 150 subjects that were divided into 5 subgroups (4 pathological, 1 healthy). For each subject systolic (see Figure 3.5a) and diastolic frames (see Figure 3.5b) are provided with corresponding segmentations (see Figure 3.5c and 3.5d), which enables direct comparison via e.g. the Dice score. The segmentations contain only four values with 0, 1, 2 and 3 representing

pixels located in the background, in the RV cavity, in the myocardium, and in the LV cavity. The frames themselves are 3D volumes with size  $216 \times 256 \times 10$ , thus ten image slices can be extracted which each have size  $216 \times 256$ . The same holds true for the segmentation, which means that both the image data and segmentations can generated in the same manner.

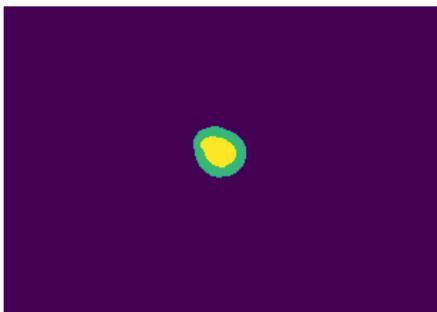
For the training data, the original 4D data was used as we do not require the segmentations for the end-systolic and end-diastolic frames. As the 4D data has size  $216 \times 256 \times 10 \times 30$  we can extract 30 frames for each of the ten slices. These can be sorted into 251376 image pairs for training. For the validation and test data the end-systolic and end-diastolic frames with their segmentations were used given us 641 image pairs each as we can only have these two frames to align.



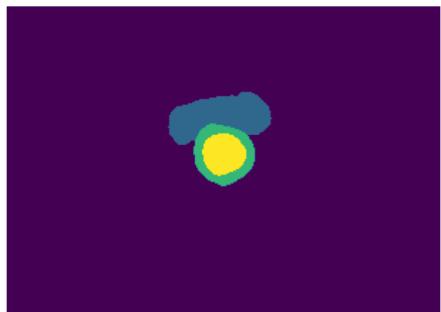
(a) Image of systolic phase.



(b) Image of diastolic phase.



(c) Segmentation of systolic phase.



(d) Segmentation of diastolic phase.

Figure 3.5: Example images for systolic (a) and diastolic frames (b) with corresponding segmentations (c),(d) taken from the ACDC dataset [35].

# 4

## Network Architectures

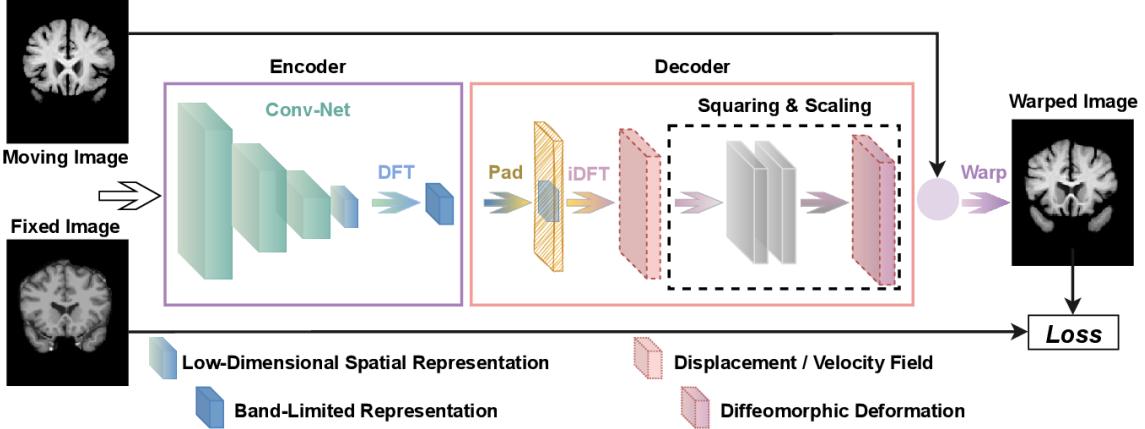
As a starting point *Fourier-Net* [7] and its successor *Fourier-Net+* [8] were used. These networks, which are explained in the following pages, enable fast and accurate registration while needing less resources compared to similar approaches. These attributes are very beneficial for a potential application like motion correction because the current networks, e.g. *LAPNet* [20], are usually supervised and require large computational resources.

### 4.1 Fourier-Net

*Fourier-Net* is a new unsupervised approach that aims to learn a low-dimensional representation of the displacement field in a band-limited Fourier domain instead of the full field in the spatial domain. This band-limited representation is then decoded by a model-driven decoder to the dense, full-resolution displacement field in the spatial domain. This allows for fewer parameters and computational operations, resulting in faster inference speeds [7]. The architecture is based on the U-Net [18], like most deep registration approaches, but replaces the expanding path with a parameter-free model-driven decoder as mentioned before. The encoder of *Fourier-Net* consists of a CNN, which takes two images (fixed and moving) as inputs. The output is a displacement field that is then converted from the spatial domain into the Fourier domain via an discrete Fourier transformation (DFT). From there, this band-limiting representation is padded with zeros to the full resolution of the original displacement field. The field is then recovered by using the inverse DFT (iDFT) to convert it back into the spatial domain. This displacement field is then used to warp the moving image into the fixed image. Additionally, squaring and scaling layers [36] can be added before warping the image in order to encourage a diffeomorphism in final deformation.

#### Encoder

The encoder of *Fourier-Net* consists of a CNN that generates the displacement field between the two inputs followed by a DFT layer that produces a band-limited representation of the full displacement field. The fully convolutional neural network (FCN) from *SYMNet* [23] was modified to function as the CNN of the encoder. Its (original) architec-

Figure 4.1: Architecture of *Fourier-Net* taken from [7].

ture (see Figure 4.2) is again based on the *U-Net* with a contracting and expanding path. The FCN concatenates the inputs images  $X$  and  $Y$  as a single 2-channel input and estimates two dense, non-linear displacement fields  $\phi_{XY}$  and  $\phi_{YX}$ , however we only need the displacement field for the moving image, denoted as  $\mathbb{S}_\phi$ , since we are not interested in transforming the fixed image. This is actually a low dimensional representation because *Fourier-Net* does not utilize the last two levels of the FCNs expanding path that would be needed to reconstruct the actual full-resolution displacement field  $\phi$ .

For each level in the contracting path of the FCN, two successive convolution layers are applied, which contain one  $3 \times 3 \times 3$  convolution layer with a stride of 1, followed by a  $3 \times 3 \times 3$  convolution layer with a stride of 2 to further compute the high-level features between the input images as well as downsample the features by half until the lowest level of the network is reached. For each level in the expanding path of the FCN, the feature maps from the contracting path are concatenated through skip connections and apply  $3 \times 3 \times 3$  convolution with a stride of 1 and  $2 \times 2 \times 2$  deconvolution layer for upsampling the feature maps to twice of its size. At the end of the expanding path, two  $5 \times 5 \times 5$  convolution layers with a stride of 1 are appended to the last convolution layer and generate the displacement fields  $\theta_{XY}$  and  $\theta_{YX}$  [23]. Each convolution layer in the FCN is followed by a rectified linear unit (ReLU) activation, except for the output convolution layer that does not have an activation function because *Fourier-Net* only uses the first two levels of the expanding path, thus leading the FCN to generate a low dimensional representation  $\mathbb{S}_\phi$  of the full-resolution displacement field  $\phi$ .

As discussed previously, the encoder aims to learn a displacement (or velocity) field in the band-limited Fourier domain. Intuitively, this may require convolutions to be able to handle complex-valued numbers, which can be done by using complex-valued CNNs [37], which are suitable when both input and output are complex values, however these complex-valued operations sacrifice computational efficiency. Other approaches like *DeepFlash* [38] tackle this problem by converting the input images to the Fourier domain and using two individual real-valued CNNs to learn the real and imaginary parts separately. This, however, increases training and inference cost. To bridge the domain gap between real-valued spatial images and complex-valued band-limited displacement fields without increasing

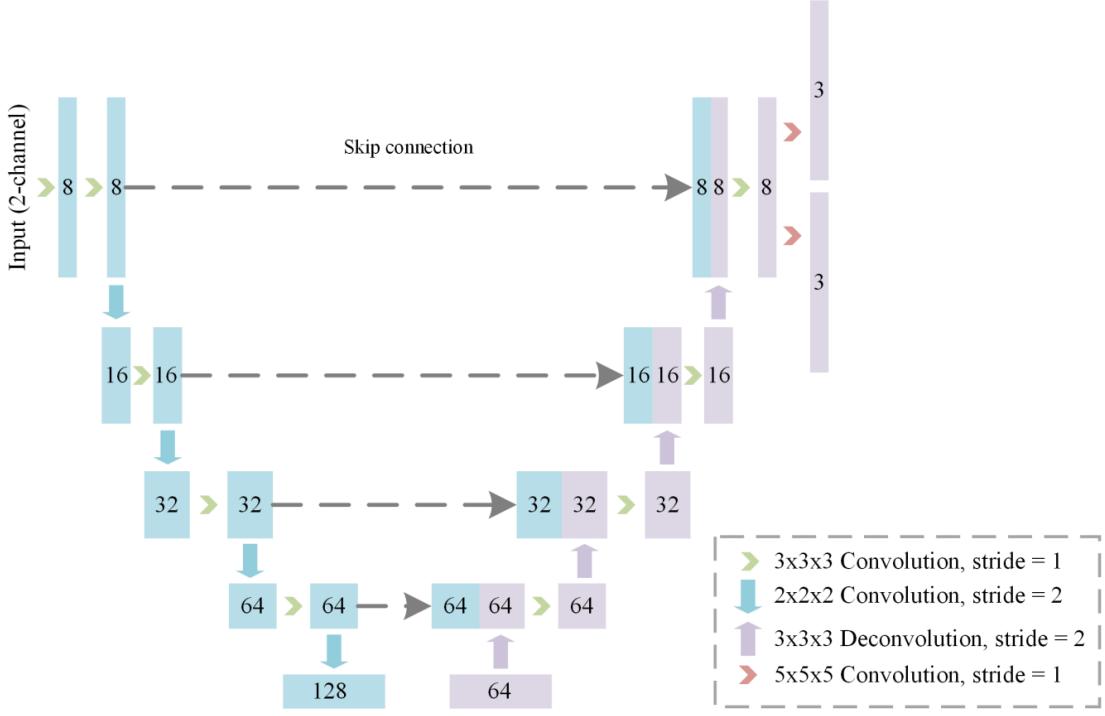


Figure 4.2: Architecture of the FCN from SYMNet [23].

complexity, *Fourier-Net* uses a DFT layer after the FCN. This is a simple and effective way to produce complex-valued band-limited displacement fields without the network needing to be able to handle complex values itself. The DFT applied to the displacement field  $\phi$  can be defined as follows:

$$[\mathcal{F}(\phi)]_{k,l} = \sum_{n=0}^{H-1} \sum_{m=0}^{W-1} \phi_{n,m} \cdot \exp\left(i \cdot \left(\frac{2\pi k}{H} \cdot n + \frac{2\pi l}{W} \cdot m\right)\right), \quad (4.1)$$

where  $\phi$  has size  $H \times W$ ,  $n \in [0, H - 1]$  and  $m \in [0, W - 1]$  are the discrete indices in the spatial domain, and  $k \in [0, H - 1]$  and  $l \in [0, W - 1]$  are the discrete indices in the frequency domain with  $i$  being the imaginary unit. However,  $\phi$  in this equation is actually the low dimensional representation of the displacement field output by the modified FCN, which can be formulated as follows:

$$\mathbb{S}_\phi = \text{FCN}(M, F; \Theta), \quad (4.2)$$

with  $M$  being the moving and  $F$  the fixed image, as well as  $\Theta$  representing the parameters of the FCN. Thus, the whole encoder can be defined as:

$$\mathbb{B}_\phi = \mathcal{F}(\mathbb{S}_\phi) = \mathcal{F}(\text{FCN}(M, F; \Theta)), \quad (4.3)$$

with the DFT layer  $\mathcal{F}$ , full-resolution spatial displacement field  $\phi$  and the complex band-limited displacement field  $\mathbb{B}_\phi$ . The low dimensional representation  $\mathbb{S}_\phi$  actually contains all the information of the band-limited Fourier coefficients in  $\mathbb{B}_\phi$ . As such, *Fourier-Net*

does not need to learn the coefficients of  $\mathbb{B}_\phi$ , but instead only the real-valued coefficients in  $\mathbb{S}_\phi$ , which is the low dimensional spatial representation of the full-resolution spatial displacement field  $\phi$ , which is then reconstructed by the decoder.

### Decoder

The decoder contains no learnable parameters, instead the usual expansive path is replaced with a zero-padding layer, an iDFT layer, and an optional squaring and scaling layer.

The output from the encoder is a band-limited representation  $\mathbb{B}_\phi$  in the frequency domain of the low dimensional displacement field  $\mathbb{S}_\phi$  in the spatial domain. To recover the full-resolution displacement field  $\phi$  in the spatial domain, we first pad the patch  $\mathbb{B}_\phi$ , containing mostly low frequency signals, to the original image resolution with zeros. We then feed the zero-padded complex-valued coefficients, denoted as  $\mathcal{F}(\phi)$ , to an iDFT layer consisting of two steps: shifting the Fourier coefficients from centers to corners and then applying the iDFT to convert them into the spatial domain:

$$\phi_{n,m} = \frac{1}{HW} \sum_{k=0}^{H-1} \sum_{l=0}^{W-1} \mathcal{D}_{k,l} [\mathcal{F}(\phi)]_{k,l} \cdot \exp \left( i \cdot \left( \frac{2\pi n}{H} \cdot k + \frac{2\pi m}{W} \cdot l \right) \right). \quad (4.4)$$

The  $H \times W$  sized sampling mask  $\mathcal{D}$  is a low-pass filter that has zeros as entries if they are on the positions of high-frequency signals in  $\phi$  and ones if they are on the low-frequency positions. Thus we can reconstruct the full spatial displacement field  $\phi$  from  $\mathbb{B}_\phi$  despite the latter being band-limited. Approaching the problem from the other side we can also think about working backwards from the final displacement. For this, after applying equation (4.1), the low-frequency signals are shifted to a center patch with size  $\frac{H}{a} \times \frac{W}{b}$  with  $a = 2 \cdot Z_a, b = 2 \cdot Z_b, Z_a, Z_b \in \mathbb{Z}^+$ , which is then center-cropped to get  $\mathbb{B}_\phi$ . This crop of  $\mathcal{F}(\phi)$  can be reconstructed using the iDFT from equation (4.4) with the cropping functioning as a kind of low-pass filtering:

$$[\mathbb{S}_\phi]_{n,m} = \frac{ab}{HW} \sum_{k=1}^{\frac{H}{a}-1} \sum_{l=1}^{\frac{W}{b}-1} [\mathbb{B}_\phi]_{k,l} \cdot \exp \left( i \cdot \left( \frac{2\pi an}{H} \cdot k + \frac{2\pi bm}{W} \cdot l \right) \right), \quad (4.5)$$

with  $n \in [0, \frac{H}{a} - 1]$  and  $m \in [0, \frac{W}{b} - 1]$  being the indices of the spatial domain, while  $k \in [0, \frac{H}{a} - 1]$  and  $l \in [0, \frac{W}{b} - 1]$  are the indices of the frequency domain with  $i$  being the imaginary unit. Thus  $\mathbb{S}_\phi$  actually contains all the necessary information from  $\phi$ , as long as they have the same low-frequency coefficients  $\mathbb{B}_\phi$ . This can be formulated as:

$$[\mathbb{S}_\phi]_{n,m} = ab \cdot \phi_{an,bm}, \quad (4.6)$$

because most entries of  $\mathcal{F}(\phi)$  are zeros, and the remaining values are exactly the same as in  $\mathbb{B}_\phi$ , which means that  $\mathbb{S}_\phi$  contains all the information  $\phi$  can provide. For this, however,  $\mathbb{B}_\phi$  needs to be padded to the original image size first in order to get the full-resolution displacement and not a low dimensional representation. This ultimately shows that there is a unique mapping between  $\mathbb{S}_\phi$  and  $\phi$ , which means that it is reasonable to use a network to learn  $\mathbb{S}_\phi$  directly from image pairs and then reconstruct the displacement field

in a very efficient manner [7]. The complete reconstruction process is visualized in Figure 4.3.

As both padding and iDFT layers are differentiable, *Fourier-Net* can be optimized via back-propagation. For *Diff-Fourier-Net* extra squaring and squaring layers [36] are needed in the decoder turning the displacement field into a stationary velocity field. Typically seven scaling and squaring layers are used to impose such a diffeomorphism [7, 36].

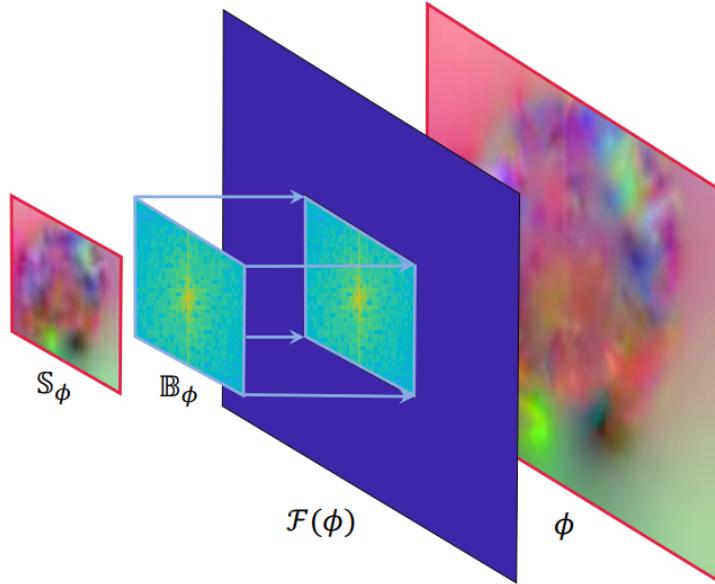


Figure 4.3: Reconstruction of the displacement field via the decoder taken from [8].

### Diffeomorphic Transforms

Diffeomorphic deformations are differentiable and invertible, thus preserving topology, which is a desirable property for transformations.  $\phi : \mathbb{R}^N \rightarrow \mathbb{R}^N$  represents the deformation that maps the coordinates from one image to coordinates in another image, as long as both images have the dimension  $N$ . When using a stationary velocity field representation like e.g. *DARTEL* [39], the deformation field is defined through the following ordinary differential equation (ODE) [36]:

$$\frac{\partial \phi^{(t)}}{\partial t} = v(\phi^{(t)}), \quad (4.7)$$

where  $\phi^{(0)} = \text{Id}$  is the identity transformation and  $t$  is time. The final registration field  $\phi^{(1)}$  can be obtained by integrating the stationary velocity field  $v$  over  $t = [0, 1]$ . This is typically done by integration numerically using scaling and squaring [40]. The integration of a stationary ODE represents a one-parameter subgroup of diffeomorphisms. In group theory,  $v$  is a member of the Lie algebra and is exponentiated to produce  $\exp(v) = \phi^{(1)}$ , which is also a member of the Lie group. From the properties of one-parameter subgroups, for any scalars  $t$  and  $t'$ ,  $\exp((t + t') \cdot v) = \exp(t \cdot v) \circ \exp(t' \cdot v)$ , where  $\circ$  is a composition map associated with the Lie group. Starting from  $\phi^{(1/2^T)} = p + v(p)$  where

$p$  is a map of spatial locations, we use the recurrence  $\phi^{(1/2)^{t-1}} = \phi^{(1/2^t)} \circ \phi^{(1/2^t)}$  to obtain  $\phi^1 = \phi^{(1/2)} \circ \phi^{(1/2)}$ .  $T$  is chosen so that  $v \approx 0$  [36].

As diffeomorphic deformations are defined as smooth and invertible deformations, the output of the iDFT layer in *Fourier-Net* can be regarded as a stationary velocity field  $v$  instead of a displacement field  $\phi$ . In Figure 4.1 scaling and squaring layers are visualized. These apply the diffeomorphic transformation in three steps [40]:

1. Scaling: Divide the velocity field  $v$  by a factor  $2^N$ , so that  $\frac{v}{2^N}$  is close to zero (depending on the desired accuracy).
2. Exponentiation: Compute  $\exp(\frac{v}{2^N}) = \phi^{(1)}(\frac{v}{2^N})$  with a first-order explicit numerical scheme.
3. Squaring:  $N$  recursive squarings of  $\phi^{(1)}(\frac{1}{2^N}) = \exp(\frac{v}{2^N})$  to yield an accurate estimation of  $\phi^{(1)}(1) = \phi^{(1)}(\frac{1}{2^N})^{2^N} = \exp(\frac{v}{2^N})^{2^N} = \exp(v)$ .

Thus, the diffeomorphic deformation can be efficiently calculated. When used, the specific version is then called *Fourier-Net Diff* to differentiate it from the baseline version.

## Spatial Transformer

The warping layer of *Fourier-Net* utilizes the *Spatial Transformer* [41], which allows for spatial image manipulation within the network. This is a differentiable and learnable module for neural networks which applies a spatial transformation to a feature map during a single forward pass. The spatial transformer mechanism is split into three parts as seen in Figure 4.4. First is the localization network, which takes the input and outputs the parameters for the transformation. These are then used to create a sample grid using the grid generator. Lastly, the sampler produces the output feature map based on the input at the grid points.

From the input feature map  $U \in \mathbb{R}^{H \times W \times C}$  with width  $W$ , height  $H$  and channels  $C$  the localization network  $f_{\text{loc}}$  computes the parameters  $\theta = f_{\text{loc}}(U)$  of the transformation  $T_\theta$  which is later applied to the feature map. Thus the size of  $\theta$  varies depending on the transformation. The localization network function can both be implemented as a fully-connected network or as a CNN, but should include a final regression layer to produce the transformation parameters.

In order to warp the input feature map, each output pixel is computed by applying a sampling kernel centered at a particular location in the input feature map. The output pixels are defined to lie on a regular grid  $G = G_i$  of pixels, forming an output feature map  $V \in \mathbb{R}^{H' \times W' \times C}$ , where  $H'$  and  $W'$  are the height and width of the grid with  $C$  again being the number of channels, which is the same for input and output.

In order to perform a spatial transformation of the input feature map  $U$ , the sampler must take the set of sampling points  $T_\theta(G)$ , along and produce the sampled output feature map  $V$ . Each coordinate  $(x_i^s, y_i^s)$  in  $T_\theta(G)$  defines the spatial location in the input where a sampling kernel is applied to get the value at a particular pixel in the output:

$$V_i^c = \sum_n^H \sum_m^W U_{nm}^c k(x_i^s - m; \Phi_x) k(y_i^s - n; \Phi_y), \quad (4.8)$$

with  $\Phi_x$  and  $\Phi_y$  being the parameters for a generic sampling kernel  $k$  that defines the image interpolation,  $U_{nm}^c$  is the value of the input feature maps at location  $(n, m)$  in the channel  $c \in [1, \dots, C]$  and  $V_i^c$  is the value for every pixel  $i \in [1, \dots, H'W']$  for the output feature map. Any sampling kernel can be used, as long as (sub-) gradients can be defined with respect to  $(x_i^s, y_i^s)$  to allow the loss gradients to flow back not only to the input feature map, but also to the sampling grid coordinates and therefore back to the transformation parameters  $\theta$  and localization network, thus enabling back-propagation [41].

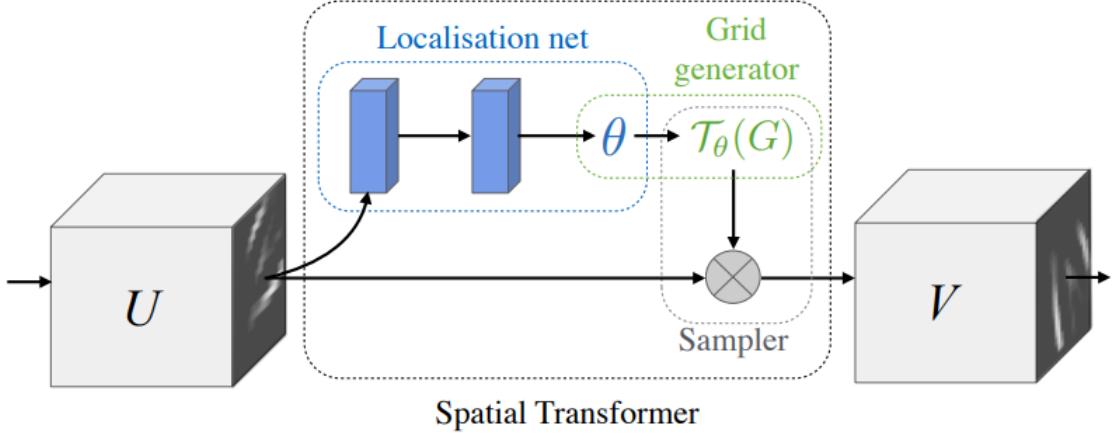


Figure 4.4: Architecture of the *Spatial Transformer* taken from [41].

### Loss Function

The loss function consists of two parts to enable unsupervised learning, which are balanced using the scalar parameter  $\lambda$ . The first,  $\mathcal{L}_1$ , measures the similarity between the fixed image and the moving image after warping, while the second,  $\mathcal{L}_2$ , ensures a smooth displacement field. Thus, the unsupervised loss  $\mathcal{L}$  can be calculated as follows:

$$\begin{aligned}\mathcal{L}(\Theta) &= \min \left( \mathcal{L}_1(\phi(\Theta)) + \lambda \cdot \mathcal{L}_2(\phi(\Theta)) \right) \\ &= \min \left( \mathcal{L}_1(v(\Theta)) + \lambda \cdot \mathcal{L}_2(v(\Theta)) \right),\end{aligned}\tag{4.9}$$

for both displacement fields  $\phi$  and velocity fields  $v$ . The first part of the loss function consists of:

$$\mathcal{L}_1(\phi(\Theta)) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}_{Sim}(M_i \circ (\phi_i(\Theta) + \text{Id}) - F_i),\tag{4.10}$$

where  $\circ$  denotes the warping operation,  $N$  the number of training pairs with moving images  $M_i$  and fixed images  $F_i$ ,  $\Theta$  the network parameters,  $\phi_i$  the displacement field,  $\text{Id}$  the identity grid.  $\mathcal{L}_{Sim}$  determines the similarity between warped moving images and fixed images via MSE or NCC, and the second term of the unsupervised loss,  $\mathcal{L}_2$ , defines the

smoothness regularization function that controls smoothness of the displacement fields:

$$\mathcal{L}_2(\phi(\Theta)) = \frac{1}{N} \sum_{i=1}^N \|\nabla \phi_i(\Theta)\|_2^2, \quad (4.11)$$

with  $\nabla$  denoting the first order gradient and  $\|\cdot\|_2^2$  denoting the squared  $L_2$ -Norm. When using the squaring and scaling layers, thus making the deformation of the moving image diffeomorphic, the loss needs to be modified by replacing the displacement field  $\theta$  with the velocity field  $v$ . Thus, both parts of the of the loss function need to be changed:

$$\mathcal{L}_1(v(\Theta)) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}_{Sim}(M_i \circ \text{Exp}(v_i(\Theta) - F_i), \quad (4.12)$$

$$\mathcal{L}_2(v(\Theta)) = \frac{1}{N} \sum_{i=1}^N \|\nabla v_i(\Theta)\|_2^2 \quad (4.13)$$

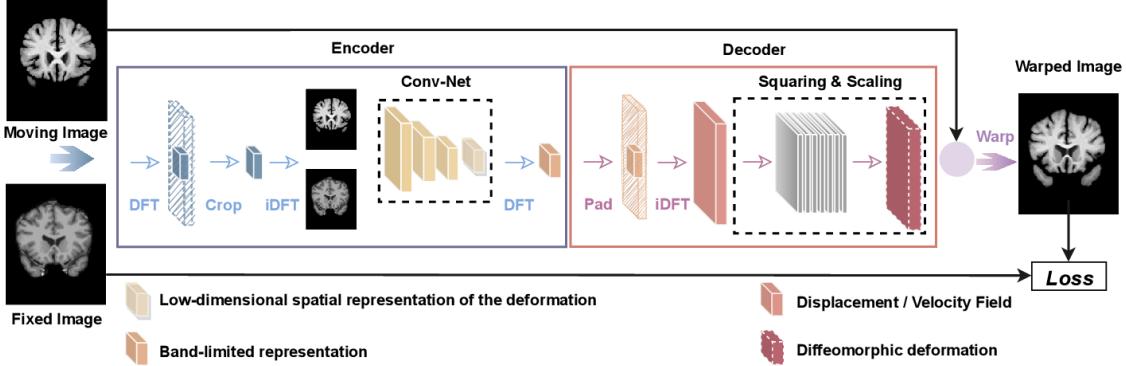
## 4.2 Fourier Net+

*Fourier-Net+*, as the name suggests, is an extension of Fourier-Net which takes the band-limited spatial representation of the images as input, instead of their original full-resolution counterparts. This leads to further reduction in the number of convolutional layers in the contracting path of the network, resulting in a decrease of parameters, memory usage, and computational operations. This makes *Fourier-Net+* even more efficient than its predecessor [8].

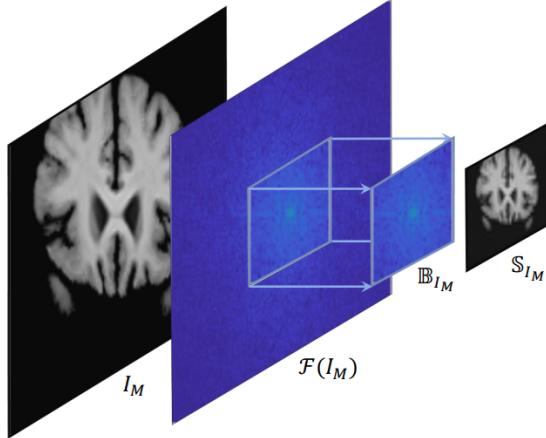
As seen in Figure 4.5, the network architecture is almost the same as for *Fourier-Net* (see Figure 4.1 for comparison). However, while the decoder, and thus the loss function, remain the same, the encoder is slightly altered to make the network even more efficient. For this, similarly to the decoder, a DFT is used, however this time the idea is applied to the input images. These are first transformed into the Fourier domain, then low-pass filtered by center-cropping and finally reconstructed from their band-limited representation back into the spatial domain via an iDFT. The two images, now compressed, are the input for the encoder of *Fourier-Net*, meaning the CNN and following DFT. However, due to the band-limiting before the CNN, the latter can be made much more light-weight, thus reducing computational cost. This is visualized in Figure 4.7. Thus, *Fourier-Net+* too is overall lighter than the baseline *Fourier-Net* in terms of the number of parameters and computations. However, such a light network may face limitations in accurately capturing complex deformations. To counter this potential weakness, the authors propose a cascaded version of *Fourier-Net+*, which uses multiple versions of *Fourier-Net+* cascaded one after the other to achieve a better overall displacement field [8]. A schematic for this can be seen in Figure 4.8.

### Changes to the Encoder

In order to further reduce the amount of computational operations, *Fourier-Net+* discards the early layers of the encoder. Instead, a DFT  $\mathcal{F}(I_M)$  followed by a center-crop to produce

Figure 4.5: Architecture of *Fourier-Net+* taken from [8].

the band-limited representation  $\mathbb{B}_{I_M}$  in the frequency domain and iDFT are used to get the spatial patch  $\mathbb{S}_{I_M}$ , while the rest of the encoder from *Fourier-Net* stays the same. The input  $I_M$  is thus compressed to a lower resolution (i.e. band-limited) using the frequency space, which reduces the computational cost. This process is visualized in Figure 4.6. The encoder of *Fourier-Net+* has several convolutional layers less in the contracting path compared to *Fourier-Net*, which leads to a further accelerated registration process while reducing the memory footprint. These advantages are visualized in Figure 4.7 where the amount of different layers between a conventional *U-Net*, *Fourier-Net* (with the smaller decoder) and *Fourier-Net+* (with a smaller encoder and decoder) are shown.

Figure 4.6: Compression in the frequency domain of the encoder used in *Fourier-Net+* taken from [8].

### Effects of Cascading

As seen in the previous section, *Fourier-Net+* is lighter than *Fourier-Net* due to the band-limited representation of both images and deformations lowering the number of parameters and computations. This, however, can lead to limitations when trying to accurately capture complex deformations. To this end, a cascaded version of *Fourier-Net+*

## 4 Network Architectures

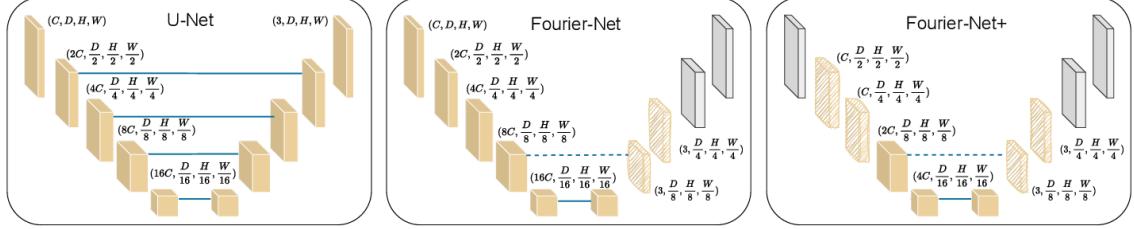


Figure 4.7: Architecture of the CNN for a typical U-Net, *Fourier-Net* and *Fourier-Net+* taken from [8].

called  $K \times \text{Fourier-Net+}$  (see Figure 4.8), with  $K$  denoting the amount of cascades, can be used where the the warped image of one cascade is used as the moving image of the next. It is important to note that the weights are not shared between cascades. The squaring and scaling layers, in case a diffeomorphic deformation is wanted, are applied after the last cascade. The same is true for the calculation of the loss function.

In order to more accurately describe this version of *Fourier-Net+* one can look at the in- and outputs of the different cascades. The first cascade of  $K \times \text{Fourier-Net+}$  has the moving image  $I_M$  and fixed image  $I_F$  as inputs. While the latter always stays the same for all cascades the warped image  $I_M^{w(1)}$  from the first cascade is used as input for the second cascade instead of the original moving image. Thus, in general  $I_M^{w(k-1)}$  and  $I_F$  are the inputs for a cascade  $k \in [1, K]$  with output  $\delta\phi^{(k)}$ . Furthermore,  $I_M^{w(k)}$  can be defined as:

$$I_M^{w(k)} = (((((I_M \circ \delta\phi^{(1)}) \circ \delta\phi^{(2)}) \circ \dots) \circ \delta\phi^{(k-1)}) \circ \delta\phi^{(k)}) = I_M \circ \phi^{(k)}, \quad (4.14)$$

with  $\phi^{(k)}$  being the displacement field computed by composing the outputs of the first cascade up to the  $k$ -th cascade:

$$\phi^{(k)} = \delta\phi^{(1)} \circ \delta\phi^{(2)} \circ \dots \circ \delta\phi^{(k-1)} \circ \delta\phi^{(k)}. \quad (4.15)$$

Thus, the output displacement field of the  $K$ -th cascade  $\phi^{(K)}$  is the final displacement which is then used to warp  $I_M$ , the original moving image, in order to compute the loss [8].

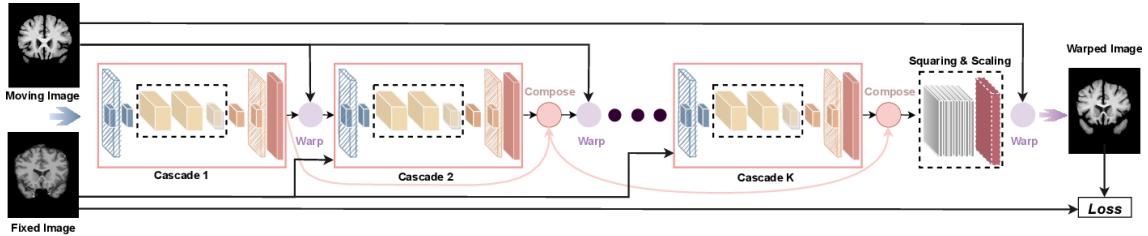


Figure 4.8: Cascaded version of *Fourier-Net+* taken from [8].

# 5

## Experiments

In the following chapter the different experiments that were conducted are described and explained. First, experiments were conducted on the *CMRxRecon* dataset to find the optimal parameters for our model. As this dataset lacks ground truth segmentations, the *ACDC* dataset was used to extend these tests.

### 5.1 Parameter Studies on the CMRxRecon Dataset

As a first step, frame-to-frame registration on the *CMRxRecon* dataset was tested. The goal of these initial experiments was to examine whether the network could learn to align the cardiac image pairs despite the movement between the frames. To start, different parameters were tested in order to find the optimal model parameters for the *CMRxRecon* dataset.

As discussed in section 3.2, the CMR images were cropped to extract the moving cardiac region. This was done for both the fully sampled and subsampled (Acc4, Acc8, Acc10) data. Image pairs for each slice from each patient of the training set were generated to train the network. In order to evaluate the training progress the mean MSE and mean SSIM metrics calculated the validation data at each checkpoint (usually every 1000 iterations). The best model parameters, as determined by the metrics, were saved in a separate folder. To further evaluate the training success the unseen test data containing new patients was used. To benchmark the networks training success the MSE and SSIM were again used, but the mean inference time and the percentage of non-positive Jacobian determinant of deformation ( $\% |J_\phi| \leq 0$ ) were also calculated. The procedure, however, was slightly altered for the subsampled data as only the displacement field of the model was used to align the corresponding fully sampled image pairs. The evaluation metrics were then calculated on these images for better comparability.

#### Fourier-Net versus Fourier-Net+

First, the baseline network *Fourier-Net* was compared to its newer version *4xFourier-Net+* once with both using the diffeomorphic version and once without the transform. Thus the effect of the diffeomorphic deformation can be compared to the normal transforma-

tion on *Fourier-Net* and *4xFourier-Net+*. All of the experiments were conducted on the fully sampled data as the subsampling should not effect the tested properties. The other parameters were also kept constant with a learning rate of 0.0001, MSE-loss and  $\lambda = 0.01$ . FT crop in *4xFourier-Net+* was set to  $24 \times 24$ .

### Starting Channel Size

Next, the impact of the starting channels was examined. For this, *4xFourier-Net+* was trained with MSE-loss, a learning rate of 0.0001,  $\lambda = 0.01$ , FT crop of  $24 \times 24$  and no diffeomorphic transform. The experiments were again only conducted on the fully sampled data. Three different channel sizes were used for training *Fourier-Net+:* 8, 16 and 32. Due to the early stopping the networks were trained for 17, 17 and 10 epochs respectively.

### Quantitative Comparison with NiftyReg

To establish a baseline, the similarity metrics were calculated between the unaligned test data. *4xFourier-Net+* was then trained with the MSE-loss, 8 channels,  $\lambda = 0.01$  and a learning rate of 0.0001. *NiftyReg*, a traditional registration algorithm that iteratively optimizes the displacement [24], was used for comparison. *NiftyReg* provides the warped image for the fully sampled data directly, but for the subsampled data the control point grid produced by *NiftyReg* was used to resample the corresponding fully sampled data for comparability. Thus only the fully sampled test data is needed for an unregistered baseline as the subsampled data is only used to produce the displacement which is then used to warp the fully sampled data.

## 5.2 Parameter Tests on the ACDC Dataset

As the evaluation on the CMRxRecon dataset is limited to image similarity measures, further testing was done on the ACDC dataset which contains cardiac data with segmentations. These can be used to better evaluate the registration performance by using e.g. the Dice score.

### Fourier-Net versus Fourier-Net+

First, *Fourier-Net*, *Fourier-Net+* and *4xFourier-Net* are compared for both normal and diffeomorphic versions. Additionally, network versions creating dense instead of band-limited displacements were used a baseline for the best achievable performance. The registration performance was evaluated on the test set using the percentage of Dice scores calculated with and without the background label, percentage of SSIM, MSE multiplied by  $10^{-3}$  (denoted by m for milli) and the percentage of negative Jacobian determinant of deformation. The mean inference time on the GPU together with memory consumption for each model are also added to the comparison with the latter containing number of trainable model parameters, mult-add operations (in millions or billion as denoted by M or G) and the total memory in Megabyte (MB). The latter is only given for the normal versions as the diffeomorphic version require the same amount of memory as well as the baseline dense

versions. The first network was trained for a maximum of 15 epochs with early stopping activating after 3 epochs without improvement of the validation Dice score. The other variants were than trained with the same number of epochs, which worked out to be only 6 epochs due to the large training set in the experiments without any data augmentation. All networks were trained on the fully sampled ACDC data with MSE as the unsupervised similarity loss, a channel size of 8, learning rate of 0.0001 and  $\lambda = 0.01$ . *Fourier-Net+* and *4xFourier-Net* used a FT crop size of  $48 \times 48$  for these experiments. All of these parameters were constant for the diffeomorphic and dense versions of the networks. The results and subsequent discussion can be found in section 6.2.

### Starting Channel Size

Next, the impact of the starting channels on the ACDC data was examined as these dictate the number of features that the network uses. For this, *Fourier-Net+* and *4xFourier-Net+* was trained with MSE-loss, a learning rate of 0.0001,  $\lambda = 0.01$ , FT crop of  $24 \times 24$  and no diffeomorphic transform. The experiments were again only conducted on the fully sampled data. Four different channel sizes were used for training: 8, 16, 32 and 64. The first network variant was trained with early stopping, which ended at 6 epochs. For better comparability, the other network variants were also trained with the same number of epochs. The registration performance was evaluated on the test set using the percentage of the Dice score calculated once on all labels, once without the background using only the cardiac labels, as well as the percentage of SSIM, MSE multiplied by  $10^{-3}$  and the percentage of negative Jacobian determinant of deformation. Additionally the number of network parameters, number of mult-add operations (in millions or billion as denoted by M or G) and memory consumption in Megabyte (MB) are listed for all network variants to better compare the relationship between registration performance and memory consumption to potentially find a optimal configuration for both performance and efficiency. The mean inference time was measured on GPU in seconds. The results and subsequent discussion can be found in section 6.2.

### Fourier-Transform Crop Size

The FT crop size, used for compressing the input images, is a parameter specific to *Fourier-Net+* and *4xFourier-Net+*. Four different sizes of the FT crop were analyzed:  $80 \times 168$ ,  $40 \times 84$ ,  $48 \times 48$  and  $24 \times 24$ . A larger crop would obviously compress the images less, thus leading to less efficiency, however, a very small crop, while very efficient, might lose some of the image details leading to a decrease in performance. The experiments were again only conducted on the fully sampled data as one would expect similar results on subsampled data. The network variants were trained with MSE as the unsupervised similarity loss, a channel size of 8, learning rate of 0.0001 and  $\lambda = 0.01$ . Again, the networks were trained for only 5 epochs, as this is were the first model ended training due to early stopping. The registration performance was evaluated on the test set using the Dice score calculated once on all labels, once without the background, as well as the SSIM, MSE and the percentage of negative Jacobian determinant of deformation. Additionally, the number of network parameters, number of mult-add operations in millions (M) and memory

consumption in Megabyte (MB) are listed for all network variants to better compare the relationship between registration performance and memory consumption to potentially find a optimal configuration for both performance and efficiency. The mean inference time was measured on GPU in seconds. The results and subsequent discussion can be found in section 6.2.

### Comparison with VoxelMorph

After finding parameters which optimize registration performance while trying to minimize memory consumption, a comparison to another commonly used registration network should be made. For this, the famous *VoxelMorph* [6] was used. It was trained with the default number of layers, MSE as the similarity loss, a learning rate of 0.0001,  $\lambda = 0.01$  and provides a dense displacement field to align the image pair. For comparison *Fourier-Net*, *Fourier-Net+* and *4xFourier-Net+* were trained for 6 epochs with MSE as the unsupervised similarity loss, a channel size of 16, learning rate of 0.0001 and  $\lambda = 0.01$  as well as a FT crop size of  $48 \times 48$ , which provide a good balance of performance and memory imprint for the networks. To evaluate registration performance, the Dice score was again computed for all labels and without the background as well as the similarity metrics SSIM, MSE and the percentage of negative Jacobian determinant of deformation. All networks were trained on the fully sampled ACDC data. The inference times of all networks were computed on the GPU and the number of network parameters, number of mult-add operations in billions (G) and memory consumption in Megabyte (MB) are also given to further compare the efficiency of the networks.

### Dense Displacement on Accelerated Data

The difference between a dense displacement field and a band-limited one was already explored in section 5.2. This, however, only includes fully sampled data, not accelerated data. This section explores potential performance changes on the latter with the hypothesis that the network variants with the dense displacement will perform worse than the variants with the band-limited displacement field for strongly subsampled data. This expectation stems from the fact that frequencies outside the center region of the k-space are dropped for the accelerated data thus reducing the advantage of the dense displacement in comparison to the band-limited (i.e. center-cropped) displacement.

The tests again included *Fourier-Net*, *Fourier-Net+* and *4xFourier-Net+*, which were all trained for 6 epochs with MSE as the unsupervised similarity loss, a channel size of 16, learning rate of 0.0001 and  $\lambda = 0.01$ . The registration performance was evaluated on the test set using the Dice score calculated once on all labels, once without the background, as well as the SSIM, MSE and the percentage of negative Jacobian determinant of deformation. Additionally, the number of network parameters, number of mult-add operations in millions (M) and memory consumption in Megabyte (MB) are given.

### Comparison on Subsampled Data

After comparing *Fourier-Net*, *Fourier-Net+* and *4xFourier-Net+* with both dense and band-limited displacement fields as well as another comparison with *VoxelMorph* on fully sam-

pled data an extension to subsampled data needs to be made. For this *Fourier-Net*, *Fourier-Net+* and *4xFourier-Net+* were compared to a baseline consisting of the unaligned test image pairs as well as *NiftyReg*, *VoxelMorph* and *LAPNet*. Thus, observations of the performance on accelerated data can be made as well as comparisons to the aligned baseline (which all methods should aim to beat), a traditional registration algorithm, a dense registration network working in the image domain and a another dense network which works in the k-space domain.

*Fourier-Net*, *Fourier-Net+* and *4xFourier-Net+* were trained for 6 epochs with MSE as the unsupervised similarity loss, a channel size of 16, learning rate of 0.0001 and  $\lambda = 0.01$  as well as a FT crop size of  $48 \times 48$  for all subsampling data. *VoxelMorph* was trained with default layers and MSE as the similarity loss, a learning rate of 0.0001 and  $\lambda = 0.01$ . Three different accelerations were tested: *Acc4* ( $R=4$ ), *Acc8* ( $R=8$ ) and *Acc10* ( $R=10$ ). For reference, the data for the fully sampled data ( $R=0$ ) is also provided. The registration performance of all methods was evaluated on the test set using the Dice score calculated once on all labels, once without the background, as well as the SSIM and MSE as image metrics. All times are computed on CPU for better comparability as *NiftyReg* currently only works on CPU as the GPU capable versions were deprecated and the other networks would gain an unfair time advantage when ran on GPU. The baseline obviously has no time associated as it only denotes the metrics on the test data before registration. The results and an comprehensive discussion are provided in section 6.2.

### 5.3 Data Consistency Experiments and Network Changes

In this section, changes to the network architecture and loss function are discussed.

#### Data Consistency Loss

As seen in the section 6.2, the performance for all methods decreases drastically for highly subsampled/accelerated data due to reconstruction artifacts in the image domain due to the missing frequencies in the k-space domain. However, since we are interested in maintaining performance on subsampled data, it might prove useful to use k-space information for the loss calculation and thus for the network training as it should not contain artifacts compared to the image domain. Thus, it would act as a kind of data consistency constraint during training. In order to achieve this, the similarity loss for *Fourier-Net*, *Fourier-Net+* and *4xFourier-Net+* was calculated not on the warped and fixed images, but instead on their corresponding masked k-spaces. This would hopefully still provide the network with useful information during training that is free of any aliasing artifacts, which might decrease performance on subsampled data. The k-spaces were masked to only compare the sampled frequencies to avoid large differences between the deleted frequencies of the fixed image with potential interpolated frequencies of the warped image. After trying the MSE as the similarity loss, it quickly became apparent that this yielded wildly inconsistent displacements which often corrupted the warped image and could no longer be reigned in by the smoothness constraint of the loss function. Even with a  $\lambda = 10$  to put a greater emphasize on the smoothness of the displacement, unrealistic

deformations could not be prevented as the training did not appear to converge.

### Changes to the Network Architecture

Next, a change to the network architecture was tested that changed the domain from image to k-space in hopes of avoiding the reconstruction artifacts present in highly accelerated images. For this we first looked at *LAPNet* which works on k-space data, but is trained in a supervised manner. As a first test, we tried to train it in an unsupervised method with the same pipeline as *Fourier-Net*. To achieve this we needed to generate  $33 \times 33$  crops of the k-spaces as *LAPNet* iteratively predicts the displacement for the center-pixel of each crop, thus producing a dense displacement for the whole image.

## 5.4 Integration into a Motion-Corrected Reconstruction Pipeline

Describe qualitative test to see whether the network can be used to improve a reconstruction pipeline for cardiac data.

### General Reconstruction Pipeline

Overview of the reconstruction pipeline using *Fourier-Net+*.

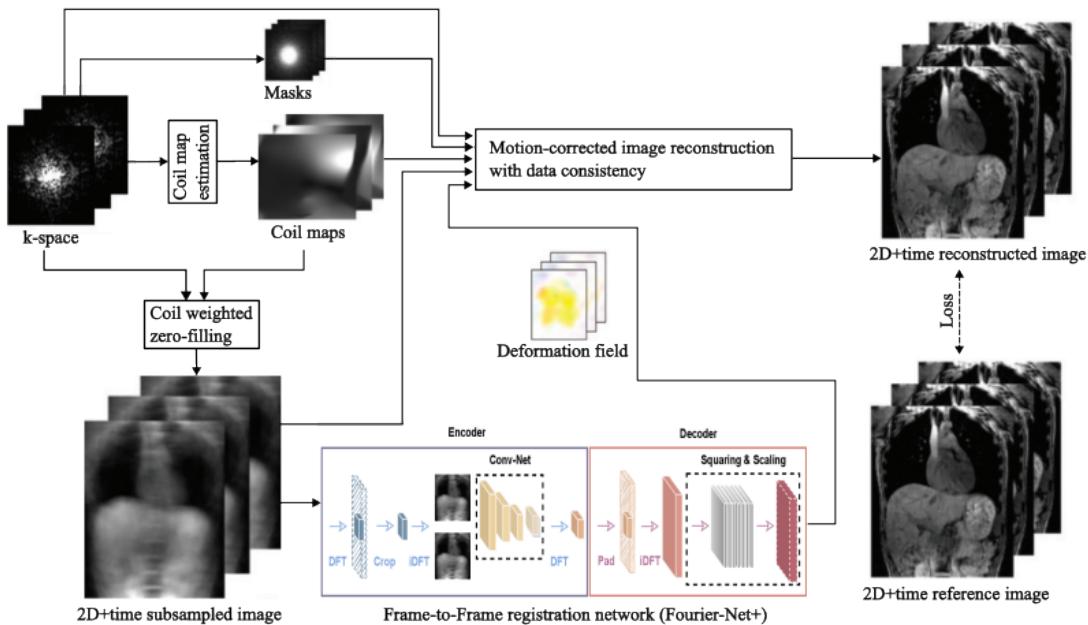


Figure 5.1: Overview of the general reconstruction pipeline using *Fourier-Net+* as a registration network.

# 6

## Results and Discussion

In this chapter, the results for the experiment described in the previous chapter will be examined. Possible implications will be discussed to determine the effectiveness of this new method.

### 6.1 Parameter Studies on the CMRxRecon Dataset

Different parameters were tested in order to find the optimal model parameters for usage on the *CMRxRecon* dataset. After some tests on the fully sampled data, a comparison with the traditional iterative algorithm *NiftyReg* was conducted on both the fully sampled and subsampled data.

#### Fourier-Net versus Fourier-Net+

First, the baseline network *Fourier-Net* was compared to its more efficient successor *4xFourier-Net+* with both using the diffeomorphic version as we can then compare the effect of the diffeomorphic deformation compared to the normal transformation on *4xFourier-Net+*. However, this effect should be similar for *Fourier-Net*. All of the experiments were conducted on the fully sampled data as the subsampling should not effect the tested properties. The baseline MSE (note that all MSE values are multiplied by  $10^{-3}$  as denoted by the m) for unregistered test images is 0.535 and the SSIM percentage (% SSIM) is 91.71. The results for *Fourier-Net* and *4xFourier-Net+* can be seen in Table 6.1.

The *Fourier-Net* model without the diffeomorphic transform was trained for 17 epochs, while the diffeomorphic version trained for 100 epochs. Despite this, the latter is only marginally better in terms of MSE and SSIM. The percentage of negative Jacobian determinants was actually better for the normal transformation compared to the diffeomorphic one. The normal *Fourier-Net+* model trained for 37 epochs, while the diffeomorphic version trained for 100 epochs. *Fourier-Net+* performs slightly worse than *Fourier-Net*, however it is more efficient.

Table 6.1: Results for *Fourier-Net* and *Fourier-Net+* on the fully sampled *CMRxRecon* test data with a diffeomorphic transform.

Model	MSE (m)	% SSIM	% $ J_\phi  \leq 0$	Time [s]
Fourier-Net	$0.319 \pm 0.352$	$93.29 \pm 3.85$	$0.0022 \pm 0.0183$	0.0025
4xFourier-Net+	$0.352 \pm 0.449$	$93.08 \pm 4.09$	$0.0003 \pm 0.0057$	0.0123
Diff-Fourier-Net	$0.318 \pm 0.388$	$93.39 \pm 3.82$	$0.0273 \pm 0.0897$	0.0028
Diff-4xFourier-Net+	$0.352 \pm 0.448$	$93.08 \pm 4.11$	$0.0029 \pm 0.0214$	0.0103

### Starting Channel Size

Next, the impact of the starting channels was examined. For this, *4xFourier-Net+* was trained with three different channel sizes: 8, 16, 32 and 64. Due to the early stopping the *4xFourier-Net+* variants were trained for 8, 17 and 10 epochs respectively. The results can be seen in Table 6.2 with the inference time for each version being calculated on the GPU. Despite training for the least amount of epochs channel size 32 performs the best for *4xFourier-Net+*.

Table 6.2: Results for three different starting channel sizes of *Fourier-Net+* on the fully sampled *CMRxRecon* test data.

Starting Channels	MSE (m)	% SSIM	% $ J_\phi  \leq 0$	Time [s]
8	$0.352 \pm 0.449$	$93.08 \pm 4.09$	$0.0003 \pm 0.0057$	0.0123
16	$0.347 \pm 0.439$	$93.12 \pm 4.04$	$0.0003 \pm 0.0059$	0.0113
32	$0.344 \pm 0.429$	$93.13 \pm 4.01$	$0.0008 \pm 0.0098$	0.0106

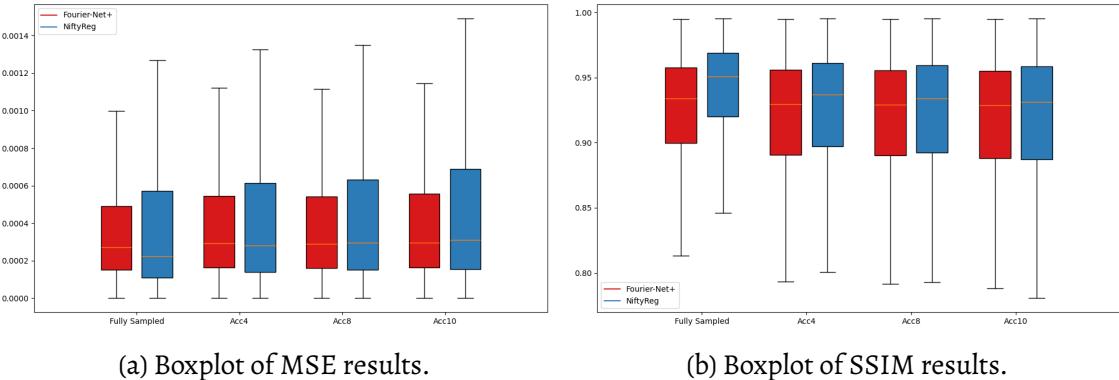
### Quantitative Comparison with NiftyReg

For this first quantitative test similarity metrics such as MSE and SSIM were used to judge image similarity as an indicator for proper image alignment. The results can be seen in Table 6.3 with the SSIM given in percent % and the MSE multiplied by  $10^{-3}$  for better readability (this is denoted by m for milli although the MSE has no unit strictly speaking). They are further visualized in Figure 6.4a for the MSE and Figure 6.4b for the SSIM results.

*Fourier-Net+* needed about 0.0104 seconds (s) per frame pair for registration with  $|J|_{<0} \% < 0.001$  for all levels of subsampling, while *NiftyReg* took up to 15 s for aligning an image pair. The results for *Fourier-Net+* look good on their own as the metrics are stable even for the subsampled data compared to the fully sampled data. However, its baseline performance is slightly worse than *NiftyReg* except for the Acc10 (R=10) data, where both perform worse than the unaligned baseline of the test data, however *4xFourier-Net+* is slightly better than *NiftyReg*. Overall, the network is much faster than the iterative registration algorithm as expected, but the latter performs better in terms of similarity metrics.

Table 6.3: Benchmark results for *Fourier-Net+* and *NiftyReg* on the *CMRxRecon* test data.

Subsampling	Method	MSE (m)	% SSIM
R=0	4xFourier-Net+	$0.435 \pm 0.554$	$92.35 \pm 4.59$
	NiftyReg	$0.514 \pm 0.771$	$93.97 \pm 3.91$
R=4	4xFourier-Net+	$0.482 \pm 0.626$	$91.81 \pm 5.01$
	NiftyReg	$0.551 \pm 0.798$	$92.38 \pm 4.94$
R=8	4xFourier-Net+	$0.479 \pm 0.620$	$91.74 \pm 5.12$
	NiftyReg	$0.565 \pm 0.803$	$92.06 \pm 5.09$
R=10	4xFourier-Net+	$0.491 \pm 0.640$	$91.59 \pm 5.23$
	NiftyReg	$0.926 \pm 2.867$	$90.66 \pm 10.32$

Figure 6.4: Boxplots for MSE and SSIM results of the test set for fully sampled and subsampled (Acc4, Acc8, Acc10) data for both *Fourier-Net+* and *NiftyReg*.

## 6.2 Further Tests on the ACDC Dataset

Due to the limitations of the *CMRxRecon* dataset in terms of segmentations, the registration performance is further evaluated on the *ACDC* dataset. Some of the parameter studies done on *CMRxRecon* were repeated, others were added to further test the networks ability.

### Fourier-Net versus Fourier-Net+

Again, *Fourier-Net*, *Fourier-Net+* and *4xFourier-Net* are compared in terms of registration performance and inference time (on GPU) together with their diffeomorphic versions, however, versions with dense instead of band-limited displacements are also added as a baseline for peak performance as these utilize all of the k-space information. The memory consumption is also added to the comparison with the number of parameters, mult-add operations and the total memory. The latter is only given for the dense and normal versions as the diffeomorphic version required roughly the same amount of memory as the normal ones, while the dense versions have a slightly larger encoder. Note that all net-

work variants were trained for 6 epochs. The results can be seen in Table 6.5. For the dense displacement variants the diffeomorphic transform increases the performance in terms of Dice score (with and without the background label) slightly for *Fourier-Net* and *4xFourier-Net*, but decreases for *Fourier-Net+*. The SSIM decreases for all models slightly, while the MSE stays about the same. Unsurprisingly, the percentage of non-positive Jacobian determinants goes down for the diffeomorphic variants. The times between the baseline and diffeomorphic versions do not vary a lot with the *Fourier-Net* being faster without the diffeomorphis, while *Fourier-Net+* and *4xFourier-Net* are faster.

The results for the band-limited versions are very similar with *Fourier-Net* and *Fourier-Net+* performing slightly worse in terms of Dice (with background) with the diffeomorphis, while *4xFourier-Net* performs better by almost 2%. When excluding the background label *Fourier-Net* again performs slightly worse, but *Fourier-Net+* and *4xFourier-Net* both improve, the latter with almost 3% more than the baseline version. These changes are not consistent with the dense displacement versions, however, the band-limited *4xFourier-Net* version might be an outlier, as all of the other models are not affected as much by the diffeomorphic transform. The SSIM values slightly decrease for *Fourier-Net* and *Fourier-Net+* with the diffeomorphism, but again improve slightly for *4xFourier-Net*. The MSE is slightly better for *Fourier-Net* and *4xFourier-Net*, but a bit worse for *Fourier-Net+*. The percentage of non-positive Jacobian determinants again decreases for the diffeomorphic variants, which is consistent with the previous observations on the dense displacement versions. The times again vary slightly with *Fourier-Net* and *Fourier-Net+* being a bit slower, while *4xFourier-Net* is quite a bit faster.

Now to the difference between the dense and band-limited displacement. Overall, the dense displacement versions perform better in terms of Dice score for all models both with and without the background label. This is to be expected as the dense displacement versions can use all of the available k-space data for the registration task, while the band-limited versions have only a limited amount of k-space data, directly impacting and limiting their performance. The dense displacement is also better in the SSIM metric, however the difference is not quite as large as with the Dice scores. As the differences between the frames is not quite as large, the weaker registration performance does not impact this metric as much as the Dice score, which is focused on the moving cardiac region. The same is true for the MSE, but the dense displacement is again far superior. Only in terms of the percentage of non-positive Jacobian determinants does the band-limited displacement actually perform better as the band-limiting probably allows for only a smaller deformation. However, the difference is still very small, especially for the diffeomorhpic versions, and it could be argued that a better registration performance at the cost of some image folding is an acceptable trade-off. In terms of inference time the dense displacement variants are surprisingly slightly faster despite having a larger encoder, however all of the model are in the low milliseconds (< 40 ms). Last but not least, the memory consumption needs to be addressed. The dense displacement variants have a larger encoder as discussed before and thus have a lot more parameters, especially for the more optimized *Fourier-Net+* and *4xFourier-Net* (about 5 times more). The number of Mult-Adds and the amount of total memory are more than doubled (almost tripled for *4xFourier-Net*) for the dense displacement versions compared to those with a band-limited displacement across all models.

## 6 Results and Discussion

Table 6.5: Results for *Fourier-Net*, *Fourier-Net+* and *4xFourier-Net+* with both dense and band-limited displacement fields as well as diffeomorphic transforms on the fully sampled ACDC test data.

Metrics	Dense Displacement		
	Fourier-Net	Fourier-Net+	4xFourier-Net+
% Dice	78.22 $\pm$ 13.84	78.43 $\pm$ 13.96	79.34 $\pm$ 14.03
% Dice*	71.03 $\pm$ 18.87	70.50 $\pm$ 19.17	72.02 $\pm$ 18.94
% SSIM	91.92 $\pm$ 3.32	89.09 $\pm$ 4.10	89.65 $\pm$ 3.85
MSE (m)	0.09 $\pm$ 0.06	0.17 $\pm$ 0.13	0.15 $\pm$ 0.12
% $ J_\phi  \leq 0$	0.53 $\pm$ 0.52	0.32 $\pm$ 0.55	0.04 $\pm$ 0.09
Time [s]	0.0070	0.0107	0.0219
Diff-Fourier-Net Diff-Fourier-Net+ Diff-4xFourier-Net+			
% Dice	78.56 $\pm$ 14.13	77.88 $\pm$ 13.81	79.49 $\pm$ 14.26
% Dice*	71.31 $\pm$ 19.13	70.19 $\pm$ 18.82	72.12 $\pm$ 19.33
% SSIM	91.79 $\pm$ 3.38	89.08 $\pm$ 4.10	89.62 $\pm$ 3.98
MSE (m)	0.09 $\pm$ 0.06	0.16 $\pm$ 0.13	0.15 $\pm$ 0.12
% $ J_\phi  \leq 0$	0.03 $\pm$ 0.06	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00
Time [s]	0.0082	0.0075	0.0203
Parameters	645,216	380,470	1,521,880
Mult-Add (G)	1.12	0.04	0.18
Memory [MB]	97.63	5.81	21.90
Band-limited Displacement			
Metrics	Fourier-Net Fourier-Net+ 4xFourier-Net+		
	Fourier-Net	Fourier-Net+	4xFourier-Net+
% Dice	77.95 $\pm$ 14.71	75.35 $\pm$ 14.28	76.59 $\pm$ 14.84
% Dice*	69.96 $\pm$ 21.55	64.74 $\pm$ 21.12	66.82 $\pm$ 21.55
% SSIM	91.35 $\pm$ 3.51	88.42 $\pm$ 3.94	88.59 $\pm$ 4.23
MSE (m)	1.00 $\pm$ 0.71	2.06 $\pm$ 1.54	1.92 $\pm$ 1.46
% $ J_\phi  \leq 0$	0.36 $\pm$ 0.38	0.13 $\pm$ 0.35	0.03 $\pm$ 0.12
Time [s]	0.0145	0.0145	0.0335
Diff-Fourier-Net Diff-Fourier-Net+ Diff-4xFourier-Net+			
% Dice	77.93 $\pm$ 14.93	75.17 $\pm$ 13.43	78.20 $\pm$ 14.01
% Dice*	69.91 $\pm$ 21.99	65.67 $\pm$ 18.61	69.61 $\pm$ 19.11
% SSIM	91.24 $\pm$ 3.56	87.81 $\pm$ 3.94	88.81 $\pm$ 4.13
MSE (m)	0.98 $\pm$ 0.67	2.29 $\pm$ 1.61	1.85 $\pm$ 1.41
% $ J_\phi  \leq 0$	0.01 $\pm$ 0.03	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00
Time [s]	0.0150	0.0207	0.0188
Parameters	434,519	75,429	301,716
Mult-Add (M)	595.04	10.98	43.91
Memory [MB]	44.69	2.25	7.66

## Starting Channel Size

Next, the impact of the starting channels was examined. For this, *Fourier-Net+* and *4xFourier-Net+* were trained with four different channel sizes: 8, 16, 32 and 64. The memory consumption as well as inference time on the GPU were calculated to extend the evaluation given by the other metrics. All network variants were trained for 6 epochs. The results can be seen in Table 6.6.

For both *Fourier-Net+* and *4xFourier-Net+* there is a clear increase in Dice (with and without the background label) with larger channel sizes as more features increase the registration performance. The SSIM and MSE metrics only increase marginally with larger channel size. The percentage of non-positive Jacobian determinants actually decreases as the registration performance increases, likely due to better/more features enabling a smoother and better displacement field to be generated by the models.

The inference time is not heavily impacted by channel size as the times remain very similar for all sizes. The memory however increases exponentially when doubling the channel sizes as all layers are effected, not just the starting one as the name suggests. Thus the number of parameters drastically increases with starting size leading to an increase in Mult-Add and thus overall memory. As *4xFourier-Net+* is just a cascaded version of *Fourier-Net+* a channel size of 16 for the latter is about the same as channel size 8 for the cascaded version and so on. Thus it seems that there is no point at which both the performance and memory consumption are maximized at the same time. This can also be seen in the direct comparison between *Fourier-Net+* and *4xFourier-Net+*, where the latter has a better performance, but also needs more memory. While an increase of channel size also increases performance in terms of Dice (with the background label) by 0.36%, 0.43% and 0.41% for the latter, which is quite consistent, it can be observed that an increase of channel size for *Fourier-Net+* has an bigger impact at the beginning (1.38%, 0.78% and 0.15%) and fades towards larger channels sizes. This effect can also be seen when excluding the background label (0.69%, 0.78% and 0.42% for *4xFourier-Net+* compared to 2.16%, 1.09% and 0.10% for *Fourier-Net+*), but is less pronounced for the SSIM and MSE as these do not capture the difference in the cardiac regions as well. Note that all versions of *Fourier-Net+* and *4xFourier-Net+* are smaller in terms of total memory than *Fourier-Net* with channel size 8 (44.69 MB), except for *4xFourier-Net+* with channel size 64.

## Fourier-Transform Crop Size

As a parameter specific to *Fourier-Net+* and *4xFourier-Net+*, the impact of the FT crop size used for compressing the images was analyzed. Four different sizes of the FT crop were used for training:  $80 \times 168$ ,  $40 \times 84$ ,  $48 \times 48$  and  $24 \times 24$ . The experiments were again only conducted on the fully sampled data and each network trained for 5 epochs. The results can be seen in Table 6.7.

The registration performance, as measured by the Dice score both with and without the background label, decreases drastically with a smaller crop size. This is as expected as the smaller FT drop compresses the images more thus making an accurate registration harder. The SSIM and MSE also get worse (SSIM decreases, MSE increases), however the amount of performance lost is hard to gauge, especially for the MSE metric. Interestingly,

## 6 Results and Discussion

Table 6.6: Results for different starting channel sizes of *Fourier-Net+* and *4xFourier-Net+* on the fully sampled ACDC test data.

Metrics	Starting Channels - Fourier-Net+			
	8	16	32	64
% Dice	75.50 $\pm$ 13.79	76.88 $\pm$ 13.86	77.66 $\pm$ 13.60	77.81 $\pm$ 13.76
% Dice*	65.70 $\pm$ 18.96	67.86 $\pm$ 19.06	68.95 $\pm$ 18.65	69.05 $\pm$ 18.80
% SSIM	88.05 $\pm$ 3.89	88.67 $\pm$ 3.88	88.83 $\pm$ 3.83	89.02 $\pm$ 3.67
MSE (m)	0.22 $\pm$ 0.16	0.20 $\pm$ 0.15	0.19 $\pm$ 0.15	0.19 $\pm$ 0.15
% $ J_\phi  \leq 0$	0.07 $\pm$ 0.25	0.04 $\pm$ 0.14	0.01 $\pm$ 0.04	0.00 $\pm$ 0.03
Time [s]	0.0072	0.0072	0.0080	0.0079
Parameters	75,429	300,477	1,199,469	4,793,037
Mult-Adds (M)	10.98	42.89	169.54	674.10
Memory [MB]	2.25	4.64	11.22	31.57
Metrics	Starting Channels - 4xFourier-Net+			
	8	16	32	64
% Dice	77.54 $\pm$ 13.73	77.90 $\pm$ 13.92	78.33 $\pm$ 14.14	78.74 $\pm$ 13.91
% Dice*	68.52 $\pm$ 18.62	69.21 $\pm$ 19.02	69.99 $\pm$ 19.21	70.41 $\pm$ 19.02
% SSIM	88.70 $\pm$ 4.17	88.89 $\pm$ 4.05	89.08 $\pm$ 4.01	89.29 $\pm$ 3.74
MSE (m)	0.19 $\pm$ 0.14	0.19 $\pm$ 0.14	0.18 $\pm$ 0.14	0.18 $\pm$ 0.14
% $ J_\phi  \leq 0$	0.02 $\pm$ 0.07	0.01 $\pm$ 0.04	0.01 $\pm$ 0.05	0.00 $\pm$ 0.00
Time [s]	0.0301	0.0244	0.0297	0.0277
Parameters	301,716	1,201,908	4,797,876	19,172,148
Mult-Adds (G)	0.04	0.17	0.68	2.70
Memory [MB]	7.66	17.23	43.56	124.94

the percentage of non-positive Jacobian determinants actually decreases for smaller crop sizes, perhaps because the displacements on the compressed images are not as extreme and thus more smooth. Note however, that the crop size of  $48 \times 48$  is a bit of an outlier in this regard breaking the trend with a higher percentage as  $40 \times 84$ . The time also seems to decrease slightly with a smaller crop, although this is quite noisy as the time from  $80 \times 168$  to  $40 \times 84$  is halved for *Fourier-Net+*, but the time needed for  $48 \times 48$  and  $24 \times 24$  increases again very slightly. For *4xFourier-Net+* the time decrease for all smaller crop sizes, except for  $24 \times 24$ .

The number of parameters does not change for different crop sizes as the networks themselves do not change, however the number of Mult-Adds and the total memory still change with the image size. Both decrease with a larger crop size as the image gets smaller. This effect is again not linear as a reduction in image size yields a larger reduction in memory going from  $80 \times 168$  to  $40 \times 84$  than from  $48 \times 48$  to  $24 \times 24$ . Again, there is no sweet-spot to be found that maximizes both memory efficiency and registration performance similar to the experiments in the previous chapter.

Table 6.7: Results for four different FT crop sizes for *Fourier-Net+* and *4xFourier-Net+* examined on the fully sampled ACDC test data.

Metrics	FT crop size - Fourier-Net+			
	80 × 168	40 × 84	48 × 48	24 × 24
% Dice	78.24 ± 14.43	76.61 ± 13.90	75.27 ± 13.49	73.77 ± 14.42
% Dice*	70.66 ± 19.70	67.60 ± 19.24	65.48 ± 18.76	63.37 ± 20.07
% SSIM	89.81 ± 4.09	88.58 ± 3.87	87.98 ± 3.91	87.00 ± 3.99
MSE (m)	0.15 ± 0.12	0.20 ± 0.15	0.22 ± 0.16	0.27 ± 0.19
% $ J_\phi  \leq 0$	0.22 ± 0.46	0.05 ± 0.15	0.09 ± 0.25	0.00 ± 0.02
Time [s]	0.0158	0.0077	0.0079	0.0081
Parameters	75,429	75,429	75,429	75,429
Mult-Adds (M)	64.03	16.37	10.98	2.74
Memory [MB]	9.51	2.97	2.25	1.12
Metrics	FT crop size - 4xFourier-Net+			
	80 × 168	40 × 84	48 × 48	24 × 24
% Dice	78.12 ± 15.01	77.49 ± 14.67	74.95 ± 14.15	72.58 ± 14.67
% Dice*	70.08 ± 21.92	68.03 ± 21.57	64.54 ± 20.75	61.19 ± 21.63
% SSIM	89.91 ± 4.01	89.03 ± 3.98	87.98 ± 3.94	87.02 ± 3.95
MSE (m)	1.45 ± 1.11	1.81 ± 1.37	2.22 ± 1.61	2.64 ± 1.84
% $ J_\phi  \leq 0$	0.12 ± 0.23	0.03 ± 0.15	0.06 ± 0.23	0.02 ± 0.15
Time [s]	0.0390	0.0301	0.0277	0.0289
Parameters	301,716	301,716	301,716	301,716
Mult-Adds (M)	256.14	65.50	43.91	10.98
Memory [MB]	36.70	10.54	7.66	3.15

### Comparison with VoxelMorph

In this test *Fourier-Net*, *Fourier-Net+* and *4xFourier-Net+* are compared to the famous *VoxelMorph* [6]. This unsupervised network brought deep learning based registration approaches into the mainstream and computes a dense displacement field in contrast to the band-limited displacement of the other networks. From the previous experiments a FT crop size of  $48 \times 48$  and channel size of 16 was chosen to strike a good balance between performance and memory efficiency. The experiment was only done on the fully sampled test data and all networks were trained for 6 epochs. The results can be seen in Table 6.8. *VoxelMorph* performs worst in terms of Dice score with the background label, however it does outperform *Fourier-Net+* in Dice without the background label. The best registration performance in terms of Dice (both with and without the background label) has *Fourier-Net* followed by *4xFourier-Net+* as expected. *VoxelMorph*, however, does perform best in terms of SSIM and MSE followed by *Fourier-Net*, so it seems that the dense network aligns the overall image well, but struggles to compensate the strong changes in the cardiac region between frames. It is conspicuous, however, that the percentage of non-positive Jacobian determinants for *VoxelMorph* is quite high with almost 90% while the other models are all under one percent. This is perhaps caused by the dense displacement which causes a better overall alignment (as indicated by the good SSIM and MSE values)

at the cost of more folding occurring. In terms of time *4xFourier-Net+* is slowest followed by *VoxelMorph*, *Fourier-Net* and *Fourier-Net+*, however all pretty fast ( $< 30ms$ ). The model with the least amount of parameters and Mult>Adds is *VoxelMorph*, however perhaps due to the dense displacement it needs quite a lot of memory with almost 40 MB. *Fourier-Net* is the largest model in terms of model paramters, Mult>Adds and total memory (90 MB). *Fourier-Net+* and *4xFourier-Net+* both have a higher number of paramters and Mult>Adds, but a lower total amount of memory needed. Overall, *4xFourier-Net+* is the only model that is truly more efficient and has better performance in terms of Dice than *VoxelMorph*, as *Fourier-Net+* is more efficient, but not necessarily better in terms of Dice score, while *Fourier-Net* is far superior in registration performance (as measured by the Dice score), but also less efficient in terms of memory consumption.

Table 6.8: Comparison of *Fourier-Net*, *Fourier-Net+*, *4xFourier-Net+* and *VoxelMorph* with similarity metrics and memory consumption on the fully sampled ACDC test data.

	Fourier-Net	Fourier-Net+	4xFourier-Net+	VoxelMorph
% Dice	$78.31 \pm 13.96$	$76.88 \pm 13.86$	$77.90 \pm 13.92$	$75.84 \pm 13.46$
% Dice*	$71.17 \pm 18.99$	$67.86 \pm 19.06$	$69.21 \pm 19.02$	$68.25 \pm 18.36$
% SSIM	$91.53 \pm 3.49$	$88.67 \pm 3.88$	$88.89 \pm 4.05$	$93.53 \pm 3.30$
MSE (m)	$0.09 \pm 0.07$	$0.20 \pm 0.15$	$0.19 \pm 0.14$	$0.06 \pm 0.04$
% $ J_\phi  \leq 0$	$0.44 \pm 0.45$	$0.04 \pm 0.14$	$0.01 \pm 0.04$	$89.96 \pm 1.32$
Time [s]	0.0099	0.0071	0.0244	0.0145
Parameters	1,735,447	300,477	1,201,908	84,322
Mult>Adds (G)	2.35	0.04289	0.17157	0.00157
Memory [MB]	90.08	4.64	17.23	39.04

### Dense Displacement on Accelerated Data

The difference between a dense displacement field and a band-limited one was already explored in section 6.2. However, the results in Table 6.5 only include fully sampled data, not accelerated data. These new tests again included *Fourier-Net*, *Fourier-Net+* and *4xFourier-Net+*, but are extended to subsampled data. Results for Acc4 data can be seen in Table 6.9, for Acc8 data in Table 6.10 and for Acc10 in Table 6.11.

For Acc4, the dense version of *Fourier-Net* performed worse than the band-limited version in terms of Dice, as we expected, however this is not the case for *Fourier-Net+* and *4xFourier-Net+* where the performance decreases for the band-limited versions as indicated by the dice score. The same is true for the SSIM and MSE values. The percentage of non-positive Jacobian determinants is higher for the band-limited *Fourier-Net*, while band-limited *Fourier-Net+* and *4xFourier-Net+* show no folding. While this at first seems like a positive metric it can also indicate a lack of complex displacements which could explain the lower performance. All band-limited versions are faster due to the smaller network size as the encoder lack some layers compared to the dense versions.

For Acc8 all band-limited networks perform worse in terms of Dice, SSIM and MSE. The percentage of non-positive Jacobian determinants is again higher for the band-limited

## 6 Results and Discussion

*Fourier-Net* compared to the dense version, while the band-limited *Fourier-Net+* and *4xFourier-Net+* show no folding. The band-limited versions surprisingly are slower compared to the dense versions except for *Fourier-Net+*.

For Acc10, similar to Acc4, band-limited *Fourier-Net* performs better in terms of Dice, SSIM and MSE, while band-limited *Fourier-Net+* and *4xFourier-Net+* are again worse compared to the dense versions. The percentage of non-positive Jacobian determinants for the band-limited *Fourier-Net* is lower than for the dense version this time, while band-limited *Fourier-Net+* and *4xFourier-Net+* again show no folding. As for the Acc4 data, the band-limited networks are faster in terms of inference time.

The results for the accelerated data are not quite consistent as the band-limited *Fourier-Net* performs better for Acc4 and Acc10, but not for Acc8. Band-limited *Fourier-Net+* and *4xFourier-Net+* perform consistently worse than their dense versions across all acceleration levels. The band-limiting of the displacement seems to help with avoiding folding and in most cases decrease inference time. Our hypothesis that the advantage of the dense displacement compared to the band-limited displacement disappears for accelerated data does not seem to hold except for *Fourier-Net* on Acc4 and Acc10 data.

Table 6.9: Results for *Fourier-Net*, *Fourier-Net+* and *4xFourier-Net+* with both dense and band-limited displacement fields on the Acc4 ACDC test data.

Metrics	Dense Displacement		
	Fourier-Net	Fourier-Net+	4xFourier-Net+
% Dice	$67.84 \pm 14.07$	$73.84 \pm 15.18$	$72.83 \pm 15.64$
% Dice*	$56.14 \pm 18.45$	$64.31 \pm 21.00$	$63.19 \pm 21.59$
% SSIM	$71.91 \pm 4.85$	$77.58 \pm 10.64$	$77.32 \pm 10.53$
MSE (m)	$0.28 \pm 0.12$	$0.20 \pm 0.14$	$0.21 \pm 0.15$
% $ J_\phi  \leq 0$	$0.37 \pm 0.37$	$0.15 \pm 0.26$	$0.15 \pm 0.27$
Time [s]	0.0082	0.0092	0.0337
Metrics	Band-limited Displacement		
	Fourier-Net	Fourier-Net+	4xFourier-Net+
% Dice	$72.30 \pm 14.04$	$70.18 \pm 17.43$	$70.36 \pm 16.63$
% Dice*	$61.87 \pm 19.32$	$59.36 \pm 23.90$	$59.47 \pm 22.80$
% SSIM	$78.77 \pm 7.22$	$75.78 \pm 10.63$	$75.91 \pm 10.72$
MSE (m)	$0.16 \pm 0.10$	$0.28 \pm 0.19$	$0.27 \pm 0.19$
% $ J_\phi  \leq 0$	$0.46 \pm 0.54$	$0.00 \pm 0.00$	$0.00 \pm 0.00$
Time [s]	0.0068	0.0079	0.0183

## 6 Results and Discussion

Table 6.10: Results for *Fourier-Net*, *Fourier-Net+* and *4xFourier-Net+* with both dense and band-limited displacement fields on the Acc8 ACDC test data.

Metrics	Dense Displacement		
	Fourier-Net	Fourier-Net+	4xFourier-Net+
% Dice	$69.12 \pm 14.05$	$74.20 \pm 14.50$	$74.36 \pm 14.97$
% Dice*	$57.99 \pm 18.66$	$64.50 \pm 20.57$	$65.01 \pm 20.82$
% SSIM	$77.27 \pm 3.74$	$86.08 \pm 4.19$	$85.53 \pm 3.99$
MSE (m)	$0.21 \pm 0.10$	$0.15 \pm 0.12$	$0.15 \pm 0.12$
% $ J_\phi  \leq 0$	$0.47 \pm 0.30$	$0.15 \pm 0.27$	$0.15 \pm 0.28$
Time [s]	0.0117	0.0169	0.0412

Metrics	Band-limited Displacement		
	Fourier-Net	Fourier-Net+	4xFourier-Net+
% Dice	$66.07 \pm 13.72$	$71.13 \pm 15.37$	$73.13 \pm 15.26$
% Dice*	$53.76 \pm 17.03$	$60.10 \pm 21.34$	$62.61 \pm 21.39$
% SSIM	$73.34 \pm 3.83$	$84.02 \pm 4.36$	$84.61 \pm 4.44$
MSE (m)	$0.26 \pm 0.09$	$0.21 \pm 0.16$	$0.19 \pm 0.15$
% $ J_\phi  \leq 0$	$0.78 \pm 0.37$	$0.00 \pm 0.00$	$0.01 \pm 0.04$
Time [s]	0.0138	0.0133	0.0418

Table 6.11: Results for *Fourier-Net*, *Fourier-Net+* and *4xFourier-Net+* with both dense and band-limited displacement fields on the Acc10 ACDC test data.

Metrics	Dense Displacement		
	Fourier-Net	Fourier-Net+	4xFourier-Net+
% Dice	$69.78 \pm 14.35$	$73.62 \pm 15.13$	$73.88 \pm 15.05$
% Dice*	$58.27 \pm 19.54$	$63.39 \pm 21.54$	$64.18 \pm 21.00$
% SSIM	$83.27 \pm 2.86$	$89.69 \pm 3.34$	$89.83 \pm 3.27$
MSE (m)	$0.16 \pm 0.10$	$0.14 \pm 0.12$	$0.13 \pm 0.12$
% $ J_\phi  \leq 0$	$0.37 \pm 0.31$	$0.21 \pm 0.39$	$0.12 \pm 0.19$
Time [s]	0.0112	0.0116	0.0433

Metrics	Band-limited Displacement		
	Fourier-Net	Fourier-Net+	4xFourier-Net+
% Dice	$72.28 \pm 14.65$	$71.83 \pm 15.91$	$72.55 \pm 15.66$
% Dice*	$62.58 \pm 19.99$	$61.10 \pm 22.07$	$62.00 \pm 21.97$
% SSIM	$90.68 \pm 2.93$	$87.77 \pm 3.50$	$89.13 \pm 3.47$
MSE (m)	$0.10 \pm 0.08$	$0.19 \pm 0.15$	$0.17 \pm 0.15$
% $ J_\phi  \leq 0$	$0.23 \pm 0.26$	$0.00 \pm 0.01$	$0.00 \pm 0.01$
Time [s]	0.0091	0.0088	0.0310

## Comparison on Subsampled Data

In section 5.2, *Fourier-Net*, *Fourier-Net+* and *4xFourier-Net+* were already compared to *VoxelMorph*, however, these comparisons were only done on fully sampled data, not accelerated ones. To further add to the comparison, *NiftyReg* was again used as a traditional registration algorithm. The unaligned test image pairs were used as a baseline for the lower bound of registration performance. Values which are worse than the baseline are marked in red, while the best results for each acceleration level is highlighted with blue. The results can be seen in Table 6.12. All times are computed on CPU due to *NiftyReg* only working on the CPU as the GPU capable versions were deprecated. Additionally, Figure 6.13 shows the performance of the methods for each segmentation label (excluding the background) while Figure 6.14 shows example images and segmentations warped by the different methods for a visual comparison.

For  $R = 0$ , *Fourier-Net* performs best in terms of Dice (both with and without the background label), closely followed by *4xFourier-Net+* while *NiftyReg* performs worse than the baseline. *VoxelMorph* performs best in terms of SSIM followed by *Fourier-Net* and *NiftyReg* as well as MSE where *NiftyReg* again performs worse than the baseline. *Fourier-Net* is the fastest method with under 0.1s per image pair, while *NiftyReg* takes over 100s making it the slowest by a large margin (all other methods are under 1s).

For  $R = 4$ , the registration performance decreases for all methods with *NiftyReg*, *Fourier-Net+* and *4xFourier-Net+* performing worse than the baseline, while *VoxelMorph* performs best in terms of Dice (both with and without the background label). *NiftyReg* performs best in terms of SSIM, while *Fourier-Net+* and *4xFourier-Net+* perform worse than the baseline. *NiftyReg* and *VoxelMorph* both perform best in terms of MSE with none of methods being worse than the baseline. In terms of time *NiftyReg* is again the worst method with about 80s, while all other methods need under 1s with *Fourier-Net+* again being the fastest. For  $R = 8$ , *4xFourier-Net+* performs best for Dice with the background label, while *VoxelMorph* performs best in term of Dice without the background label. *NiftyReg* and *Fourier-Net* perform worse than the baseline in terms of Dice (both with and without the background label), while *Fourier-Net+* performs a bit better then the baseline for Dice with the background, but slightly worse than the baseline for Dice without the background. *NiftyReg* again performs best for SSIM with *Fourier-Net*, *Fourier-Net+* and *4xFourier-Net+* performing worse than the baseline, however *VoxelMorph* is best in terms of MSE with only *4xFourier-Net+* performing worse than the baseline. *NiftyReg* is again the worst method in terms of time with over 80s while *Fourier-Net+* again is the fastest.

For  $R = 10$ , *4xFourier-Net+* performs best in terms of Dice with the background label, while *VoxelMorph* performs best in terms of Dice without the background. *NiftyReg* again performs worse than the baseline for Dice (both with and without background label), but is the best in terms of SSIM, while *VoxelMorph*, *Fourier-Net+* and *4xFourier-Net+* perform worse than the baseline. Both *VoxelMorph* and *Fourier-Net* perform best in terms of MSE with no method being worse than the baseline. *Fourier-Net+* is again the fastest method with under 0.1s, while *NiftyReg* is again the slowest despite its best time yet with about 47s.

To summarize, *Fourier-Net+* is the fastest method for all acceleration levels with *NiftyReg* being far off in terms of execution time as it is not machine learning based. *VoxelMorph*

## 6 Results and Discussion

is best in terms of MSE for all acceleration levels and is best in terms of Dice without the background label for subsampled data. *4xFourier-Net+* performs best in terms of Dice with the background label for  $R = 8$  and  $R = 10$ . *NiftyReg* is best in terms of SSIM for subsampled data, but performs badly in all other metrics. In general the performance decreases for all methods for higher acceleration levels as seen in the Dice metric, but, counter-intuitively, the SSIM is highest and the MSE the lowest for the  $R = 10$  baseline showing once again that these image similarity metrics are not very useful for evaluating registration performance across different acceleration levels.

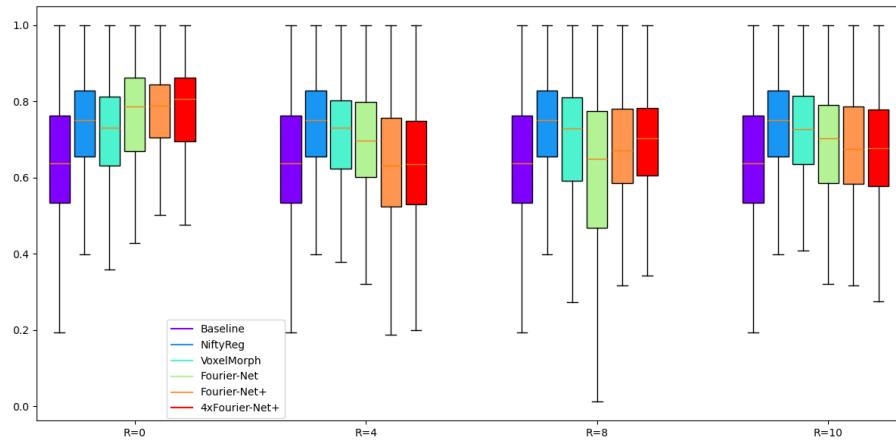
To further evaluate the difference between the methods performance one can look at the Dice scores of the three individual cardiac labels. The results can be seen as boxplots in Figure 6.13a for the left ventricle, Figure 6.13b for the myocardium and Figure 6.13c for the right ventricle.

As seen in Figure 6.13, the Dice scores for each label vary widely as the performance of all methods is best for the right ventricle and worst for the myocardium. *NiftyReg* performs well for the left ventricle, even having the best performance for the subsampled data, while it is the worst method for the myocardium (being only slightly better than the baseline) and the right ventricle (being worse than the baseline). *VoxelMorph* performs best for the subsampled data on the myocardium and is still very good on the left and right ventricle.

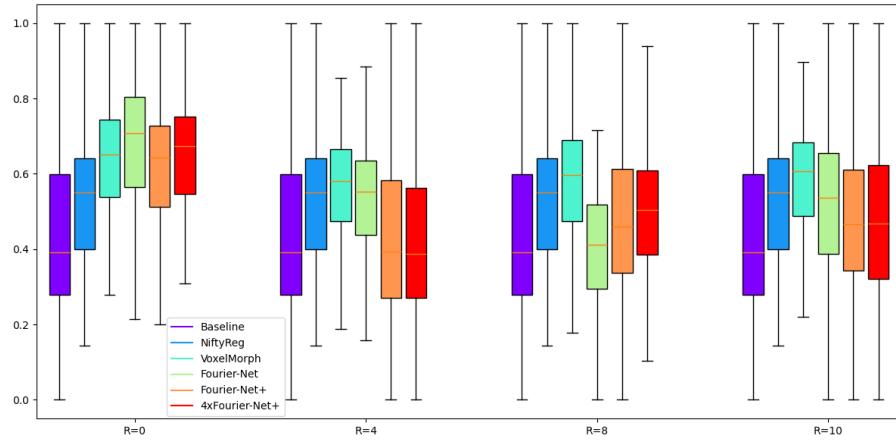
Table 6.12: Test results for *NiftiReg*, *VoxelMorph*, *Fourier-Net*, *Fourier-Net+* and *4xFourier-Net+* on the ACDC test data with an unaligned baseline for comparison. The best results for each metric and subsampling are highlighted in blue, while values worse than the unaligned baseline are marked with red.

	Method	% DICE	% DICE*	% SSIM	MSE (m)	Time [s]
$R=0$	Baseline	$70.85 \pm 18.27$	$60.35 \pm 25.24$	$86.39 \pm 4.08$	$0.33 \pm 0.23$	-
	NiftyReg	<b><math>70.74 \pm 13.77</math></b>	<b><math>57.56 \pm 18.86</math></b>	$91.26 \pm 3.18$	<b><math>1.94 \pm 1.61</math></b>	122.52
	VoxelMorph	$75.84 \pm 13.46$	$68.25 \pm 18.36$	<b><math>93.53 \pm 3.30</math></b>	<b><math>0.06 \pm 0.04</math></b>	0.1845
	Fourier-Net	<b><math>78.31 \pm 13.96</math></b>	<b><math>71.17 \pm 18.99</math></b>	$91.53 \pm 3.49$	$0.09 \pm 0.07$	0.1918
	Fourier-Net+	$76.88 \pm 13.86$	$67.86 \pm 19.06$	$88.67 \pm 3.88$	$0.20 \pm 0.15$	<b><math>0.0893</math></b>
	4xFourier-Net+	$77.90 \pm 13.92$	$69.21 \pm 19.02$	$88.89 \pm 4.05$	$0.19 \pm 0.14$	0.3262
$R=4$	Baseline	$70.85 \pm 18.27$	$60.35 \pm 25.24$	$76.80 \pm 11.02$	$0.28 \pm 0.19$	-
	NiftyReg	<b><math>69.89 \pm 13.73</math></b>	<b><math>56.47 \pm 17.86</math></b>	<b><math>86.03 \pm 6.41</math></b>	<b><math>0.15 \pm 0.14</math></b>	80.08
	VoxelMorph	<b><math>73.45 \pm 13.40</math></b>	<b><math>63.75 \pm 18.51</math></b>	$78.56 \pm 7.32$	<b><math>0.15 \pm 0.08</math></b>	0.1858
	Fourier-Net	$72.30 \pm 14.04$	$61.87 \pm 19.32$	$78.77 \pm 7.22$	$0.16 \pm 0.10$	0.1890
	Fourier-Net+	<b><math>70.18 \pm 17.43</math></b>	<b><math>59.36 \pm 23.90</math></b>	<b><math>75.78 \pm 10.63</math></b>	$0.28 \pm 0.19$	<b><math>0.0935</math></b>
	4xFourier-Net+	<b><math>70.36 \pm 16.63</math></b>	<b><math>59.47 \pm 22.80</math></b>	<b><math>75.91 \pm 10.72</math></b>	$0.27 \pm 0.19$	0.3462
$R=8$	Baseline	$70.85 \pm 18.27$	$60.35 \pm 25.24$	$85.35 \pm 4.43$	$0.22 \pm 0.17$	-
	NiftyReg	<b><math>70.04 \pm 13.42</math></b>	<b><math>56.37 \pm 17.63</math></b>	<b><math>91.07 \pm 2.80</math></b>	$0.12 \pm 0.13$	88.36
	VoxelMorph	$72.94 \pm 14.20$	<b><math>64.27 \pm 18.97</math></b>	$88.32 \pm 3.01$	<b><math>0.09 \pm 0.06</math></b>	0.1725
	Fourier-Net	<b><math>66.07 \pm 13.72</math></b>	<b><math>53.76 \pm 17.03</math></b>	<b><math>73.34 \pm 3.83</math></b>	<b><math>0.26 \pm 0.09</math></b>	0.2109
	Fourier-Net+	$71.13 \pm 15.37$	<b><math>60.10 \pm 21.34</math></b>	<b><math>84.02 \pm 4.36</math></b>	$0.21 \pm 0.16$	<b><math>0.1487</math></b>
	4xFourier-Net+	<b><math>73.13 \pm 15.26</math></b>	$62.61 \pm 21.39$	<b><math>84.61 \pm 4.44</math></b>	$0.19 \pm 0.15$	0.5373
$R=10$	Baseline	$70.85 \pm 18.27$	$60.35 \pm 25.24$	$89.17 \pm 3.58$	$0.20 \pm 0.17$	-
	NiftyReg	<b><math>70.40 \pm 13.34</math></b>	<b><math>56.61 \pm 17.71</math></b>	<b><math>93.47 \pm 2.37</math></b>	$0.11 \pm 0.13$	47.44
	VoxelMorph	$72.46 \pm 13.20$	<b><math>62.82 \pm 18.08</math></b>	<b><math>88.27 \pm 2.36</math></b>	<b><math>0.10 \pm 0.06</math></b>	0.1173
	Fourier-Net	$72.28 \pm 14.65$	$62.58 \pm 19.99$	$90.68 \pm 2.93$	<b><math>0.10 \pm 0.08</math></b>	0.1088
	Fourier-Net+	$71.83 \pm 15.91$	$61.10 \pm 22.07$	<b><math>87.77 \pm 3.50</math></b>	$0.19 \pm 0.15$	<b><math>0.0438</math></b>
	4xFourier-Net+	<b><math>72.55 \pm 15.66</math></b>	$62.00 \pm 21.97$	<b><math>89.13 \pm 3.47</math></b>	$0.17 \pm 0.15$	0.1539

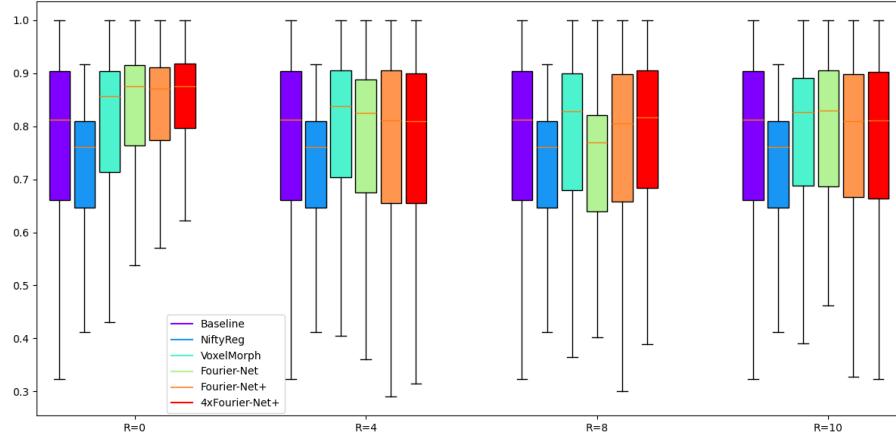
## 6 Results and Discussion



(a) Boxplot of the Dice scores for the left ventricle cavity.



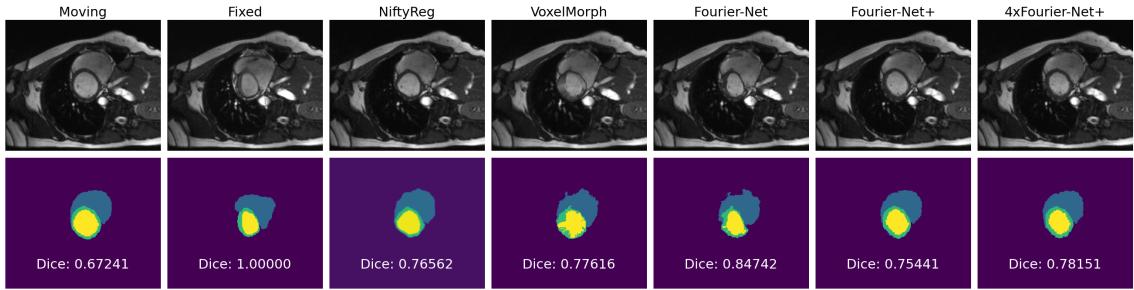
(b) Boxplot of the Dice scores for the myocardium.



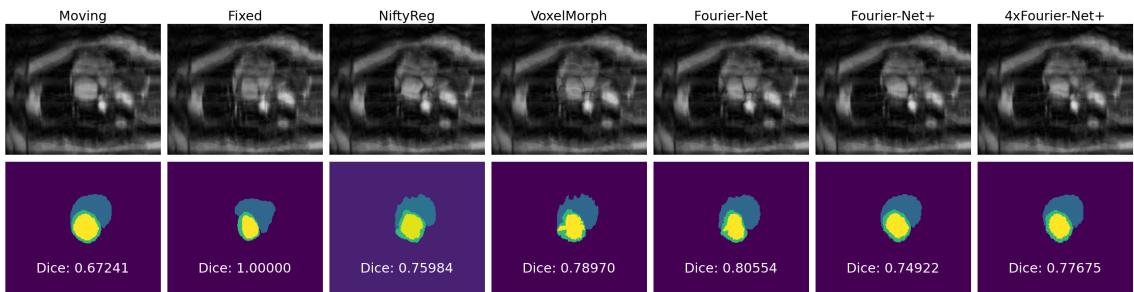
(c) Boxplot of the Dice scores for right ventricle cavity.

Figure 6.13: Boxplots of Dice scores (split by label excluding the background) for all models on fully sampled ( $R = 0$ ), Acc4 ( $R = 4$ ), Acc8 ( $R = 8$ ) and Acc10 ( $R = 10$ ) ACDC test data.

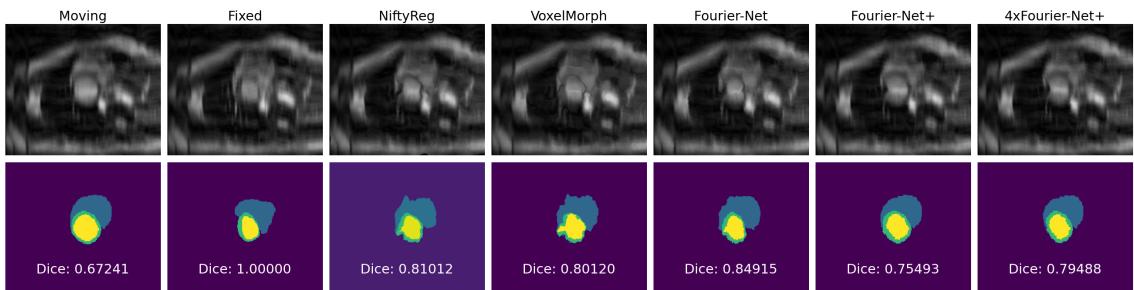
## 6 Results and Discussion



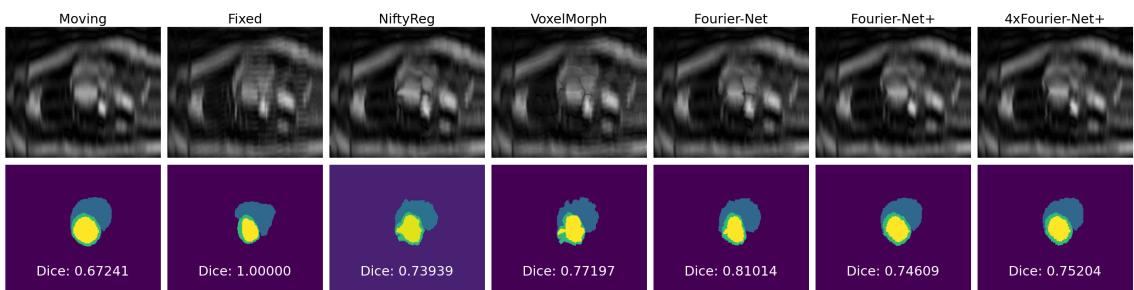
(a) Example images and segmentations for fully sampled ( $R = 0$ ) data.



(b) Example images and segmentations for Acc4 ( $R = 4$ ) data.



(c) Example images and segmentations for Acc8 ( $R = 8$ ) data.



(d) Example images and segmentations for Acc10 ( $R = 10$ ) data.

Figure 6.14: Examples of warped images and segmentations for *NiftyReg*, *VoxelMorph*, *Fourier-Net*, *Fourier-Net+* and *4xFourier-Net+* together with the original image pair from the ACDC test data.

# 7

## Conclusion

Summery of all stuff...

## Bibliography

- [1] Chen, X., Diaz-Pinto, A., Ravikumar, N., and Frangi, A. Deep learning in medical image registration. In: *Progress in Biomedical Engineering*, Dec. 2020. ISSN: 2516-1091. DOI: 10.1088/2516-1091/abd37c.
- [2] Haskins, G., Kruger, U., and Yan, P. Deep learning in medical image registration: a survey. In: *Machine Vision and Applications* 31(1–2), Jan. 2020. ISSN: 1432-1769. DOI: 10.1007/s00138-020-01060-x.
- [3] Fu, Y., Lei, Y., Wang, T., Curran, W. J., Liu, T., and Yang, X. Deep learning in medical image registration: a review. In: *Physics in Medicine & Biology* 65(20):20TR01, Oct. 2020. ISSN: 1361-6560. DOI: 10.1088/1361-6560/ab843e.
- [4] Zou, J., Gao, B., Song, Y., and Qin, J. A review of deep learning-based deformable medical image registration. In: *Frontiers in Oncology* 12, Dec. 2022. ISSN: 2234-943X. DOI: 10.3389/fonc.2022.1047215.
- [5] Chen, J., Liu, Y., Wei, S., Bian, Z., Subramanian, S., Carass, A., Prince, J. L., and Du, Y. A survey on deep learning in medical image registration: new technologies, uncertainty, evaluation metrics, and beyond. 2023. DOI: 10.48550/ARXIV.2307.15615.
- [6] Balakrishnan, G., Zhao, A., Sabuncu, M. R., Guttag, J., and Dalca, A. V. Voxel-Morph: A Learning Framework for Deformable Medical Image Registration. In: *IEEE Transactions on Medical Imaging* 38(8):1788–1800, Aug. 2019. ISSN: 1558-254X. DOI: 10.1109/tmi.2019.2897538.
- [7] Jia, X., Bartlett, J., Chen, W., Song, S., Zhang, T., Cheng, X., Lu, W., Qiu, Z., and Duan, J. Fourier-Net: Fast Image Registration with Band-Limited Deformation. In: *Proceedings of the AAAI Conference on Artificial Intelligence*. Vol. 37. 1. 2023, pp. 1015–1023.
- [8] Jia, X., Thorley, A., Gomez, A., Lu, W., Kotecha, D., and Duan, J. *Fourier-Net+: Leveraging Band-Limited Representation for Efficient 3D Medical Image Registration*. 2023. DOI: 10.48550/ARXIV.2307.02997.
- [9] Serai, S. D. Basics of magnetic resonance imaging and quantitative parameters T<sub>1</sub>, T<sub>2</sub>, T<sub>2</sub><sup>\*</sup>, T<sub>1</sub>rho and diffusion-weighted imaging. In: *Pediatric Radiology* 52(2):217–227, Apr. 2021. ISSN: 1432-1998. DOI: 10.1007/s00247-021-05042-7.
- [10] Hennig, J. K-space sampling strategies. In: *European Radiology* 9(6):1020–1031, July 1999. ISSN: 1432-1084. DOI: 10.1007/s003300050788.
- [11] Mansfield, P. Real-Time Echo-Planar Imaging By NMR. In: *British Medical Bulletin* 40(2):187–190, 1984. ISSN: 0007-1420. DOI: 10.1093/oxfordjournals.bmb.a071970.
- [12] Tam, A. L., Lim, H. J., Wistuba, I. I., Tamrazi, A., Kuo, M. D., Ziv, E., Wong, S., Shih, A. J., Webster, R. J., Fischer, G. S., et al. Image-Guided Biopsy in the Era of Personalized Cancer Care: Proceedings from the Society of Interventional Radi-

## Bibliography

- ology Research Consensus Panel. In: *Journal of Vascular and Interventional Radiology* 27(1):8–19, Jan. 2016. ISSN: 1051-0443. DOI: 10.1016/j.jvir.2015.10.019.
- [13] Chen, A. M., Hsu, S., Lamb, J., Yang, Y., Agazaryan, N., Steinberg, M. L., Low, D. A., and Cao, M. MRI-guided radiotherapy for head and neck cancer: initial clinical experience. In: *Clinical and Translational Oncology* 20(2):160–168, June 2017. ISSN: 1699-3055. DOI: 10.1007/s12094-017-1704-4.
- [14] Rigaud, B., Simon, A., Castelli, J., Lafond, C., Acosta, O., Haigron, P., Cazoulat, G., and Crevoisier, R. de Deformable image registration for radiation therapy: principle, methods, applications and evaluation. In: *Acta Oncologica* 58(9):1225–1237, June 2019. ISSN: 1651-226X. DOI: 10.1080/0284186x.2019.1620331.
- [15] Yang, D., Li, H., Low, D. A., Deasy, J. O., and Naqa, I. E. A fast inverse consistent deformable image registration method based on symmetric optical flow computation. In: *Physics in Medicine and Biology* 53(21):6143–6165, Oct. 2008. ISSN: 1361-6560. DOI: 10.1088/0031-9155/53/21/017.
- [16] Vercauteren, T., Pennec, X., Perchant, A., and Ayache, N. Diffeomorphic demons: Efficient non-parametric image registration. In: *NeuroImage* 45(1):S61–S72, Mar. 2009. ISSN: 1053-8119. DOI: 10.1016/j.neuroimage.2008.10.040.
- [17] Kingma, D. P. and Ba, J. *Adam: A Method for Stochastic Optimization*. 2014. DOI: 10.48550/ARXIV.1412.6980.
- [18] Ronneberger, O., Fischer, P., and Brox, T. U-Net: Convolutional Networks for Biomedical Image Segmentation. In: *Medical Image Computing and Computer-Assisted Intervention – MICCAI 2015*. Springer International Publishing, 2015, pp. 234–241. ISBN: 9783319245744. DOI: 10.1007/978-3-319-24574-4\_28.
- [19] Fan, J., Cao, X., Yap, P.-T., and Shen, D. BIRNet: Brain image registration using dual-supervised fully convolutional networks. In: *Medical Image Analysis* 54:193–206, May 2019. ISSN: 1361-8415. DOI: 10.1016/j.media.2019.03.006.
- [20] Küstner, T., Pan, J., Qi, H., Cruz, G., Gilliam, C., Blu, T., Yang, B., Gatidis, S., Botnar, R., and Prieto, C. LAPNet: Non-Rigid Registration Derived in k-Space for Magnetic Resonance Imaging. In: *IEEE Transactions on Medical Imaging* 40(12):3686–3697, Dec. 2021. ISSN: 1558-254X. DOI: 10.1109/tmi.2021.3096131.
- [21] Zhang, J. *Inverse-Consistent Deep Networks for Unsupervised Deformable Image Registration*. 2018. DOI: 10.48550/ARXIV.1809.03443.
- [22] Chen, J., Frey, E. C., He, Y., Segars, W. P., Li, Y., and Du, Y. TransMorph: Transformer for unsupervised medical image registration. In: *Medical Image Analysis* 82:102615, Nov. 2022. ISSN: 1361-8415. DOI: 10.1016/j.media.2022.102615.
- [23] Mok, T. C. and Chung, A. C. Fast Symmetric Diffeomorphic Image Registration with Convolutional Neural Networks. In: *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, June 2020. DOI: 10.1109/cvpr42600.2020.00470.
- [24] Modat, M., Ridgway, G. R., Taylor, Z. A., Lehmann, M., Barnes, J., Hawkes, D. J., Fox, N. C., and Ourselin, S. Fast free-form deformation using graphics processing units. In: *Computer Methods and Programs in Biomedicine* 98(3):278–284, June 2010. ISSN: 0169-2607. DOI: 10.1016/j.cmpb.2009.09.002.
- [25] Odille, F., Vuissoz, P., Marie, P., and Felblinger, J. Generalized Reconstruction by Inversion of Coupled Systems (GRICS) applied to free-breathing MRI. In: *Mag-*

## Bibliography

- Magnetic Resonance in Medicine* 60(1):146–157, June 2008. ISSN: 1522-2594. DOI: 10.1002/mrm.21623.
- [26] Küstner, T., Pan, J., Gilliam, C., Qi, H., Cruz, G., Hammernik, K., Blu, T., Rueckert, D., Botnar, R., Prieto, C., et al. Self-Supervised Motion-Corrected Image Reconstruction Network for 4D Magnetic Resonance Imaging of the Body Trunk. In: *APSIPA Transactions on Signal and Information Processing* 11(1), 2022. ISSN: 2048-7703. DOI: 10.1561/116.00000039.
- [27] Xuan, K., Sun, S., Xue, Z., Wang, Q., and Liao, S. Learning MRI k-Space Subsampling Pattern Using Progressive Weight Pruning. In: *Lecture Notes in Computer Science*. Springer International Publishing, 2020, pp. 178–187. ISBN: 9783030597139. DOI: 10.1007/978-3-030-59713-9\_18.
- [28] Yiasemis, G., Zhang, C., Sánchez, C. I., Sonke, J.-J., and Teuwen, J. Deep MRI reconstruction with radial subsampling. In: *Medical Imaging 2022: Physics of Medical Imaging*. Ed. by W. Zhao and L. Yu. SPIE, Apr. 2022. DOI: 10.1117/12.2609876.
- [29] Yiasemis, G., Sánchez, C. I., Sonke, J.-J., and Teuwen, J. On Retrospective k-space Subsampling schemes For Deep MRI Reconstruction. 2023. DOI: 10.48550/ARXIV.2301.08365.
- [30] Marcus, D. S., Wang, T. H., Parker, J., Csernansky, J. G., Morris, J. C., and Buckner, R. L. Open Access Series of Imaging Studies (OASIS): Cross-sectional MRI Data in Young, Middle Aged, Nondemented, and Demented Older Adults. In: *Journal of Cognitive Neuroscience* 19(9):1498–1507, Sept. 2007. ISSN: 1530-8898. DOI: 10.1162/jocn.2007.19.9.1498.
- [31] Hoopes, A., Hoffmann, M., Fischl, B., Guttag, J., and Dalca, A. V. HyperMorph: Amortized Hyperparameter Learning for Image Registration. In: *Information Processing in Medical Imaging*. Springer International Publishing, 2021, pp. 3–17. ISBN: 9783030781910. DOI: 10.1007/978-3-030-78191-0\_1.
- [32] Hering, A., Hansen, L., Mok, T. C. W., Chung, A. C. S., Siebert, H., Hager, S., Lange, A., Kuckertz, S., Heldmann, S., Shao, W., et al. Learn2Reg: Comprehensive Multi-Task Medical Image Registration Challenge, Dataset and Evaluation in the Era of Deep Learning. In: *IEEE Transactions on Medical Imaging* 42(3):697–712, Mar. 2023. ISSN: 1558-254X. DOI: 10.1109/tmi.2022.3213983.
- [33] Wang, C., Lyu, J., Wang, S., Qin, C., Guo, K., Zhang, X., Yu, X., Li, Y., Wang, F., Jin, J., et al. CMRxRecon: An open cardiac MRI dataset for the competition of accelerated image reconstruction. 2023. DOI: 10.48550/ARXIV.2309.10836.
- [34] Pruessmann, K. P., Weiger, M., Scheidegger, M. B., and Boesiger, P. SENSE: Sensitivity encoding for fast MRI. In: *Magnetic Resonance in Medicine* 42(5):952–962, Nov. 1999. ISSN: 1522-2594. DOI: 10.1002/(sici)1522-2594(199911)42:5<952::aid-mrm16>3.0.co;2-s.
- [35] Bernard, O., Lalande, A., Zotti, C., Cervenansky, F., Yang, X., Heng, P.-A., Cetin, I., Lekadir, K., Camara, O., Gonzalez Ballester, M. A., et al. Deep Learning Techniques for Automatic MRI Cardiac Multi-Structures Segmentation and Diagnosis: Is the Problem Solved? In: *IEEE Transactions on Medical Imaging* 37(11):2514–2525, Nov. 2018. ISSN: 1558-254X. DOI: 10.1109/tmi.2018.2837502.
- [36] Dalca, A. V., Balakrishnan, G., Guttag, J., and Sabuncu, M. R. Unsupervised Learning for Fast Probabilistic Diffeomorphic Registration. In: *Lecture Notes in*

## Bibliography

- Computer Science*. Springer International Publishing, 2018, pp. 729–738. ISBN: 9783030009281. DOI: 10.1007/978-3-030-00928-1\_82.
- [37] Trabelsi, C., Bilaniuk, O., Zhang, Y., Serdyuk, D., Subramanian, S., Santos, J. F., Mehri, S., Rostamzadeh, N., Bengio, Y., and Pal, C. J. *Deep Complex Networks*. 2017. DOI: 10.48550/ARXIV.1705.09792.
- [38] Wang, J. and Zhang, M. *DeepFLASH: An Efficient Network for Learning-based Medical Image Registration*. 2020. DOI: 10.48550/ARXIV.2004.02097.
- [39] Ashburner, J. A fast diffeomorphic image registration algorithm. In: *NeuroImage* 38(1):95–113, Oct. 2007. ISSN: 1053-8119. DOI: 10.1016/j.neuroimage.2007.07.007.
- [40] Arsigny, V., Commowick, O., Pennec, X., and Ayache, N. A Log-Euclidean Framework for Statistics on Diffeomorphisms. In: *Lecture Notes in Computer Science*. Springer Berlin Heidelberg, 2006, pp. 924–931. ISBN: 9783540447085. DOI: 10.1007/11866565\_113.
- [41] Jaderberg, M., Simonyan, K., Zisserman, A., and Kavukcuoglu, K. Spatial Transformer Networks. In: *ArXiv* abs/1506.02025, 2015. URL: <https://api.semanticscholar.org/CorpusID:6099034>.