1The Optimisation of Deep Neural Networks for2Segmenting Multiple Knee Joint Tissues from MRIs

4	
5	Dimitri A Kessler ^{a‡} , James W MacKay ^{a,b} , Victoria Crowe ^c , Frances Henson ^d ,
6	Martin J Graves ^c , Fiona J Gilbert ^{a*} and Joshua D Kaggie ^{a*‡}
7	
8	^a Department of Radiology, University of Cambridge, Cambridge, United Kingdom
9	^b Norwich Medical School, University of East Anglia, Norwich, United Kingdom
10 11	^c Cambridge University Hospitals NHS Foundation Trust, Addenbrooke's Hospital, Cambridge, United Kingdom.
12	^d Department of Veterinary Medicine, University of Cambridge, United Kingdom
13	
14	* JDK and FJG share senior authorship to this work.
15	
16	‡Corresponding Authors:
17	Dimitri A Kessler, M.Sc., Department of Radiology, University of Cambridge, School of
18	Clinical Medicine, Box 218, Cambridge Biomedical Campus, Cambridge, United Kingdom,
19	CB2 0QQ.
20	Tel: +44-7711266011, E-Mail: dak50@cam.ac.uk
21	
22	Joshua Kaggie, Ph.D., Department of Radiology, University of Cambridge, School of
23	Clinical Medicine, Box 218, Cambridge Biomedical Campus, Cambridge, United Kingdom,
24 25	CB2 UQQ. Tal: ± 44.1222746444 E Mail: $ik626@aam aa uk$
25 26	Tel. +44-1225/40444, E-Mail. JK050@call.ac.uk
27	
28	Submitted to Computerized Medical Imaging and Graphics as a Full Paper
29	
30	
31	
32	
33	
34	
35	
36	

The Optimisation of Deep Neural Networks for Segmenting Multiple Knee Joint Tissues from MRIs

3

4 Abstract

Automated semantic segmentation of multiple knee joint tissues is desirable to allow faster
and more reliable analysis of large datasets and to enable further downstream processing e.g.
automated diagnosis.

8

9 In this work, we evaluate the use of conditional Generative Adversarial Networks (cGANs) 10 as a robust and potentially improved method for semantic segmentation compared to other 11 extensively used convolutional neural network, such as the U-Net. As cGANs have not yet 12 been widely explored for semantic medical image segmentation, we analysed the effect of 13 training with different objective functions and discriminator receptive field sizes on the 14 segmentation performance of the cGAN. Additionally, we evaluated the possibility of using transfer learning to improve the segmentation accuracy. The networks were trained on i) the 15 16 SKI10 dataset which comes from the MICCAI grand challenge "Segmentation of Knee 17 Images 2010", ii) the OAI ZIB dataset containing femoral and tibial bone and cartilage segmentations of the Osteoarthritis Initiative cohort and iii) a small locally acquired dataset 18 (Advanced MRI of Osteoarthritis (AMROA) study) consisting of 3D fat-saturated spoiled 19 gradient recalled-echo knee MRIs with manual segmentations of the femoral, tibial and 20 21 patellar bone and cartilage, as well as the cruciate ligaments and selected peri-articular 22 muscles. The Sørensen–Dice Similarity Coefficient (DSC), volumetric overlap error (VOE) 23 and average surface distance (ASD) were calculated for segmentation performance evaluation. 24

25

DSC \geq 0.95 were achieved for all segmented bone structures, DSC \geq 0.83 for cartilage and muscle tissues and DSC of \approx 0.66 were achieved for cruciate ligament segmentations with both cGAN and U-Net on the in-house AMROA dataset. Reducing the receptive field size of the cGAN discriminator network improved the networks segmentation performance and resulted in segmentation accuracies equivalent to those of the U-Net. Pretraining not only increased segmentation accuracy of a few knee joint tissues of the fine-tuned dataset, but also

1	increased the network's capacity to preserve segmentation capabilities for the pretrained
2	dataset.
3	
4	cGAN machine learning can generate automated semantic maps of multiple tissues within the
5	knee joint which could increase the accuracy and efficiency for evaluating joint health.
6	
7	Key Words: magnetic resonance imaging (MRI); musculoskeletal; image segmentation;
8	convolutional neural network (CNN); generative adversarial network (GAN)
9	
10	
11	
12	
13	
14	
15	
16	
17	
18	
19	
20	
21	
22	
23	
24	
25	
26	
27	
28	
29	
30	
31	
32	
33	
34	

1 **1. Introduction**

2 Osteoarthritis (OA) is a degenerative disease involving the entire synovial joint (Goldring et al., 2017; Hunter and Eckstein, 2009; Martel-Pelletier et al., 2016). 3 Important risk factors for the development of OA include age, muscle weakness, 4 5 abnormal joint loading due to joint malalignment or overloading (obesity, high impact sport), and injury to the menisci and ligaments (Ismail and Vincent, 2017; Lohmander et 6 7 al., 2007; Martel-Pelletier et al., 2016). Distinctive hallmarks of OA include the 8 progressive destruction of articular cartilage structure and alterations in the surrounding 9 joint tissues, including bone, meniscus, ligament and peri-articular muscle. Magnetic Resonance Imaging (MRI) is a commonly used tool to evaluate clinical abnormalities of the 10 11 knee (Blumenkrantz and Majumdar, 2016). Morphological changes due to OA are well demonstrated with MRI (Benhamou et al., 2001; Hunter et al., 2015; MacKay et al., 2018; 12 13 Neogi et al., 2013; Wise et al., 2018). Tissue specific masks of the knee joint can be useful 14 for the analysis of OA, especially as automated tools continue to be developed and validated 15 (Bindernagel et al., 2011; Deniz et al., 2018; Lee et al., 2014; Liu et al., 2017; Ng et al., 2006; Patel and Singh, 2018; Seim et al., 2010; Shan et al., 2014; Shrivastava et al., 2014; Swanson 16 17 et al., 2010; Xia et al., 2013; Zhou et al., 2016).

18

For both clinical and research usage, a significant amount of time is spent manually 19 20 segmenting images to designate tissue-specific regional masks, also known as regions-of-21 interest (ROIs). Image masking remains a very significant challenge within medical imaging 22 due to heterogeneity in organ appearance and disease progression and presentation. The 23 segmentation of neighbouring soft tissues such as the cruciate ligaments, cartilages and muscles in the knee joint which have similar image intensities (and therefore poor contrast 24 25 resolution) is an especially demanding task. ROIs can be generated through manual or semi-26 manual delineation by a trained reader, or they may be generated automatically using signal 27 thresholding (Swanson et al., 2010), shape (Bindernagel et al., 2011; Seim et al., 2010), atlas 28 (Lee et al., 2014; Shan et al., 2014), or derive from region based (Ng et al., 2006; Patel and 29 Singh, 2018; Shrivastava et al., 2014) approaches, as well as with machine learning 30 approaches (Deniz et al., 2018; Liu et al., 2017; Xia et al., 2013; Zhou et al., 2016). Machine learning methods include unsupervised learning, such as k-means clustering, which segments 31 32 based on spatial clusters of similar signal intensities in an image (Ng et al., 2006; Patel and Singh, 2018; Shrivastava et al., 2014), or supervised learning by training the algorithm on 33

image masks that have been obtained from any previous masking technique (Deniz et al., 2018; Liu et al., 2017; Xia et al., 2013; Zhou et al., 2016). The number of high-quality label maps for supervised learning is typically very small, and the performance of a machine learning network trained on a low number of data is limited due to the lack of heterogeneity of images presented during training. Transfer learning may be used to mitigate this by pretraining a network on a large dataset with different but related similarities to the actual task, followed by network refinement on the small dataset (Shie et al., 2015).

8

9 Convolutional neural networks (CNNs), in particular U-Nets (Ronneberger et al., 2015), have 10 demonstrated their capability to automate the segmentation of musculoskeletal MRIs (Liu et al., 2017; Norman et al., 2018). Nevertheless, a drawback of this approach with CNNs is that 11 they usually use pixel-wise measures such as the absolute (L1) or square (L2) error loss 12 which can be non-optimal for image data, and, in the case of L2, result in blurry boundaries 13 14 (Pathak et al., 2016). In contrast, generative adversarial networks (GANs) (Goodfellow et al., 15 2014) learn a similarity measure (feature-wise metric) that adapts to the training task by 16 implementing two competing, or adversarial, neural networks. During adversarial training, 17 one network focusses on image discrimination and guides a second network which focusses 18 on image generation to create "real" images that have a data distribution indistinguishable from the training data distribution. The generator and discriminator are trained 19 20 simultaneously and competitively in a mini-max game while convergence is achieved when 21 the Nash equilibrium is reached, i.e. no network can improve through further training if one 22 remains unchanged (Zhao et al., 2017).

23

24 Conditional GANs (cGANs) modify the GAN approach to learn image-to-image mappings 25 (Goodfellow et al., 2014; Isola et al., 2017). In comparison to traditional GANs that learn a 26 mapping from random noise to a generated output, cGANs learn a mapping from an observed 27 variable, for example an image to generate an output, such as a label map (Goodfellow et al., 2014; Isola et al., 2017). cGANs have been used to produce image labels for neurological 28 (Rezaei et al., 2017), cardiac (Dou et al., 2018), abdominal (Huo et al., 2018), respiratory 29 30 (Chen et al., 2018) and musculoskeletal imaging (Liu, 2018, Gaj et al., 2019). (Liu, 2019) used unpaired image-to-image translation with a method called cycle-consistent generative 31 32 adversarial network (CycleGAN) to perform semantic image segmentation of femorotibial 33 cartilage and bone of the knee joint of unlabelled MRI datasets. The "pix2pix" framework is 34 one cGAN approach that has demonstrated segmentation capability (Isola et al., 2017).

Semantic segmentation with cGANs, particularly those combining U-Net generators and
 Markov Random Field discriminators (patch-based discriminators), is relatively unexplored.
 The method has previously been performed for semantic segmentation of the brain (Rezaei et al., 2017). In (Gaj et al., 2019), a cGAN was used for semantic segmentation of knee cartilage
 and meniscus but with an image-wise discriminator rather than a patch-wise discriminator.

7 The aim of this study was to implement and evaluate a cGAN for automated semantic 8 segmentation of multiple joint tissues from MR images: the femoral, tibial and patellar bones 9 and cartilage surfaces; the cruciate ligaments; and two selective muscles, the medial vastus 10 and gastrocnemius. Our essential contributions are summarised as followed:

- 11
- Implementation of a cGAN based on the "pix2pix" framework introduced by (Isola et al., 2017) using a U-Net generator and a patch-based discriminator for automatic segmentation of multiple knee joint tissues. As far as we know, cGANs have not previously been used for semantic segmentation of the patellar bone and cruciate ligaments, as well as muscles of the knee joint.
- Evaluating the segmentation performance of the cGAN with different objective
 functions by combining the cGAN loss with different pixel-wise error losses and
 modifying the weighting hyperparameter between the cGAN loss and pixel-wise
 error loss.
- 22

25

28

- 3. Assessing the choice of the generator depth and discriminator receptive field size on
 the performance of the cGAN for multi-tissue segmentation.
- 26 4. Quantitative comparison of the cGAN approach with the well-known U-Net27 approach.
- 5. Exploring the use of transfer learning for improved segmentation performance of
 both cGAN and U-Net.
- 31
- 32
- 33
- 34

1 2. Material and Methods

2 2.1 Image datasets

Three image datasets were used for network training and testing; the publicly available SKI10 and OAI ZIB datasets, consisting of 100 and 507 labelled knee MRs, respectively, and a locally acquired dataset of ten segmented knee MRs (Advanced MRI of Osteoarthritis (AMROA) study).

7

8 2.1.1 SKI10

9 The "Segmentation of Knee Images 2010" (SKI10) dataset (Heimann et al., 2010), consists of 10 approximately 90% 1.5T and 10% 3.0T sagittal MR images using multiple system vendors -GE, Siemens, Philips, Toshiba, and Hitachi. The sequences were varied and included both 11 12 gradient echo and spoiled gradient echo sequences, commonly with fat suppression. The 13 images were segmented on a slice-by-slice basis by experts from Biomet, Inc., initially 14 through intensity thresholds and thereafter with manual editing. One hundred 3D image 15 datasets of the SKI10 challenge were provided with semi-manual masks of femoral and tibial cartilage and bone. In our study, 70 datasets were used for network training and 30 for 16 17 network testing.

18

19 2.1.2 OAI ZIB

The OAI ZIB dataset (Ambellan et al., 2019) is comprised of segmentations of femoral and 20 tibial cartilage and bone of 507 MR imaging volumes from the publicly available 21 22 Osteoarthritis Initiative dataset ("The Osteoarthritis Initiative," n.d.). The MR images were acquired on Siemens 3T Trio systems using a 3D double echo steady state (DESS) sequence 23 24 with water excitation. Outlines of femoral and tibial bone and cartilage were generated using 25 a statistical shape model (Seim et al., 2010) with manual adjustments performed by experts at Zuse Institute Berlin. The OAI ZIB data covers all degrees of OA (KL 0 - 4), with more 26 27 cases having severe OA (KL \geq 3) (Ambellan et al., 2019). As with the SKI10 dataset, we split 28 the dataset in 70% (355) for network training and 30% (152) for testing.

29

30 2.13 AMROA

The locally acquired participant cohort consisted of ten subjects: five healthy volunteers and five patients with mild-to-moderate OA. The patients followed at least one subset of 1 American College of Rheumatology criteria for OA and were recruited between April 2017 2 to April 2018 (Table 1). The healthy volunteers were approximately matched to OA patients 3 for age, sex, and body mass. Network training was performed on data from four subjects with 4 OA and four healthy subjects. Two individuals (one with OA and one healthy) were used as a 5 unique set for test measurements. The number of test individuals was chosen such that roughly 80% of the data could be used for training. Ethical approval was obtained from the 6 7 National Research Ethics Service, and all subjects provided written informed consent before 8 participation.

9

The source images (Fig. 2A) for each subject were 3D fat-saturated spoiled gradient recalled-10 echo (3D-FS SPGR) images and were acquired on a 3.0T MRI system (MR750, GE 11 Healthcare, Waukesha, WI, USA) using an 8-channel transmit/receive knee coil (InVivo, 12 13 Gainesville, FL, USA). The 3D-FS SPGR sequence parameters were: field-ofview=150x128x136 mm³, matrix size=512x380x136 zero-fill interpolated to 512x512x136, 14 voxel size= $0.29 \times 0.29 \times 1.0 \text{ mm}^3$, TR = 12.5 ms, TE = 2.4 ms, flip angle = 25°, coil 15 acceleration factor (ASSET) = 2, partial Fourier phase encoding = 0.5 (half-NEX), bandwidth 16 17 $= \pm 11.9$ kHz, with fat-suppression.

18

Semi-manual segmented masks (Fig. 2A) of the patella, tibia, and femur bones as well as of 19 20 their respective surrounding patellar, tibial and femoral cartilages (Fig. 2b) were created from the 3D-FS SPGR images by a musculoskeletal radiologist with 8 years' experience, using the 21 22 Stradwin software v5.4a (University of Cambridge Department of Engineering, Cambridge, UK, now freely available as 'StradView' at http://mi.eng.cam.ac.uk/Main/StradView/) 23 24 (MacKay et al., 2020). Additionally, masks of the vastus medialis and medial head of 25 gastrocnemius muscles were created. This semi-manual segmentation pipeline consists of sparse manual contour generation (every 2nd-5th sagittal image/2-5 mm) followed by 26 27 automatic surface triangulation using the regularised marching tetrahedra method. Volume preserving surface smoothing allows creation of an accurate segmentation from relatively 28 sparse manual contours (Treece et al., 1999). Manual segmentations of the anterior cruciate 29 ligament (ACL) and posterior cruciate ligament (PCL) were created on the 3D-FS SPGR 30 images using ITK SNAP (Yushkevich et al., 2006) by a radiologist with 3 years' experience. 31

1 2.2 Training Data and Masking

2 Each of the major structures were given a separate image value, i.e., colour, in the 3 segmentation mask, such that the network determined the unique weights to generate a similar regional colour-value from an MR image. On a 256-bit colour-scale, the three bones 4 5 were stored in the blue colour channel where the femur colour code was 50, tibia was 100, 6 and patella was 150. The cartilages were stored in the green colour channel where the femoral 7 cartilage colour code was 50, the tibial was 100 and the patellar was 150. Additionally, for 8 the AMROA dataset, the muscles were stored in the red colour channel with the medial 9 vastus muscle code set to 100 and the medial gastrocnemius muscle colour code set to 200. 10 The ACL mask was stored in the blue colour channel and the PCL in the green colour 11 channel with both colour codes set to 200.

12

The MRIs and image masks were converted from the DICOM and NIFTI formats (Larobina and Murino, 2014), respectively, to a common image format (Portable Network Graphics, PNG) before training. Noise-only images were not used for training or testing, as training a network to fit against zero-valued masks results in a poor constraint. After network training, a tissue- / region-specific Boolean mask was created on the predicted test images by removing prediction values outside of ± 20 colour scale units of the tissue specific value. 3D mask predictions were obtained by iterating over the 2D segmented slices.

20

21 **2.3 Network Specifications**

This work uses the "pix2pix" framework of a conditional GAN (cGAN) described by Nvidia (Isola et al., 2017). The cGAN consists of two deep neural networks, a generator (*G*) and a discriminator (*D*). For our task, *G* learns to translate sagittal MR images of the knee joint (source images *x*) to semantic segmentation maps (G(x)), while *D* aims to differentiate between the real segmentation map (*y*) and the synthetically generated.

27

28 The structure of a cGAN is illustrated in Figure 1. The loss function for this cGAN is

29
$$\mathcal{L}_{cGAN}(G,D) = \mathbb{E}_{x,y}[log D(x,y)] + \mathbb{E}_{x}[log(1-D(x,G(x))]]$$

The loss function describes how G is minimized against a maximised D. Since both optimisation processes are dependent on each other, convergence is achieved by reaching a saddle point (simultaneously minimum / maximum for both networks' cost) rather than a

(1)

minimum. The loss also incorporates a L1 distance to reduce image blurring and ensure that the generated image from G(x) are not significantly different from the target image y (Isola et al., 2017; Regmi and Borji, 2018). This L1 loss is given by

$$\mathcal{L}_{L1}(G) = \mathbb{E}_{x,y}[||y - G(x)||_1]$$
(2)

5 The overall objective of the cGAN is to find the optimal solution to

4

$$G * = \arg\min_{G} \max_{D} \mathcal{L}_{cGAN}(G, D) + \lambda \mathcal{L}_{L1}(G)$$
(3)

7 with λ being a hyper-parameter used for balancing the two losses (Regmi and Borji, 2018).
8

9 The cGAN used in this work utilises the U-Net encoder-decoder architecture for the 10 generator, which is frequently used for image segmentation problems (Ronneberger et al., 11 2015). The generator was trained to generate images that are indistinguishable from a target 12 image (i.e., the segmented map). Spatial consistency of the data is not guaranteed with a U-13 Net segmented map, which can cause inaccurate boundaries (Ronneberger et al., 2015). 14 However, adversarial losses in the discriminator regulate and therefore increase the accuracy 15 to higher order shapes (Yang et al., 2017).

16

We modified the U-Net generator from the "pix2pix" network by increasing the input layer to
be able to train on 512 x 512 resolution images. For this an additional ConvolutionBatchNorm-leakyReLU layer was inserted in the encoding and a Convolution-BatchNormReLU layer in the decoding network part.

21

The discriminator is a patch-based fully convolutional neural network, PatchGAN (Li and Wand, 2016; Long et al., 2018), which models the image as a Markov random field. It performs a convolutional patch-wise (N x N) classification with all the outputs in the patch averaged and taken as the output of D. D is therefore less dependent on distant pixels/voxels beyond a "patch diameter" and is a form of neighbouring texture loss. The PatchGAN can be applied to arbitrarily large images, due to a fixed size of the patch.

28

To analyse the cGANs performance we compared it to the performance of a U-Net network, which is widely used for image segmentation processes. We used the cGAN generator network as the U-Net network to maintain an effective comparison. 1

The networks were implemented using PyTorch (Torch v1.0.1) and all training was performed on a Nvidia P6000 GPU card (3840 CUDA cores, 24 GB GDDR5X). The training phase of optimisation was performed as described by the "pix2pix" network, using stochastic gradient descent to minimise D(x,y) and stochastic gradient ascent to maximise D(x,G(x)). The Adam solver was used with a learning rate 0.0002 and momentum parameters, $\beta_1 = 0.5$, $\beta_2 = 0.999$. We introduced random noise (jitter) during training by resizing the input images to 542 x 542 using bi-cubic interpolation followed by random cropping back to 512 x 512.

- 10 A detailed description of the network architectures can be found in the Appendix.
- 11

12 **2.4 Segmentation Evaluation Metrics**

The Sørensen–Dice Similarity Coefficient (DSC) (Dice, 1945; Sørensen, 1948) was used to evaluate the overlap between the generated segmentation and the manual segmentation. The DSC ranges between 0 and 1, with 0 representing no overlap and 1 complete overlap between the two sets. DSC is defined as twice the size of the intersect divided by the sum of the sizes of two sample sets, given as

18
$$DSC = \frac{2|X \cap Y|}{|X| + |Y|}$$
 (4)

for Boolean metrics. For the experiments involving the SKI10 and OAI ZIB datasets, the volumetric overlap error (VOE) and the boundary distance-based metric average surfaces distance (ASD) were determined to assess segmentation accuracy and allow an appropriate comparison with previous studies using these datasets.. The VOE can be calculated as

$$VOE = 1 - \frac{|X \cap Y|}{|X \cup Y|}$$
(5)

25 with small values for VOE expressing greater accuracy.

26

27 The ASD is expressed in *mm* and is defined as

$$ASD = \frac{1}{N_X + N_Y} \left(\sum_{i=1}^{N_X} D_X(y) + \sum_{i=1}^{N_Y} D_Y(x) \right)$$
(6)

where D_X(y) = min_{x∈X} ||y - x||is the distance of a voxel y to a surface X and ||·|| denotes the
Euclidean norm.

4

5 **2.5 Evaluation of Network Characteristics**

6 This section aims at evaluating and adjusting specific network characteristics towards 7 improving overall network performance, for both cGAN and U-Net. All networks in this 8 section were trained for 100 epochs and all cGANs with a 70 x 70 PatchGAN discriminator 9 unless otherwise stated.

10

11 **2.5.1 Evaluation of Network Objective Function:**

We evaluated the cGANs performance with different objective functions by combining the cGAN loss with different pixel-wise error losses. In this work the cGAN is tasked to output a segmentation map of multiple tissues having different features and locations in the input MR image. We assessed the shortcomings and strengths of including the \mathcal{L}_{L1} , \mathcal{L}_{L2} and Smooth L1 (\mathcal{L}_{SmL1}) (Girshick, 2015) loss functions in the cGAN objective. The \mathcal{L}_{L2} loss and \mathcal{L}_{SmL1} loss are given by

$$\mathcal{L}_{L2}(G) = \mathbb{E}_{x,y}[||y - G(x)||_2^2]$$
(7)

19
$$\mathcal{L}_{SmL1}(G) = \begin{cases} 0.5 \cdot \mathbb{E}_{x,y}[||y - G(x)||_2^2], & \text{if } |y - G(x)| < 1\\ \mathbb{E}_{x,y}[||y - G(x)||_1] - 0.5, & \text{otherwise} \end{cases}$$
(8)

Furthermore, the weighting hyperparameter λ between the cGAN loss and pixel-wise error loss was changed to vary the balance between the two task losses. $\lambda = 0.01$, 1, 100 and 10000 were investigated. Network training with the cGAN loss alone ($\lambda = 0$) was additionally performed and evaluated.

24

We also trained the U-Net with the same three different pixel-wise error losses (\mathcal{L}_{L1} , \mathcal{L}_{L2} and \mathcal{L}_{SmL1}) as the cGAN to maintain an effective comparison.

1 2.5.2 Evaluation of Altering the Loss Objective during Training:

After obtaining initial results, we observed that the cGAN was unable to segment muscle tissues, independent of the objective function trained on. Therefore, we decided to explore the effect of varying the loss objective during training. For this, we trained a cGAN with $\mathcal{L}_{cGAN} + \lambda \mathcal{L}_{L2}$ loss and a U-Net with \mathcal{L}_{L2} loss for 50 epochs and then changed the loss functions for the ensuing 50 epochs to $\mathcal{L}_{cGAN} + \lambda \mathcal{L}_{L1}$ and \mathcal{L}_{L1} , respectively.

7

8 **2.5.3 Evaluation of the Generator Depth:**

9 We analysed the effect of changing the depth of the generator network on the cGANs and U-10 Nets quantitative performance. In addition to the generator down-sampling the input through 11 nine convolutional networks, we tested a generator consisting of seven and five convolutions during down-sampling. Furthermore, we assessed the quantitative performance of the 12 13 generator network with different numbers feature channels. We compared networks starting 14 with different minimum number of feature channels (16, 32, 64 and 128) and thus end at 15 different maximum numbers of feature channels (128, 256, 512 and 1024). All cGANs were trained with $\mathcal{L}_{cGAN} + \lambda \mathcal{L}_{L1}$ loss with $\lambda = 100$ and all U-Nets with the \mathcal{L}_{L1} loss. Detailed 16 descriptions of the generator network architectures can be found in the Appendix. 17

18

19 **2.5.4 Evaluation of the PatchGAN Receptive Field Size:**

We evaluated the effect of changing the PatchGAN receptive field size on the cGANs qualitative (artefact emergence) and quantitative (segmentation accuracy) performance. In addition to the 70 x 70 PatchGAN, we tested a 1 x 1 (PixelGAN), 34 x 34 and 286 x 286 PatchGAN. All cGANs were trained with $\mathcal{L}_{cGAN} + \lambda \mathcal{L}_{L1}$ loss with $\lambda = 100$. Detailed descriptions of the discriminator network architectures can be found in the Appendix.

25

26 **2.5.5 Evaluation of Transfer Learning:**

Since the AMROA dataset only comprises of a low number of subjects (N=8) for training, we assess the influence of transfer learning on network performance, by initially training both a cGAN ($\mathcal{L}_{cGAN} + \lambda \mathcal{L}_{L1}$) and a U-Net (\mathcal{L}_{L1}) for 20 epochs on the larger SKI10 and OAI ZIB training datasets separately followed by network fine-tuning for 80 epochs on the smaller AMROA training set. Additionally, a cGAN and a U-Net were trained for 20 epochs on the AMROA training dataset followed by network refinement training for 80 epochs on either the SKI10 or OAI ZIB training set to analyse the potential segmentation improvement of SKI10 and OAI ZIB. Network performance evaluations were performed using AMROA, SKI10 and OAI ZIB testing datasets. As determined from the previous sections, the cGAN trained with the $\mathcal{L}_{cGAN} + \lambda \mathcal{L}_{L1}$ loss objective (λ =100) and a 1 x 1 PixelGAN as well as the U-Net trained with the \mathcal{L}_{L1} loss objective achieved the highest segmentation accuracies for most knee joint tissues segmented in the AMROA dataset and were used in this section.

6 7

8 3. Results and Discussion

9 **3.1 Network Training and Testing**

Semi-manual segmentation of the AMROA images by the reader required \sim 30 minutes per subject-volume. cGAN training was performed in 80 seconds/epoch for the AMROA training dataset, and 390 seconds/epoch for the SKI10 dataset. U-Net training was performed in 45 seconds/epoch for the AMROA training dataset, and 185 seconds/epoch for the SKI10 dataset. Segmentation post-training on a single slice was processed in \approx 0.13s. The highlights of the upcoming sections are:

- 16 3.2 The U-Net trained with \mathcal{L}_{L1} loss objective outperformed the cGANs and the U-17 Nets trained with different loss objectives in the segmentation performance of 18 most knee joint tissues.
- 3.3 Altering the network objective function midway through cGAN and U-Net
 training lead to unanticipated but advantageous results. This variation resulted in
 improved segmentation performances of several tissues and the cGANs capability
 to segment muscle tissue, which previously had not been possible with non-altered
 objective function training.
- 3.4 The cGAN and U-Net trained with nine convolutions/transpose convolutions in
 the networks encoding/decoding parts and a minimum feature channel change of
 64 achieved the highest segmentation accuracies for most knee joint tissues
 annotated.
- 3.5 The greatest improvements in segmentation performance of the cGAN was
 achieved by reducing the receptive field size of the discriminator network. This
 resulted in segmentation accuracies equivalent to those of the U-Net.
- 3.6 Transfer learning not only increased segmentation accuracy of some tissues of the
 fine-tuned dataset, but also increased the network's capacity to maintain
 segmentation capabilities for the pretrained dataset.

- 1 3.7 Overall, the cGAN trained with the $\mathcal{L}_{cGAN} + \lambda \mathcal{L}_{L1}$ loss objective (λ =100) and a 1 x 2 1 PixelGAN as well as the U-Net trained with the \mathcal{L}_{L1} loss objective achieved 3 comparable and the highest segmentation accuracies for most knee joint tissues 4 segmented.
- 5

6 **3.2 Evaluation of Network Objective Function**

The quantitative results of assessing the impact of combining the cGAN objective with three 7 different pixel error losses with varying weightings λ on the cGANs segmentation 8 9 performance are in Table 2, with the qualitative results depicted in Figure 2B. The cGANs 10 trained with larger values for λ (λ =100 and 10000) achieved the highest segmentation 11 performance for all tissues and the produced segmentation maps were less affected by artefacts compared to the cGANs trained with $\lambda = 0.01$ and 1. For instance, the images from 12 the networks trained with $\mathcal{L}_{cGAN} + \lambda \mathcal{L}_{L1}$ ($\lambda = 0.01$), $\mathcal{L}_{cGAN} + \lambda \mathcal{L}_{L2}$ ($\lambda = 1$) and $\mathcal{L}_{cGAN} + \lambda \mathcal{L}_{L2}$ 13 $\lambda \mathcal{L}_{Sml,1}$ (1) had artefacts where the networks seem to detect bone or cartilage structures 14 where there were none in the original MR input image. By increasing the weighting 15 16 hyperparameter λ , more emphasis is put on the pixel error losses to guide the network to produce more accurate representations of the ground truth segmentation map and reduces 17 18 these artefacts. However, the influence of GAN loss diminishes with very large values for λ with the discriminator having minimal effect on generator training. 19

20

The qualitative results of training a U-Net with different pixel error losses are presented in 21 Figure 2C while the quantitative results are listed in Table 3. The U-Net trained with \mathcal{L}_{L1} loss 22 objective achieves the highest accuracy for all tissues compared to \mathcal{L}_{L2} and \mathcal{L}_{SmL1} loss except 23 for the muscle tissues. Muscle tissues appeared on the majority of 2D MR knee images seen 24 25 by the network during training, however we only segmented two selective medial muscles in the AMROA dataset due to time constraints. It is interesting to note that although the U-Net 26 trained with \mathcal{L}_{L1} was not able to capture the medial head of gastrocnemius and vastus 27 medialis muscles, the cGAN trained with the $\mathcal{L}_{cGAN} + \lambda \mathcal{L}_{L1}$ objective ($\lambda = 10000$) was. 28 Simple absolute difference (\mathcal{L}_{L1}) was not capable of differentiating lateral muscle textures 29 from medial. The U-Nets trained with \mathcal{L}_{L2} and \mathcal{L}_{SmL1} losses were capable of segmenting the 30 selective muscles with high accuracies as they are penalised more by the squaring term in 31 their loss objectives when the difference between ground truth and model predictions are 32 large. Interestingly, although the patella bone and cartilage only appear on very few slices in 33

1 a 3D dataset, and ACL and PCL on even fewer, the U-Net with \mathcal{L}_{L1} segmented these tissues better than the \mathcal{L}_{L2} and \mathcal{L}_{SmL1} (\mathcal{L}_{L2} : DSC_{P Bone} < 0.2%, DSC_{P Cartilage} < 5.3%, DSC_{ACL} < 2 15.2%, $DSC_{PCL} < 21.3\%$; \mathcal{L}_{SmL1} : $DSC_{P Bone} < 0.4\%$, $DSC_{P Cartilage} < 6.0\%$, $DSC_{ACL} < 6.9\%$, 3 4 $DSC_{PCL} < 17.8\%$). This could be explained by the cruciate ligament and patellar tissues either 5 being present or not on a 2D training image and the network is not being constrained to only 6 segment medial tissues. Overall, the U-Net with \mathcal{L}_{L1} produced sharper boundaries, especially 7 for the smaller ligament structures, as compared to the segmentation maps produced by U-8 Nets trained with \mathcal{L}_{L2} and \mathcal{L}_{SmL1} , in which the boundaries are more diffused.

9

10 We decided to assess the model's performance when including noise-only images in the 11 testing dataset as we excluded them during model training, and this might limit the models' 12 use in a clinical setting. This effect was only evaluated for a the cGAN trained with the $\mathcal{L}_{cGAN} + \lambda \mathcal{L}_{L1}$ ($\lambda = 100$) objective function and the U-Net trained with the \mathcal{L}_{L1} loss objective. 13 14 The quantitative results are listed in Table 4 with qualitative results displayed in Figure 3. 15 Both networks showed comparable segmentation performances after testing with noise-only images with percentage differences (%-Diff) of the DSC for all segmented tissues $\leq 2.3\%$. 16 17 Including noise-only images into the testing set had greater effects on the cGAN DSC of the medial vastus muscle (VM muscle) (%-Diff = 1.5%), the ACL (%-Diff = 1.6%) and the PCL 18 19 (%-Diff = 1.9%) as well as on the U-Net DSC of the ACL (%-Diff = 2.3%). These higher 20 differences could be explained by the lower segmentation capability of these structures by the cGAN and U-Net models to begin with (cGAN: DSC_{VM muscle}: 0.113 vs 0.098, DSC_{ACL}: 0.577 21 vs 0.593; DSC_{PCL}: 0.073 vs 0.092; U-Net: DSC_{ACL}: 0.643 vs 0.620). Furthermore, the larger 22 %-Diff in the DSC of the VM muscle is caused by the cGAN model irregularly segmenting 23 24 VM muscle tissues on noise only images (Figure 3B).

25

3.3 Evaluation of Altering Loss Objective during Training

Figure 4 compares the qualitative results and Table 5 compares the DSCs obtained from a cGAN and a U-Net, in which the objective functions were changed midway through training to the cGANs and U-Nets trained with non-altered objective functions. Training a cGAN with varied loss objective ($\mathcal{L}_{cGAN} + \lambda \mathcal{L}_{L2} \rightarrow \mathcal{L}_{cGAN} + \lambda \mathcal{L}_{L1}$) notably reduced its ability to segment the ACL, however considerably improved its segmentation performance on the medial vastus and gastrocnemius muscles, as well as PCL, compared to the other cGANs ($\mathcal{L}_{cGAN} + \lambda \mathcal{L}_{L1}$) and $\mathcal{L}_{cGAN} + \lambda \mathcal{L}_{L2}$). The images in Figure 4B show the improvements in muscle 1 segmentation with the cGAN trained with varied loss objective. This was a surprising result as neither the cGAN trained with $\mathcal{L}_{cGAN} + \lambda \mathcal{L}_{L1}$ nor with $\mathcal{L}_{cGAN} + \lambda \mathcal{L}_{L2}$ alone were able to 2 segment muscle. Looking at the different training epochs of the cGAN trained with varied 3 4 loss, during $\mathcal{L}_{cGAN} + \lambda \mathcal{L}_{L2}$ no muscle tissue was being semantically segmented. However, 5 when changing to $\mathcal{L}_{cGAN} + \lambda \mathcal{L}_{L1}$ and between training epochs 50 and 60, the network started 6 segmenting muscle tissue (Figure 5). After the initial 50 epochs of $\mathcal{L}_{cGAN} + \lambda \mathcal{L}_{L2}$ training, the 7 cGANs weights must have been favourable for continuing training with $\mathcal{L}_{cGAN} + \lambda \mathcal{L}_{L1}$ to 8 additionally semantically segment muscle tissue.

9

The U-Net trained with altered objective function $(\mathcal{L}_{L2} \rightarrow \mathcal{L}_{L1})$ also showed notable 10 11 improvements in the segmentation performance of the medial vastus and gastrocnemius muscles while the segmentation scores of the other knee tissues remained comparable with 12 those of the other U-Nets (\mathcal{L}_{L1} and \mathcal{L}_{L2}). Figure 4C qualitatively compares the results of a U-13 14 Net trained with altered loss objective to those of the U-Nets trained with a single, non-15 altered loss objective. As mentioned in the corresponding method section, this idea came after reviewing a few initial training results. While the U-Net trained with the \mathcal{L}_{L1} objective 16 17 was not able to segment the medial vastus and gastrocnemius muscles after training, the U-Net with the \mathcal{L}_{L2} loss objective was. However, these images were slightly blurrier, and the 18 segmentation accuracy of the remaining tissues was poorer than compared to \mathcal{L}_{L1} . By varying 19 the loss objective during training, the strengths of \mathcal{L}_{L2} and \mathcal{L}_{L1} were combined. We decided to 20 21 first train the network with \mathcal{L}_{L2} loss to capture all tissues and then to change to \mathcal{L}_{L1} halfway 22 through training to make the images sharper and increase segmentation accuracy. This 23 method created a more proficient network capable of segmenting all tissues with higher or 24 comparable accuracies to the networks trained with non-altered loss objectives.

25

3.4 Evaluation of the Generator Depth

The quantitative results of assessing the impact of generator network depth on the cGANs and U-Nets segmentation performances are in Tables 6 and 7.

29

The cGAN with a generator down-sampling the input through nine convolutional networks achieved the highest DSC scores for tibial and patellar bone, as well as for femoral and patellar cartilage. Femoral bone and tibial cartilage were best segmented by the cGAN with five convolutions/transpose convolutions in the generator encoding/decoding parts. The medial vastus and gastrocnemius muscles, as well as ACL and PCL were best segmented by the cGAN with seven convolutions. Training the cGAN with a minimum feature channel change of 64 resulted in the highest segmentation scores for most tissues except for femoral bone, tibial cartilage and the medial vastus muscle.

5

6 The U-Net trained with nine convolutions/transpose convolutions in the networks 7 encoding/decoding parts achieved the highest segmentation accuracies for all but one tissue 8 (femoral cartilage), which was slightly better segmented by the U-Net with five 9 convolutions/transpose convolutions. Training the U-Net with a minimum feature channel 10 change of 64 resulted in the highest DSC scores for most tissues apart from patella cartilage 11 and ACL which were segmented best by the U-Net trained with a minimum feature channel 12 change of 128.

13

14 It is important to note for this section that increasing the number of convolutions and feature 15 channels in the generator network substantially increases the overall number of parameters in 16 the network and the time per epoch required to train the network (see network architectures in 17 the Appendix for details). A considered decision between increase in learning time and 18 significant improvement in segmentation accuracy has to be made.

19

20 **3.5 Evaluation of PatchGAN Receptive Field Size**

21 Figure 6 shows the qualitative comparison of the effect of using different patch sizes in the 22 discriminator network, while the corresponding DSCs are listed in Table 6. The cGAN 23 trained with the 1 x 1 PatchGAN (PixelGAN) achieved the highest segmentation accuracy for 24 most tissues except for femoral and tibial cartilage and both muscle tissues, which were best 25 segmented by the 34 x 34 PatchGAN. Increasing the receptive field size increases the number 26 of parameters in the discriminator network and therefore may be more difficult to train. 27 Additionally, as in the 'pix2pix' paper (Isola et al., 2017), we also noticed the repetitive tiling / checkerboard artefact (Figure 7). However, in our instance, the artefacts become more 28 29 pronounced with every increase in patch size instead of the inverse tendency as seen by (Isola 30 et al., 2017). This could be a result of us assigning the cGANs with the reverse task (image to 31 label) compared to the one performed by (Isola et al., 2017) (label to image).

Figure 8 depicts the loss evolution during network training of the cGAN trained with the 1 x PatchGAN discriminator. The loss evolutions of the cGAN generator (\mathcal{L}_{cGAN} and \mathcal{L}_{L1}) and discriminator (\mathcal{L}_{real} and \mathcal{L}_{fake}) are shown in Figure 8A and 8B, respectively. Figure 8B highlights how the Nash equilibrium was reached for the discriminator network during cGAN training.

6

7 **3.6 Evaluation of Transfer Learning**

8 The quantitative results of this section are presented in Tables 9 and 10 with qualitative 9 comparisons between single step (one dataset) and two step training (transfer learning) 10 displayed in Figures 9 and 10.

11

When comparing the segmentation performances of the proposed cGAN and U-Net without 12 13 and with transfer learning and testing on the SKI10 testing dataset (Table 9, Figures 9A-C), the AMROA-pretrained / SKI10-retrained (AMROA \rightarrow SKI10) U-Net showed the highest 14 15 DSC scores for femoral and tibial bone and the highest boundary accuracy (i.e. smallest 16 ASDs) for femoral bone, while the SKI10-only trained U-Net segmented the tibial bone with the highest boundary accuracy. Femoral cartilage was best segmented by the AMROA-17 18 pretrained / SKI10-retrained (AMROA \rightarrow SKI10) cGAN and tibial cartilage by the SKI10only trained cGAN. 19

20

Testing the OAI ZIB testing dataset on the proposed cGAN and U-Net without and with transfer learning (Table 9, Figures 9D-F), the AMROA-pretrained / OAI ZIB-retrained (AMROA \rightarrow OAI ZIB) cGAN showed the highest accuracies for tibial bone and femoral cartilage, while the OAI ZIB-only trained cGAN segmented the femoral bone and tibial cartilage with the highest accuracies.

26

When testing the cGANs and U-Nets on the AMROA testing dataset (Table 10, Figure 10), the SKI10-pretrained / AMROA-retrained (SKI10 \rightarrow AMROA) U-Net had the highest DSCs for femoral and tibial bone as well as the ACL. Femoral cartilage as well as patellar bone and cartilage was segmented most accurately by the OAI ZIB-pretrained / AMROA-retrained (OAI ZIB \rightarrow AMROA) U-Net. The AMROA only trained U-Net showed the best segmentation accuracy for tibial cartilages. The SKI10-pretrained / AMROA-retrained (SKI10 \rightarrow AMROA) cGAN provided the highest segmentation score for the vastus medialis 1 muscle while the medial head of gastrocnemius muscle and the PCL was best segmented by the OAI ZIB-pretrained / AMROA-retrained (OAI ZIB \rightarrow AMROA) cGAN. Compared to 2 3 the U-Net, the cGAN could successfully segment both medial muscles which could promote 4 a strength of the cGAN. A further note is that, although the SKI10 and OAI ZIB datasets only 5 comprised of segmentations of femoral and tibial bone and cartilage, the cGANs and U-Nets 6 initialised with the respective SKI10- and OAI ZIB-pretrained network weights and retrained 7 on the AMROA dataset were able to recuperate and capture patellar, ligament and muscle 8 tissues.

9

A challenge of any machine learning technique is obtaining a training set that optimises the 10 amount of variation from the rare morphology of pathological conditions or image artefacts. 11 12 The AMROA dataset was highly controlled, with the patients and imaging occurring with a single imaging protocol on a single MRI system. The images showed a clear bone-cartilage 13 separation and enabled better cartilage segmentation scores after training than the SKI10 14 15 dataset. The OAI ZIB dataset highlights the benefits of training on a very large number of 16 images with the cGAN and U-Net (OAI ZIB-only trained) achieving $DSC \ge 0.984$ for bone 17 and DSC ≥ 0.837 for cartilage segmentations.

18

The ability for the network to be used under variable conditions was simulated by using three 19 20 knee datasets (AMROA, SKI10 and OAI ZIB). Even without transfer learning, the AMROA 21 training enabled SKI10 and OAI ZIB segmentation and vice versa, albeit not with high 22 accuracy, but nonetheless indicating the robustness of deep learning methods. Transfer 23 learning not only improved the segmentation accuracy for some tissues of the local dataset 24 but also enhanced the networks ability to segment the SKI10 / OIA ZIB test dataset by 25 introducing more heterogeneity into the model. Even though the SKI10- and OAI ZIB-26 pretrained networks were then fine-tuned to segment the local AMROA dataset, it could 27 segment the SKI10 and OAI ZIB testing dataset with an improved performance compared to the AMROA-only trained network without pretraining. This effect was seen for both cGANs 28 29 and U-Nets.

30

31 **3.7 AMROA: Comparison to Previous Studies**

In this subsection, the results obtained for the different tissues semantically segmented in this study are compared to those of previous studies. The cGAN and U-Net achieving the highest segmentation accuracy on the AMROA dataset for each respective tissue is chosen for this
 purpose.

3

4 **Bone:**

5 While cartilage has been traditionally studied for OA, bone shape has been under increasing 6 investigations (Ambellan et al., 2019; Felson and Neogi, 2004). Bone shape has been linked 7 to radiographic OA (Hunter et al., 2015; Neogi et al., 2013; Wise et al., 2018) and associated 8 with longitudinal pain progression (Hunter et al., 2015). Segmented bone can be used to 9 separate out bone-specific diseases, such as osteochondral defects.

10

11 The OAI ZIB-pretrained / AMROA-retrained cGAN trained with the $\mathcal{L}_{cGAN} + \lambda \mathcal{L}_{L1}$ loss 12 objective (λ =100) and a 1 x 1 PixelGAN generated segmentations of femoral (DSC = 0.972), tibial (DSC = 0.962) and patellar (DSC = 0.947) bone with the highest accuracy. The SKI10-13 pretrained / AMROA-retrained U-Net (\mathcal{L}_{L1} loss objective) achieved slightly higher 14 15 segmentation accuracies for femoral and tibial bone tissues (femoral: DSC = 0.974; tibial: 16 DSC = 0.965) and the OAI ZIB-pretrained / AMROA-retrained U-Net for patellar bone (DSC = 0.948), compared to the cGANs. The boundaries of the images, near the top and bottom of 17 18 any 2D slice, did not always segment all bone, which is where the MRI radiofrequency (RF) transmit and receive uniformity was poor due to characteristics of the MRI coil. Traditional 19 20 semi-automatic approaches involving signal threshold, region-based or clustering segmentation can be similarly sensitive to image non-uniformities (Swanson et al., 2010). 21 22 These non-uniformities are shown as a change in signal-to-noise or darkening of the 23 surrounding muscle tissues (see lower regions of Figure 2). These effects from RF transmit or 24 receive non-uniformity could be mitigated with a larger training population, as more complex 25 modelling of data is possible. Nevertheless, segmentation of the patella achieved the lowest 26 accuracy. The patella has the widest range of inter-subject variability when compared to the 27 larger tibial and femoral bones. The patella bone can vary in both shape and position, shifting 28 due to the orientation and bend of the knee. Additionally, due to its smaller volume, fewer 29 training images are used for the patella segmentation.

30

The cGAN and U-Net bone segmentation scores achieved in this study are similar to those achieved by a CycleGAN method using unannotated knee MR images for femoral (DSC = 0.95 - 0.97) and tibial (DSC = 0.93 - 0.95) bone segmentation (Liu, 2019), and a convolutional encoder-decoder network combined with a 3D fully connected conditional
random field and simplex deformable modelling for femoral (DSC = 0.970), tibial (DSC =
0.962) and patellar (DSC = 0.898) bone segmentation (Zhou et al., 2018).

4

5 **Cartilage**

For a long time, OA was considered a disease primarily involving variations in articular cartilage composition and morphology. Therefore, the attention was predominantly placed on the extraction of OA biomarkers from quantitative MR imaging techniques using manual or semi-manual segmentation techniques that suffer from intra- and inter-observer variability (Pedoia et al., 2016). Deep learning methods can provide a fast and repeatable alternative to overcome these time-consuming and operator-dependent procedures.

12

The OAI ZIB-pretrained / AMROA-retrained cGAN trained with the $\mathcal{L}_{cGAN} + \lambda \mathcal{L}_{L1}$ loss 13 14 objective (λ =100) and a 1 x 1 PixelGAN generated segmentations of femoral (DSC = 0.875), tibial (DSC = 0.811) and patellar (DSC = 0.879) cartilage with the highest accuracy from all 15 16 cGAN trainings. The OAI ZIB-pretrained / AMROA-retrained U-Net (\mathcal{L}_{L1} loss objective) achieved marginally higher accuracies for femoral (DSC = 0.893) and patellar (DSC = 0.898) 17 cartilage segmentations and the AMROA-only trained U-Net (\mathcal{L}_{L1} loss objective) achieved a 18 slightly higher segmentation accuracy for tibial cartilage (DSC = 0.834) compared to the 19 20 cGAN results.

21

The cartilage segmentation performances of both cGAN and U-Net are comparable to those attained by a 2D U-Net for femoral, tibial and patellar cartilage segmentations on T1pweighted (DSC=0.632 - 0.702) and DESS MR images (DSC=0.767 - 0.878) (Norman et al., 2018), a CycleGAN method for femoral and tibial cartilage segmentation on PD-weighted (DSC=0.65 - 0.66) and T2-weighted FSE images (DSC=0.81 - 0.75) (Liu, 2019), as well as the recently investigated cGAN for femoral, tibial and patellar segmentation on DESS MR images (DSC=0.843 - 0.918) (Gaj et al., 2019).

29

30 Muscle

31 As muscle weakness and atrophy can be regarded as preceding risk factors and resulting pain-

32 related consequences for the development and progression of OA, studying morphological

changes in knee joint muscles has become increasingly important (Fink et al., 2007;
 Slemenda et al., 1997).

3

The SKI10-pretrained / AMROA-retrained cGAN and the OAI ZIB-pretrained / AMROAretrained cGAN trained with the $\mathcal{L}_{cGAN} + \lambda \mathcal{L}_{L1}$ loss objective (λ =100) and a 1 x 1 PixelGAN segmented the medial gastrocnemius muscle (DSC = 0.909) and medial vastus muscle (DSC = 0.922) with the highest accuracies, respectively. The U-Net trained with altered loss objective ($\mathcal{L}_{L2} \rightarrow \mathcal{L}_{L1}$) achieved the highest segmentation accuracies for both the medial gastrocnemius (DSC = 0.933) and vastus (DSC = 0.914) muscles.

10

11 Our results are comparatively lower compared to those of a semi-automatic single-atlas (DSC 12 = 0.95 - 0.96) and fully-automatic multi-atlas (DSC = 0.91 - 0.94) based approach for medial vastus segmentation (Le Troter et al., 2016), and a 2D U-Net for quadriceps (DSC = 0.98) 13 14 segmentation (Kemnitz et al., 2019). A crucial difference between these studies and ours is 15 the plane in which segmentation was performed. While muscles are typically segmented on 16 axial images as this provides a more straightforward task with clearer separation between different muscles, our multi-class tissue segmentation approach was performed on sagittal 17 images. Segmenting different muscles in the sagittal plane is a demanding task, especially in 18 19 areas of the calf muscles where the two-headed gastrocnemius muscle overlaps (medial and 20 lateral) while also overlaying the soleus muscle.

21

22 Cruciate Ligament

23 There has been a growing interest in investigating and understanding the mechanism 24 responsible for the post-traumatic development of OA following injury to the cruciate 25 ligaments, especially the ACL (Chaudhari et al., 2008; Messer et al., 2019; Monu et al., 26 2017). Although ACL reconstruction and rehabilitation can help restore patients to normal life and previous activities, it cannot prevent the long-term risk of developing OA (Paschos, 27 2017). Accurate and repeatable segmentations of the cruciate ligaments are required when 28 aiming at evaluating longitudinal changes in the cruciate ligaments following reconstructive 29 30 surgery.

31

In our study, the OAI ZIB-pretrained / AMROA-retrained cGAN trained with the 1 x 1 PixelGAN and $\mathcal{L}_{cGAN} + \lambda \mathcal{L}_{L1}$ loss objective (λ =100) achieved the highest accuracy for ACL 1 (DSC = 0.664) and PCL segmentation (DSC = 0.652). The SKI10-pretrained / AMROA-2 retrained U-Net (\mathcal{L}_{L1} loss objective) achieved a similar accuracy for ACL segmentation 3 (DSC = 0.665) and the AMROA-only trained U-Net (\mathcal{L}_{L1} loss objective) achieved a 4 marginally lower accuracy for PCL segmentation (DSC = 0.641), compared to the best 5 performing cGANs.

6

7 (Lee et al., 2013) proposed a graph cut method for automatic ACL segmentation and attained a DSC score of 0.672, while (Paproki et al., 2016) used a patch-based method for PCL 8 9 segmentation to achieve a DSC score of 0.744. Using a 3D convolutional neural network 10 (CNN), (Mallya et al., 2019) achieved DSC scores of 0.40 and 0.61 for ACL and PCL segmentations, respectively. When combining their 3D CNN with a deformable atlas-based 11 12 segmentation method, their ACL (DSC = 0.84) and PCL (0.85) segmentation accuracies increased substantially. In general, 3D networks could provide higher segmentation 13 14 accuracies especially for fine structures such as the cruciate ligaments that only appear on a few 2D slices in a 3D dataset. However, 2D segmentation techniques are useful for broader 15 16 applicability, as 2D imaging is often faster and currently still more clinically employed than 3D imaging. 17

18

19 The lower similarity scores achieved in our study compared to the other studies could arise 20 from the use of 3D-FS SPGR images as source images during training as these are non-21 optimal for the segmentation of the cruciate ligaments due to their less than ideal soft tissue 22 separation with surrounding structures and fluid. Fat-saturated proton-density-weighted fast 23 spin echo or T2-weighted fast spin echo images are more suitable for segmentation purposes 24 as shown by (Mallya et al., 2019) and (Paproki et al., 2016), respectively. These sequences 25 are clinically used for cruciate ligament assessment due to their dark appearance and clear 26 separation from fluid and other surrounding tissues.

27

3.8 SKI10 and OAI ZIB: Comparison to Previous Studies

In this subsection, the segmentation results of the SKI10 and OAI ZIB datasets in this study are compared to those of previous studies. The cGAN and U-Net achieving the highest segmentation accuracy on these datasets is chosen for this purpose.

1 **SKI10**

2 The AMROA-pretrained / SKI10-retrained U-Net (\mathcal{L}_{L1} loss objective) achieved a comparable

ASD score for femoral bone (ASD = 0.44 mm) and an improved ASD score for tibial bone (ASD = 0.26 mm) to those reported by (Liu et al., 2017) and (Ambellan et al., 2019). However, the segmentation accuracies for femoral (VOE \geq 42.2%) and tibial (VOE \geq 47.6%) cartilage achieved by our models were substantially lower.

7

8 OAI ZIB

The OAI ZIB-only trained cGAN trained with the $\mathcal{L}_{cGAN} + \lambda \mathcal{L}_{L1}$ loss objective (λ =100) and a 9 10 1 x 1 PixelGAN generated segmentations of femoral bone (DSC = 0.985) and tibial cartilage (DSC = 0.839) with the highest accuracy. AMROA-pretrained / OAI ZIB-retrained cGAN 11 trained with the 1 x 1 PixelGAN and $\mathcal{L}_{cGAN} + \lambda \mathcal{L}_{L1}$ loss objective (λ =100) achieved the 12 13 highest accuracy for tibial bone (DSC = 0.985) and femoral cartilage (DSC = 0.897) 14 segmentation. The ASD of both the femoral (ASD = 0.33 mm) and tibial (ASD = 0.29 mm) bones were smaller than image resolution of the OAI DESS images $(0.36 \times 0.36 \times 0.7 \text{ mm}^3)$. 15 Although we achieve similar DSC scores for all tissues on the OAI ZIB dataset compared to 16 those presented in (Ambellan et al., 2019), our ASD scores were larger. The pixel-wise error 17 losses (\mathcal{L}_{L1} . \mathcal{L}_{L2} and \mathcal{L}_{SmL1}) used to train the networks in our work were chosen to maintain 18 an effective comparison between the cGAN and the U-Net. However, training our models 19 20 with loss functions more traditionally used for segmentation purposes such as multi-class 21 Dice similarity or cross entropy might lead to more comparable results for boundary-22 distance-based metrics.

23

24 **3.9 Limitations**

25 The network performances are depended on the accuracy of the ground truth segmentations. 26 Inaccuracies or errors in the segmentation maps could result in a less accurate network, 27 especially when trained on a low number of image volumes, as done in this study. Additionally, training a network on a low number of high-quality images restricts the 28 networks applicability to only highly controlled studies with homogeneous data. Therefore, 29 30 the networks trained in this study might be limited in their application in clinical settings where high image quality is not always achievable due to patient conditions and operator 31 32 variabilities.

1 Network training on 2D MR image slices is considerably less computationally demanding 2 than on 3D volumes. For the purposes of this study such as investigating the effects of 3 training with different loss objectives and cGAN discriminator networks, it was sufficient to 4 train on 2D images. Nevertheless, the segmentation of small knee joint structures, such as the 5 cruciate ligaments, could benefit from 3D networks that should add spatial continuity along 6 the slice dimension.

7

8 Furthermore, the segmentation results presented in this study are from standalone networks 9 without further processing within a pipeline. Therefore, the obtained results, especially for 10 cartilage segmentation, are not comparable to those from current state-of-the-art pipeline 11 methods such as described by (Liu et al., 2017) and (Ambellan et al., 2019) that initially 12 perform automated segmentation using a CNN followed by further refinement using 13 deformable or statistical shape models, respectively.

14

Lastly, additional investigations into varying the network architectures and optimisation strategies are warranted, with ever more loss functions as well as layer combination and optimisation strategies continuously being developed.

18

19 **4. Conclusion**

20 This work demonstrated the usage of a cGAN, using a U-Net generator with a PatchGAN 21 discriminator, for the purpose of automatically segmenting multiple knee joint tissues on MR 22 images. While DSC > 0.9 were achieved for all segmented bone structures and DSC > 0.7523 for cartilage and muscle tissues, DSC of only ≈ 0.64 were achieved for cruciate ligament segmentations. Nevertheless, this segmentation performance was attained despite the low 24 25 number of subjects (N=8) for training on the local dataset. Although the U-Net outperformed 26 the cGAN in most knee joint tissue segmentations, this study provides an optimal platform 27 for future technical developments for utilising cGANs for segmentation tasks. By enabling 28 automated and simultaneous segmentation of multiple tissues we hope to increase the 29 accuracy and time efficiency for evaluating joint health in osteoarthritis.

- 30
- 31
- 32
- 33

1 Appendix

2 Network Description

Generator: The encoding part of the generator network consists of the repeated application 3 of nine 4x4 convolutions with stride 2, down-sampling the input by a factor of 2 at each 4 5 layer. Each convolution is followed by a batch normalisation layer (except the first layer) and 6 a leaky rectified linear unit (leaky ReLU) with slope 0.2. During the first encoding step the 7 number of feature channels is changed from 3 to 64. At the subsequent three encoding steps, the number of feature channels is doubled (64 - 512), while the following five are kept at 8 512. In the ensuing decoding part, the input is repeatedly up-sampled by a factor of 2 by nine 9 4x4 transpose convolutional layers with stride 2 and additional skip connections 10 (concatenations) between each layer i and 9-i, changing the number of feature channels at 11 12 each step. The first four decoder convolutions are followed by batch normalisation, dropout 13 (50%) and a ReLU. The next four decoder convolutions are followed by batch normalisation 14 and a ReLU without dropout. After the final layer a convolution followed by a Tanh 15 activation layer is applied to generate the segmentation map.

16 Total number of parameters: 66.999 M

17	Training time (s/epoch):	AMROA:	135 (cGAN with 1x1 PixelGAN)
18			130 (cGAN with 70x70 PatchGAN)
19			100 (U-Net)
20		SKI10:	380 (cGAN with 1x1 PixelGAN)
21			210 (U-Net)
22		OAI ZIB:	2710 (cGAN with 1x1 PixelGAN)
23			1530 (U-Net)

Generator with five convolutions in encoder/decoder: In this generator network, the encoding part consists of the repeated application of five 4x4 convolutions with stride 2, down-sampling the input by a factor of 2 at each layer. In the ensuing decoding part, the input is repeatedly up-sampled by a factor of 2 by five 4x4 transpose convolutional layers with stride 2 and additional skip connections between each layer i and 5-i.

29 Total number of parameters: 16.659 M

30	Training time (s/epoch):	AMROA:	110 (cGAN with 70x70 PatchGAN)
31			90 (U-Net)

1	Generator with seven conv	olutions in en	coder/decoder: The encoding part consists of the						
2	repeated application of seven 4x4 convolutions with stride 2, down-sampling the input by a								
3	factor of 2 at each layer. In the subsequent decoding part, the input is repeatedly up-sampled								
4	by a factor of 2 by seven 4x4	4 transpose cor	volutional layers with stride 2 and additional skip						
5	connections between each layer i and 7-i.								
6	Total number of parameters:	41.829 M							
7	Training time (s/epoch):	AMROA:	120 (cGAN with 70x70 PatchGAN)						
8			100 (U-Net)						
9	Generator with 16 as mini	mum number	of feature channels: In this network, the number						
10	of feature channels is chan	nged from 3 to	o 16 during the first encoding step. During the						
11	following three encoding ste	eps, the number	r of feature channels is doubled $(16 - 128)$, while						
12	the subsequent five are kept	at 128.							
13	Total number of parameters:	4.191 M							
14	Training time (s/epoch):	AMROA:	105 (cGAN with 70x70 PatchGAN)						
15			70 (U-Net)						
16	Generator with 32 as min	nimum numb	er of feature channels: The number of feature						
17	channels is changed from 3	3 to 32 during	g the first encoding step. In the following three						
18	encoding steps, the number	of feature char	nnels is doubled $(32 - 256)$, while the subsequent						
19	five are kept at 256.								
20	Total number of parameters:	16.755 M							
21	Training time (s/epoch):	AMROA:	100 (cGAN with 70x70 PatchGAN)						
22			75 (U-Net)						
23	Generator with 128 as min	nimum numbe	er of feature channels: In the first encoding step						
24	the number of feature channel	nels is changed	from 3 to 128. In the following three encoding						
25	steps, the number of feature	channels is do	oubled (128-1024), while the subsequent five are						
26	kept at 1024.								
27	Total number of parameters:	267.953 M							
28	Training time (s/epoch):	AMROA:	245 (cGAN with 70x70 PatchGAN)						
29			220 (U-Net)						
30									

Discriminator:

70 x 70 PatchGAN: The discriminator network repeatedly down-samples the input by
applying three 4x4 convolutions with stride 2 followed by two 4x4 convolutions with stride
1. Each convolution during down-sampling is followed by a batch normalisation layer
(except the first and last layer) and a leaky ReLU (slope 0.2) (except for the last layer). The
number of feature channels are doubled (64 – 512) during the first four convolutional steps.
The final convolutional layer is proceeded by a Sigmoid activation layer.

7 Total number of parameters: 2.769 M

8 1 x 1 PatchGAN (PixelGAN): This PixelGAN discriminator network applies three 1 x 1 9 convolutions with stride 1, where the first convolution is followed by a leaky ReLU (slope 10 0.2), the second convolution by a batch normalisation layer and a leaky ReLU (slope 0.2) 11 and the final convolution by a Sigmoid activation function. The number of feature channels 12 are doubled (64 – 128) during the first two convolutions.

13 Total number of parameters: 0.009 M

34 x 34 PatchGAN: This network repetitively down-samples the input by using two 4x4 convolutions with stride 2 followed by two 4x4 convolutions with stride 1. Each convolution is followed by a batch normalisation layer (except the first and last layer) and a leaky ReLU (slope 0.2) (except for the last layer). The number of feature channels are doubled (64 – 256) during the first three convolutional steps. The final layer is ensued by a Sigmoid activation layer.

20 Total number of parameters: 0.666 M

286 x 286 PatchGAN: This discriminator network consists of eight convolutional layers with
4x4 spatial filters. The first 6 convolutions have stride 2 while the last two have stride 1. Each
convolutional layer is followed by a batch normalisation layer (except the first and last layer)
and a leaky ReLU (slope 0.2) (except for the last layer). The number of feature channels are
doubled (64 – 512) during the first four convolutions and kept at 512 for the ensuing layers.
A Sigmoid activation layer succeeds the final convolution.
Total number of parameters: 11.159 M

28

20

29

30

1 Acknowledgements

2 This work acknowledges funding support by the European Union's Horizon 2020 research 3 and innovation programme under grant agreement No 761214, Addenbrooke's Charitable 4 Trust (ACT) and the NIHR comprehensive Biomedical Research Centre (BRC) award to 5 Cambridge University Hospitals (CUH) NHS Foundation Trust. FJG receives funding from the National Institute for Health Research Senior Investigator award. DAK, JWM, and JDD 6 7 acknowledge funding support from GlaxoSmithKline for their studentships and fellowships, 8 respectively. We would also like to acknowledge the support of Robert L Janiczek 9 (Experimental Medicine Imaging, GlaxoSmithKline, Philadelphia, United States) and 10 Alexandra R Roberts (Corelab Manager and Imaging Director, Antaros Medical, Uppsala, Sweden). The views expressed are those of the authors and not necessarily those of the 11 12 funding bodies.

- 13
- 14

15 **References**

Ambellan, F., Tack, A., Ehlke, M., Zachow, S., 2019. Automated segmentation of knee bone
and cartilage combining statistical shape knowledge and convolutional neural networks:
Data from the Osteoarthritis Initiative. Med. Image Anal. 52, 109–118.
https://doi.org/10.1016/j.media.2018.11.009

- Benhamou, C.L., Poupon, S., Lespessailles, E., Loiseau, S., Jennane, R., Siroux, V., Ohley,
 W., Pothuaud, L., 2001. Fractal analysis of radiographic trabecular bone texture and
 bone mineral density: Two complementary parameters related to osteoporotic fractures.
 J. Bone Miner. Res. 16, 697–704. https://doi.org/10.1359/jbmr.2001.16.4.697
- Bindernagel, M., Kainmueller, D., Seim, H., Lamecker, H., Zachow, S., Hege, H.C., 2011.
 An articulated statistical shape model of the human knee. Inform. aktuell 59–63.
 https://doi.org/10.1007/978-3-642-19335-4 14
- Blumenkrantz, G., Majumdar, S., 2016. Quantitative magnetic resonance imaging of
 articular. Eur. Cells Mater. 13, 76–86. https://doi.org/10.22203/ecm.v013a08
- Chaudhari, A.M.W., Briant, P.L., Bevill, S.L., Koo, S., Andriacchi, T.P., 2008. Knee
 kinematics, cartilage morphology, and osteoarthritis after ACL injury. Med. Sci. Sports
 Exerc. 40, 215–222. https://doi.org/10.1249/mss.0b013e31815cbb0e
- 32 Chen, C., Dou, Q., Chen, H., Heng, P.A., 2018. Semantic-aware generative adversarial nets
- 33 for unsupervised domain adaptation in chest X-ray segmentation. Lect. Notes Comput.

- Sci. (including Subser. Lect. Notes Artif. Intell. Lect. Notes Bioinformatics) 11046
 LNCS, 143–151. https://doi.org/10.1007/978-3-030-00919-9_17
- Deniz, C.M., Xiang, S., Hallyburton, R.S., Welbeck, A., Babb, J.S., Honig, S., Cho, K.,
 Chang, G., 2018. Segmentation of the Proximal Femur from MR Images using Deep
 Convolutional Neural Networks. Sci. Rep. 8, 1–14. https://doi.org/10.1038/s41598-01834817-6
- Dice, L.R., 1945. Measures of the Amount of Ecologic Association Between Species.
 Ecology 26, 297–302. https://doi.org/10.2307/1932409
- Dou, Q., Ouyang, C., Chen, C., Chen, H., Heng, P.A., 2018. Unsupervised cross-modality
 domain adaptation of convnets for biomedical image segmentations with adversarial
 loss. IJCAI Int. Jt. Conf. Artif. Intell. 2018-July, 691–697.
- Felson, D.T., Neogi, T., 2004. Osteoarthritis: Is It a Disease of Cartilage or of Bone? Arthritis
 Rheum. 50, 341–344. https://doi.org/10.1002/art.20051
- Fink, B., Egl, M., Singer, J., Fuerst, M., Bubenheim, M., Neuen-Jacob, E., 2007.
 Morphologic changes in the vastus medialis muscle in patients with osteoarthritis of the
 knee. Arthritis Rheum. 56, 3626–3633. https://doi.org/10.1002/art.22960
- Gaj, S., Yang, M., Nakamura, K., Li, X., 2019. Automated cartilage and meniscus
 segmentation of knee MRI with conditional generative adversarial networks. Magn.
 Reson. Med. 1–13. https://doi.org/10.1002/mrm.28111
- Girshick, R., 2015. Fast R-CNN. Proc. IEEE Int. Conf. Comput. Vis. 2015 Inter, 1440–1448.
 https://doi.org/10.1109/ICCV.2015.169
- Goldring, M.B., Culley, K.L., Otero, M., 2017. Pathogenesis of Osteoarthritis in General, in:
 Grässel, S., Aszódi, A. (Eds.), Cartilage: Volume 2: Pathophysiology. Springer
 International Publishing. https://doi.org/10.1007/978-3-319-45803-8
- Goodfellow, I.J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S.,
 Courville, A., Bengio, Y., 2014. Generative Adversarial Networks. arXiv Prepr.
 arXiv1406.2661v1 1–9. https://doi.org/10.1001/jamainternmed.2016.8245
- Heimann, T., Styner, M., Warfield, S.K., 2010. Segmentation of Knee Images : A Grand
 Challenge Segmentation of Knee Images : [WWW Document]. URL
 http://www.ski10.org/ski10.pdf
- Hunter, D., Nevitt, M., Lynch, J., Kraus, V.B., Katz, J.N., Collins, J.E., Bowes, M.,
 Guermazi, A., Roemer, F.W., Losina, E., 2015. Longitudinal validation of periarticular
 bone area and 3D shape as biomarkers for knee OA progression? Data from the FNIH
 OA Biomarkers Consortium. Ann. Rheum. Dis. annrheumdis-2015-207602.

- Hunter, D.J., Eckstein, F., 2009. Exercise and osteoarthritis. J. Anat. 214, 197–207.
 https://doi.org/10.1111/j.1469-7580.2008.01013.x
- Huo, Y., Xu, Z., Bao, S., Assad, A., Abramson, R.G., Landman, B.A., 2018. Adversarial
 synthesis learning enables segmentation without target modality ground truth. Proc. Int. Symp. Biomed. Imaging 2018-April, 1217–1220.
 https://doi.org/10.1109/ISBI.2018.8363790
- Ismail, H.M., Vincent, T.L., 2017. Cartilage Injury and Osteoarthritis, in: Grässel, S., Aszódi,
 A. (Eds.), Cartilage: Volume 2: Pathophysiology. Springer International Publishing.
 https://doi.org/10.1007/978-3-319-45803-8
- Isola, P., Zhu, J.Y., Zhou, T., Efros, A.A., 2017. Image-to-image translation with conditional adversarial networks. Proc. - 30th IEEE Conf. Comput. Vis. Pattern Recognition, CVPR
 2017 5967–5976. https://doi.org/10.1109/CVPR.2017.632
- Kemnitz, J., Baumgartner, C.F., Eckstein, F., Chaudhari, A., Ruhdorfer, A., Wirth, W., Eder, 14 15 S.K., Konukoglu, E., 2019. Clinical evaluation of fully automated thigh muscle and adipose tissue segmentation using a U-Net deep learning architecture in context of 16 osteoarthritic knee Magn. Reson. Physics, Biol. Med. 17 pain. Mater. 18 https://doi.org/10.1007/s10334-019-00816-5
- Larobina, M., Murino, L., 2014. Medical image file formats. J. Digit. Imaging 27, 200–206.
 https://doi.org/10.1007/s10278-013-9657-9
- Le Troter, A., Fouré, A., Guye, M., Confort-Gouny, S., Mattei, J.P., Gondin, J., SalortCampana, E., Bendahan, D., 2016. Volume measurements of individual muscles in
 human quadriceps femoris using atlas-based segmentation approaches. Magn. Reson.
 Mater. Physics, Biol. Med. 29, 245–257. https://doi.org/10.1007/s10334-016-0535-6
- Lee, H., Hong, H., Kim, J., 2013. Anterior Cruciate Ligament Segmentation from Knee MR
 Images Using Graph Cuts with Geometric and Probabilistic Shape Constraints, in: Lee
- K.M., Matsushita Y., Rehg J.M., Hu Z. (Eds) Computer Vision ACCV 2012. Lecture
 Notes in Computer Science, Vol 7725. Springer, Berlin, Heidelberg. pp. 305–315.
- 29 https://doi.org/10.1007/978-3-642-37444-9_24
- Lee, J.G., Gumus, S., Moon, C.H., Kwoh, C.K., Bae, K.T., 2014. Fully automated
 segmentation of cartilage from the MR images of knee using a multi-atlas and local
 structural analysis method. Med. Phys. 41. https://doi.org/10.1118/1.4893533
- Li, C., Wand, M., 2016. Precomputed Real-Time Texture Synthesis with Markovian
 Generative Adversarial Networks. arXiv Prepr. arXiv1604.04382v1 1–17.

¹ https://doi.org/10.1136/annrheumdis-2015-207602

- Liu, F., 2019. SUSAN: segment unannotated image structure using adversarial network.
 Magn. Reson. Med. 81, 3330–3345. https://doi.org/10.1002/mrm.27627
- Liu, F., Zhou, Z., Jang, H., Samsonov, A., Zhao, G., Kijowski, R., 2017. Deep convolutional
 neural network and 3D deformable approach for tissue segmentation in musculoskeletal
 magnetic resonance imaging. Magn. Reson. Med. 79, 2379–2391.
 https://doi.org/10.1002/mrm.26841
- Lohmander, L.S., Englund, P.M., Dahl, L.L., Roos, E.M., 2007. The long-term consequence
 of anterior cruciate ligament and meniscus injuries: Osteoarthritis. Am. J. Sports Med.
 35, 1756–1769. https://doi.org/10.1177/0363546507307396
- Long, J., Shelhamer, E., Darrell, T., 2018. Fully Convolutional Adaptation Networks for
 Semantic Segmentation. Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit.
 6810–6818. https://doi.org/10.1109/CVPR.2018.00712
- MacKay, J.W., Kaggie, J.D., Treece, G.M., McDonnell, S.M., Khan, W., Roberts, A.R.,
 Janiczek, R.L., Graves, M.J., Turmezei, T.D., McCaskie, A.W., Gilbert, F.J., 2020.
 Three-Dimensional Surface-Based Analysis of Cartilage MRI Data in Knee
 Osteoarthritis: Validation and Initial Clinical Application. J. Magn. Reson. Imaging 1–
 13. https://doi.org/10.1002/jmri.27193
- MacKay, J.W., Kapoor, G., Driban, J.B., Lo, G.H., McAlindon, T.E., Toms, A.P., McCaskie,
 A.W., Gilbert, F.J., 2018. Association of subchondral bone texture on magnetic
 resonance imaging with radiographic knee osteoarthritis progression: data from the
 Osteoarthritis Initiative Bone Ancillary Study. Eur. Radiol. 28, 4687–4695.
 https://doi.org/10.1007/s00330-018-5444-9
- Mallya, Y., J., V., M. S., V., Venugopal, V.K., Mahajan, V., 2019. Automatic delineation of
 anterior and posterior cruciate ligaments by combining deep learning and deformable
 atlas based segmentation. Med. Imaging 2019 Biomed. Appl. Mol. Struct. Funct.
 Imaging. 10953. https://doi.org/10.1117/12.2512431
- Martel-Pelletier, J., Barr, A.J., Cicuttini, F.M., Conaghan, P.G., Cooper, C., Goldring, M.B.,
 Goldring, S.R., Jones, G., Teichtahl, A.J., Pelletier, J.P., 2016. Osteoarthritis. Nat. Rev.
 Dis. Prim. 2. https://doi.org/10.1038/nrdp.2016.72
- Messer, D.J., Shield, A.J., Williams, M.D., Timmins, R.G., Bourne, M.N., 2019. Hamstring
 muscle activation and morphology are significantly altered 1–6 years after anterior
 cruciate ligament reconstruction with semitendinosus graft. Knee Surgery, Sport.
 Traumatol. Arthrosc. 0, 0. https://doi.org/10.1007/s00167-019-05374-w
- Monu, U.D., Jordan, C.D., Samuelson, B.L., Hargreaves, B.A., Gold, G.E., McWalter, E.J.,

- 2017. Cluster analysis of quantitative MRI T 2 and T 1p relaxation times of cartilage
 identifies differences between healthy and ACL-injured individuals at 3T. Osteoarthr.
 Cartil. 25, 513–520. https://doi.org/10.1016/j.joca.2016.09.015
- Neogi, T., Bowes, M.A., Niu, J., De Souza, K.M., Vincent, G.R., Goggins, J., Zhang, Y.,
 Felson, D.T., 2013. Magnetic resonance imaging-based three-dimensional bone shape of
 the knee predicts onset of knee osteoarthritis: Data from the osteoarthritis initiative.
 Arthritis Rheum. 65, 2048–2058. https://doi.org/10.1002/art.37987
- Ng, H.P., Ong, S.H., Foong, K.W.C., Goh, P.S., Nowinski, W.L., 2006. Medical Image
 Segmentation Using K-Means Clustering and Improved Watershed Algorithm. IEEE
 Southwest Symp. Image Anal. Interpret. 61–65.
 https://doi.org/10.1109/SSIAI.2006.1633722
- Norman, B., Pedoia, V., Majumdar, S., 2018. Use of 2D U-Net Convolutional Neural
 Networks for Automated Cartilage and Meniscus Segmentation of Knee MR Imaging
 Data to Determine Relaxometry and Morphometry. Radiology 288, 177–185.
 https://doi.org/10.1148/radiol.2018172322
- 16 Paproki, A., Wilson, K.J., Surowiec, R.K., Ho, C.P., Pant, A., Bourgeat, P., Engstrom, C., Crozier, S., Fripp, J., 2016. Automated segmentation and T2-mapping of the posterior 17 18 cruciate ligament from MRI of the knee: Data from the osteoarthritis initiative. Proc. -IEEE Biomed. 19 2016 13th Int. Symp. Imaging 424-427. 20 https://doi.org/10.1109/ISBI.2016.7493298
- Paschos, N.K., 2017. Anterior cruciate ligament reconstruction and knee osteoarthritis. World
 J. Orthop. 8, 212–217. https://doi.org/10.5312/wjo.v8.i3.212
- Patel, F.K., Singh, M., 2018. Segmentation of cartilage from knee MRI images using the
 watershed algorithm. Int. J. Adv. Res. Ideas Innov. Technol. 4, 1727–1730.
- 25 Pathak, D., Krahenbuhl, P., Donahue, J., Darrell, T., Efros, A.A., 2016. Context Encoders:
- Feature Learning by Inpainting. arXiv Prepr. arXiv1604.07379v2 1–12.
 https://doi.org/10.1109/CVPR.2016.278
- Pedoia, V., Majumdar, S., Link, T.M., 2016. Segmentation of joint and musculoskeletal
 tissue in the study of arthritis. Magn. Reson. Mater. Physics, Biol. Med. 29, 207–221.
 https://doi.org/10.1007/s10334-016-0532-9
- Regmi, K., Borji, A., 2018. Cross-View Image Synthesis using Conditional GANs. arXiv
 Prepr. arXiv1803.03396v2.
- 33 Rezaei, M., Harmuth, K., Gierke, W., Kellermeier, T., Fischer, M., Yang, H., Meinel, C.,
- 34 2017. A Conditional Adversarial Network for Semantic Segmentation of Brain Tumor.

- BrainLes 2017 Brainlesion Glioma, Mult. Sclerosis, Stroke Trauma. Brain Inj. Springer
 10670, 241–252.
- Ronneberger, O., Fischer, P., Brox, T., 2015. U-Net: Convolutional Networks for Biomedical
 Image Segmentation. arXiv Prepr. arXiv1505.04597v1 1–8.
- Seim, H., Kainmueller, D., Lamecker, H., Bindernagel, M., Malinowski, J., Zachow, S.,
 2010. Model-based auto-segmentation of knee bones and cartilage in MRI data. Med.
 Image Anal. Clin. A Gd. Chall. 215–223.
- 8 Shan, L., Zach, C., Charles, C., Niethammer, M., 2014. Automatic atlas-based three-label
 9 cartilage segmentation from MR knee images. Med. Image Anal. 18, 1233–1246.
 10 https://doi.org/10.1016/j.media.2014.05.008
- Shie, C.K., Chuang, C.H., Chou, C.N., Wu, M.H., Chang, E.Y., 2015. Transfer representation
 learning for medical image analysis. Proc. Annu. Int. Conf. IEEE Eng. Med. Biol. Soc.
 EMBS 711–714. https://doi.org/10.1109/EMBC.2015.7318461
- Shrivastava, K., Gupta, N., Sharma, N., 2014. Medical Image Segmentation using Modified
 K Means Clustering. Int. J. Comput. Appl. 103, 12–16.
- Slemenda, C., Brandt, K.D., Heilman, D.K., Mazzuca, S., Braunstein, E.M., Katz, B.P.,
 Wolinsky, F.D., 1997. Quadriceps weakness and osteoarthritis of the knee. Ann. Intern.
 Med. 127, 97–104. https://doi.org/10.7326/0003-4819-127-2-199707150-00001
- Sørensen, T.J., 1948. A method of establishing groups of equal amplitude in plant sociology
 based on similarity of species and its application to analyses of the vegetation on Danish
 commons. Biol. Skr. 5, 1–34.
- Swanson, M.S., Prescott, J.W., Best, T.M., Powell, K., Jackson, R.D., Haq, F., Gurcan, M.N.,
 2010. Semi-automated Segmentation to Assess the Lateral Meniscus in Normal and
 Osteoarthritic Knees. Osteoarthr. Cartil. 18, 344–353.
- 25 https://doi.org/10.1016/j.joca.2009.10.004
- 26 The Osteoarthritis Initiative [WWW Document], n.d. . https://nda.nih.gov/oai/.
- Treece, G.M., Prager, R.W., Gee, A.H., 1999. Regularised marching tetrahedra: improved
 iso-surface extraction. Comuters Graph. 23, 583–598.
- Wise, B.L., Niu, J., Zhang, Y., Liu, F., Pang, J., Lynch, J.A., Lane, N.E., 2018. Bone shape
 mediates the relationship between sex and incident knee osteoarthritis. BMC
 Musculoskelet. Disord. 19, 1–9. https://doi.org/10.1186/s12891-018-2251-z
- 32 Xia, Y., Fripp, J., Chandra, S.S., Schwarz, R., Engstrom, C., Crozier, S., 2013. Automated
- bone segmentation from large field of view 3D MR images of the hip joint. Phys. Med.
- Biol. 58, 7375–7390. https://doi.org/10.1088/0031-9155/58/20/7375

1	Yang, D., Xu, D., Zhou, K., Georgescu, B., Chen, M., Grbic, S., Metaxas, D., Comaniciu, D.,
2	2017. Automatic Liver Segmentation Using an Adversarial Image-to-Image Network.
3	Med. Image Comput. Comput. Assist. Interv MICCAI 2017, Springer 10435, 507-
4	515.
5	Yushkevich, P.A., Piven, J., Hazlett, H.C., Smith, R.G., Ho, S., Gee, J.C., Gerig, G., 2006.
6	User-guided 3D active contour segmentation of anatomical structures: Significantly
7	improved efficiency and reliability. Neuroimage 31, 1116–1128.
8	https://doi.org/10.1016/j.neuroimage.2006.01.015
9	Zhao, J.J., Mathieu, M., LeCun, Y., 2017. Energy-based Generative Adversarial Network.
10	arXiv Prepr. arXiv 1609.03126v4.
11	Zhou, L., Chav, R., Cresson, T., Chartrand, G., De Guise, J., 2016. 3D knee segmentation
12	based on three MRI sequences from different planes. Proc. Annu. Int. Conf. IEEE Eng.
13	Med. Biol. Soc. EMBS 2016-Octob, 1042–1045.
14	https://doi.org/10.1109/EMBC.2016.7590881
15	Zhou, Z., Zhao, G., Kijowski, R., Liu, F., 2018. Deep convolutional neural network for
16	segmentation of knee joint anatomy. Magn. Reson. Med. 80, 2759–2770.
17	https://doi.org/10.1002/mrm.27229
18	
19	
20	
21	
22	
23	
24	
25	
26	
27	
28	
29	
30	
31	
32	
33	
34	

1 Tables

Table 1 - Participant characteristics showing the mean age, number of males/females (M/F),

3 average body-mass-index (BMI), Kellgren-Lawrence (KL) osteoarthritis score and the 4 number of training/testing set images of the locally acquired dataset. Additionally, the

5 number of participants (N) and training/testing set images of the SKI10 and OAI ZIB datasets

6 are given.

Dataset	Variable	Training Set	Testing Set
Local	Ν	8	2
	Images	806	171
	Mean Age (years)	53	52
	Sex (M/F)	5/3	0/2
	Mean BMI (kg/m ²)	27.8	27.7
	KL (0/2/3)	4/1/3	1/1/0
SKI10	N	70	30
	Images	6133	2626
OAI ZIB	N	355	152
	Images	43814	18517

1 **Table 2** – **Results of the Network Objective Function: cGAN.** The influence of mixing the cGAN objective with different pixel-wise error losses and varying 2 their significance by changing the weighting hyperparameter λ on the segmentation performance of the proposed cGAN was assessed. Highest network scores 3 achieved for each tissue are highlighted grey and in bold.

- 4 Training and testing were performed on the AMROA training and testing datasets, respectively.
- 5 Results are presented as mean \pm standard deviation.
- 6 Abbreviations: F Bone femoral bone, T Bone tibial bone, P Bone patellar bone, F Cartilage femoral cartilage, T Cartilage tibial cartilage, P Cartilage –
- patellar cartilage, VM Muscle vastus medialis muscle, GM Muscle medial head of gastrocnemius medialis muscle, ACL anterior cruciate ligament, PCL –
 posterior cruciate ligament, DSC Sørensen–Dice similarity coefficient

	Network Objective Function Results										
	cGAN										
Pixel Loss	λ	F Bone	T Bone	P Bone	F Cartilage	T Cartilage	P Cartilage	VM Muscle	GM Muscle	ACL	PCL
		DSC	DSC	DSC	DSC	DSC	DSC	DSC	DSC	DSC	DSC
L1	0	0.931 ± 0.020	0.864 ± 0.008	0.911 ± 0.036	0.774 ± 0.030	0.717 ± 0.108	0.872 ± 0.030	0.000 ± 0.000	0.000 ± 0.000	0.000 ± 0.000	0.000 ± 0.000
	0.01	0.900 ± 0.018	0.890 ± 0.031	0.912 ± 0.002	0.727 ± 0.023	0.715 ± 0.060	0.850 ± 0.048	0.000 ± 0.000	0.000 ± 0.000	0.509 ± 0.009	0.171 ± 0.208
	1	0.899 ± 0.014	0.856 ± 0.010	0.807 ± 0.060	0.465 ± 0.037	0.666 ± 0.022	0.426 ± 0.098	0.611 ± 0.181	0.595 ± 0.054	0.000 ± 0.000	0.000 ± 0.000
	100	0.918 ± 0.011	0.948 ± 0.018	0.928 ± 0.002	0.812 ± 0.002	0.748 ± 0.042	$\textbf{0.863} \pm \textbf{0.043}$	0.113 ± 0.085	0.000 ± 0.000	0.577 ± 0.020	0.073 ± 0.103
	10000	$\textbf{0.968} \pm \textbf{0.006}$	0.944 ± 0.026	0.917 ± 0.008	$\textbf{0.875} \pm \textbf{0.021}$	$\textbf{0.810} \pm \textbf{0.036}$	0.840 ± 0.065	0.879 ± 0.036	0.793 ± 0.080	0.432 ± 0.237	0.338 ± 0.386
L2	0.01	0.902 ± 0.004	0.915 ± 0.003	0.923 ± 0.005	0.750 ± 0.002	0.740 ± 0.079	0.834 ± 0.077	0.000 ± 0.000	0.000 ± 0.000	0.000 ± 0.000	0.000 ± 0.000
	1	0.902 ± 0.046	0.902 ± 0.008	0.902 ± 0.044	0.741 ± 0.004	0.736 ± 0.033	0.838 ± 0.041	0.000 ± 0.000	0.000 ± 0.000	0.149 ± 0.104	0.002 ± 0.002
	100	0.928 ± 0.015	0.939 ± 0.007	0.921 ± 0.022	0.768 ± 0.016	0.752 ± 0.049	0.862 ± 0.039	0.001 ± 0.001	0.000 ± 0.000	$\textbf{0.652} \pm \textbf{0.094}$	0.101 ± 0.074
	10000	0.952 ± 0.000	$\textbf{0.950} \pm \textbf{0.015}$	0.923 ± 0.001	0.828 ± 0.043	0.684 ± 0.092	0.832 ± 0.054	0.814 ± 0.145	0.856 ± 0.121	0.440 ± 0.084	0.293 ± 0.358
SmL1	0.01	0.914 ± 0.034	0.902 ± 0.003	0.920 ± 0.011	0.726 ± 0.007	0.729 ± 0.042	0.762 ± 0.068	0.000 ± 0.000	0.000 ± 0.000	0.343 ± 0.066	0.000 ± 0.000
	1	0.884 ± 0.044	0.912 ± 0.006	0.926 ± 0.013	0.740 ± 0.014	0.732 ± 0.044	0.829 ± 0.067	0.055 ± 0.007	0.000 ± 0.000	0.000 ± 0.000	0.000 ± 0.000
	100	0.903 ± 0.019	0.944 ± 0.006	0.936 ± 0.003	0.776 ± 0.035	0.741 ± 0.066	0.857 ± 0.029	0.031 ± 0.044	0.070 ± 0.100	0.578 ± 0.053	0.044 ± 0.052
	10000	0.951 ± 0.002	0.946 ± 0.018	0.935 ± 0.015	0.825 ± 0.035	0.738 ± 0.047	0.797 ± 0.088	0.914 ± 0.001	$\textbf{0.837} \pm \textbf{0.146}$	0.261 ± 0.073	$\textbf{0.374} \pm \textbf{0.341}$

1 Table 3 – Results of the Network Objective Function: U-Net. The influence of different pixel-wise error losses on the segmentation performance of the U-Net 2 was assessed. Highest network scores achieved for each tissue are highlighted grey and in bold.

- 3 Training and testing were performed on the AMROA training and testing datasets, respectively.
- 4 Results are presented as mean \pm standard deviation.
- 5 Abbreviations: F Bone femoral bone, T Bone tibial bone, P Bone patellar bone, F Cartilage femoral cartilage, T Cartilage tibial cartilage, P Cartilage –
- 6 patellar cartilage, VM Muscle vastus medialis muscle, GM Muscle medial head of gastrocnemius medialis muscle, ACL anterior cruciate ligament, PCL -
- 7 posterior cruciate ligament, DSC Sørensen–Dice similarity coefficient

	Network Objective Function Results										
U-Net											
Pixel Loss	F Bone	T Bone	P Bone	F Cartilage	T Cartilage	P Cartilage	VM Muscle	GM Muscle	ACL	PCL	
	DSC	DSC	DSC	DSC	DSC	DSC	DSC	DSC	DSC	DSC	
L1	$\textbf{0.972} \pm \textbf{0.006}$	$\textbf{0.960} \pm \textbf{0.001}$	$\textbf{0.941} \pm \textbf{0.010}$	$\textbf{0.886} \pm \textbf{0.007}$	$\textbf{0.834} \pm \textbf{0.010}$	$\textbf{0.890} \pm \textbf{0.034}$	0.000 ± 0.000	0.000 ± 0.000	$\textbf{0.643} \pm \textbf{0.153}$	$\textbf{0.641} \pm \textbf{0.008}$	
L2	0.950 ± 0.007	0.957 ± 0.009	0.939 ± 0.003	0.831 ± 0.020	0.723 ± 0.068	0.837 ± 0.051	0.888 ± 0.000	0.881 ± 0.021	0.491 ± 0.136	0.428 ± 0.196	
SmL1	0.953 ± 0.001	0.953 ± 0.009	0.937 ± 0.004	0.843 ± 0.021	0.771 ± 0.036	0.830 ± 0.088	0.894 ± 0.002	0.910 ± 0.045	0.574 ± 0.230	0.463 ± 0.174	

11

12

13

14

15

16

1 **Table 4 - Results of additionally testing on noise only images.** The influence of including noise only images in the testing set on the overall segmentation 2 performance of a cGAN trained with $\mathcal{L}_{cGAN} + \lambda \mathcal{L}_{L1}$ ($\lambda = 100$) loss objective and a U-Net trained with \mathcal{L}_{L1} objective. Training was performed on the AMROA 3 training dataset without noise only images.

4 Abbreviations: F Bone – femoral bone, T Bone – tibial bone, P Bone – patellar bone, F Cartilage – femoral cartilage, T Cartilage – tibial cartilage, P Cartilage –

5 patellar cartilage, VM Muscle - vastus medialis muscle, GM Muscle – medial head of gastrocnemius medialis muscle, ACL – anterior cruciate ligament, PCL – 6 posterior cruciate ligament, DSC – Sørensen, Dice similarity coefficient % Diff – absolute percentage difference

6 posterior cruciate ligament, DSC - Sørensen–Dice similarity coefficient, %-Diff – absolute percentage difference

		Influence of Noise Only Images										
	cGAN											
_	Testing	F Bone DSC	T Bone DSC	P Bone DSC	F Cartilage DSC	T Cartilage DSC	P Cartilage DSC	VM Muscle DSC	GM Muscle DSC	ACL DSC	PCL DSC	
	No Noise With Noise	$\begin{array}{c} 0.918 \pm 0.011 \\ 0.925 \pm 0.012 \end{array}$	$\begin{array}{c} 0.948 \pm 0.018 \\ 0.946 \pm 0.017 \end{array}$	$\begin{array}{c} 0.928 \pm 0.002 \\ 0.928 \pm 0.004 \end{array}$	$\begin{array}{c} 0.812 \pm 0.002 \\ 0.810 \pm 0.003 \end{array}$	$\begin{array}{c} 0.748 \pm 0.042 \\ 0.752 \pm 0.045 \end{array}$	$\begin{array}{c} 0.863 \pm 0.043 \\ 0.858 \pm 0.054 \end{array}$	0.113 ± 0.085 0.098 ± 0.114	0.000 ± 0.000 0.000 ± 0.000	$\begin{array}{c} 0.577 \pm 0.020 \\ 0.593 \pm 0.028 \end{array}$	$\begin{array}{c} 0.073 \pm 0.103 \\ 0.092 \pm 0.131 \end{array}$	
_	%-Diff	0.7	0.2	0.0	0.2	0.4	0.5	1.5	0.0	1.6	1.9	
	U-Net											
_	Testing	F Bone DSC	T Bone DSC	P Bone DSC	F Cartilage DSC	T Cartilage DSC	P Cartilage DSC	VM Muscle DSC	GM Muscle DSC	ACL DSC	PCL DSC	
	No Noise With Noise	0.972 ± 0.006 0.968 ± 0.001	0.960 ± 0.001 0.957 ± 0.009	$0.941 \pm 0.010 \\ 0.938 \pm 0.016 \\ 0.2$	$0.886 \pm 0.007 \\ 0.885 \pm 0.004 \\ 0.1$	0.834 ± 0.010 0.833 ± 0.010	0.890 ± 0.034 0.894 ± 0.026	$\begin{array}{c} 0.000 \pm 0.000 \\ 0.000 \pm 0.000 \end{array}$	0.000 ± 0.000 0.000 ± 0.000	$0.643 \pm 0.153 \\ 0.620 \pm 0.156 \\ 2.3$	$0.641 \pm 0.008 \\ 0.643 \pm 0.025 \\ 0.2$	
7	%-D1II	0.4	0.5	0.5	0.1	0.1	0.4	0.0	0.0	2.3	0.2	
, 8												
9												
10												
11												
12												

1

11

Table 5 – Results of Altering the Loss Objective during Training. Assessing the influence of altering the loss objective function during training on the segmentation performance of the proposed cGAN and U-Net. A cGAN was trained with $\mathcal{L}_{cGAN} + \lambda \mathcal{L}_{L2}$ objective and a U-Net with \mathcal{L}_{L2} objective for 50 epochs followed by a further 50 epochs training with $\mathcal{L}_{cGAN} + \lambda \mathcal{L}_{L1}$ and \mathcal{L}_{L1} objectives, respectively. Segmentation performances are compared with the previously trained cGANs ($\mathcal{L}_{cGAN} + \lambda \mathcal{L}_{L1}$ and $\mathcal{L}_{cGAN} + \lambda \mathcal{L}_{L2}$; $\lambda = 100$; 100 epochs) and U-Nets (\mathcal{L}_{L1} and \mathcal{L}_{L2} ;100 epochs). Highest network scores achieved for each tissue are highlighted grey and in bold.

7 Training and testing were performed on the AMROA training and testing datasets, respectively.

- Abbreviations: FB femoral bone, TB tibial bone, PB patellar bone, FC femoral cartilage, TC tibial cartilage, PC patellar cartilage, VM Muscle vastus
 medialis muscle, GM Muscle medial head of gastrocnemius medialis muscle, ACL anterior cruciate ligament, PCL posterior cruciate ligament, DSC -
- 10 Sørensen–Dice similarity coefficient

	Altering the Loss Objective during Training Results										
	cGAN										
Network Loss	F Bone	T Bone	P Bone	F Cartilage	T Cartilage	P Cartilage	VM Muscle	GM Muscle	ACL	PCL	
Objective	DSC	DSC	DSC	DSC	DSC	DSC	DSC	DSC	DSC	DSC	
$\mathcal{L}_{cGAN} + \lambda \mathcal{L}_{L1}$	0.918 ± 0.011	$\textbf{0.948} \pm \textbf{0.018}$	0.928 ± 0.002	0.812 ± 0.002	0.748 ± 0.042	0.863 ± 0.043	0.113 ± 0.085	0.000 ± 0.000	0.577 ± 0.020	0.073 ± 0.103	
$\mathcal{L}_{cGAN} + \lambda \mathcal{L}_{L2}$	0.928 ± 0.015	0.939 ± 0.007	0.921 ± 0.022	0.768 ± 0.016	0.752 ± 0.049	0.862 ± 0.039	0.001 ± 0.001	0.000 ± 0.000	0.652 ± 0.094	0.101 ± 0.074	
$ \begin{array}{c} \mathcal{L}_{cGAN} + \lambda \mathcal{L}_{L2} \rightarrow \\ \mathcal{L}_{cGAN} + \lambda \mathcal{L}_{L1} \end{array} $	0.936 ± 0.007	0.938 ± 0.021	0.884 ± 0.078	0.800 ± 0.021	$\textbf{0.760} \pm \textbf{0.035}$	0.855 ± 0.031	$\boldsymbol{0.739 \pm 0.010}$	$\boldsymbol{0.772 \pm 0.005}$	0.115 ± 0.032	0.392 ± 0.128	
					U-Net						
Network Loss	F Bone	T Bone	P Bone	F Cartilage	T Cartilage	P Cartilage	VM Muscle	GM Muscle	ACL	PCL	
Objective	DSC	DSC	DSC	DSC	DSC	DSC	DSC	DSC	DSC	DSC	
\mathcal{L}_{L1}	0.972 ± 0.006	0.960 ± 0.001	0.941 ± 0.010	$\textbf{0.886} \pm \textbf{0.007}$	0.834 ± 0.010	$\textbf{0.890} \pm \textbf{0.034}$	0.000 ± 0.000	0.000 ± 0.000	0.643 ± 0.153	0.641 ± 0.008	
\mathcal{L}_{L2}	0.950 ± 0.007	0.957 ± 0.009	0.939 ± 0.003	0.831 ± 0.020	0.723 ± 0.068	0.837 ± 0.051	0.888 ± 0.000	0.881 ± 0.021	0.491 ± 0.136	0.428 ± 0.196	
$\mathcal{L}_{L2} \rightarrow \mathcal{L}_{L1}$	0.970 ± 0.006	$\textbf{0.961} \pm \textbf{0.007}$	$\textbf{0.941} \pm \textbf{0.003}$	0.869 ± 0.016	0.793 ± 0.021	0.886 ± 0.027	$\textbf{0.914} \pm \textbf{0.008}$	$\textbf{0.933} \pm \textbf{0.010}$	0.632 ± 0.170	0.567 ± 0.094	

1 Table 6 – Results of Varying Generator Network Depth: Number of Convolutions. The influence of varying the number of convolutions during down2 sampling in the generator networks of both the cGAN and U-Net was assessed. Highest network scores achieved for each tissue are highlighted grey and in bold.

3 Training and testing were performed on the AMROA training and testing datasets, respectively.

4 Results are presented as mean \pm standard deviation.

5 Abbreviations: F Bone – femoral bone, T Bone – tibial bone, P Bone – patellar bone, F Cartilage – femoral cartilage, T Cartilage – tibial cartilage, P Cartilage –

6 patellar cartilage, VM Muscle - vastus medialis muscle, GM Muscle - medial head of gastrocnemius medialis muscle, ACL - anterior cruciate ligament, PCL -

7 posterior cruciate ligament, DSC - Sørensen–Dice similarity coefficient

	Generator Network Depth Results – Number of Convolutions during Down-Sampling										
	cGAN										
Number	F Bone	T Bone	P Bone	F Cartilage	T Cartilage	P Cartilage	VM Muscle	GM Muscle	ACL	PCL	
Down Convs	DSC	DSC	DSC	DSC	DSC	DSC	DSC	DSC	DSC	DSC	
5	$\textbf{0.928} \pm \textbf{0.006}$	0.929 ± 0.006	0.893 ± 0.029	0.721 ± 0.029	0.751 ± 0.039	0.838 ± 0.042	0.049 ± 0.069	0.000 ± 0.000	0.622 ± 0.042	0.286 ± 0.189	
7	0.889 ± 0.023	0.921 ± 0.026	0.928 ± 0.002	0.764 ± 0.047	0.624 ± 0.039	0.846 ± 0.057	0.171 ± 0.226	$\textbf{0.167} \pm \textbf{0.236}$	$\textbf{0.626} \pm \textbf{0.041}$	$\textbf{0.289} \pm \textbf{0.408}$	
9	0.918 ± 0.011	$\textbf{0.948} \pm \textbf{0.018}$	0.928 ± 0.002	$\textbf{0.812} \pm \textbf{0.002}$	0.748 ± 0.042	0.863 ± 0.043	0.113 ± 0.085	0.000 ± 0.000	0.577 ± 0.020	0.073 ± 0.103	
					U-Net						
5	0.969 ± 0.002	0.952 ± 0.016	0.919 ± 0.022	$\textbf{0.887} \pm \textbf{0.018}$	0.823 ± 0.001	0.888 ± 0.031	0.000 ± 0.000	0.000 ± 0.000	0.631 ± 0.125	0.544 ± 0.249	
7	0.964 ± 0.003	0.956 ± 0.005	0.921 ± 0.008	0.874 ± 0.032	0.787 ± 0.044	0.869 ± 0.029	0.000 ± 0.000	0.000 ± 0.000	0.539 ± 0.160	0.592 ± 0.120	
9	$\boldsymbol{0.972 \pm 0.006}$	0.960 ± 0.001	0.941 ± 0.010	0.886 ± 0.007	$\textbf{0.834} \pm \textbf{0.010}$	$\textbf{0.890} \pm \textbf{0.034}$	0.000 ± 0.000	0.000 ± 0.000	$\textbf{0.643} \pm \textbf{0.153}$	$\textbf{0.641} \pm \textbf{0.008}$	

8

Table 7 - Results of Varying Generator Network Depth: Number of Minimum Feature Maps. The influence of starting with different numbers of minimum 1

- feature channel maps in the generator networks of both the cGAN and U-Net was assessed. Highest network scores achieved for each tissue are highlighted grey 2 and in bold. 3
- Training and testing were performed on the AMROA training and testing datasets, respectively. 4
- Results are presented as mean \pm standard deviation. 5
- Abbreviations: F Bone femoral bone, T Bone tibial bone, P Bone patellar bone, F Cartilage femoral cartilage, T Cartilage tibial cartilage, P Cartilage -6
- patellar cartilage, VM Muscle vastus medialis muscle, GM Muscle medial head of gastrocnemius medialis muscle, ACL anterior cruciate ligament, PCL -7 posterior cruciate ligament, DSC - Sørensen-Dice similarity coefficient
- 8

Generator Network Depth Results – Number of Minimum Feature Channel Maps											
cGAN											
Feature Maps	F Bone DSC	T Bone DSC	P Bone DSC	F Cartilage DSC	T Cartilage DSC	P Cartilage DSC	VM Muscle DSC	GM Muscle DSC	ACL DSC	PCL DSC	
16	0.774 ± 0.059	0.903 ± 0.040	0.858 ± 0.003	0.547 ± 0.236	0.473 ± 0.269	0.771 ± 0.070	0.000 ± 0.000	0.000 ± 0.000	0.000 ± 0.000	0.000 ± 0.000	
32	0.899 ± 0.004	0.937 ± 0.001	0.875 ± 0.027	0.750 ± 0.028	0.720 ± 0.038	0.831 ± 0.030	$\textbf{0.414} \pm \textbf{0.260}$	0.000 ± 0.000	0.000 ± 0.000	0.000 ± 0.000	
64	0.918 ± 0.011	$\textbf{0.948} \pm \textbf{0.018}$	$\textbf{0.928} \pm \textbf{0.002}$	$\textbf{0.812} \pm \textbf{0.002}$	0.748 ± 0.042	0.863 ± 0.043	0.113 ± 0.085	0.000 ± 0.000	$\textbf{0.577} \pm \textbf{0.020}$	$\textbf{0.073} \pm \textbf{0.103}$	
128	$\textbf{0.925} \pm \textbf{0.006}$	0.935 ± 0.021	0.831 ± 0.032	0.805 ± 0.010	$\textbf{0.773} \pm \textbf{0.081}$	0.784 ± 0.061	0.341 ± 0.256	0.000 ± 0.000	0.336 ± 0.219	0.011 ± 0.016	
					U-Net						
16	0.966 ± 0.000	0.950 ± 0.021	0.912 ± 0.028	0.868 ± 0.011	0.795 ± 0.001	0.864 ± 0.028	0.000 ± 0.000	0.000 ± 0.000	0.000 ± 0.000	0.202 ± 0.110	
32	0.969 ± 0.006	0.946 ± 0.016	0.914 ± 0.005	0.875 ± 0.026	0.795 ± 0.051	0.878 ± 0.032	0.000 ± 0.000	0.000 ± 0.000	0.000 ± 0.000	0.453 ± 0.039	
64	$\textbf{0.972} \pm \textbf{0.006}$	$\textbf{0.960} \pm \textbf{0.001}$	$\textbf{0.941} \pm \textbf{0.010}$	$\textbf{0.886} \pm \textbf{0.007}$	$\textbf{0.834} \pm \textbf{0.010}$	0.890 ± 0.034	0.000 ± 0.000	0.000 ± 0.000	0.643 ± 0.153	$\textbf{0.641} \pm \textbf{0.008}$	
128	0.968 ± 0.006	0.960 ± 0.004	0.929 ± 0.014	0.884 ± 0.022	0.823 ± 0.010	$\textbf{0.897} \pm \textbf{0.013}$	0.000 ± 0.000	0.000 ± 0.000	0.645 ± 0.053	0.597 ± 0.025	

Table 8 – **Results of PatchGAN Receptive Field Size**. Comparison of segmentation performance of the proposed cGAN with different N x N receptive field sizes of the PatchGAN discriminator network. Highest network scores achieved for each tissue are highlighted grey and in bold.

3 The cGANs were trained with the $\mathcal{L}_{cGAN} + \lambda \mathcal{L}_{L1}$ objective with $\lambda = 100$ with training and testing being performed on the AMROA dataset.

7 8

9

10

11

12

13

14

15

16

4 Abbreviations: FB - femoral bone, TB - tibial bone, PB - patellar bone, FC - femoral cartilage, TC - tibial cartilage, PC - patellar cartilage, VM Muscle -

vastus medialis muscle, GM Muscle – medial head of gastrocnemius medialis muscle, ACL – anterior cruciate ligament, PCL – posterior cruciate ligament, DSC
 Sørensen–Dice similarity coefficient

PatchGAN Receptive Field Size Results T Cartilage Receptive F Bone P Bone F Cartilage P Cartilage VM Muscle GM Muscle ACL PCL T Bone Field Size DSC 0.613 ± 0.143 0.849 ± 0.046 0.869 ± 0.069 1 x 1 0.971 ± 0.005 $\textbf{0.953} \pm \textbf{0.012}$ 0.947 ± 0.007 0.804 ± 0.024 0.869 ± 0.053 0.812 ± 0.066 0.618 ± 0.140 34 x 34 0.968 ± 0.007 0.952 ± 0.015 0.941 ± 0.013 0.849 ± 0.002 0.795 ± 0.013 0.868 ± 0.023 $\textbf{0.883} \pm \textbf{0.007}$ 0.876 ± 0.009 0.621 ± 0.096 0.594 ± 0.118 70 x 70 0.918 ± 0.011 0.948 ± 0.018 0.928 ± 0.002 0.812 ± 0.002 0.748 ± 0.042 0.863 ± 0.043 0.113 ± 0.085 0.577 ± 0.020 0.073 ± 0.103 0.000 ± 0.000 286 x 286 0.941 ± 0.000 0.938 ± 0.008 0.920 ± 0.012 0.766 ± 0.020 0.731 ± 0.003 0.767 ± 0.049 0.702 ± 0.022 0.597 ± 0.078 0.383 ± 0.090 0.070 ± 0.022

1 Table 9 – Results of Transfer Learning. Comparison of segmentation performance of the proposed cGAN and U-Net without and with transfer learning and testing on the

2 SKI10 and OAI ZIB testing dataset. Highest network scores achieved for each tissue are highlighted grey and in bold.

3 SKI10/OAI ZIB \rightarrow AMROA: Pretraining the network for 20 epochs on the SKI10/OAI ZIB dataset followed by network fine-tuning for 80 epochs on the AMROA dataset.

4 AMROA \rightarrow SKI10/OAI ZIB: Pretraining the network for 20 epochs on the AMROA dataset followed by network fine-tuning for 80 epochs on the SKI10/OAI ZIB dataset.

5 Abbreviations: FB – femoral bone, TB – tibial bone, FC – femoral cartilage, TC – tibial cartilage, DSC - Sørensen–Dice similarity coefficient, ASD – average surface distance,

6 VOE – volumetric overlap error

Transfer Learning Results											
SKI10 Testing											
Network	Training	F Bone		T Bone		F Cartilage		T Cartilage			
		DSC	ASD	DSC	ASD	DSC	VOE	DSC	VOE		
	AMROA	0.929 ± 0.040	3.726 ± 1.758	0.893 ± 0.069	3.368 ± 1.935	0.488 ± 0.093	67.19 ± 8.36	0.465 ± 0.114	69.01 ± 10.00		
	SKI10	0.974 ± 0.013	1.445 ± 1.918	0.979 ± 0.007	0.527 ± 0.403	0.736 ± 0.058	41.49 ± 6.99	$\textbf{0.684} \pm \textbf{0.070}$	$\textbf{47.58} \pm \textbf{7.98}$		
cGAN	SKI10 → AMROA	0.938 ± 0.039	3.229 ± 1.776	0.929 ± 0.041	2.696 ± 2.326	0.544 ± 0.077	62.23 ± 7.45	0.480 ± 0.100	67.86 ± 8.89		
	AMROA → SKI10	0.974 ± 0.012	1.280 ± 1.484	0.977 ± 0.010	0.802 ± 1.139	$\textbf{0.738} \pm \textbf{0.059}$	41.19 ± 7.08	0.675 ± 0.071	48.65 ± 7.94		
	AMROA	0.925 ± 0.038	1.856 ± 0.997	0.907 ± 0.055	1.868 ± 1.336	0.545 ± 0.082	62.16 ± 7.62	0.462 ± 0.112	69.26 ± 9.86		
	SKI10	0.973 ± 0.015	0.756 ± 0.995	0.978 ± 0.008	0.254 ± 0.340	0.728 ± 0.058	42.42 ± 6.88	0.674 ± 0.066	48.85 ± 7.55		
U-Net	SKI10 → AMROA	0.943 ± 0.032	1.071 ± 0.682	0.936 ± 0.038	1.436 ± 1.083	0.576 ± 0.078	59.18 ± 7.86	0.456 ± 0.115	69.76 ± 9.93		
	AMROA → SKI10	0.975 ± 0.013	$\textbf{0.440} \pm \textbf{0.492}$	$\boldsymbol{0.979 \pm 0.007}$	0.258 ± 0.288	0.731 ± 0.056	42.08 ± 6.74	0.670 ± 0.070	49.19 ± 7.84		
	OAI ZIB Testing										
	AMROA	0.939 ± 0.016	4.153 ± 1.962	0.914 ± 0.080	4.681 ± 3.197	0.611 ± 0.068	55.66 ± 7.10	0.601 ± 0.089	56.44 ± 9.14		
	OAI ZIB	0.985 ± 0.002	0.328 ± 0.123	0.985 ± 0.003	0.293 ± 0.072	0.895 ± 0.023	18.92 ± 3.64	0.839 ± 0.040	$\textbf{27.55} \pm \textbf{5.90}$		
cGAN	OAI ZIB → AMROA	0.961 ± 0.009	1.786 ± 1.202	0.961 ± 0.018	4.426 ± 2.902	0.641 ± 0.071	52.41 ± 7.87	0.738 ± 0.055	41.23 ± 6.70		
	AMROA → OAI ZIB	0.985 ± 0.002	0.403 ± 0.268	$\textbf{0.985} \pm \textbf{0.003}$	$\textbf{0.293} \pm \textbf{0.068}$	$\textbf{0.897} \pm \textbf{0.022}$	18.68 ± 3.57	0.837 ± 0.042	27.82 ± 6.19		
	AMROA	0.934 ± 0.015	5.424 ± 2.799	0.915 ± 0.094	6.282 ± 3.647	0.643 ± 0.065	52.26 ± 7.03	0.626 ± 0.063	54.12 ± 6.74		
	OAI ZIB	0.985 ± 0.002	0.388 ± 0.169	0.984 ± 0.003	0.304 ± 0.079	0.896 ± 0.020	18.83 ± 3.19	0.837 ± 0.038	27.80 ± 5.57		
U-Net	OAI ZIB → AMROA	0.966 ± 0.006	1.244 ± 0.791	0.961 ± 0.017	1.880 ± 1.133	0.734 ± 0.046	41.83 ± 5.82	0.741 ± 0.058	40.83 ± 6.97		
	AMROA → OAI ZIB	0.985 ± 0.002	0.390 ± 0.361	0.985 ± 0.003	0.327 ± 0.127	0.893 ± 0.023	19.24 ± 3.64	0.838 ± 0.037	27.75 ± 5.50		

1 Table 10 – Results of Transfer Learning. Comparison of segmentation performance of the proposed cGAN and U-Net without and with transfer learning and testing on the 2 AMROA testing dataset. Highest network scores achieved for each tissue are highlighted grey and in bold.

SKI10/OAI ZIB \rightarrow AMROA: Pretraining the network for 20 epochs on the SKI10/OAI ZIB dataset followed by network fine-tuning for 80 epochs on the AMROA dataset.

AMROA \rightarrow SKI10/OAI ZIB: Pretraining the network for 20 epochs on the AMROA dataset followed by network fine-tuning for 80 epochs on the SKI10/OAI ZIB dataset.

Abbreviations: FB – femoral bone, TB – tibial bone, PB – patellar bone, FC – femoral cartilage, TC – tibial cartilage, PC – patellar cartilage, VM Muscle - vastus medialis muscle, GM Muscle – medial head of gastrocnemius medialis muscle, ACL – anterior cruciate ligament, PCL – posterior cruciate ligament, DSC - Sørensen–Dice similarity

7 coefficient

3

4 5

6

Transfer Learning Results - AMROA Testing											
Network	Training	F Bone	T Bone	P Bone	F Cartilage	T Cartilage	P Cartilage	VM Muscle	GM Muscle	ACL	PCL
		DSC	DSC	DSC	DSC	DSC	DSC	DSC	DSC	DSC	DSC
	AMROA	0.971 ± 0.005	0.953 ± 0.012	0.947 ± 0.007	0.849 ± 0.046	0.804 ± 0.024	0.869 ± 0.053	0.812 ± 0.066	0.869 ± 0.069	0.618 ± 0.140	0.613 ± 0.143
	SKI10	0.940 ± 0.024	0.947 ± 0.013		0.735 ± 0.005	0.561 ± 0.190					
	OAI ZIB	0.962 ± 0.009	0.951 ± 0.010		0.817 ± 0.032	0.790 ± 0.014					
cGAN	SKI10 → AMROA	0.970 ± 0.008	0.961 ± 0.004	0.940 ± 0.001	0.871 ± 0.029	0.774 ± 0.039	0.858 ± 0.038	0.922 ± 0.037	0.897 ± 0.057	0.586 ± 0.043	0.468 ± 0.186
	OAI ZIB → AMROA	0.972 ± 0.003	0.962 ± 0.001	0.947 ± 0.001	0.875 ± 0.026	0.811 ± 0.042	0.879 ± 0.022	0.908 ± 0.053	$\boldsymbol{0.909 \pm 0.077}$	0.664 ± 0.058	0.652 ± 0.112
	AMROA → SKI10	0.954 ± 0.015	0.949 ± 0.005		0.761 ± 0.025	0.544 ± 0.085					
	AMROA → OAI ZIB	0.960 ± 0.007	0.951 ± 0.012		0.821 ± 0.042	0.815 ± 0.015					
	AMROA	0.972 ± 0.006	0.960 ± 0.001	0.941 ± 0.010	0.886 ± 0.007	$\textbf{0.834} \pm \textbf{0.010}$	0.890 ± 0.034	0.000 ± 0.000	0.000 ± 0.000	0.643 ± 0.153	0.641 ± 0.008
	SKI10	0.937 ± 0.031	0.944 ± 0.026		0.754 ± 0.009	0.637 ± 0.044					
	OAI ZIB	0.959 ± 0.003	0.953 ± 0.010		0.820 ± 0.026	0.798 ± 0.012					
U-Net	SKI10 → AMROA	$\textbf{0.974} \pm \textbf{0.003}$	0.965 ± 0.000	0.947 ± 0.004	0.879 ± 0.012	0.815 ± 0.016	0.896 ± 0.031	0.000 ± 0.000	0.000 ± 0.000	0.665 ± 0.114	0.000 ± 0.000
	OAI ZIB → AMROA	0.973 ± 0.004	0.964 ± 0.005	$\textbf{0.948} \pm \textbf{0.005}$	$\textbf{0.893} \pm \textbf{0.010}$	0.817 ± 0.043	$\boldsymbol{0.898 \pm 0.011}$	0.000 ± 0.000	0.000 ± 0.000	0.648 ± 0.104	0.000 ± 0.000
	AMROA → SKI10	0.950 ± 0.031	0.959 ± 0.002		0.758 ± 0.010	0.681 ± 0.009					
	AMROA → OAI ZIB	0.962 ± 0.006	0.951 ± 0.010		0.813 ± 0.032	0.790 ± 0.039					

1 Figures

Figure 1 – Conditional GAN structure. The generator is a U-Net that progressively downsamples / encodes and then up-samples / decodes an input by a series of convolutional layers, with additional skip-connections between each major layer. The generated, 'fake' segmentation image is then fed together with the ground truth segmentation image into a discriminator network (PatchGAN (Isola et al., 2017)) that gives its prediction of whether the generated image is a 'real' representation of the ground truth image, or not. A detailed description of the network architecture can be found in the Appendix.



Figure 2 – Results of Network Objective Function. Qualitative results of B) training a cGAN with different objective functions by combining the cGAN loss with different pixelwise error losses with varying weightings and C) training a U-Net with different pixel-wise error losses.



Figure 3 - Results of testing on noise only images. Assessing the segmentation performance of a cGAN trained with $\mathcal{L}_{cGAN} + \lambda \mathcal{L}_{L1}$ ($\lambda = 100$) loss objective and a U-Net trained with \mathcal{L}_{L1} objective and tested on noise only images. Training was performed on the AMROA training dataset without noise only images. A) and B) are two example results of testing the models on noise only source images and comparing to ground truth segmentation maps.



Figure 4 – Results of Altering the Loss Objective during Training. Assessing the influence of varying the objective function halfway during cGAN and U-Net training on their segmentation performance with comparison to the respective cGANs and U-Nets trained with constant loss function.



- 1 Figure 5 Influence of altering the loss objective during cGAN training on the segmentation
- 2 performance of the medial gastrocnemius and vastus muscles.
- 3 The cGAN was trained with a $\mathcal{L}_{cGAN} + \lambda \mathcal{L}_{L2}$ loss objective for 50 epochs followed by a
- 4 further 50 epochs training with $\mathcal{L}_{cGAN} + \lambda \mathcal{L}_{L1}$.
- 5 Abbreviations: VMM vastus medialis muscle, GMM medial head of gastrocnemius
- 6 muscle, DSC Dice Similarity Coefficient



1 Figure 6 – Results of PatchGAN Receptive Field Size. Assessing the influence of varying

2 the discriminator receptive field size on segmentation performance of cGAN when trained

3 and tested on the AMROA dataset.



6

5

4

7 Figure 7 – Image Artefact due to the choice of PatchGAN Receptive Field Size. Influence

8 of discriminator receptive field size on checkerboard artefact emergence of a cGAN trained

9 and tested on the AMROA dataset.





PatchGAN Receptive Field Size Variation



- 1 Figure 8 Loss Evolution during cGAN Training. The loss evolutions of the A) generator
- $(\mathcal{L}_{cGAN} \text{ and } \mathcal{L}_{L1}) \text{ and } \mathbf{B})$ discriminator $(\mathcal{L}_{real} \text{ and } \mathcal{L}_{fake})$ are shown for a cGAN trained with a
- 3 U-Net generator and a 1x1 PatchGAN discriminator for 100 epochs.



- 1 Figure 9 Results of Transfer Learning: SKI10 and OAI ZIB. Assessing the influence of
- 2 transfer learning on segmentation performance of cGAN and U-Net when tested on the
- 3 SKI10 and OAI ZIB test datasets.
- 4 SKI10 / OAI ZIB \rightarrow AMROA: Pretraining the network for 20 epochs on the SKI10 / OAI
- 5 ZIB training dataset followed by network fine-tuning for 80 epochs on the AMROA training
- 6 dataset.
- 7 AMROA \rightarrow SKI10 / OAI ZIB: Pretraining the network for 20 epochs on the AMROA
- 8 training dataset followed by network fine-tuning for 80 epochs on the SKI10 / OAI ZIB
- 9 training dataset.

SKI10 Testing



- 1 Figure 10 Results of Transfer Learning: AMROA. Assessing the influence of transfer
- 2 learning on segmentation performance of cGAN and U-Net when tested on the AMROA test
- 3 datasets.
- 4 SKI10 / OAI ZIB \rightarrow AMROA: Pretraining the network for 20 epochs on the SKI10 / OAI
- 5 ZIB training dataset followed by network fine-tuning for 80 epochs on the AMROA training
- 6 dataset.
- 7 AMROA \rightarrow SKI10 / OAI ZIB: Pretraining the network for 20 epochs on the AMROA
- 8 training dataset followed by network fine-tuning for 80 epochs on the SKI10 / OAI ZIB
- 9 training dataset.

AMROA Testing



Transfer Learning with SKI10



Transfer Learning with OAI ZIB

