Fast Segmentation of 4D Microtomography Volumes from Core-flooding Experiments in Porous Rock using Convolutional Neural Network

YILI YANG¹, Sohan Seth¹, Ian B. Butler¹, and Florian Fusseis¹

¹University of Edinburgh

November 21, 2022

Abstract

Multiphase fluid flow in porous media has been extensively studied for its applications in carbon capture and storage, hydrocarbon recovery, aquifer contamination, soil hydrology and subsurface energy resources. Fluid displacement in porous media can be investigated using highly time-resolved synchrotron X-ray microtomography. One consequence of extremely fast imaging can be compromised image quality, including noise and decreased contrast, which makes the images hard to be segmented. We trained an established convolutional neural network (CNN) architecture (U-Net) with 18,072 images from multiphase flow experiments generated by synchrotron μ CT. The trained neural network can segment synchrotron μ CT images of core-flooding experiments rapidly and accurately without any pre-processing of the raw image. Segmenting one μ CT scan volume of size 1004x496x496 takes 5.6 minutes on an Nvidia Quadro K5200 GPU, while a conventional segmentation pipeline using CPU for the same size data takes 50.2 minutes. On the test dataset, the AUC-ROC score of individual class reached above 0.99 and the mean accuracy of the three segmentation classes reached 99%. The average IoU of the three classes is 0.98. The accuracy of the CNN segmentation is of the same order as conventional methods but it is significantly faster.

Fast Segmentation of 4D Microtomography Volumes from Core-flooding Experiments in Porous Rock using Convolutional Neural Network

Y. Yang, S. Seth, I. B. Butler, F. Fusseis

University of Edinburgh

Key Points:

1

2

3

4

5

6

7	- Machine learning segmentation enables fast segmentation of μCT data and can
8	reduce processing time by ten-fold.
9	- Machine learning segmentation quality is insensitive to μCT noise and ring arte-
10	facts.
11	• Machine learning segmentation is applicable to the fast processing of large datasets
12	such as those from time resolved synchrotron microtomography.

Corresponding author: Yili Yang, yili.yang@ed.ac.uk

13 Abstract

Multiphase fluid flow in porous media has been extensively studied for its applications 14 in carbon capture and storage, hydrocarbon recovery, aquifer contamination, soil hydrol-15 ogy and subsurface energy resources. Fluid displacement in porous media can be inves-16 tigated using highly time-resolved synchrotron X-ray microtomography. One consequence 17 of extremely fast imaging can be compromised image quality, including noise and decreased 18 contrast, which makes the images hard to be segmented. We trained an established con-19 volutional neural network (CNN) architecture (U-Net) with 18,072 images from multi-20 phase flow experiments generated by synchrotron μ CT. The trained neural network can 21 segment synchrotron µCT images of core-flooding experiments rapidly and accurately 22 without any pre-processing of the raw image. Segmenting one µCT scan volume of size 23 $1004 \times 496 \times 496$ takes 5.6 minutes on an Nvidia Quadro K5200 GPU, while a conven-24 tional segmentation pipeline using CPU for the same size data takes 50.2 minutes. On 25 the test dataset, the AUC-ROC score of individual class reached above 0.99 and the mean 26 accuracy of the three segmentation classes reached 99%. The average IoU of the three 27 classes is 0.98. The accuracy of the CNN segmentation is of the same order as conven-28 tional methods but it is significantly faster. 20

30 1 Introduction

Experimental studies of fluid flow related processes in porous media, including mul-31 tiphase flow and reactive flow, have increasingly made use of X-ray microtomography 32 (µCT) techniques to image fluid distributions and changes in porous media in-situ. These 33 studies range from those which employ laboratory µCT instruments (e.g. Pak et al., 2015; 34 AlRatrout et al., 2018) to those which employ synchrotron μ CT to generate 4D time re-35 solved data that reveals dynamic fluid flow processes (e.g. Berg et al., 2014; Reynolds 36 et al., 2017; Berg et al., 2013). The high photon flux at synchrotron X-ray sources en-37 ables experimenters to acquire multiple µCT scan volumes, each scan lasting a few sec-38 onds or less, over hours of experimentation. With 4D datasets often approaching a few 39 terabytes and comprising many 10s to 100s of discrete data volumes, the efficiency of im-40 age processing and segmentation can present a bottleneck to downstream data analy-41 sis. Furthermore, fast scans with brief exposure times and few projections, combined with 42 progressive damage to scintillators resulting from high X-ray fluxes can lead to increased 43

-2-

⁴⁴ noise, image artifacts and low contrast, making image segmentation both difficult and⁴⁵ time consuming.

 μ CT image segmentation is part of a μ CT image processing workflow. By segmen-46 tation, reconstructed µCT images are classified into labelled images such that objects 47 of interest in the images are separated for further analysis. Conventionally, the image 48 quality of the reconstructed µCT images is evaluated at first. Then a selection of denois-49 ing filters are tested based on the type of noise and artefacts present in the reconstructed 50 images. The denoised images are evaluated again for the difficulty of segmentation, and 51 then tested with different segmentation algorithms such as global thresholding, water-52 shed (Neubert & Protzel, 2014), random walker (Grady, 2006) or adaptive threshold-53 ing (Sauvola & Pietikäinen, 2000). Due to the, often, extremely large volume of data, 54 the optimal processing workflow is decided based on both effectiveness and efficiency. 55

Although the conventional segmentation workflow can produce accurate and re-56 liable results, there are two factors that make it suboptimal for large datasets of 10s to 57 100s of μ CT scan volumes. First, it is slow, and in the case of the data used in this study 58 taking approximately 50 minutes per µCT scan volume on average (for computing spec-59 ifications see section 2). Second, the conventional segmentation workflow involves sig-60 nificant human supervision and decision making, and due to variations in the degree of 61 noise and artefacts between µCT scan volumes, the parameters of each processing step 62 often need to be hand-tuned for each µCT scan volume to achieve reliable segmentation 63 results. 64

Recently, deep learning segmentation methods have been increasingly used on med-65 ical CT data (e.g. Skourt et al., 2018; Weston et al., 2019). Although X-ray µCT has 66 been a conventional imaging technique in multiphase fluid flow study, deep learning meth-67 ods are mainly used for predicting fluid flow and porous media properties (e.g. Mo et 68 al., 2019; Alqahtani et al., 2020; Kanin et al., 2019). For example, Karimpouli and Tah-69 masebi (2019) segmented different mineral phases of Berea sandstone μ CT images us-70 ing a deep learning method. However, to our knowledge, there is no published work on 71 segmenting µCT images of multiphase fluid flow in porous rocks using a deep learning 72 method yet. 73

Unlike conventional segmentation methods that only take a single, or a simple com bination of features into account, deep neural networks allow abstraction of multiple lev-

-3-

els of features, such as curves, lines and even complex patterns of data, and use such fea-76 tures to perform segmentation (Rawat & Wang, 2017). CNNs (LeCun et al., 2015) are 77 a class of deep neural network that have been proved effective in visual recognition tasks 78 (e.g. Krizhevsky et al., 2012; Long et al., 2015; Girshick et al., 2014). Ronneberger et 79 al. (2015) proved the reliability of CNN in segmenting low-contrast gray scale biomed-80 ical images, and introduced the U-Net architecture. Since experimental multiphase flow 81 data may also show the characteristics of gray scale images with low-contrast objects, 82 we have chosen to use the this architecture to train a segmentation model to yield ac-83 curate and precise segmentation of multiple tomographic volumes. 84

CNN is a supervised learning method that requires access to an adequate amount 85 of manually annotated training samples, i.e., it requires training images where a human 86 has established the boundaries between different classes. Although time consuming, man-87 ual annotation is sensible when objects are easily and readily identifiable. However, un-88 like conventional images with relatively smooth and simple boundaries, pore scale ex-89 perimental images have highly irregular, multi-scale, and even fractal boundaries which 90 make manual annotation difficult and imprecise (Figure 1). For the experimental μ CT 91 data in our study, manual annotation has three main problems: 1) the highly irregular 92 and multi-scale boundaries, 2) similar grey scale for different objects, i.e., brine and rock, 93 and 3) noise and artefacts. Therefore, we acquired the ground truth segmentation of the 94 input images using a weakly supervised algorithm, with a few pre-processing steps such 95 as denoising and masking (details of this approach are in section 2.2). 96

⁹⁷ Our goal is to reproduce, or emulate, the segmentation generated by the weakly ⁹⁸ supervised approach, using machine learning to achieve an accurate, fast and scaleable ⁹⁹ result. By training the CNN segmentation model, it learns the 'relation' between the raw ¹⁰⁰ μ CT image and the optimal segmentation regardless of the processing steps leading to ¹⁰¹ the segmentation. Given an input of raw μ CT image, the trained model can produce the ¹⁰² desired segmentation result that has the same 'relation' to the input image as it learned ¹⁰³ from the training.

In this contribution we trained an U-Net architecture to efficiently segment synchrotron μCT images obtained during core-flooding experiments. The images comprise three classes of objects: rock, oil and water. The trained CNN network is able to segment the three classes directly from reconstructed images accurately and rapidly with-

-4-



1646 pixels

Figure 1. Cross section of a cylinder core of Indiana limestone imaged by synchrotron μ CT. The core is fully saturated with H₂O, the lighter grey colour is oolitic limestone grains, the darker grey colour is H₂O in the pore space. The green box shows the cropping area which is the region of interest.

¹⁰⁸ out pre-processing steps such as applying of noise filters. It can reduce processing times

- ¹⁰⁹ 10-fold compared to conventional segmentation methods. We have tested the accuracy
- of the trained CNN model on a similar dataset with brief retraining, and we also tested
- the robustness of the model against artificial ring artefacts. We show the model perfor-
- mance to be robust and accurate. We provide both the trained network and a template
- training code (Python) so users can fine-tune the model to their own data (https://github.com/yclipse/CNN-

¹¹⁴ core-flooding-muCT).

- ¹¹⁵ 2 Materials and Methods
- 116

2.1 Multiphase Fluid Flow Data Acquisition

We used the beam line 2-BM at the Advanced Photon Source in Argonne National 117 Laboratory in Chicago, U.S. A carbonate rock, Indiana limestone, was used as porous 118 medium in the experiments. The rock was cored $(21.2mm \text{ length} \times 3mm \text{ diameter})$ and 119 installed in an X-ray transparent cell (Fusseis et al., 2014). For fluid injection we used 120 n-dodecane and potassium iodide solution (2.4M KI as a contrast agent) as oil and aque-121 ous phase respectively. The fluids are immiscible, and the processes involve no chemi-122 cal reaction between the fluids and the rock. Synchrotron X-ray images were acquired 123 using pink beam by one second acquisition time in every 20 seconds during fluid injec-124 tions (600 projections per scan, $\frac{1}{600}s$ exposure time per projection, referred to as fast scan). 125 The μ CT image resolution is 2.2 micron. A total of 140 time stamps of μ CT scan vol-126 umes each sized $1004 \times 1646 \times 1646$ were acquired during the experimentation (Fig-127 ure 2). 128

The 140 µCT scans recorded three individual experimental processes: brine injec-129 tion (0-40), oil injection (40-80) and simultaneous injection of both fluids (80-140) (Fig-130 ure 2). Although the physical processes during the experiments are complex and entirely 131 different, the differences shown on the μ CT images are simple and repetitive: pores switch-132 ing between oil-filled, partially-filled and brine-filled states. Therefore, in terms of im-133 age segmentation, even a small segment of the full data is highly representative. In this 134 study, we used the first half of the brine injection process to train, validate and test the 135 CNN model, and we used the second half of the oil injection process for extended test-136 ing. These two sets of data cover all of the different attributes of the µCT images. 137

-6-



Synchrotron µCT datasets of 140 time-stamps

Figure 2. Experimental design. A total of 140 time stamps of µCT scan volumes were acquired during the experiment. The first 40 time stamps of scans recorded the process of brine displacing oil in the Indiana limestone core. The 40th-80th scans recorded the process of oil displacing brine. The last 60 scans recorded the fluid displacement process by injecting oil and brine simultaneously. The first 18 scans were used to train, validate and test the initial model, while 20 more scans from time stamp 60-80 were used to extend the initial model and test it further. An independent reference scan was acquired before the experimentation, which is used to generate a high quality binary mask of the rock.

Before the core-flooding experiment, we imaged the pore structure of the Indiana 138 limestone core with a high-resolution reference scan (Figure 2 Reference Scan). The ref-139 erence scan had $2.5 \times$ more exposure time than the fast scans and each scan involved 1500 140 projections, therefore the reference scan μCT images are of significantly higher image 141 quality than the images acquired from the fast scans. On the reference scan images (Fig-142 ure 1) there are fewer artefacts and less noise, the object edges are sharper, and the pore 143 space is fully saturated with a single fluid that contrasts well with the rock. We segmented 144 the high quality reference image into a binary image labelled with rock and pore space 145 using the seeded random walker algorithm (Grady, 2006). The pore space binary image 146 is used as a mask to further separate the two fluid phases (for details of the ground truth 147 segmentation see section 2.2). All processing and computation was done on a HP Z820 148 workstation. It has 192Gb usable memory and 2 Intel(R) Xeon(R) E5-2640 V3 proces-149 sors (32 cores). The GPU is an Nvidia[™] Quadro K5200 (8Gb, 8 cores). 150



Figure 3. The image segmentation work flow for generating ground truth. (a) raw μ CT image containing three object classes: oil, brine and rock. The fluid classes are hardly distinguishable because the rock has similar intensity characteristics as the brine. Green box shows the cropped area that was used as network input. (b) the raw μ CT image with rock masked as black using the rock segmentation acquired from the dry reference scan which makes further segmentation of the two fluid classes possible. (c) masked image filtered by the non-local-means algorithm to reduce noise. (d) ground truth segmentation produced by seeded random-walker algorithm.

Figure 4. Gray value profile along a line in (yellow) raw reconstructed image. The three different classes, i.e., rock, brine and oil are all included in the profile line. These classes are not separable by their gray values, especially between brine and rock. The range of the grey values of the three classes are highlighted in blue, purple and yellow for rock, oil and brine respectively. The kernel density estimation also shows no sign of distinguishable difference between classes.

151

2.2 Ground Truth Segmentation

The μ CT projections (raw data acquired from fast μ CT scans) were reconstructed 152 using filtered back projection (Octopus8.6[™](Dierick et al., 2004)). The reconstructed im-153 ages are referred to as the 'raw µCT images' that are a sequence of cross-sectional im-154 age slices of the scanned object normal to the rotation axis. The raw μ CT images are 155 the input images for the CNN segmentation. Shown in Figure 3(a), the raw μ CT image 156 has three target classes (inside the cylindrical rock core) marked in red boxes as exam-157 ples: the appearance of oil is dark grey shade, the appearance of brine is light grey shade, 158 and the appearance of rock is also light grey but lighter grey than brine. Figure 4 shows 159 the grey values plotted along a profile line containing all three phases of an example raw 160 reconstructed image. The figure shows that the classes are similar in terms of their gray 161 value distribution, and therefore, they cannot be easily separated by their grey values 162 alone. We acquired the ground truth images by processing the images as follows. 163

We registered and applied the binary pore space mask using Avizo9[™] with the fast scan reconstructed µCT images. Figure 3(b) shows the raw CT images after applying the mask. The rock-brine-oil system now is more distinguishable, but the brine and oil are still very noisy. We used a non-local means filter (Buades et al., 2005) to denoise the two classes in the masked images (Figure 3(c)). We used the ImageJ implementation of the filter and used a sigma value of 30 and the default window size. This reduces the noise in brine and oil, but at the cost of making the class boundary less pronounced. We implemented the seeded random-walker algorithm (Grady, 2006) to separate brine and oil (Figure 3(d)). Seeded random walker is a region growing method which is capable of inferring ambiguous boundaries. Implementation in Scikit-image was used and took 50.2 minutes to run, on average, on the same size data volume as above.

We chose a total of 18 consecutive μ CT scan volumes from the start to the end of brine invasion, these volumes recorded the on-going process of brine displacing oil. These volumes of size $1004 \times 1646 \times 1646$ were segmented in this weakly supervised fashion to generate target annotation for the CNN training dataset.

For the ground truth segmentation, an internal validation of the conservation of 179 pore space (i.e. pore space is fixed over all time stamps) was used to estimate the error 180 in segmentation of the fluid volumes. The experimental procedures ensure that the pore 181 space is filled with two fluids, i.e., either oil or brine. The proportion of the fluids can 182 change through the experiment but the volume of oil and brine together should always 183 sum up to the pore space volume. we measured the total fluid volume of each scan. We 184 also measured the total pore space volume from the reference scan. We assume the av-185 erage total fluid volume should be equal to the total pore space volume. We found that 186 there is 1.86% of random error, determined from segmenting different fluids in replicate 187 volumes, this is taken as the random error of the ground truth annotation. This amount 188 of error is unlikely to affect any qualitative measurements in 3D such as connectivity. 189 The amount of error also has negligible impact on quantitative measurements such as 190 bulk volume. For such pore-scale images with very irregular and blurry class boundaries, 191 manual segmentation of eighteen thousand images will be extremely time-consuming, and 192 may not guarantee a reduction in error. In addition to determining the error based on 193 volumetric measurements, we visually compared the results of the ground-truth segmen-194 tation with greyscale image data and found the differences in the fine details of pore struc-195 ture to be insignificant. 196

This segmentation workflow is not scalable partly because of the high computational requirements, but also because segmentation parameters cannot be applied across different batches of scans. We acquired a total of 140 µCT scan volumes during the ex-

-10-

periment, and used only the first 18 μCT scan volumes (containing 18,072 images) to
 train the CNN model. This means that 90% of the data were processed using the more
 efficient CNN model.

203

2.3 Convolutional Neural Network Segmentation

204

2.3.1 Training Methods

Before training, we prepared the training dataset. Every raw reconstructed image 205 was paired with the corresponding ground truth segmentation image to make a set of 206 input-target image pairs to train the network. A total of 18 µCT scan volumes were ran-207 domly divided into 14 training sets, 2 test sets and 2 validation sets (Figure 2 Training, 208 Val and Test). Each µCT scan volume has 1004 reconstructed slices. The raw reconstructed 209 images $(1646 \times 1646 \text{ pixels})$ were down-sampled $(2 \times 2 \text{ mean pooling})$ to size 823×823 210 pixels. A rectangular area of size 496×496 pixels was cropped as the training area (shown 211 as the green box in Figure 3). The down-sampling of the original images was: 1) to fit 212 the images into the GPU memory. 2) To exclude pores adjacent to the perimeter which 213 do not contain useful information because the perimeter has different wettability than 214 the rock (for a solid material, wettability indicates the tendency of one particular fluid 215 phase to spread over the solid surface in presence of another fluid). The choice of the par-216 ticular size 496×496 was based on the architecture of the CNN network, which has four 217 max-pooling layers which each will divide the input size by two, so the input size needs 218 to be four times divisible by 2. 219

We used the CNN architecture U-Net introduced by Ronneberger et al. (2015) and 220 implemented by van Vugt (2017) using the open-source deep learning platform Pytorch 221 (Paszke et al., 2017). The network comprises a contracting path (Figure 5 left box) and 222 an expanding path (Figure 5 right box). The contracting path, like conventional CNN, 223 uses 3×3 convolution, followed by a rectified linear unit (ReLU), and then 2×2 max-224 pooling for down-sampling. The expanding path up-samples the feature map with $2\times$ 225 2 up-convolution followed by 3×3 transposed convolution with ReLU. A concatenation 226 of the feature map of the corresponding contracting path is applied to the up-convolved 227 feature map (Figure 5 grey arrow). In total the network has 23 convolutional layers. 228

Figure 6 shows the overall workflow of training, validating and applying the CNN model. In training, the model takes a training image as its input and produces a pre-

-11-

Figure 5. Modified from Ronneberger et al. (2015), this figure shows the architecture of the U-Net. It consists of a contracting path that extracts different levels of features, and an expanding path that up-convolves the image to produce the final segmentation. The two paths form a U-shaped network. The original paper uses VALID padding (i.e. no padding), so the height and width of each feature map decreases after each convolution. In this implementation SAME padding (i.e. zero padding by 1 on each side) is used so the height and width of the feature map will stay the same (e.g. Silburt et al., 2019).

diction of the probability of each pixel belonging to a class (i.e. rock, oil and brine). Then 231

the prediction is compared with the ground truth image to yield a training loss value which 232

quantifies the difference between the prediction and ground truth. The model iteratively 233

- updates the internal adjustable parameters (i.e. weights) as training proceeds, to min-234
- imize the loss and, thus, to improve the prediction in a step-wise manner. We trained 235
- the model for 34 epochs, where a training epoch is a full traverse of all training data. 236

Figure 6. CNN training, testing and application work flow. The training process uses training datasets and validation datasets to adjust the best fit of a CNN model representing the relation between input images and ground-truth images. The testing process verify the model with a held-out dataset. The application process uses the trained CNN model to segment new μ CT data.

237 238

239

241

242

Cross-entropy loss was used as the loss function during training. The loss function calculates the distance between the two probability distributions, i.e., the output prediction and the corresponding ground-truth. Cross-entropy decreases when the output and ground-truth have a higher resemblance. We used the Pytorch built-in cross-entropy 240 loss function which is written as:

$$\operatorname{Loss} = -\sum_{i \in I} \sum_{c \in C} y_{i,c} \log \hat{y}_{i,c} \tag{1}$$

where $y = (y_1, y_2, \dots, y_n)$ and $\hat{y} = (\hat{y}_1, \hat{y}_2, \dots, \hat{y}_n)$ are the ground truth values and the 243 network prediction values of all pixels of a training image flattened to 1D arrays. C is 244

the set of all classes and I is the set of all pixels. $y_{i,c}$ and $\hat{y}_{i,c}$ are values of the *i*th pixel in class c.

Hyper-parameters are user-defined parameters that control the overall learning be-247 haviours of the network during training. Learning rate is one dimensionless hyper-parameter 248 that defines the size of steps to adjust the network weights to minimize loss; a higher learn-249 ing rate implies a larger adjustment in each iteration. Multiple learning rates were tested 250 and we found the value of 5×10^{-5} to be optimal for loss to converge quickly and smoothly. 251 During training, we divided the training datasets into mini batches each containing two 252 images. This can improve the efficiency of training. An ℓ_2 regularization parameter of 253 1×10^{-5} was applied to alleviate over-fitting (ℓ_2 regularization is also know as Ridge 254 regression, it adds squared magnitude of coefficient as penalty term to the loss function 255 to avoid over-fitting). The sequence of the input data was shuffled at the start of each 256 epoch. This can eliminate bias that is produced by the high similarity of two adjacent 257 μ CT slices. The optimizer used to train was the Adam algorithm (Kingma & Ba, 2014) 258 which is built into Pytorch. 259

260

2.3.2 Validation

At the end of each training epoch, a validation loss value is computed on the val-261 idation set (Figure 6 Training). The validation loss is not back-propagated and the net-262 work weights are not updated. Thus, the model occasionally 'sees' this data but never 263 learns from it, and thus the process of validation provides an unbiased, on-going eval-264 uation of a model fit to the training dataset. The validation loss value is plotted through-265 out training and compared with the training loss value. The model is best-trained at the 266 epoch when the validation loss starts to increase, or the validation accuracy ceases to 267 decrease. After this the model starts to overfit. 268

For the pixels in the network prediction of an input image, a pixel is assigned to a class if the corresponding probability is the highest among the three classes. There is a chance for a rare exception where a pixel has exactly equal probability for all three classes (33% oil, 33% brine, 33% rock), or equal probability for any two classes (e.g. 50% brine and 50% rock). These pixels are assigned to the rock class for the reason that such pixels are often on the contact surface between the fluids and the rock.

-14-

For each class of a prediction, we calculated the number of true positive (TP), false 275 positive (FP), true negative (TN) and false negative (FN) pixels with respect to the ground 276 truth. For example, for the oil class in a CNN segmented image, true positive pixels are 277 the pixels which are assigned to oil, and the corresponding ground truth for the same 278 pixels are also oil. False positive pixels are the pixels which are assigned to oil, but the 279 corresponding ground truth for the same pixels are not oil. True negative pixels are the 280 pixels which not assigned to oil, and which agree with the ground truth. False negative 281 pixels are the pixels which are not assigned to oil, but the ground truth for these pix-282 els is oil. 283

We used three metrics to measure the segmentation performance, they are the Rand index (accuracy, (Rand, 1971)), IoU (intersection over union) and AUC-ROC (area-undercurve of the receiver operating characteristic (Bradley, 1997)). The Rand index is the ratio of correctly classified pixels to total pixels. The Rand index was calculated using Eq. (2)

$$accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$
(2)

The intersection over union measures how well the prediction is spatially overlapped with the ground truth. It is the ratio between two values: the area of overlap and the area of union (Eq. (3)).

289

293

$$IoU = \frac{Area \text{ of intersection}}{Area \text{ of union}} = \frac{TP}{TP + FN + FP}$$
(3)

The AUC-ROC is a common metric for assessing the performance of a classifier. 294 ROC is a curve plotted in the coordinate system with true positive rate $(TPR = \frac{TP}{TP+FN})$ 295 against false positive rate (FPR = $\frac{FP}{FP+TN}$), over various thresholds. An ROC curve 296 essentially represents the trade-off between sensitivity and specificity of a classification 297 model. AUC is the area under the ROC curve which measures the model's separability 298 of different classes. The AUC-ROC value varies between 0.5 and 1. When AUC equals 299 1, it represents a perfect classification model that classifies all samples accurately. When 300 the AUC value is around 0.5, it represents the case where all classifications are random 301 due to a completely overlapped predictive probability distribution of true and false val-302 ues. 303

After training, the model is tested on 'unseen' data to check its generalizability and assess the segmentation quality. The separate test dataset was only used once the model had been trained (Figure 6 Testing). The test dataset provides a standard to evaluate the model.

308 **3 Results**

309

3.1 Training Result

Figure 7 shows the training process by plotting the validation accuracy and loss 310 over training epochs. The top graph shows that as training progresses the validation ac-311 curacy increases before finally stabilizing for all three classes. The bottom graph shows 312 that the training loss drops sharply at the outset of training but became stable as the 313 training progresses while the validation error starts at a relatively low value and decreases 314 slowly. Theoretically it is assumed that a validation loss curve is U-shaped since it first 315 drops as the algorithm learns to generalize, then stabilizes, and finally starts to rise as 316 the model tends to overfit. Therefore the criterion to stop training is set as a continual 317 increase of validation loss over 5 consecutive epochs. The network was trained for 34 epochs 318 and the best fit was identified at the 29th epoch. The total training time was 205.4 hours. 319

The segmentation by the model during training and after training are shown in Figure 8. The first row shows the raw reconstructed image and the ground truth. The second row shows the segmentation of the same image during the 11th epoch of training and after 29 epochs of full training. This result illustrates an improvement of the model during training as at the 11th epoch the segmentation result is sub-standard since it only partially captures the pore space structure and the boundary decision is not well defined. However, at the 29th epoch the segmentation result is very similar to the ground truth.

327

3.2 Testing Result

The model trained for 29 epochs was identified as the best segmentation model and tested on the full testing dataset. The performance measured on the test set shows high accuracy for all three classes. The test result shows that the accuracy of rock, brine, oil and average accuracy over all three classes are 98.5%, 99.3%, 99.1% and 99.0% respectively. The AUC-ROC score for brine, rock and oil are 0.99, 0.99 and 0.99 respectively. It shows very high separability of the model. Table 1 shows that the CNN segmentation

-16-

Figure 7. Top: The plots shows the variation of accuracy and IoU score over training epochs for three classes separately. Bottom: the image shows the variation of training loss and validation loss over training epochs. The IoU curve completely overlaps with the curve of rock accuracy.

- result correctly labelled 93.6% oil phase as oil, 96.3% brine as brine, and 99.5% rock as rock. The total oil phase and brine phase were slightly underestimated with regard to the ground truth. We measured an IoU score of 0.98 showing that the CNN segmentation is spatially overlapping with the ground truth correctly.
- The probability map (Figure 9) visualises the probability distribution of all three 338 classes. It is generated by the softmax function that transforms the model output into 339 a probability distribution such that the probability of three classes at the same pixel add 340 up to 1. The colour map indicates that the likelihood of a pixel belonging to one par-341 ticular phase, with red implying extremely likely and blue implying extremely unlikely 342 while white implying ambiguous likelihood. The three classes are well separated since 343 the ambiguous likelihood pixels are of minor amounts and are mostly located at the bound-344 aries of two phases. It also illustrates a very robust performance of the model in segmen-345 tation. 346

(a) Raw reconstructed image

(b) Ground truth

(c) CNN prediction at the 11th training epoch

(d) CNN prediction at the end of training

Figure 8. CNN prediction during training and after training. (a) Raw reconstructed image before segmentation. (b) Ground truth image. (c) CNN prediction of 3-phase probability at the 11th training epoch. (d) CNN probability prediction after 29 epochs of training. The segmentation was on the same image from the testing dataset. This illustrates the improvement of segmentation quality during training.

(d) CNN predicted probability of brine

Figure 9. Probability distribution of each phase at the network output. The color bar shows probability of a pixel belonging to that phase. This illustrates that the majority of segmentation was certain (probability close to either 0 or 1) with very few uncertain pixels (probability equals 0.5). The uncertain pixels were classified into the rock class as they are always presented between the contact between the rock and two fluids.

 Table 1. Contingency table of each class segmented by the CNN model compared with the ground truth.

			Ground truth	
		Oil	Brine	Rock
	Oil	0.93	0.02	0.06
Estimation	Brine	0.04	0.96	0.03
	Rock	0.04	0.01	0.99

347 **4** Discussion

348

4.1 Segmentation Robustness

Ring artefacts are concentric rings caused by mis-calibration or failure of a detec-349 tor element. They are one of the major and most common type of artefacts that hinders 350 μ CT image segmentation. We added artificial ring artefacts to an image selected from 351 the test set (scan 17) to examine the robustness of the model to image quality degra-352 dation due to ring artefacts. The chosen test image is visually separable and does not 353 suffer from original µCT ring artefacts. Artificial ring artefacts were generated by adding 354 concentric circles with random radii and intensities at a fixed coordinate on the image. 355 We measured the severity of this artefact in terms of the number of rings present (Fig-356 ure 10). The direct impact of the artificial added ring artefacts on the image quality is 357 statistically shown in Figure 11: compared with the histogram curve of the original im-358 age (blue), the ring artefacts caused spikes on a histogram curve of the same distribu-359 tion (orange). 360

Figure 12 shows the variation of AUC-ROC score and mean accuracy over degree of severity. The CNN segmentation performed surprisingly well for artificial ring artefacts. Only the brine phase of the severely affected images shows some minor mis-classification. It is also surprising because the training images do not have any significant ring artefacts. This test provides an insight on the insensitivity of the CNN segmentation with ring artefacts, which is one of the most intrusive feature for most conventional segmentation methods.

(a) Image with 10 artificial ring artefacts

(c) Image with 15 artificial ring artefacts

(b) CNN segmentation of (a)

Ring=15

(d) CNN segmentation of (c)

Ring=25

(f) CNN segmentation of (e)

Figure 10. CNN segmentation model tested on artificial ring artefact. Left: Same μ CT image with increasing severity of artificial ring artefacts where the severity is measured in terms of the number of rings added. Right: Corresponding CNN segmentation results where major mis-classifications are marked by green circles.

Figure 11. Image histogram of the three datasets segmented by the CNN model. All three histograms show unimodal distribution with the maximum indicating the rock class. The ring artefact image histogram shows the same distribution with the training/validation/testing image. The extended test image histogram shows a shifted distribution which is due to different μ CT reconstruction method and parameters.

Figure 12. The variation of accuracy and AUC-ROC of CNN segmentation over varying degree of ring artefact. The segmentation quality is stable with a minor decrease of performance at heavily degraded images.

4.2 Extended Test on Similar Data With Extra Training

368

To exclude data leakage that may occur when using highly similar training and testing data, we also tested the trained model on 20 µCT scan volumes which belong to the same experiment as the training dataset but from later timestamps with the same rock and fluids (Figure 2 Extended Test, 60th-80th time stamps). This dataset is of the same size as the training dataset, and uses the same approach for determining the ground-truth. This dataset has more noise compared to the training data set, and has horizontal stripe artefacts which the training dataset does not have.

The difference of image quality is because this dataset was reconstructed using a 376 different reconstruction tool, TomoPy (Gürsoy et al., 2014) and, the reconstructed im-377 age appearance has a small, but systematic variation due to the different reconstruction 378 parameters. The systematic variation is shown in the green curve on Figure 11: images 379 from the extended dataset have a gray value distribution the maximum of which is shifted 380 to the left compared with other datasets. Therefore, the trained model might not work 381 accurately on these data right away (Figure 13(c)). To mitigate this, the model was briefly 382 trained for five epochs using only the first volume of the 20, and tested with the remain-383 ing 19 volumes. 384

Figure 13 shows the CNN segmentation result by directly using the trained model 385 did not capture most of the pore structures and did not separate the two fluids (Figure 13(c)). 386 After five epochs of additional training, the CNN segmentation converged towards the 387 ground truth image (Figure 13(d)). The extra test results are: accuracy of rock is 93.5%, 388 accuracy of oil is 95.2%, accuracy of brine is 99.1%, and the mean accuracy is 95.9%. The 389 average IoU of the three classes is 0.98. The performance remains good with only a mod-390 est decrease in accuracy. This confirms that the CNN model can quickly be adapted to 391 handle similar datasets with brief re-training. 392

We tested whether the prior training was essential, or whether 5 epochs alone were sufficient. The result (Figure 13(e)) shows that the result is inconsistent with the ground truth and confirms that the prior training is necessary and the CNN requires it in order to effectively adapt to new data, but needing relatively few epochs of additional training to achieve that adaptation.

To address the impact of CNN segmentation and conventional segmentation on experimental measurements of the fluids, we compared the fluid saturation measured from CNN segmentation results and conventional segmentation results. Fluid saturation measures how much of a fluid is present in the pore space of a rock. In this comparison, we expect the fluid saturation measured from both CNN segmentation and random-walker segmentation to give similar results.

Figure 14 shows the saturation of oil measured from both CNN segmentation and 404 conventional segmentation are very close. The plot shows that measurements taken from 405 the random-walker segmentation is overall higher than the measurements taken from the 406 CNN segmentation. A one-sided Wilcoxon test was performed with the two saturation 407 measurements to examine if the measurement taken from random-walker is statistically 408 larger. The Wilcoxon test returned a p-value of less than 0.002 and confirms that the 409 difference is statistically significant. The median of the Wilcoxon test values is 2, indi-410 cating, however, that the amount amount of difference is small. The CNN segmentation 411 model is able to produce reliable segmentation that is close to the result obtained us-412 ing conventional segmentation method. The CNN model has the ability to adapt to vari-413 ations in the dataset with minimal retraining. 414

-24-

(a) Raw reconstructed image

(b) Ground truth

(c) CNN prediction with pre-trained model

(d) CNN prediction after additional training

(e) CNN prediction trained from scratch

Figure 13. Segmentation result of the model before and after extra training. (a) Raw reconstructed image before segmentation. (b) Ground truth image. (c) CNN prediction of 3-class probability using the trained network without additional training. (d) CNN segmentation result after 5 epochs of additional training with one μ CP-scan volume. (e) CNN prediction trained from scratch. This illustrates that the CNN model can be adapted to similar dataset with brief re-training.

Figure 14. Comparison of oil saturation measured from random-walker and CNN segmentation. Oil saturation of consecutive 19 scans from timestamp 60 to 80 was measured from both CNN segmentation and conventional random-walker segmentation. The measurements of saturation from the segmentation using two methods are highly similar.

415

4.2.1 Impact of Different Image Quality Disturbance

Comparing the CNN model performance on the extended test and the ring arte-416 fact test, the CNN model can perform reliably with the impact of low to moderate de-417 gree of ring artefacts without additional training. But it needs a small amount of addi-418 tional training to perform well on the extended test. Figure 11 shows that the distur-419 bance in the ring artefact images are local spikes on the histogram, while the disturbance 420 in the extended test images is a systematic shift of the image histogram. The model has 421 better resistance to local disturbance produced by ring artefacts than to a systematic 422 change of the image histogram. The results indicate that 1) the CNN segmentation can 423 be applied to datasets that are affected by low to moderate degrees of ring artefacts with-424 out additional training, and 2) the CNN segmentation needs minor additional training 425 to perform well on datasets that have a systematic shift of the histogram. It also sug-426 gests that for heterogeneous datasets, the training data set should be sampled across the 427 entire range to decrease bias and reduce additional training. However it may require a 428 greater amount of training to generalize the model. 429

430 4.2.2 Temporal Efficiency

For a dataset of size $1004 \times 496 \times 496$, the CNN model took 5 minutes 40 seconds 431 to segment the dataset. The conventional segmentation pipeline using the random-walker 432 algorithm took 50 minutes 13 seconds to segment the same dataset. The comparison shows 433 a ten-fold speed advantage using the CNN method. This comparison was done on a CPU-434 dedicated workstation which has 32 CPU cores, without which the processing time of 435 the conventional segmentation can be up to 15 hours and 32 minutes if using a single 436 thread. The GPU used in this study was released in 2014. The GPU has a compute ca-437 pability of 3.5, while current GPUs have a compute capability of 8.6, meaning the micro-438 architecture is five-generations behind and lack of many powerful computational features. 439 Although it is difficult to give a precise number of improvement, we estimate a further 440 improvement in temporal efficiency at least ten times of magnitude if the latest GPU 441 is used. 442

443

4.3 Future Improvements

The current training dataset can be augmented by introducing noise, mirroring, rotating, scaling, cropping or translating the original training images. Data augmentation allows amplification of the training dataset based on existing dataset and therefore can further generalize the CNN model to achieve a more robust segmentation of different datasets. This data augmentation strategy was tested with the U-Net architecture and was found to produced excellent segmentation (Ronneberger et al., 2015).

As a working CNN model of segmenting µCT data of core-flooding experiments, 450 the current model learned the segmentation process, rather than applying an intrinsic 451 knowledge of 'what is rock' and 'what are fluids'. The reason is that the training data 452 itself is highly biased towards the particular rock and fluid type of this experiment. Dif-453 ferent rocks have different internal structures. The appearance of a rock on μCT is highly 454 dependent on the imaging settings, and even the same rock can look different under dif-455 ferent conditions such as X-ray energy and sources. Different fluids can vary in μCT im-456 ages too. It is beneficial to keep training on future experiments to generalise the model 457 with more kinds of rocks and fluids. As more type of imaging conditions, noises, litholo-458 gies and fluid phases are collected in the training data, the model will be increasingly 459 generalized. The future training needed will eventually decrease as the CNN network has 460

'seen' more rocks, fluids and experiments. It becomes a more mature model that can segment multiphase fluid flow experiment data regardless of noise, rock type, fluid type and
beam line condition.

Comparison of different network architectures can be further tested. Apart from
the U-Net, there are other powerful CNN architectures such as ResNet (He et al., 2016),
GoogLeNet (Szegedy et al., 2015), VGGNet (Simonyan & Zisserman, 2014) etc. with different advantages and specialities. The approach introduced in this paper can be implemented on different CNN architectures for a variety of segmentation demands.

469 5 Conclusion

Use of a convolutional neural network can significantly improve the segmentation 470 workflow for multiphase flow µCT images. The advantages are five-fold. First, once a 471 network for segmenting multiphase flow images is trained, it can be applied to future data 472 without retraining or only with fine-tuning. Second, the segmentation is directly per-473 formed on the reconstructed image, and so considerable time spent on the pre-processing 474 (tuning of filtering, registration, masking etc.) can be avoided. This is significant because 475 faster segmentation means that more time can be devoted to analysis and interpretation 476 of data. Third, the CNN method can segment a dataset of size $1004 \times 496 \times 496$ in 5 477 minutes, which is ten-fold faster than a conventional segmentation pipeline using the random-478 walker method. The speed advantage is obvious even on our CPU-dedicated worksta-479 tion, and the speed of CNN may be significantly increased by using a more recent GPU. 480 Fourth, the algorithm is capable of segmenting images that are highly affected by ring 481 artefacts, which often require additional correction or removal steps for conventional pro-482 cessing paths (e.g. Rashid et al., 2012; Brun et al., 2009). Finally, the performance of 483 the CNN network improves as more data is available through re-training. With little re-484 training it can be easily adapted to new datasets when the previous training is biased. 485 Overall the CNN segmentation is a powerful and efficient tool for μ CT image segmen-486 tation, especially for large datasets. 487

488 Acknowledgments

We acknowledge the support of Petrobras and Royal Dutch Shell for financial support
of the International Centre for Carbonate Reservoirs (ICCR). This study comes from the
SatuTrack II sub-project of ICCR.

-28-

492 Data Availability Statement

⁴⁹³ The data that support the findings of this study are available from Petrobras and Shell.

- ⁴⁹⁴ Restrictions apply to the availability of these data, which were used under license for this
- 495 study. Data are available from the authors with the permission of Petrobras and Shell.

496 References

- Alqahtani, N., Alzubaidi, F., Armstrong, R. T., Swietojanski, P., & Mostaghimi, P.
 (2020). Machine learning for predicting properties of porous media from 2D
 X-ray images. Journal of Petroleum Science and Engineering, 184, 106514.
- AlRatrout, A., Blunt, M. J., & Bijeljic, B. (2018). Wettability in complex porous
 materials, the mixed-wet state, and its relationship to surface roughness. *Proceedings of the National Academy of Sciences*, 115(36), 8901–8906.
- Berg, S., Armstrong, R., Ott, H., Georgiadis, A., Klapp, S. A., Schwing, A., ... Leu,
 L. (2014). Multiphase flow in porous rock imaged under dynamic flow con ditions with fast X-ray computed microtomography. *Petrophysics*, 55(04),
 304–312.
- Berg, S., Ott, H., Klapp, S. A., Schwing, A., Neiteler, R., Brussee, N., ... Schwarz,
 J.-O. (2013). Real-time 3D imaging of haines jumps in porous media flow. *Proceedings of the National Academy of Sciences*, 110(10), 3755–3759.
- Bradley, A. P. (1997). The use of the area under the ROC curve in the evaluation of machine learning algorithms. *Pattern recognition*, 30(7), 1145–1159.
- Brun, F., Kourousias, G., Dreossi, D., & Mancini, L. (2009). An improved method
 for ring artifacts removing in reconstructed tomographic images. , 926–929.
- Buades, A., Coll, B., & Morel, J.-M. (2005). A non-local algorithm for image denoising., 2, 60–65.
- Dierick, M., Masschaele, B., & Van Hoorebeke, L. (2004). Octopus, a fast and user friendly tomographic reconstruction package developed in labview(R). Measure ment Science and Technology, 15(7), 1366.
- ⁵¹⁹ Fusseis, F., Steeb, H., Xiao, X., Zhu, W.-l., Butler, I. B., Elphick, S., & Mäder, U.
- (2014). A low-cost X-ray-transparent experimental cell for synchrotron-based
 X-ray microtomography studies under geological reservoir conditions. Journal
 of synchrotron radiation, 21(1), 251–253.
- Girshick, R., Donahue, J., Darrell, T., & Malik, J. (2014). Rich feature hierarchies

524	for accurate object detection and semantic segmentation. , $580{-}587.$
525	Grady, L. (2006). Random walks for image segmentation. <i>IEEE Transactions on</i>
526	Pattern Analysis & Machine Intelligence(11), 1768–1783.
527	Gürsoy, D., De Carlo, F., Xiao, X., & Jacobsen, C. (2014). TomoPy: a framework
528	for the analysis of synchrotron tomographic data. Journal of synchrotron radi-
529	ation, 21(5), 1188–1193.
530	He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image
531	recognition., 770–778.
532	Kanin, E., Osiptsov, A., Vainshtein, A., & Burnaev, E. (2019). A predictive model
533	for steady-state multiphase pipe flow: Machine learning on lab data. Journal
534	of Petroleum Science and Engineering, 180, 727–746.
535	Karimpouli, S., & Tahmasebi, P. (2019). Segmentation of digital rock images using
536	deep convolutional autoencoder networks. Computers & geosciences, 126, 142–
537	150.
538	Kingma, D. P., & Ba, J. (2014). Adam: A method for stochastic optimization.
539	Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). Imagenet classification with
540	deep convolutional neural networks. , 1097–1105.
541	LeCun, Y., Bengio, Y., & Hinton, G. (2015). Deep learning. nature, 521(7553), 436.
542	Long, J., Shelhamer, E., & Darrell, T. (2015). Fully convolutional networks for se-
543	mantic segmentation. , 3431–3440.
544	Mo, S., Zhu, Y., Zabaras, N., Shi, X., & Wu, J. (2019). Deep convolutional encoder-
545	decoder networks for uncertainty quantification of dynamic multiphase flow in
546	heterogeneous media. Water Resources Research, $55(1)$, 703–728.
547	Neubert, P., & Protzel, P. (2014). Compact watershed and preemptive slic: On im-
548	proving trade-offs of superpixel segmentation algorithms. , $996{-}1001.$
549	Pak, T., Butler, I. B., Geiger, S., van Dijke, M. I., & Sorbie, K. S. (2015). Droplet
550	fragmentation: 3D imaging of a previously unidentified pore-scale process dur-
551	ing multiphase flow in porous media. Proceedings of the National Academy of
552	$Sciences, \ 112(7), \ 1947-1952.$
553	Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., Lerer, A.
554	(2017). Automatic differentiation in PyTorch.
555	Rand, W. M. (1971). Objective criteria for the evaluation of clustering methods.
556	Journal of the American Statistical association, 66(336), 846–850.

- Rashid, S., Lee, S. Y., & Hasan, M. K. (2012). An improved method for the removal
 of ring artifacts in high resolution ct imaging. *EURASIP Journal on Advances in Signal Processing*, 2012(1), 93.
- Rawat, W., & Wang, Z. (2017). Deep convolutional neural networks for image classi fication: A comprehensive review. Neural computation, 29(9), 2352–2449.
- Reynolds, C. A., Menke, H., Andrew, M., Blunt, M. J., & Krevor, S. (2017). Dy namic fluid connectivity during steady-state multiphase flow in a sandstone.

Proceedings of the National Academy of Sciences, 114 (31), 8187–8192.

564

- Ronneberger, O., Fischer, P., & Brox, T. (2015). U-net: Convolutional networks for
 biomedical image segmentation., 234–241.
- Sauvola, J., & Pietikäinen, M. (2000). Adaptive document image binarization. Pat tern recognition, 33(2), 225–236.
- Silburt, A., Ali-Dib, M., Zhu, C., Jackson, A., Valencia, D., Kissin, Y., ... Menou,
- 570 K. (2019). Lunar crater identification via deep learning. *Icarus*, 317, 27–38.
- Simonyan, K., & Zisserman, A. (2014). Very deep convolutional networks for large scale image recognition. arXiv preprint arXiv:1409.1556.
- Skourt, B. A., El Hassani, A., & Majda, A. (2018). Lung CT image segmentation
 using deep neural networks. *Procedia Computer Science*, 127, 109–113.
- Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., ... Rabinovich,
 A. (2015). Going deeper with convolutions. , 1–9.
- van Vugt, J. (2017). *PyTorch U-Net.* https://github.com/jvanvugt/pytorch -unet. GitHub.
- ⁵⁷⁹ Weston, A. D., Korfiatis, P., Kline, T. L., Philbrick, K. A., Kostandy, P., Sakinis,
- T., ... Erickson, B. J. (2019). Automated abdominal segmentation of CT scans for body composition analysis using deep learning. *Radiology*, 290(3), 669–679.