

1 **Neural network-based CO₂ interpretation from 4D**
2 **Sleipner seismic images**

3 **Bei Li¹and Yunyue Elita Li¹**

4 ¹Department of Civil and Environmental Engineering, National University of Singapore, 117576 Singapore

5 **Key Points:**

- 6 • We train a 3D U-Net for an automatic end-to-end mapping from 4D seismic im-
7 ages to 3D CO₂ distribution.
- 8 • Successful applications of the trained neural network demonstrate its robustness
9 and consistency.
- 10 • We analyze the neural network interpretation standards and provide training strat-
11 egy for 2D sparse labels.

Corresponding author: Yunyue Elita Li, elita.li@nus.edu.sg

12 **Abstract**

13 Time-lapse or 4D seismic survey is a crucial monitoring tool for CO₂ geological sequestration.
 14 Conventional time-lapse interpretation provides detailed characterization of CO₂
 15 distribution in the storage unit. However, the process is labour-intensive and hard to keep
 16 interpretation consistency throughout the monitoring history, due to the inevitable in-
 17 terpreter subjectivity and the long term nature of CCS project. Moreover, reliable in-
 18 terpretation is based on adequate match between the baseline and time-lapse datasets,
 19 requiring much data processing effort. We propose a neural network (NN)-based inter-
 20 pretation directly from 4D seismic images to the corresponding 3D probability of CO₂
 21 distribution. We use a simplified 3D U-net, whose training, validation and test are all
 22 based on the publicly available datasets from the Sleipner CO₂ storage project. The lim-
 23 ited labels for training are derived from the interpreted CO₂ plume outlines within the
 24 internal sandstone layers for 2010. Then we apply the trained NN on different time-lapse
 25 seismic datasets from 1999 to 2010. The results suggest that our NN-based CO₂ inter-
 26 pretation has the following advantages: (1) high interpretation efficiency by providing
 27 an automatic end-to-end mapping; (2) robustness against the processing-caused mismatch
 28 between the baseline and time-lapse inputs, which can potentially relax the baseline re-
 29 processing demands for newly acquired or reprocessed time-lapse datasets; and (3) in-
 30 herent interpretation consistency throughout multi-vintage datasets. Qualitative anal-
 31 ysis for the NN interpretation standards shows that both amplitude difference and event
 32 similarity contribute to the determination of CO₂ distribution. We also compare 2D and
 33 3D U-Nets under scenario where only sparse 2D slice labels are available for training.
 34 The results suggest that 3D U-Net provide high-quality interpretation at the cost of large
 35 computational resources for training and application, in comparison with 2D U-Net.

36 **Plain Language Summary**

37 To reduce the greenhouse effect, captured and compressed CO₂ can be injected into
 38 subsurface area with special geological settings that permanently accommodate the green-
 39 house gas. Many pilot projects for such geological CO₂ sequestration have been ran for
 40 decades, during which various monitoring techniques have been experimented to study
 41 the interactions between injected CO₂ and the storage unit mainly for the purpose of
 42 evaluating storage safety. Such studies suggest that time-lapse 3D seismic survey is the
 43 key tool for detailed understanding of CO₂ behavior along time. However, the recorded
 44 seismic data must go through rigorous processing route to match the baseline and time-
 45 lapse datasets, so that the differences can be detected and interpreted by human labours.
 46 In this study, we propose to use the neural network (NN) to facilitate the human inter-
 47 pretation. It offers an efficient end-to-end mapping directly from the 4D seismic images
 48 to 3D CO₂ distribution. By incorporating differently processed data during training, the
 49 NN gains robustness against moderate mismatch between the baseline and time-lapse
 50 images. The generalized applications of the trained NN on different time-lapse data show
 51 great consistency throughout the monitoring history, which provides reliable analysis for
 52 CO₂ plume development as a function of time.

53 **1 Introduction**

54 CO₂ capture and storage (CCS) is an important measure for greenhouse gas mit-
 55 igation, among which geological CO₂ sequestration aims at keeping the CO₂ underground
 56 so that they are permanently removed from the atmosphere. During CCS injection projects,
 57 long-term monitoring is necessary for the purposes of understanding CO₂ behaviour in
 58 the reservoir, detecting CO₂ leakage from the storage unit, and assessing effects of con-
 59 tingency measures in case of leakage (Furre et al., 2017). Time-lapse or 4D seismic sur-
 60 vey is the key tool for detailed and quantitative CO₂ characterization varying along time
 61 (Arts et al., 2008; Boait et al., 2012).

For CCS monitoring, the time-lapse seismic analysis estimates the subsurface parameter changes between two separate seismic experiments due to CO₂ injection. The density and bulk modulus difference between the injected supercritical CO₂ and the originally saturated brine leads to dramatic velocity decrease in the storage reservoir, and consequently creates large reflectivities at interfaces between the CO₂ and brine saturated rocks (R. Chadwick et al., 2005). Hence, conventional quantitative interpretations for CO₂ plume thickness and saturation are mainly based on velocity pushdown and reflection amplitude (A. Chadwick et al., 2010). Although detailed and informative, such interpretation depends on reliable processing with relative amplitude preservation and satisfying match between the baseline and time-lapse datasets. Furthermore, it also involves tedious human interpretations, such as horizon picking in multiple 3D data volumes. Hence, it is hard to keep the whole process efficient and consistent throughout the long-term CCS project. To partially automate this process, stratigraphic inversion (Clochard et al., 2010) and full waveform inversion (FWI) (Romdhane & Querendez, 2014) have been utilized to build high-resolution models of elastic impedance or velocity, based on which the CO₂ plume can be interpreted more easily and accurately. However, the inversion-based interpretation requires satisfying initial models, which still rely on much human interpretation, and the computational cost for 3D inversion is expensive.

Different from conventional interpretation, machine learning-based seismic interpretation, trained by labels obtained from experienced interpreter or realistic model building, can deliver satisfying results with much higher efficiency, e.g., in fault and horizon detection (Wu et al., 2019; Geng et al., 2020). For CO₂ interpretation in CCS projects, machine learning is more attractive due to its inherent interpretation consistency throughout the long-term monitoring history. Using synthetic data generated by flow simulation, rock physics modeling and acoustic wave equation modeling, (Wang et al., 2020) applied different machine learning algorithms to predict CO₂ saturation from seismic attributes and other downhole measurements. (Sinha et al., 2020) treated the CO₂ leakage detection as an anomaly detection problem from CO₂ injection rates and pressure data using various neural networks (NNs).

Despite of these successful machine learning applications on these frequently repeated measurements with relatively small data size, direct applications of NNs on the time-lapsed 3D surface seismic data are rarely reported due to the difficulty of acquiring or simulating such large volume datasets and corresponding labels for supervised learning. In this study, we propose a NN-based CO₂ interpretation that offers an end-to-end mapping from 4D seismic images to 3D CO₂ distribution probability. We employ a simplified 3D U-net which has been successfully utilized for fault detection from stacked 3D seismic images (Wu et al., 2019). We train it by the publicly available seismic datasets and labels from the Sleipner CCS project, which is the world first industrial offshore CCS project starting injection from 1996, and storing around 18.5 million tonnes CO₂ by 2020 (Williams & Chadwick, 2021). The trained NN is applied to time-lapse datasets acquired and processed in different years. With a single TITAN RTX GPU (24 G), the runtime for training and validation is around 3 hours, while the application takes only few seconds. The NN interpretation results show high resolution with valid consistency throughout the injection history. In addition, the NN also present robustness against processing-caused moderate mismatch between the baseline and time-lapse images for the input, which can potentially alleviate the reprocessing demands for newly acquired or reprocessed time-lapse datasets. To understand the interpretation standards of the trained NN, we perform tests based on the primary reflections in the reservoir area and its corresponding surface-related multiples below the reservoir. The outcome suggests our trained NN considers both amplitude difference and the image similarity for detecting CO₂ distribution. Incidentally, we discuss the situation when only sparse 2D labels exist. The comparison between 2D U-Net and 3D U-Net with spare training weights reveals that the 3D U-Net is advantageous in terms of interpretation resolution and continuity, but requires much more computational resources.

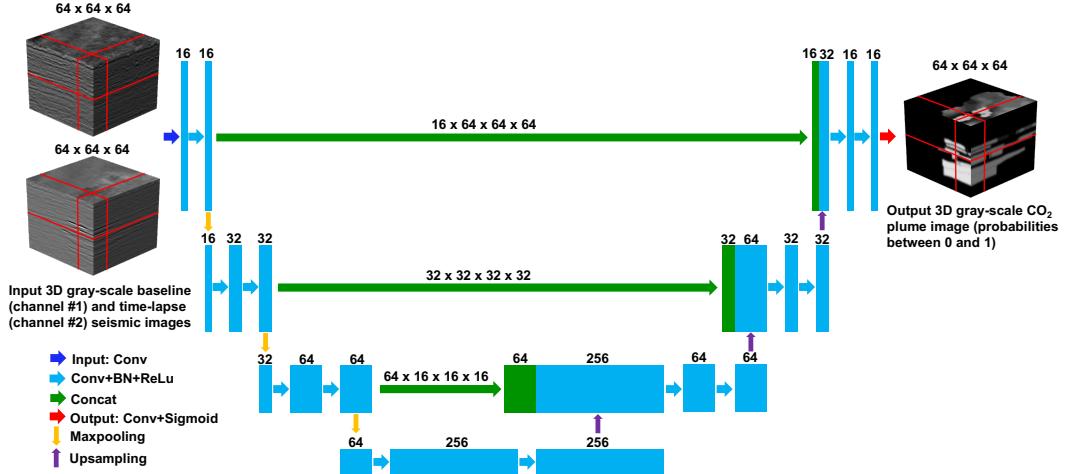


Figure 1. Simplified 3D U-Net for CO₂ interpretation from 4D seismic images.

116 The structure of this abstract is as following: first, we introduce the NN architecture;
 117 second, we address the issue of training dataset generation with attested interpre-
 118 tation labels, and present the training and validation processes; third, we test the trained
 119 NN on datasets acquired and processed in different years; finally, we discuss the NN in-
 120 terpretation standards and provide a guideline for the case of sparse 2D labels, before
 121 we conclude.

122 2 Simplified 3D U-Net for CO₂ interpretation

123 The CO₂ interpretation in this study aims at depicting 3D CO₂ distribution in the
 124 reservoir from 4D seismic images. Therefore, we consider it as an image segmentation
 125 task, which assigns ones to CO₂ saturated parts, while zeros to other parts, in the cor-
 126 responding seismic image domain. Similar image segmentation has been achieved in fault
 127 identification from seismic images (Wu et al., 2019), using a simplified version of the orig-
 128 inal U-Net (Ronneberger et al., 2015). We make some adjustments on the fault-detection
 129 network for CO₂ interpretation as shown in Figure 1. Firstly, our network input has two
 130 channels accommodating both baseline and time-lapse seismic images for CO₂ interpre-
 131 tation. Such input structure is trying to mimic the conventional interpretation, where
 132 we identify CO₂ according to amplitude changes and velocity pushdown in the time-lapse
 133 image w.r.t the baseline image. The analysis (downward) and synthesis (upward) paths
 134 are almost identical to those used for fault segmentation. However, our network uses half
 135 the number of feature maps at the bottom of the U-Net, to reduce memory and com-
 136 putational cost while preserving the performance for CO₂ interpretation. In addition, we
 137 also add an extra layer of batch normalization (BN) (Ioffe & Szegedy, 2015) before each
 138 rectified linear unit (ReLU), for faster convergence, as suggested by (Çiçek et al., 2016).
 139 The kernel size for max pooling in the downward layer is $2 \times 2 \times 2$ with a stride of 2
 140 along each dimension. Correspondingly, we use trilinear algorithm for upsampling with
 141 the scale factor of 2. The kernel size for the 3D convolutional layer is $3 \times 3 \times 3$, except
 142 for the output layer, where convolution kernel size becomes $1 \times 1 \times 1$. The final Sig-
 143 moid function outputs a 3D probability volume for CO₂ presence in the corresponding
 144 seismic image domain.

145 The dimension of 3D cubes for both input and output is $64 \times 64 \times 64$, with the
 146 voxel size of $25\text{m} \times 25\text{m} \times 8\text{ms}$ along inline, crossline and traveltime directions, respec-
 147 tively. Hence, the receptive field of each voxel in the bottom feature map is $200\text{m} \times 200\text{m} \times$

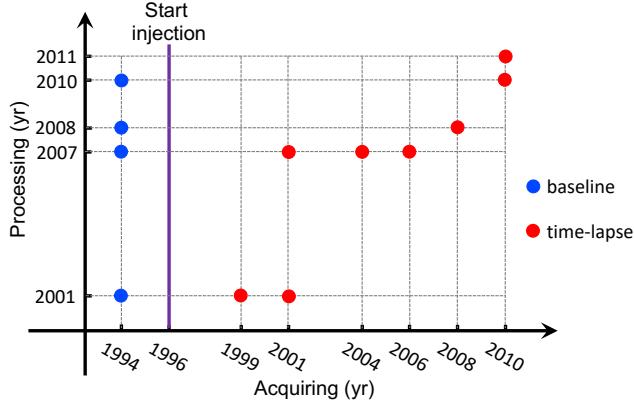


Figure 2. Acquiring and processing years for available seismic images in the 4D Sleipner seismic dataset.

64ms, which provides a satisfying balance between the computational cost and feature detection for CO₂ interpretation. To avoid the influence of disparate amplitude ranges between multi-vintages with different processing routes, we normalize the input cubes of baseline and time-lapse images, respectively, i.e., each image cube is subtracted by its mean value and divided by its standard deviation before being concatenated and fed to the NN.

3 Training and validation

3.1 Datasets generation

For the training input, publicly shared Sleipner 4D seismic dataset (Equinor, 2020a) includes abundant pairs of baseline and time-lapse datasets that have gone through time-lapse processing in different years. Figure 2 shows the acquiring and processing years of different datasets. The baseline 3D seismic data is acquired in 1994, before the CO₂ injection started in 1996 (Baklid et al., 1996). Then multiple time-lapse 3D seismic surveys have been repeated. In the shared datasets, the time-lapse data acquired in 1999, 2001, 2004, 2006, 2008 and 2010, are available. We refer to certain dataset as *xxpyy*, if it is acquired in the year of *xx* and processed in the year of *yy*. Each time-lapse dataset *xxpyy* has a correspondingly reprocessed baseline dataset *94pyy*, except for 10p11, which only has gone through image processing route without a matched baseline. Although all shared datasets contain near-, middle-, far-, and full-offset stacked images, we only utilize the near-offset part, which is adequate for the CO₂ interpretation in this study.

For the training output or label, it requires attested CO₂ interpretation corresponding to the input data. Here, we utilize the interpreted CO₂ plume boundaries of 2010 in nine internal sandstone layers provided in Sleipner 2019 Benchmark Model (Equinor, 2020b). Due to the lack of information for the CO₂ layers' thickness, we simply fill ones to the whole corresponding sandstone layer vertically within the plume lateral boundaries. Figure 3a illustrates the labeled CO₂ distribution in the model domain, along with the available interval velocity and reservoir interfaces from the benchmark model as well. Using these information, we convert the CO₂ label along depth into the image domain along traveltimes (Figure 3b). The generated label can be seen as 3D probability volume of CO₂ distribution with lateral resolution of 100 m (as indicated in the benchmark model), and vertical resolution of corresponding sandstone layers' thickness (in terms of corresponding two-way traveltimes). Considering the interpretation uncertainties in the interval velocity and interface positions for the depth to traveltimes conversion, we slightly smooth

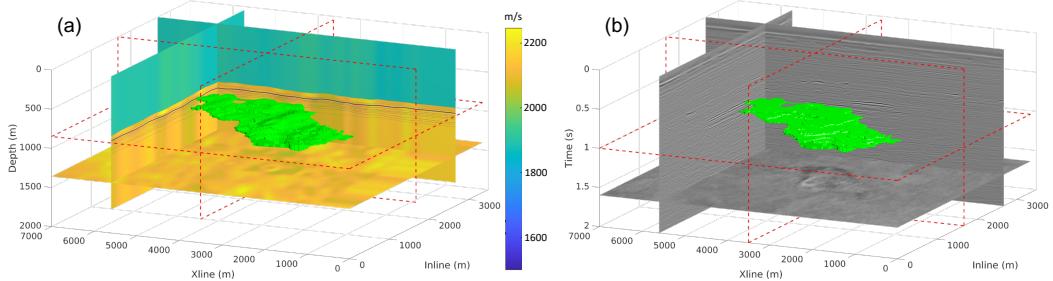


Figure 3. CO₂ plume labels for 2010 in (a) the model domain and (b) the image domain. The green blobs represent the CO₂ plume mask in 3D model or image domain. The red dashed squares indicate the sampled slice positions in 3D. The background color of the slices in (a) represent the interval velocity, while the wiggles represent the depth of interfaces in Utsira formation. The slices shown in (b) are stacked near-offset images from 10p10.

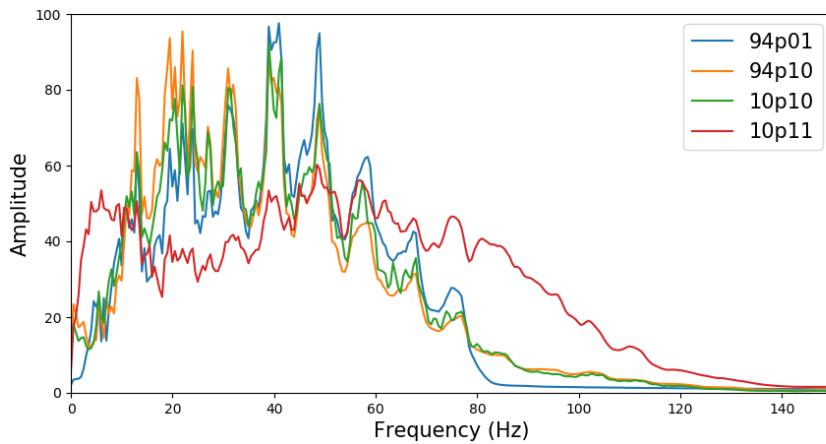


Figure 4. Comparison among amplitude spectra of 94p01, 94p10, 10p10, and 10p11.

the generated image-domain label by a 3D Gaussian filter, so that the label margins are adjusted to lower confidence.

Since the labels are limited to 2010, it is natural to choose the baseline and time-lapse pair of 2010 datasets, 94p10 and 10p10, to generate training inputs. However, we also include the originally processed baseline 94p01 and the newly processed 10p11 in training dataset generation. Consequently, there are four types of combinations between baseline and time-lapse images as 94p01 vs. 10p10, 94p01 vs. 10p11, 94p10 vs. 10p10, and 94p10 vs. 10p11. These various combinations create sufficient processing discrepancy between the baseline and time-lapse inputs. Such discrepancy in the training dataset is necessary to improve the NN's tolerance against processing-caused mismatch between the baseline and the time-lapse inputs. Figure 4 shows the comparison among amplitude spectra of the four datasets. Reasonably, 94p10 and 10p10 are mostly identical to each other, and 94p01 is slightly different from them. However, 10p11 shows marked difference in the bandwidth compared to the others, since it has not gone through the time-lapse processing route.

We randomly sample 3D cubes correspondingly in one baseline, one time-lapse and the label volumes, which covers around 3.35×7.0 km on the surface and 2.0 s along travelttime. Originally, we have 720 cubes whose centers distributing randomly in the entire

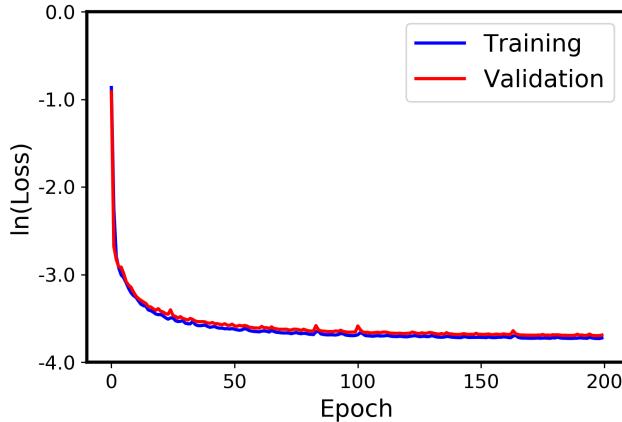


Figure 5. Training and validation losses in natural log scale.

volume. For each cube, the four combinations of baseline and time-lapse inputs corresponding to the same label are generated. Therefore, total 2880 samples are composed initially. Among these samples, around 60% show absolute zero CO₂ probability for their labels, which can significantly slow down the training convergence. Therefore, we reduce the number of zero CO₂ samples by randomly discarding some of them. Eventually, we have 1576 samples left, in which only 480 (around 1/3) of them are zero CO₂ samples. We divide these samples randomly into training and validation datasets with 1500 and 76 samples, respectively.

207 3.2 Training and validation

208 We use the binary cross-entropy (BCE) loss function since our labels are CO₂ prob-
 209 ability between zero and one. The total training epoch is 200. We use the Adam opti-
 210 mizer (Kingma & Ba, 2014) with a learning rate of 0.0002. The batch size is 30 to di-
 211 minish overfitting and improve the training efficiency with a single TITAN RTX GPU
 212 (24G). We use pytorch to implement the whole training and validation process and it
 213 takes approximately 3 hours.

214 Figure 5 shows the training and validation loss varying along the epoch number.
 215 Eventually, the decrease in the order of magnitude for both training and validation loss
 216 are around 2.8. In Figure 6, we show the NN predictions for two different samples from
 217 the training datasets. Comparing the baseline and time-lapse images in the samples, we
 218 observe good similarity for the area without CO₂ distribution. On the contrary, some
 219 amplitude anomalies indicate the CO₂ area in the corresponding labels. The NN pre-
 220 dictions for both samples are quite consistent with the labels, despite of slightly lower
 221 resolution. More epochs or smaller batch size could further improve the resolution, but
 222 the robustness and generalization of the NN will probably be compromised.

223 4 Applications

224 4.1 Robustness test

225 After our NN is trained based on random samples from the seismic images and cor-
 226 responding CO₂ labels for 2010, we firstly apply it to the same 4D seismic datasets, but
 227 with regularly sampled cubes on the entire 3D volume. The sampling number is 3×5×
 228 6 along inline, crossline and traveltime directions, resulting in 90 samples with approx-
 229 imately 40% overlapping in 3D. The runtime for such 90-sample test takes only few sec-

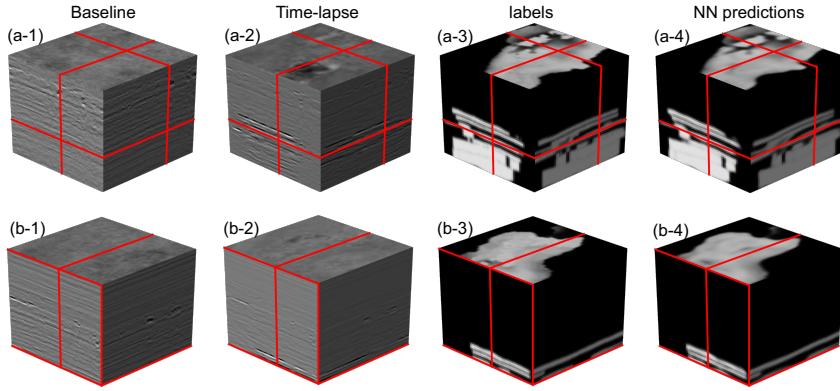


Figure 6. Trained NN predictions for two samples from the training datasets. The first two columns are the baseline and time-lapse seismic images as the NN input, and the last two columns are the human interpreted labels and NN predictions. The red lines indicate the slice positions shown by the cube surfaces.

230 seconds. The obtained NN predictions for the 90 samples can be reconstructed to the original seismic image dimension through weighted summation. The 3D weighting function
231 is:

$$w(x, y, t) = f(x)f(y)f(t), \quad (1)$$

235 where x , y and t represent inline, crossline and travelttime directions, respectively; f denotes a 1D weighting function with ones in the middle and gradually decaying to zero
236 at the edges using Hanning window as follows:
237

$$f(a) = \begin{cases} 1, & |a| \leq \frac{1}{2}\alpha L \\ \cos^2\left(\frac{\pi}{L}(|a| - \frac{1}{2}\alpha L)\right), & \frac{1}{2}\alpha L < |a| \leq \frac{L}{2} \\ 0, & |a| > \frac{L}{2}, \end{cases} \quad (2)$$

$$(3)$$

241 where $L = 64$ is the valid length of the weighting function which is consistent to the
242 NN output size along each dimension, $\alpha = 0.6$ indicates the portion of ones w.r.t L in
243 the middle of the weighting function. Thus, the weighted summation of the 90 NN pre-
244 dictions is

$$M(x, y, t) = \frac{\sum_{i=1}^{90} p_i w(x - x_i, y - y_i, t - t_i)}{\sum_{i=1}^{90} w(x - x_i, y - y_i, t - t_i)}, \quad (4)$$

245 where M is the reconstructed 3D CO₂ prediction, p_i is the NN prediction for the i th sam-
246 ple, whose cube center is (x_i, y_i, t_i) .
247

248 To test the NN robustness against processing-caused mismatch between baseline
249 and time-lapse inputs, we apply our NN on the four different sets of 90 samples corre-
250 sponding to the four different combinations as 94p01 vs. 10p10, 94p01 vs. 10p11, 94p10
251 vs. 10p10, and 94p10 vs. 10p11, respectively. Figure 7 shows the reconstructed CO₂ dis-
252 tributions based on different NN predictions. We can see that the predicted CO₂ dis-
253 tributions for different combinations of baseline and time-lapse images are almost iden-
254 tical to the label (Figure 7e). Further calculating the BCE loss between each reconstruc-
255 ted prediction and the label, Table 1 shows that 94p10 vs. 10p10 provides the best result,
256 and using 10p11 as time-lapse input always leads to larger losses in comparison with 10p10.
257 Such observations are consistent with the dataset similarity represented by the spectra
258 comparison shown in Figure 4. Regardless of the insignificant differences in the BCE losses,
259

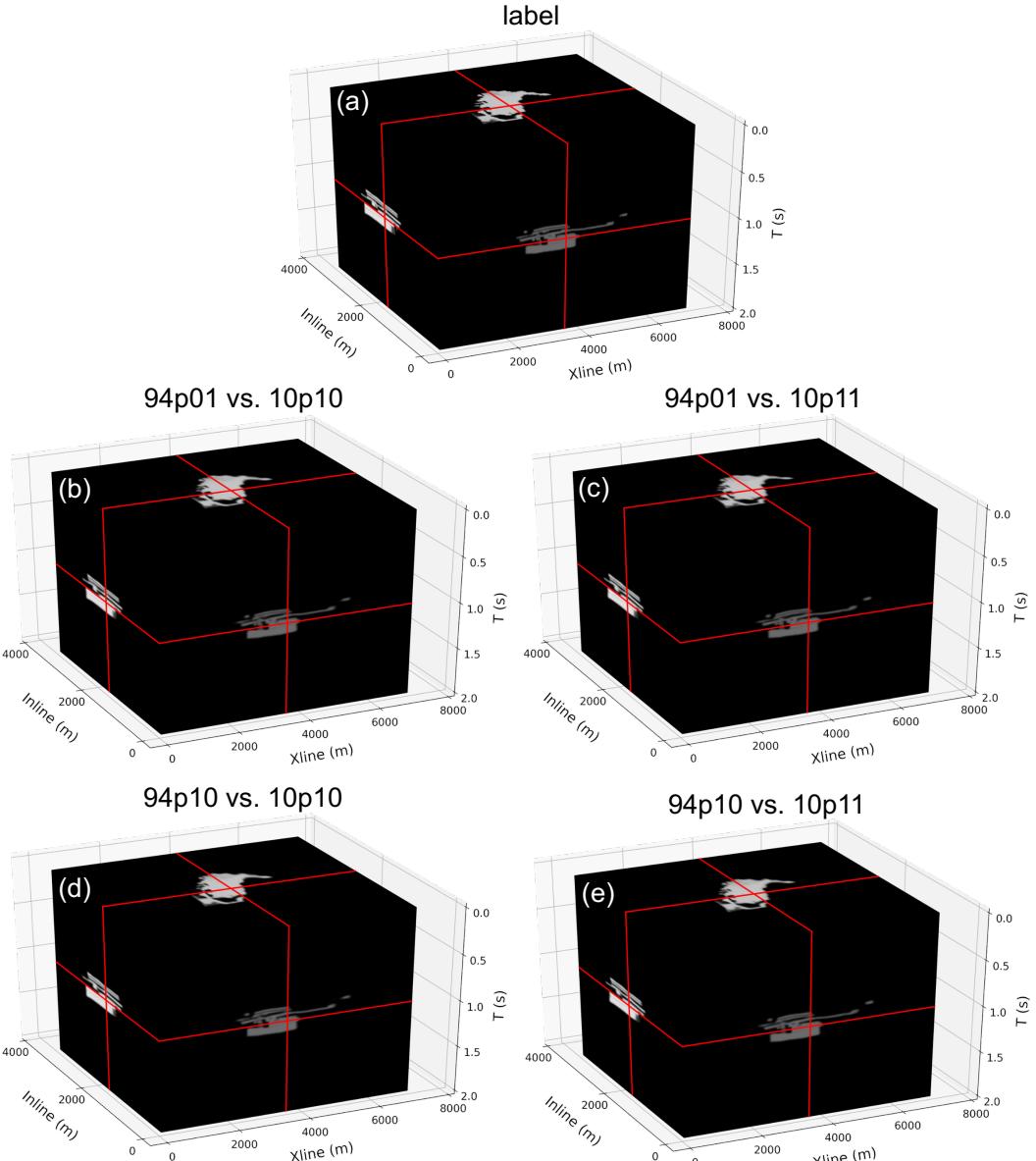


Figure 7. Comparison between (a) the label and (b,c,d,e) the reconstructed CO₂ distributions based on NN predictions using four different combinations of baseline and time-lapse images.

the NN predictions are visually undifferentiated, reflecting strong robustness of the trained NN against moderate processing mismatch between the baseline and time-lapse images.

In Figure 8, we further compare the absolute amplitude anomaly with the human interpreted and NN predicted CO₂ distributions in the top sand wedge layer above the Utsira Formation. The 2D amplitude anomaly is obtained by calculating the vertical mean of the absolute amplitude difference between 94p10 and 10p10 within the designated layer, whereas the interpreted CO₂ distributions are the vertical means of the 3D CO₂ probabilities within the same layer. We can see that the amplitude anomaly does provide a rough clue for the CO₂ distribution. However, it requires further processing and more detailed analysis to obtain the high-resolution CO₂ plume as shown in the label (Figure 8b). Contrarily, the trained NN achieves accurate CO₂ depiction with high resolution.

Table 1. BCE loss of the reconstructed CO₂ distributions w.r.t the label.

	BCE loss	baseline	94p01	94p10
time-lapse				
10p10			4.86e ⁻³	4.85e ⁻³
10p11			4.91e ⁻³	4.93e ⁻³

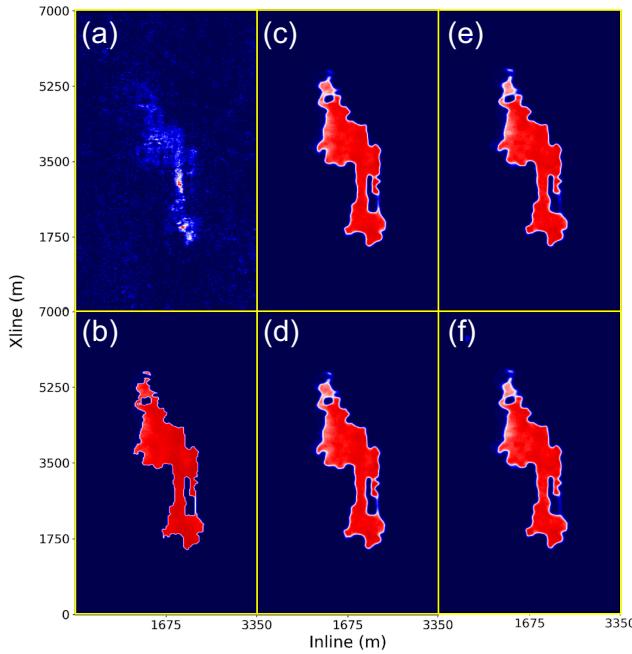


Figure 8. Comparison between (a) the absolute amplitude anomaly with (b) the human interpreted and (c,d,e,f) NN predicted CO₂ distributions in the top sand wedge layer above the Utsira Formation. For baseline input, (c) and (e) use 94p01, while (d) and (f) use 94p10; for time-lapse input, (c) and (d) use 10p10, while (e) and (f) use 10p11.

271 tion directly from the baseline and time-lapse images, even when noticeable processing
272 mismatch exists.

273 4.2 Consistency test

274 To further generalize our trained NN, we apply it to other available 4D seismic vintages
275 shared by the Sleipner CO₂ storage project. Since we have proved the robustness
276 of our trained NN, we use the same originally processed 94p01 as the NN baseline in-
277 put for all time-lapse inputs: 99p01, 01p01, 04p07, 06p07, 08p08, and 10p10. Figure 9
278 displays the NN interpreted CO₂ distribution in the top (L9), middle (L5) and base (L1)
279 of the internal sandstone layers in Utsira Formation, developing from 1999 to 2010. In
280 all displayed layers, the NN predictions are reasonably compacted with clear and con-
281 tinuous boundaries. Moreover, they are growing steadily throughout the decade, although
282 the CO₂ plume in L1 expands noticeably slower, due to the buoyancy of supercritical
283 CO₂ in the saline aquifer (Arts et al., 2008). In 1999, the top layer result (Figure9(L9-
284 1999)) shows a singularity on the probability map, indicating the injected CO₂ has just
285 reached the top of the formation (A. Chadwick et al., 2010). Similar singularities are also

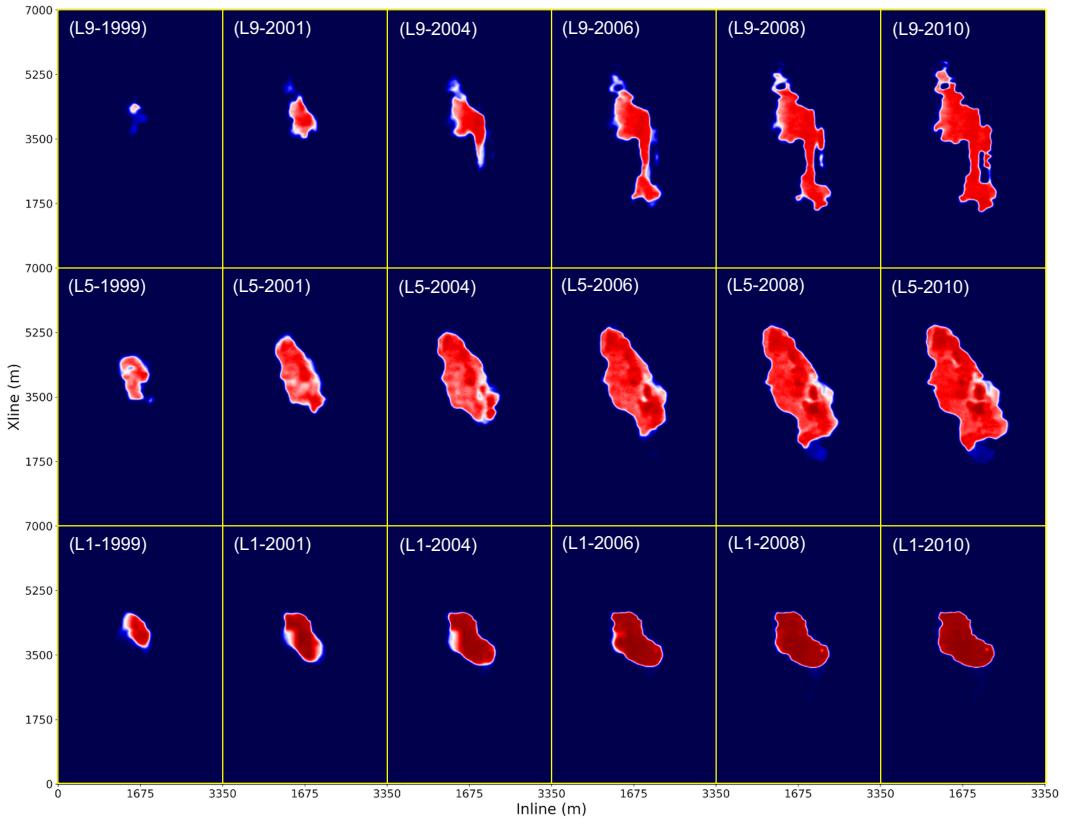


Figure 9. NN interpreted CO₂ plume expanding along L9 (top), L5 (middle) and L1 (base) of the Utsira Formation, from 1999 to 2010.

visible in Figures 9(L9-2004), (L9-2008) and (L5-1999), suggesting that our NN interpretation has the potential for high-resolution leakage detection or feeder recovery. Finally, we can also identify the migration directions of CO₂ plume in different layers, e.g., in L5, the CO₂ plume is mainly lengthened along SW-NE directions, and specifically, towards the NE direction since 2004. Generally, the NN interpretations along time are reasonably consistent in terms of CO₂ migration and plume expansion in the storage unit.

5 Discussion

5.1 Analysis for NN interpretation standards

We design a test to qualitatively explain how does the trained NN determine the CO₂ distribution given baseline and time-lapse inputs. Direct observation from Figure 6 indicates that large amplitude anomaly in the time-lapse image w.r.t the baseline image is the apparent key. In view of this intuitive hypothesis, we sample two cubes whose centers are both in the inline assemble at 1625 m from 94p10 and 10p10 as shown in Figure 10. The first sample includes the primary reflections caused by CO₂ accumulation in the storage unit, and the second sample contains the corresponding surface-related multiples right below the first sampling position. Hence, there are amplitude anomalies in both samples, despite that their labels are totally different as shown in Figure 10c. We feed these two samples of baseline and time-lapse cubes (Figures 11a and b) into the trained NN, and the predictions are shown in Figure 11c, which are consistent with the labels (Figure 11d). This implies that the NN does not solely rely on the amplitude dif-

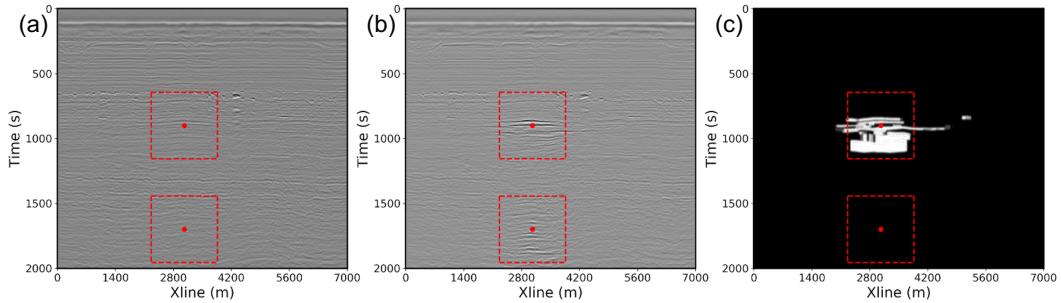


Figure 10. The inline assembles at 1625 m for (a) baseline and (b) time-lapse images from 94p10 and 10p10, along with the (c) CO₂ distribution label. The red dots and dashed lines indicate the centers and boundaries of the sampled cubes within the inline assemble.

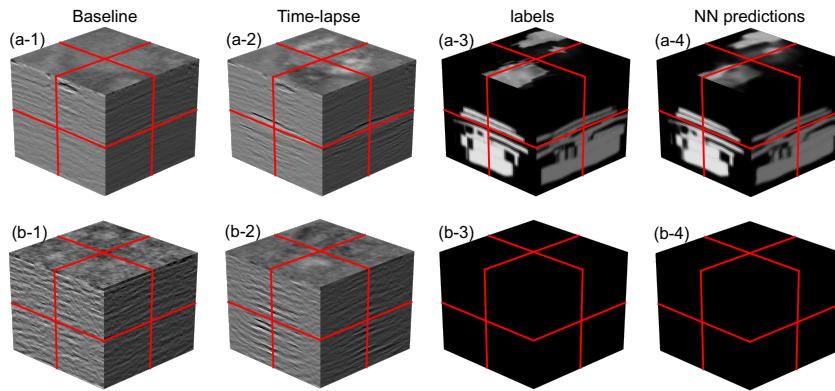


Figure 11. Trained NN predictions for the two samples shown in Figure 10. The first two columns are the baseline and time-lapse seismic images as the NN input, and the last two columns are the human interpreted labels and NN predictions.

ference anomaly between the baseline and time-lapse images to determine the CO₂ distribution.

To further explore why the NN does not misinterpret the surface-related multiples, we modify the time-lapse image for this sample. The modifications and corresponding NN predictions are shown in Figure 12. In the first modified sample, we scale up the center area of the original time-lapse image (Figure 11b-2) as the new time-lapse image (Figure 12c), whereas in the second modified sample, we scale up the center area of the original baseline image (Figure 11b-1) as the new time-lapse image (Figure 12e). Both modified samples now share the same baseline image as shown in Figure 12a and the locally scaling multiplier is displayed in Figure 12b, whose center area has the value of 4, whereas the outer area is 1. By feeding the modified samples to the trained NN, we obtain the predictions shown in Figures 12d and f. It appears that further increasing the amplitude of the multiples does not create significant CO₂ probability in the prediction. However, when we directly scale up the baseline amplitude in specific area as the time-lapse image, significant CO₂ probability emerges in the scaled up area. Hence, it reveals that the trained NN also considers the image similarity between the baseline and time-lapse images, in addition to the amplitude difference anomaly between them. This indicates that our trained NN performs similarly with human interpretation.

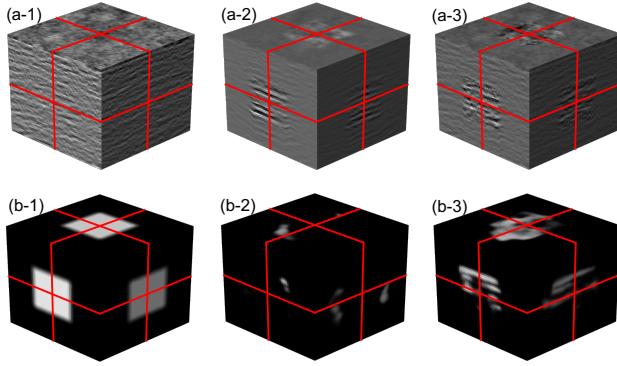


Figure 12. Modifications of the time-lapse images in the second sample shown in Figure 11. (a) is the common baseline; (b) is the locally scaling multiplier; (c) and (e) are new time-lapse images; (d) and (f) are NN predictions corresponding to (c) and (e), respectively.

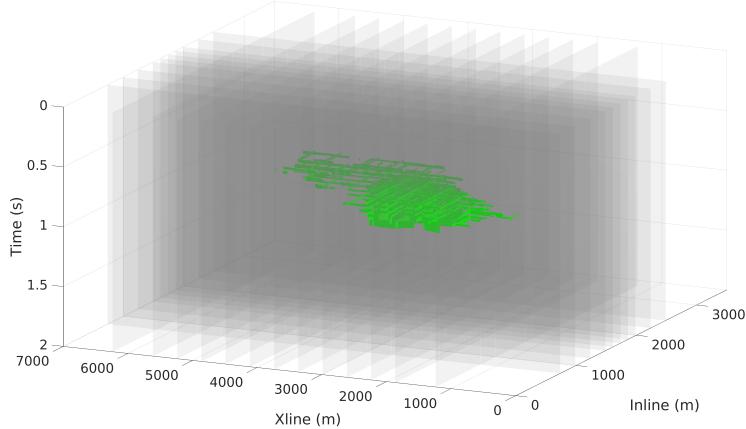


Figure 13. 2D slice labels for CO₂ distribution, indicated by the green patches on each slice represented by translucent surfaces.

324 5.2 2D vs. 3D NN using sparse 2D labels

325 The presented NN based on the Sleipner 4D seismic dataset is trained by a com-
 326 plete 3D label generated from CO₂ plume boundary interpretation for 2010. However,
 327 in more general cases, the interpreted labels for CO₂ distribution are often available in
 328 sparse 2D slices of inline and/or crossline assemblies, due to human perception limita-
 329 tions. We create an example of interpreted slices along both inline and crossline direc-
 330 tions shown in Figure 13. We sample 10 slices along inline direction, and 13 slices along
 331 crossline direction. The slice interval is smaller (125 m along inline direction and 375 m
 332 along crossline direction) near the storage unit, and larger (250 m along inline direction
 333 and 625 m along crossline direction) away from the storage unit, for a better represen-
 334 tation of the label target.

335 One way of utilizing such sparse 2D labels is to directly train a 2D U-Net. Another
 336 way is to train the 3D U-Net with corresponding weight on the sparse labels during loss
 337 evaluation (Çiçek et al., 2016). Here, we compare these two strategies to offer a guide-
 338 line under such realistic scenario of sparse interpretation labels.

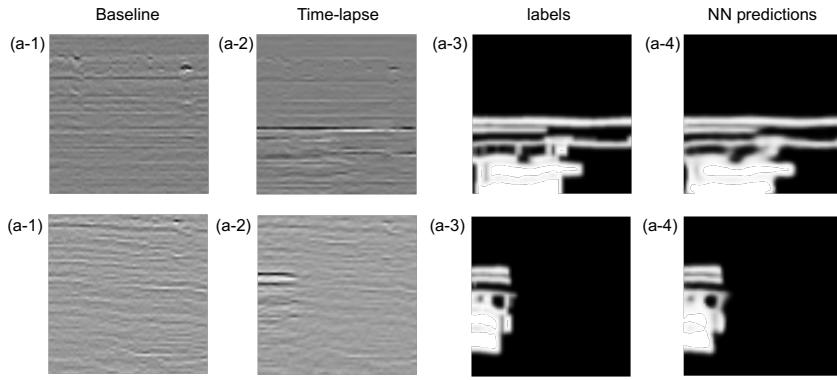


Figure 14. Trained 2D NN predictions for two samples from the training dataset with sparse 2D slice labels. The first two columns are the baseline and time-lapse seismic images as the NN input, and the last two columns are the human interpreted labels and NN predictions.

339 5.2.1 2D U-Net with sparse 2D labels

340 We use the same architecture shown in Figure 1 for the 2D U-Net, except we re-
 341 duce the dimensionality from 3D to 2D for the convolution, max pooling and upsampling.
 342 To generate training dataset, we also utilize both 94p01 and 94p10 as the baseline in-
 343 put, along with 10p10 and 10p11 as the time-lapse input. Since trace intervals along both
 344 inline and crossline directions are the same, we can train a 2D U-Net applicable on both
 345 inline and crossline assemblies. The 2D patch size is 64×64 , with the grid size of $25\text{m} \times$
 346 8ms along inline/crossline and traveltime directions, respectively. We sample 2D patches
 347 on all available slices shown in Figure 13 correspondingly in one baseline image, one time-
 348 lapse image and the label slice. The number of initially sampled patch centers is 1320,
 349 resulting in 5280 2D training samples. We also reduce the number of zero CO_2 samples,
 350 and eventually keep 1500 samples for training, of which around 25% are zero CO_2 sam-
 351 ples. The 2-D U-Net training uses the same training parameters as for the 3D U-Net.
 352 Notice that due to the dimensionality reduction, the training time has reduced signif-
 353 icantly from 3 hours to 5 minutes for 200 epochs using the same TITAN RTX GPU. Fig-
 354 ure 14 shows the results of two training samples. It appears that the NN predictions are
 355 consistent with the corresponding 2D labels.

356 To test the trained 2D U-Net, we apply it on 94p01 vs. 04p07 and 94p01 vs. 10p10,
 357 respectively, by regularly sampling 2D patches along all inline and crossline assemblies
 358 for the 3D dataset volume. The sampling numbers are also $(3, 5, 6)$ along inline, crossline
 359 and traveltime directions. Traversing all the inline and crossline assemblies respectively,
 360 we have the total sampling number as $269 \times (5 \times 6) + 561 \times (3 \times 6) = 18168$, in which
 361 269 and 561 are the number of inline and crossline assemblies, respectively. We combine
 362 these 2D NN predictions together by the weighted summation as shown in Equation 4,
 363 only now with 2D Hanning-windowed weighting functions. Figures 16a and b display the
 364 reconstructed 3D CO_2 distributions for the two tests. Compared to the label and the
 365 predictions from the original 3-D U-Net shown in Figure 7, the 2010 result predicted by
 366 2D U-Net shows lower resolution and more artifacts around the edges of the volume. Sim-
 367 ilar defects are also visible in the 2004 result obtained from 2D U-Net. We further dis-
 368 play the reconstructed 2D CO_2 distributions along the top sand wedge layer for both tests
 369 in Figures 17a and b. Distinct scratch-like artifacts are visible in comparison with the
 370 corresponding 3D U-Net predictions shown in Figure 9. This is because the 2D U-Net
 371 cannot preserve the continuity along the 3rd dimension vertical to the plane where the
 372 2D U-Net is applied. Although we apply the 2D U-Net along both inline and crossline
 373 assemblies then combine the predictions by weighted summation, the outcomes simply

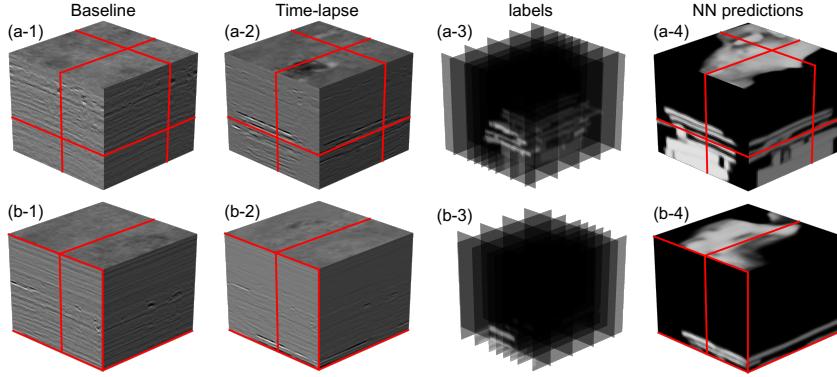


Figure 15. Trained 3D NN predictions for two samples from the training dataset with sparse 2D slice labels. The first two columns are the baseline and time-lapse seismic images as the NN input, and the last two columns are the human interpreted labels and NN predictions.

present discontinuities along both directions (the horizontal and vertical "scratches") as shown in Figures 17a and b.

376 5.2.2 3D U-Net with sparse 2D labels

377 We train a new 3D U-Net with sparse 2D labels using the same training dataset
 378 and parameters as those used in the original 3D U-Net with full 3D labels. Thus, for each
 379 sampled cube, the baseline and time-lapses images are unchanged, but the label now has
 380 a corresponding weight, in which the sampled 2D slice position is assigned as one, whereas
 381 other part is zero. The weight is implemented during the BCE loss calculation and then
 382 backpropagated to influence the 3D U-Net update. The runtime for 200 epochs are ba-
 383 sically the same as for the original U-Net training. Figure 15 shows the same samples
 384 displayed in Figure 6. We can see the NN inputs are exactly the same, whereas the la-
 385 bels are vastly different, since now we only have certain vertical slices (Figures 15a-3 and
 386 b-3) instead of the whole cube (Figures 6a-3 and b-3). However, the NN predictions are
 387 reasonably consistent. The horizontal slices shown in Figures 15a-4 and b-4 have been
 388 retrieved with satisfying resolution and continuity, even though they are not exactly the
 389 same as predictions from NN trained by full 3D labels shown in Figures 6a-4 and b-4.

390 We also test the sparsely trained 3D U-Net on 94p01 vs. 04p07 and 94p01 vs. 10p10,
 391 respectively. The 3D reconstructed CO₂ distributions are shown in Figures 16c and d.
 392 It appears that the sparsely trained 3D U-Net results provide much higher resolution than
 393 those obtained from 2D U-Net (Figures 16a and b). Figures 17c and d further display
 394 the top sand wedge layer CO₂ distributions from the sparsely trained 3D U-Net. Com-
 395 pared to the 2D U-Net results (Figures 17a and b), there are no scratch-like artifacts and
 396 the plume boundaries exhibit more continuity.

397 In summary, the 3D U-Net trained by weighted sparse labels generally present higher-
 398 quality interpretation than the 2D U-Net trained by the same labels, in terms of reso-
 399 lution, boundary continuity and artifacts. However, the 3D U-Net training and appli-
 400 cations require much more computational resources than the 2D U-Net. Hence, we sug-
 401 gest to use 3D U-Net even for sparse 2D labels as long as necessary computational power
 402 (GPU with large enough memory) is accessible.

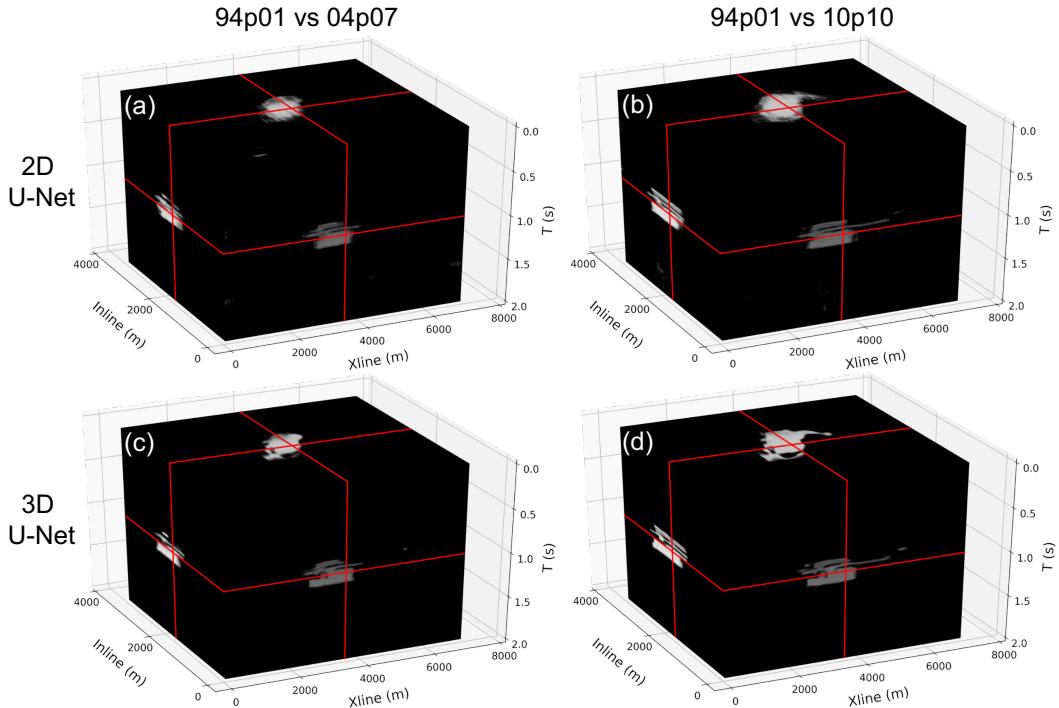


Figure 16. Reconstructed 3D CO₂ distributions from NN predictions for the tests of 94p01 vs. 04p07 and 94p01 vs. 10p10, using trained 2-D U-Net and 3-D U-Net, respectively.

6 Conclusion

We utilized a simplified 3D U-Net to interpret the 3D CO₂ distribution from large 4D seismic images. The NN is trained on Sleipner datasets acquired from 1994 and 2010, but processed in 2001, 2010 and 2011. Hence, the trained NN shows reasonable robustness against processing-caused mismatch between the baseline and time-lapse images. Moreover, the generalized applications on other time-lapse images acquired from 1999 to 2008 also achieve satisfying results with great interpretation consistency. We also analyzed the NN interpretation standards and provide NN training strategy under more realistic scenario where only sparse 2D labels are available. Overall, the studied 3D U-Net is proved to be an efficient, robust and flexible tool for CO₂ interpretation from 4D seismic monitoring during long-term CO₂ injection projects.

Acknowledgments

The authors acknowledge the Singapore Economic Development Board for its financial support through the Petroleum Engineering Professorship. The python package and all supplementary materials related to this article can be found online https://github.com/nusbei/CO2_Sleipner.

References

- Arts, R., Chadwick, A., Eiken, O., Thibeau, S., & Nooner, S. (2008). Ten years' experience of monitoring CO₂ injection in the Utsira Sand at Sleipner, offshore Norway. *First Break*, 26(1).
- Baklid, A., Korbol, R., Owren, G., et al. (1996). Sleipner Vest CO₂ disposal, CO₂ injection into a shallow underground aquifer. In *Spe annual technical confer-*

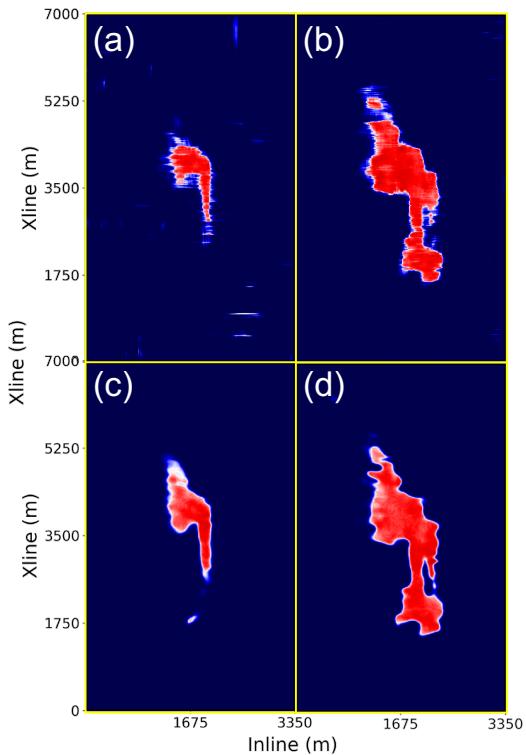


Figure 17. NN interpreted CO₂ distributions in the top sand wedge layer above the Utsira Formation. (a) and (b) are the 2D U-Net results from 94p01 vs. 04p07 and 94p01 vs. 10p10, respectively, while (c) and (d) are the 3D U-Net results from 94p01 vs. 04p07 and 94p01 vs. 10p10, respectively.

- ence and exhibition.
- Boait, F., White, N., Bickle, M., Chadwick, R., Neufeld, J., & Huppert, H. (2012). Spatial and temporal evolution of injected CO₂ at the Sleipner Field, North Sea. *Journal of Geophysical Research: Solid Earth*, 117(B3), B03309.
- Chadwick, A., Williams, G., Delepine, N., Clochard, V., Labat, K., Sturton, S., ... others (2010). Quantitative analysis of time-lapse seismic monitoring data at the Sleipner CO₂ storage operation. *The Leading Edge*, 29(2), 170–177.
- Chadwick, R., Arts, R., & Eiken, O. (2005). 4D seismic quantification of a growing CO₂ plume at Sleipner, North Sea. In *Geological society, london, petroleum geology conference series* (Vol. 6, pp. 1385–1399).
- Çiçek, Ö., Abdulkadir, A., Lienkamp, S. S., Brox, T., & Ronneberger, O. (2016). 3D U-Net: Learning Dense Volumetric Segmentation from Sparse Annotation. In *International conference on medical image computing and computer-assisted intervention* (pp. 424–432).
- Clochard, V., Delépine, N., Labat, K., & Ricarte, P. (2010). CO₂ plume imaging using 3D pre-stack stratigraphic inversion: A case study on the Sleipner field. *First Break*, 28(1).
- Equinor. (2020a). *Sleipner 2019 Benchmark Model*. Retrieved from <https://co2datashare.org/dataset/e6f67cbd-abf3-4d85-a118-ed386a994c2c> doi: 10.11582/2020.00004
- Equinor. (2020b). *Sleipner 4D Seismic Dataset*. Retrieved from <https://co2datashare.org/dataset/cbdc354c-fa61-4ab4-a0b4-134e1350a82b> doi: 10.11582/2020.00005
- Furre, A.-K., Eiken, O., Alnes, H., Vevatne, J. N., & Kiær, A. F. (2017). 20 years of Monitoring CO₂-injection at Sleipner. *Energy Procedia*, 114, 3916–3926.
- Geng, Z., Wu, X., Shi, Y., & Fomel, S. (2020). Deep learning for relative geologic time and seismic horizons. *Geophysics*, 85(4), WA87–WA100.
- Ioffe, S., & Szegedy, C. (2015). Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning* (pp. 448–456).
- Kingma, D. P., & Ba, J. (2014). Adam: A Method for Stochastic Optimization. *arXiv preprint arXiv:1412.6980*.
- Romdhane, A., & Querendez, E. (2014). CO₂ Characterization at the Sleipner Field with Full Waveform Inversion: Application to Synthetic and Real Data. *Energy Procedia*, 63, 4358–4365.
- Ronneberger, O., Fischer, P., & Brox, T. (2015). U-net: Convolutional networks for biomedical image segmentation. In *International conference on medical image computing and computer-assisted intervention* (pp. 234–241).
- Sinha, S., de Lima, R. P., Lin, Y., Sun, A. Y., Symons, N., Pawar, R., & Guthrie, G. (2020). Normal or abnormal? Machine learning for the leakage detection in carbon sequestration projects using pressure field data. *International Journal of Greenhouse Gas Control*, 103, 103189.
- Wang, Z., Dilmore, R. M., & Harbert, W. (2020). Inferring CO₂ saturation from synthetic surface seismic and downhole monitoring data using machine learning for leakage detection at CO₂ sequestration sites. *International Journal of Greenhouse Gas Control*, 100, 103115.
- Williams, G., & Chadwick, R. (2021). Influence of reservoir-scale heterogeneities on the growth, evolution and migration of a CO₂ plume at the Sleipner Field, Norwegian North Sea. *International Journal of Greenhouse Gas Control*, 106, 103260.
- Wu, X., Liang, L., Shi, Y., & Fomel, S. (2019). FaultSeg3D: Using synthetic data sets to train an end-to-end convolutional neural network for 3D seismic fault segmentation. *Geophysics*, 84(3), IM35–IM45. Retrieved from <https://doi.org/10.1190/geo2018-0646.1> doi: 10.1190/geo2018-0646.1