

1      Stochastic pix2vid: A new spatiotemporal method for  
2      image-to-video synthesis in geologic CO<sub>2</sub> storage  
3      prediction

4      Misael M. Morales<sup>1\*</sup>, Carlos Torres-Verdin<sup>1,2</sup>, and Michael J. Pyrcz<sup>1,2</sup>

5      1. Hildebrand Department of Petroleum and Geosystems Engineering, The University of Texas at Austin

6      2. Jackson School of Geosciences, The University of Texas at Austin

7      \*Corresponding author; email: [misaelmorales@utexas.edu](mailto:misaelmorales@utexas.edu)

8      **Abstract**

9      Numerical simulation of multiphase flow in porous media is an important step in understanding the dynamic  
10     behavior of geologic CO<sub>2</sub> storage (GCS). Scaling up GCS requires fast and accurate high-resolution modeling  
11     of the storage reservoir pressure building and saturation plume migration; however, such modeling is chal-  
12     lenging due to the high computational costs of traditional physics-based simulations. Deep learning models  
13     trained with numerical simulation data can provide a fast and reliable alternative to expensive physics-based  
14     numerical simulations. We present a pix2vid neural network architecture for solving multiphase fluid flow  
15     problems with superior speed, accuracy, and efficiency. The pix2vid model is designed based on the principles  
16     of computer vision and video synthesis and is able to generate dynamic spatiotemporal predictions of fluid  
17     flow from static reservoir models. We apply the pix2vid model to a highly-complex CO<sub>2</sub>-water multiphase  
18     problem with a wide range of reservoir models in terms of porosity and permeability heterogeneity, facies  
19     distribution, and injection configurations. The pix2vid method is first-of-its-kind in static-to-dynamic pre-  
20     diction of reservoir behavior, where a single static input is mapped to its dynamic response. The pix2vid  
21     method provides superior performance in highly heterogeneous geologic formations and complex estimation  
22     such as gas saturation and pressure buildup plume determination. The trained model can serve as a general-  
23     purpose, static-to-dynamic alternative to traditional numerical reservoir simulation of 2D CO<sub>2</sub> injection  
24     problems with significant speedups compared to traditional methods.

25      *Keywords:* Image-to-video synthesis, Spatiotemporal prediction, Convolutional neural network, Recur-  
26      rent neural network, Proxy model

## 27 1 Introduction

28 Geologic CO<sub>2</sub> sequestration (GCS) has emerged as a proven technology to reduce anthropogenic greenhouse  
29 gas emissions to the atmosphere [citation]. This has become increasingly popular worldwide due to the need  
30 to meet international climate protection agreements [citation]. However, there are several technical challenges  
31 associated with the modeling of large-scale GCS operations. In order to accurately forecast and monitor  
32 subsurface multiphase flow, physics-based high-fidelity numerical simulation is required [citation]. These nu-  
33 mercial simulations are computationally intensive and time-consuming since they require iterative solutions  
34 of large-scale nonlinear systems of equations [citation]. Similarly, due to the large degree of uncertainty in  
35 subsurface data collection, inherent uncertainty in the spatial distribution of the properties of heterogeneous  
36 porous media require a robust probabilistic assessment for improved engineering decision-making [citation].  
37 In order to capture the fine-scale multiphase flow behavior given an uncertain spatial distribution of sub-  
38 surface properties, a large number of forward numerical simulation runs are required, leading to very high  
39 computational costs [citation]. To overcome this, machine learning techniques have emerged as candidate  
40 reduced-order models (ROMs) for efficient parameterization and prediction of subsurface flow and transport  
41 behavior [citation].

42 Recent advancements in computing power, specifically GPU-enabled neural network models, have accel-  
43 erated the fields of forward and inverse modeling [citation]. Classical techniques are often hindered by the  
44 size of the models and data, specifically the volume, velocity, variety, value, and veracity encountered in big  
45 data [citation]. By analyzing extensive data sets, machine learning techniques can uncover complex latent  
46 patterns and relationships that may not be discernible through traditional methods [citation]. When com-  
47 bined with a reduced-order modeling framework, machine learning approaches can efficiently and accurately  
48 exploit latent or salient features hidden in the data, removing redundancies or noise, and decreasing the or-  
49 der of the problem significantly [citation]. These approaches can often be divided into two main categories,  
50 namely purely data-driven mapping operators or physics-informed neural networks (PINNs). Typically, the  
51 training process for PINNs is done by the minimization of the (physical) loss from the residual of the gov-  
52 erning partial differential equations (PDEs) that govern the system along with the losses associated with  
53 the initial and boundary conditions [citation]. However, over variants of PINNs such as physics-guided or  
54 physics-constrained neural networks have also proven useful for subsurface energy resource engineering appli-  
55 cations [citation]. On the other hand, data-driven mapping operators, or proxy models, are neural network  
56 architectures trained with labeled data that produce a mapping from input features to output parameters  
57 [citation]. This procedure requires significant amounts of training data but can be applied to a wide variety  
58 of settings and conditions [citation] but suffer from lack of generalization and struggle to provide accurate

59 predictions away from the domain of the training data. For both approaches, typically, spatial relationships  
60 are captured through convolutional neural networks (CNNs) and the temporal relationships through  
61 recurrent neural networks (RNNs) [citation], but recent advancements in transformer-based architectures are  
62 showing improved performance compared to the aforementioned techniques [citation]. In general, efficient  
63 compression of the input features into a representative latent space is proven as an effective approach for  
64 spatial and temporal parameterization of the forward or inverse problem.

65 A number of machine learning-based proxy (or surrogate) models have been developed to estimate the  
66 reservoir behavior in subsurface energy resource applications. Most techniques rely on the concept of image  
67 translation, or pix2pix, where a target image is predicted from an input image [citation]. Maldonado and  
68 Pyrcz [citation] developed a convolutional U-net model to predict pressure and saturation states given an  
69 uncertain geologic realization. This work is an example of image-to-image static forecasting, where the time  
70 state is given as an input, and the proxy model will predict a single response state of pressure and saturation  
71 at the given time. Wen and Benson [citation] developed a Fourier Neural Operator (FNO) architecture  
72 to predict image-to-image response states of pressure and saturation from an uncertain geologic realization  
73 and was further extended for multi-scale and nested domains [citation]. Moreover, numerous other proxy  
74 models have been developed for subsurface applications using more complex architectures such as generative  
75 adversarial networks (GANs) [citation] and transformers [citation]. However, most of these formulations are  
76 presented as an even-determined or sometimes over-determined estimation problem, with equal or greater  
77 number of input features compared to the output parameters since they are based on the pix2pix, or image-  
78 to-image formulation.

79 Moving beyond image-to-image predictions, Kim and Durlofsky [citation] developed a convolutional-  
80 recurrent proxy for image-to-series forecasting and discussed its advantages for closed-loop reservoir man-  
81 agement under geologic uncertainty. This method moves beyond the image-to-image forecasting and exploits  
82 a spatiotemporal latent space in the encoder-decoder neural network architecture to obtain well flow rates  
83 and pressures over time from a static geologic realization. The image-to-series formulation can still be an  
84 even- or over-determined estimation problem, where we have equal or more inputs than outputs. Further-  
85 more, Tang et al. [citation] and Jiang and Durlofsky [citation] developed a recurrent residual U-net (R-U-net)  
86 proxy for the prediction of dynamic pressure- and saturation-over-time from uncertain geologic realizations.  
87 This method aim to obtain dynamic response states over time from a single static input. This proxy is  
88 formulated as a more interesting under-determined estimation problem, where the number of input features  
89 is a fraction of the number of output parameters. However, the recurrent R-U-net proxy is limited by the  
90 fact that only the latent space receives spatiotemporal processing, while the model reconstruction is done via  
91 time-distributed deconvolutions, treating time as an additional “spatial” dimension, and not fully exploiting

92 the spatiotemporal relations in the data and latent space as an image-to-video forecasting formulation.

93 The problem of image-to-video forecasting, also known as video synthesis, has been approached previously  
94 by researchers in the field of computer vision. Iliadis et al. [citation] were the first to develop a deep learning-  
95 based framework for video compressive sensing to reconstruct a video sequence from a single measured frame  
96 using a deep fully-connected neural network, or artificial neural network (ANN). Despite excellent accuracy  
97 in the video predictions, this method is still limited by time-distributed fully-connected layers in the encoder  
98 and decoder portions of the network, thus not exploiting the spatiotemporal relationships in the data. Xu  
99 and Ren [citation] developed a three-part encoder-recurrent-decoder network for video reconstruction from  
100 the estimated motion fields of the encoded frames. The implementation is similar to that of Tang et al.  
101 [citation] and Jiang and Durlofsky [citation] in that it applies a recurrent update in the latent space but  
102 relies on time-distributed deconvolutions for the video frames reconstruction. Dorkenwald et al. [citation]  
103 developed a conditional invertible neural network (cINN) as a bijective mapping between image and video  
104 domains using a dynamic latent representation. The cINN architecture allowed for video-to-image and image-  
105 to-video predictions, proving possible the generation of video frames from a static input image. Finally,  
106 Holynski et al. [citation] implemented the idea of Eulerian motion fields to define the moving portions of the  
107 image and thus were able to accurately reconstruct a series of video frames from a static image using a  
108 spatiotemporal latent space parameterization. These advancements in the field of computer vision and video  
109 compressed sensing serve as a foundation for our image-to-video spatiotemporal proxy model.

110 In this work, we propose a novel image-to-video spatiotemporal proxy model for the prediction of dynamic  
111 reservoir behavior over time from an uncertain static geologic realization. In this work, we apply the  
112 spatiotemporal proxy to a large-scale GCS operation. Our model exploits the spatial and temporal structures  
113 in latent space to dynamically reconstruct the time-dependent pressure and saturation states from a static  
114 geologic realization. The encoder portion of the network receives as inputs the static geologic realization with  
115 channels representing the porosity, permeability, and facies distributions, and the location of CO<sub>2</sub> injection  
116 well(s). The uncertain geologic realizations are generated from a wide array of possible geologic scenarios  
117 (e.g., fluvial, turbidite, and deepwater lobe systems), and the number and location of CO<sub>2</sub> injection wells is  
118 also considered uncertain. The model then reconstructs the dynamic pressure and saturation distributions  
119 using a spatiotemporal decoder network with convolutional long short-term memory (ConvLSTM) layers,  
120 which are concatenated with the residuals of the spatial latent parameterizations from the encoder network.  
121 Thus, it is not an encoder-recurrent-decoder architecture, but instead a fully spatiotemporal convolutional-  
122 recurrent image-to-video model. Our proxy model shows significant advantages compared to image-to-image  
123 and encoder-recurrent-decoder models in terms of computational efficiency and prediction accuracy and can  
124 be used as a replacement for high-fidelity simulations (HFS) in GCS projects as an image-to-video mapping

125 operator.

126 In the methodology section, we discuss the proposed spatiotemporal proxy model architecture as well as  
127 the geologic modeling and numerical reservoir simulation steps required to generate the training data. In  
128 the results and discussion sections, we evaluate the training and performance of the proposed proxy model  
129 and compare its efficiency and accuracy to high-fidelity numerical simulations using a 2D synthetic case for  
130 large-scale GCS operations.

## 131 2 Methodology

132 This section describes the governing equations, reservoir model and simulation specifications, model archi-  
133 tecture, and training strategy of the pix2vid model.

134 **2.1 Governing equations** For the CO<sub>2</sub>-water multiphase flow problem, the general form of the mass  
135 accumulation for component  $\kappa = \text{CO}_2$  or water is given by [citation]:

$$\frac{\partial M^k}{\partial t} = -\nabla \bullet F^\kappa + q^\kappa. \quad (1)$$

136 For each component  $\kappa$ , the mass accumulation term  $M^\kappa$  is summed over all phases  $p$ ,

$$M^k = \phi \sum_p S_p \rho_p X_p^\kappa \quad (2)$$

137 where  $\phi$  is the porosity,  $S_p$  is the saturation of phase  $p$ ,  $\rho_p$  is the density of phase  $p$ , and  $X_p^\kappa$  is the mass  
138 fraction of component  $\kappa$  present in phase  $p$ . For each component  $\kappa$ , there is also the advective mass flux  
139  $F^\kappa|_{adv}$  obtained by summing over all phases  $p$ ,

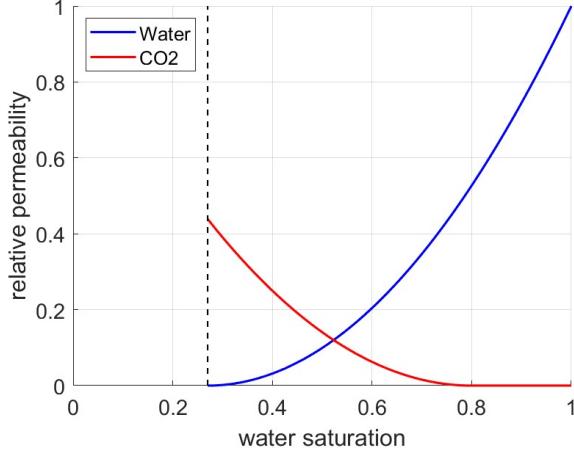
$$F^\kappa|_{adv} = \sum_p X_p^\kappa F_p \quad (3)$$

140 where each individual phase flux  $F_p$  is given by Darcy's equation:

$$F_p = \rho_p u_p = -k \frac{k_{r,p} \rho_p}{\mu_p} (\nabla P_p - \rho_p g). \quad (4)$$

141 Here,  $u_p$  is the Darcy velocity of phase  $p$ ,  $k$  is the absolute permeability,  $k_{r,p}$  is the relative permeability  
142 of phase  $p$ ,  $\mu_p$  is the viscosity of phase  $p$ , and  $g$  is the gravitational acceleration constant. The relative  
143 permeability curves for the CO<sub>2</sub>-water system are shown in Figure 1. The fluid pressure of phase  $p$ ,

$$P_p = P + P_c \quad (5)$$



**Figure 1:** Relative permeability curves for the CO<sub>2</sub>-water system. The residual saturations are 0.27 and 0.2 for water and CO<sub>2</sub>, respectively.

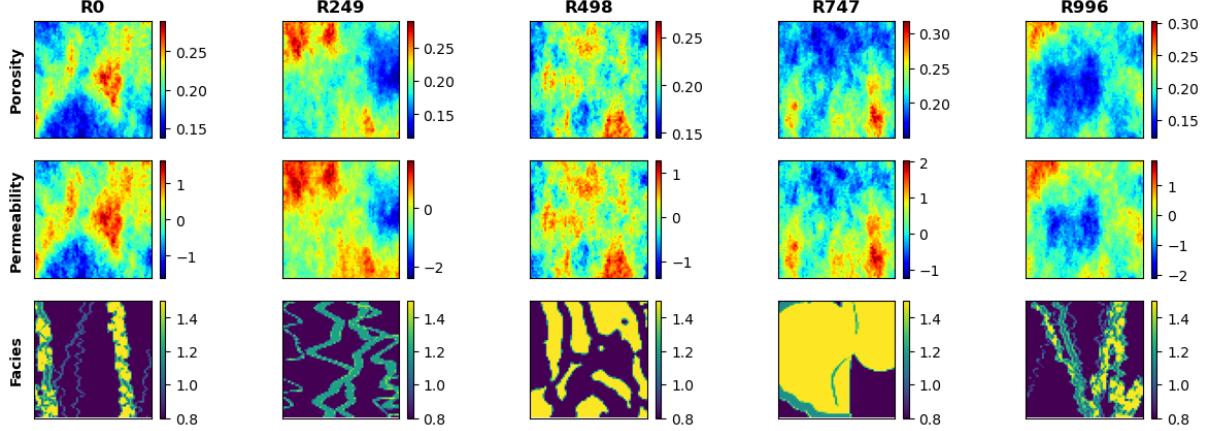
is given by the sum of the reference phase pressure  $P$  and the capillary pressure  $P_c$ . The numerical simulation does not include molecular diffusion or hydrodynamic dispersion for practical purposes.

## 2.2 Reservoir Model and Simulation

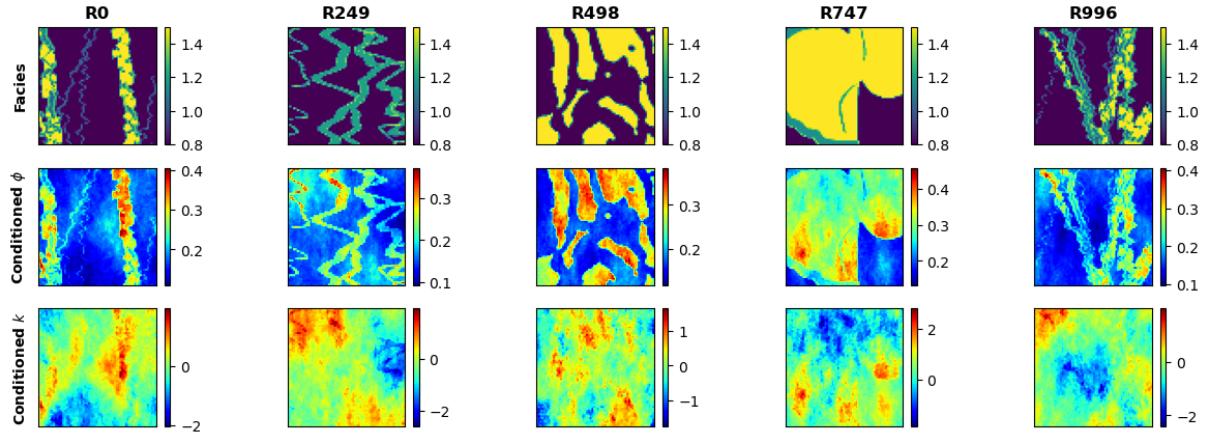
We use SGeMS [citation] to construct an ensemble of realizations that is representative of various potential geologic scenarios for CO<sub>2</sub> storage in deep geological formations. Using sequential Gaussian co-simulation [citation], we generate a set of 1,000 random porosity ( $\phi$ ) and permeability ( $k$ ) distributions with a wide range of values, as shown in Figure 2. Facies distributions are obtained from a library of deepwater fluvial training images [citation]. These encompass a wide range of possible geologic scenarios including marked point (lobe, ellipse, and bar), FluvSim (channel, channel-levee, and channel-levee-splay), surface based (compensational cycles of lobes), and bank retreat (channel complex). To generate consistent porosity and permeability distributions with the facies-based geologic scenarios, we conditionally multiply the original porosity and permeability distributions with the facies distributions. The resulting fluvial distributions are shown in Figure 3.

The conditioned fluvial porosity and permeability distributions simulated for the problem of geologic CO<sub>2</sub> storage using MRST [1]. Specifically, the MRST-co2lab module is used as an automatic-differentiation framework for the compositional simulation of the two-phase CO<sub>2</sub>-water problem. The reservoir is initialized as a fully water saturated zone (i.e., aquifer) with an initial pressure of 4,000 psi. The reservoir has constant isothermal conditions and pressure boundary conditions and represents a large-scale geologic CO<sub>2</sub> storage project with negligible dip, such as found in the Illinois Basin and parts of the North Sea and Gulf Coast.

The model has dimensions of 1km-1km-100m in the x-, y-, and z-directions, respectively. We use 64 uniform grid cells in the x- and y-directions. The grid design is sufficiently refined to resolve the pressure



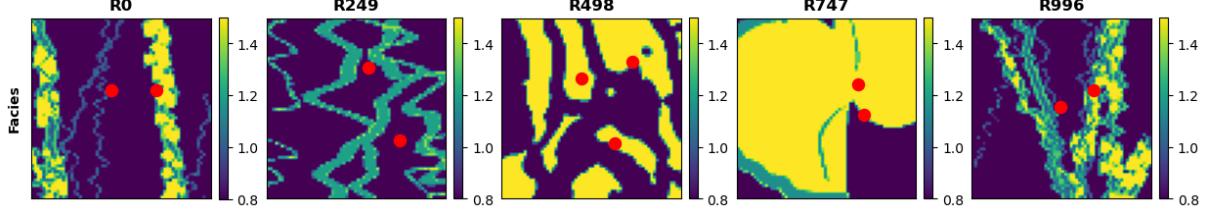
**Figure 2:** Spatial distribution of porosity (top), permeability (middle), and facies (bottom) for 5 random realizations.



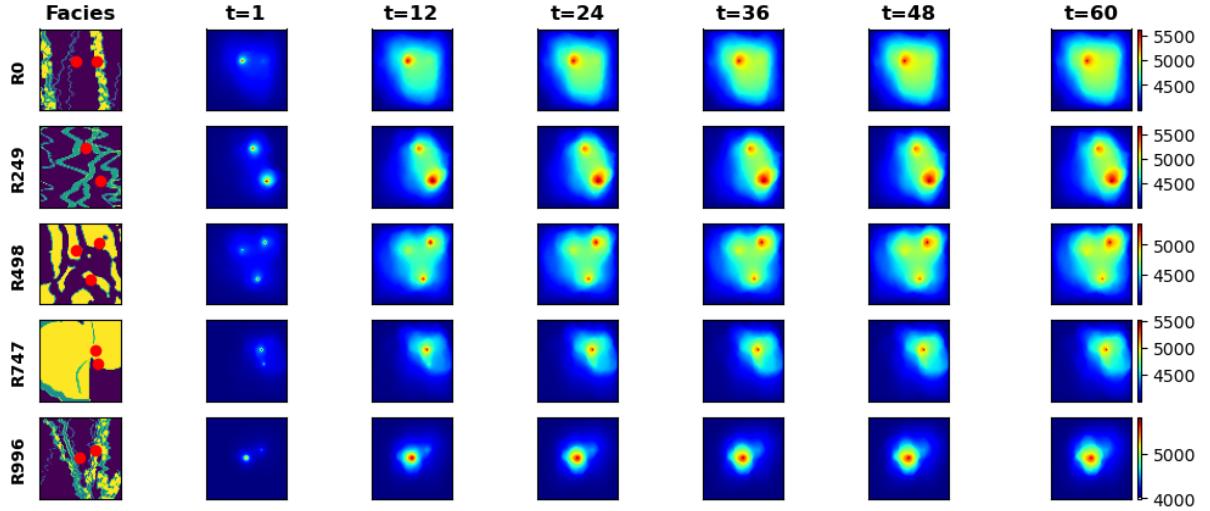
**Figure 3:** Spatial distribution conditioned to facies (top) for porosity (middle) and permeability (bottom) for 5 random realizations.

and saturation plumes in highly heterogeneous reservoirs while remaining computationally tractable for the purpose of training deep learning models. A random number of injection wells,  $w \in [1, 3]$ , are placed randomly along the reservoir for each of the 1,000 realizations, no closer than 250m from the boundaries, as shown in Figure 4. The injection well(s) are randomly placed and not conditioned to zones of preferential porosity, permeability, nor facies. Each injection well has a constant radius of 0.1m and a single and continuous perforation that injects pure supercritical CO<sub>2</sub> at a constant rate such that the total injection rate of the  $w$  well(s) is 0.5 megatons per year.

The numerical simulation is run for 5 years, monitored monthly, for a total of 60 timesteps. At each grid cell and for each time step, we resolve the implicit pressure, explicit saturation (IMPES) formulation of Eq. (1) to obtain the corresponding dynamic pressure and saturation distributions over time (videos) from the static geologic realizations of porosity and permeability conditioned to the fluvial facies (images).



**Figure 4:** CO<sub>2</sub> injection well(s) location (red) overlaid over facies distributions for 5 random realizations.



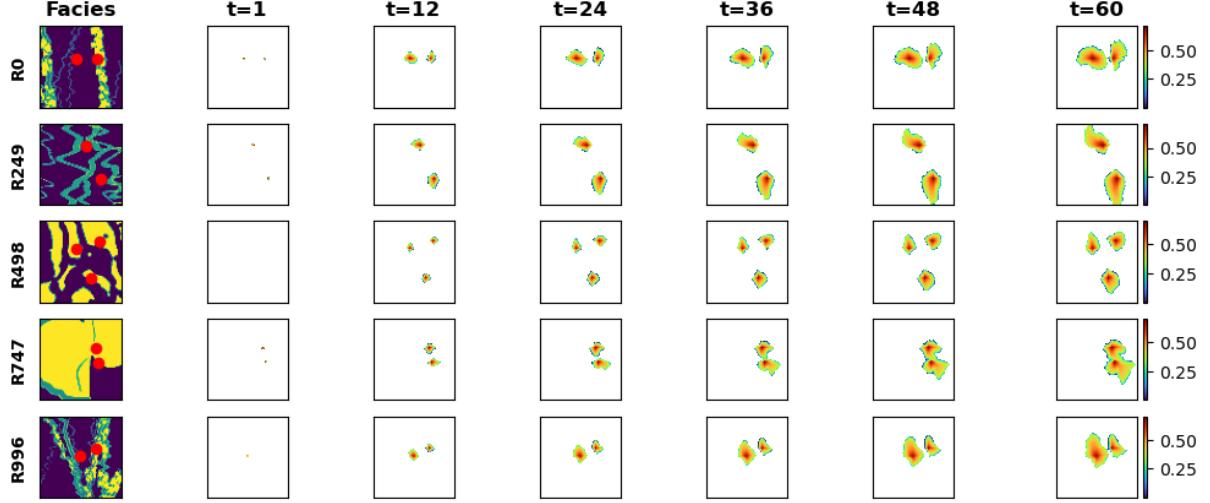
**Figure 5:** Pressure response distributions over time for 5 random realizations obtained from HFS (in psia).

176 The pressure and saturation responses corresponding to the previously-shown geologic model realizations  
177 are shown in Figures 5 and 6, respectively.

### 178 **2.3 Proxy Model Architecture**

179 The pix2vid model is designed as an image-to-video data mapping operator from the static realizations  
180 of geologic distributions of porosity, permeability and facies as well as the injector well(s) distribution, to the  
181 dynamic responses of pressure and saturation over time. A single model is trained to predict both pressure  
182 and saturation distributions over time as a multi-channel output.

183 Let  $m$  represent a geologic model realization of porosity, permeability, facies, and injector well(s) distri-  
184 butions, such that  $m = \{\phi, k, facies, w\}$ . The dynamic responses of pressure and saturation over time are  
185 given by  $d = f(m)$ , such that  $d = \{P(t), S(t)\}$  and  $f$  is the physics-based numerical reservoir simulation  
186 mapping operator. Our aim is to replace  $f$  with a more efficient data mapping operator by training the  
187 Stochastic pix2vid model. For this purpose, we exploit the latent structures in space and time of the static  
188 inputs and dynamic outputs through a spatiotemporal encoder-decoder architecture. The encoder portion  
189 of the network is comprised of sequential convolutional layers to compress the spatial features of the model  
190 realizations into a latent parameterization  $z_m$ , given by  $z_m = Enc(m)$ . In their compressed representation,



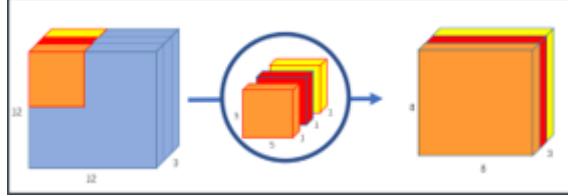
**Figure 6:** Saturation response distributions over time for 5 random realizations obtained from HFS.

191 these features represent the salient characteristics of the geologic distributions. The decoder portion of the  
 192 network is designed as a series of recursive residual convolutional-recurrent layers, such that the latent space  
 193  $z_m$  is recursively decoded into the dynamic distributions of pressure and saturation over time. The previous  
 194 timestep latent representations,  $z_d^t$ , are used in the subsequent timestep to refine the outputs and reduce  
 195 systematic error propagation in time. Thus, the full architecture is represented as

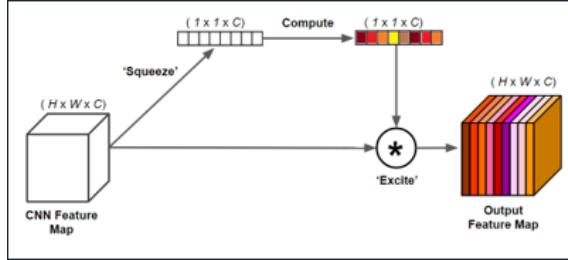
$$d = Dec^t([Enc(m); z_d^t]) \quad (6)$$

196 The encoder portion compresses the geologic realizations,  $m$ , into a latent representation  $z_m$  through the  
 197 usage of separable convolutions [citation]. This type of convolution learns the parameters for each channel in  
 198 the image separately, avoiding mixing of variables or loss of resolution, as shown in Figure 7. This is especially  
 199 important when dealing with Gaussian-distributed permeability and porosity in combination with binomial-  
 200 distributed facies and binary well(s) location distributions. Each separable convolution layer is regularized  
 201 with an  $l_1$ -norm weight of  $1 \times 10^{-6}$ . Moreover, we use a squeeze-and-excite layer to improve channel  
 202 interdependence, also avoiding mixing and loss of resolution [citation]. Each squeeze-and-excite layer will  
 203 provide the optimal network weights for each channel independent of the other channels by passing the feature  
 204 maps through a global pooling layer (squeeze) and a dense layer (excite), adding content aware mechanisms  
 205 to weight each channel adaptively, as shown in Figure 8. Furthermore, by applying instance normalization  
 206 as opposed to the more common batch normalization, we achieve channel-independent normalization of the  
 207 convolved features [citation]. Instance normalization is a special case of group normalization, where each  
 208 the numbers of channels per group is exactly 1, such that each channels gets its own normalization scheme,

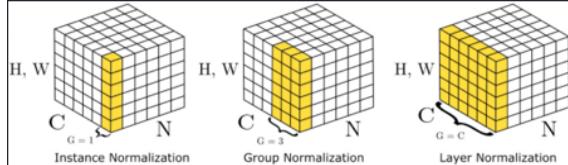
as shown in Figure 9. Parametric rectified linear units (PReLU) is used as the activation function, where at each minibatch iteration, the network learns the optimal leaky slope for activation in each layer, as shown in Figure 10. Finally, pooling and dropout are applied and the resulting feature map is reduced to half the input dimension. Through 3 convolutional encoding layers with filter size  $3 \times 3$ , we obtain the latent parameterizations  $z_m^1$ ,  $z_m^2$ , and  $z_m^3$ . Table 1 summarizes the architecture and dimensions of each layer.



**Figure 7:** Schematic for a separable convolutional layer. Each channel is convolved with its own set of convolutional filters to obtain the best representation, as opposed to traditional convolutions where the same filter is applied to all channels in the data.

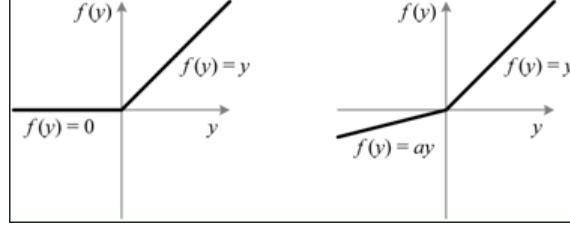


**Figure 8:** Schematic for a squeeze-and-excite layer. The "squeeze" layer takes the global average of the data for each channel, and the "excite" layer is a fully-connected layer with nonlinear activation to estimate the optimal weights for each channel in the data. The result is a weighted representation of the data based on their intrinsic global characteristics.



**Figure 9:** Schematic for instance normalization (left) compared to group normalization (center) and batch normalization (right). In an instance normalization layer, each channel will be normalized by themselves rather than normalizing the entire batch or a subset of channels (groups).

The decoder portion of the pix2vid model extracts the spatiotemporal relationships from the latent representations of  $m$  to reconstruct the dynamic pressure and saturation responses over time,  $d$ . To accurately reconstruct the spatiotemporal structure from the static latent space  $z_m$ , we employ a series of convolutional-recurrent layers, namely a convolutional long-short term memory layer (ConvLSTM). The general form of a 2D ConvLSTM layer is shown in Figure 11. Through 3 convolutional-recurrent layers, we obtain the dynamic



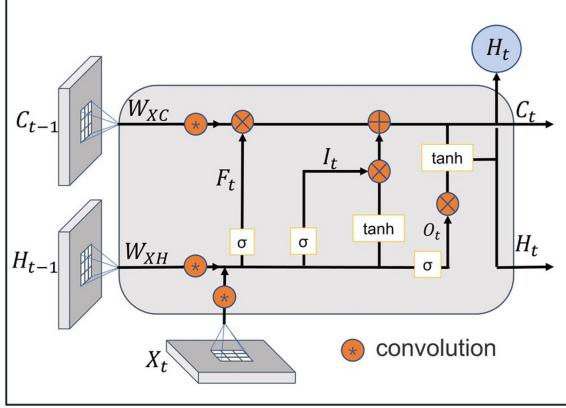
**Figure 10:** Schematic for the Parametric Rectified Linear Unit (PReLU) activation function (right) compared to the traditional ReLU activation function (left). The slope of the negative portion of the data,  $\alpha$ , is learned for each batch.

**Table 1:** Encoder network architecture

Layer Number	Architecture	Shape in (h,w,c)	Shape out (h,w,c)
1	SeparableConv2D	$64 \times 64 \times 4 (m)$	
	Squeeze-and-Excite		
	Instance Norm		
	PReLU + Pooling		
2	Spatial Dropout		$32 \times 32 \times 64 (z_m^1)$
	SeparableConv2D	$32 \times 32 \times 64$	
	Squeeze-and-Excite		
	Instance Norm		
3	PReLU + Pooling		
	Spatial Dropout		$16 \times 16 \times 128 (z_m^2)$
	SeparableConv2D	$16 \times 16 \times 128$	
	Squeeze-and-Excite		
3	Instance Norm		
	PReLU + Pooling		
	Spatial Dropout		$8 \times 8 \times 256 (z_m^3)$

219 prediction of  $d$  as follows:

- 220 Step 1: **Spatiotemporal decoding of  $z_m^3$ :** The first ConvLSTM layer takes the smallest latent represen-  
221 tation,  $z_m^3$ , and reconstructs the first decoded timestep  $z_d^3$ .
- 222 Step 2: **Residual concatenation of  $z_m^2$ :** The first decoded timestep,  $z_d^3$ , is concatenated with the inter-  
223 mediate static encoding  $z_m^2$  to retain multi-scale features and improve prediction performance and  
224 resolution.
- 225 Step 3: **Intermediate spatiotemporal decoding:** The second ConvLSTM layer takes the residual con-  
226 catenation of the intermediate latent representations,  $[z_m^2, z_d^3]$ , to predict the next spatiotemporal  
227 representation  $z_d^2$ .
- 228 Step 4: **Residual concatenation of  $z_m^1$ :** The intermediate decoded timestep,  $z_d^2$ , is concatenated with the  
229 largest static encoding  $z_m^1$ .



**Figure 11:** Schematic of a convolutional-LSTM (ConvLSTM) layer. The layer applies convolutional operations to the input data using a set of learnable filters to capture the spatial patterns. The recurrent part is a long short-term memory layer with memory and forget gates to capture the temporal patterns. LSTM units are applied to each spatial location separately allowing to capture both spatial and temporal dependencies in the data.

230 Step 5: **Final spatiotemporal decoding:** The third ConvLSTM layer takes the residual concatenation of  
 231 the larger latent representations,  $[z_m^1, z_d^2]$ , to predict the full-scale dynamic output,  $d$ .

232 To enhance the performance of the spatiotemporal decoding, each ConvLSTM layer is followed by a batch  
 233 normalization, activation, and a transpose convolutional layer, the latter for downscaling the latent space to  
 234 twice its dimension. Spatial dropout is then applied, and the concatenated features are once more convolved  
 235 and activated to obtain the layer prediction. Table 2 shows the architecture of the decoder network.

236 This process yields the first video frame prediction,  $d^1$ , from the latent representation of the geologic  
 237 realizations  $z_m$ . Each subsequent video frame prediction is obtained by another set of residual concate-  
 238 nation of the previous timestep dynamic decoded representation. The static latent representation  $z_m$  is  
 239 concatenated at each timestep with the previous dynamic decoded representation for each layer such that  
 240 we have  $[z_m, z_{d_t}^i]$ , where  $i$  is the decoding layer number and  $t$  is the timestep. By recursively implementing  
 241 spatiotemporal decoding to the latent representation  $z_m$ , we obtain the prediction of the dynamic response  
 242 at times  $[t_0, t_1, \dots, t_n]$  for each iteration  $t = 1, \dots, n$ .

243 The complete Stochastic pix2vid architecture is shown in Figure 12. Here we observe the spatial com-  
 244 pression of the geologic models  $m$  through the decoding portion of the network, and the spatiotemporal  
 245 decoding and residual multi-scale concatenations through the decoder portion of the network. The resulting  
 246 architecture provides a data mapping operator from static geologic models (images) to dynamic reservoir  
 247 response (videos).

#### 248 **2.4 Training Strategy**

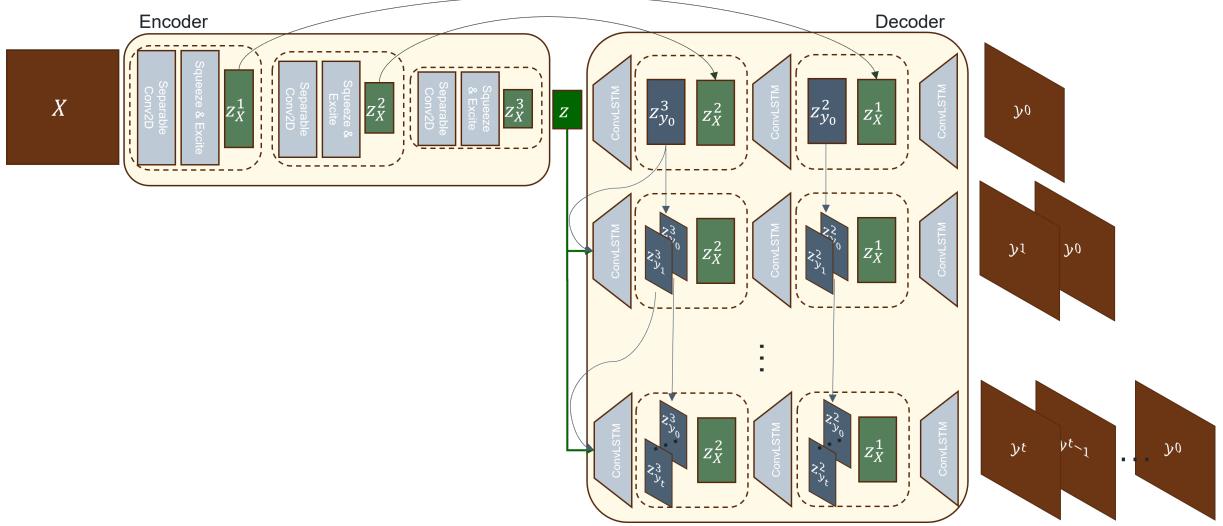
249 The inputs to the Stochastic pix2vid are the geologic realizations, comprised of the distributions of

**Table 2:** Decoder network architecture

Layer Number	Architecture	Shape in (t,h,w,c)	Shape out (t,h,w,c)
1	ConvLSTM2D	$1 \times 8 \times 8 \times 256$	
	BatchNorm + LeakyReLU		
	Conv2DTranspose		
	Spatial Dropout		
	Concatenate ( $z_m^3$ )		
2	Conv2D + Sigmoid		$t \times 16 \times 16 \times 128 (z_{d_t}^3)$
	ConvLSTM2D	$t \times 16 \times 16 \times 128$	
	BatchNorm + LeakyReLU		
	Conv2DTranspose		
	Spatial Dropout		
3	Concatenate ( $z_m^2$ )		
	Conv2D + Sigmoid		$t \times 32 \times 32 \times 64 (z_{d_t}^2)$
	ConvLSTM2D	$t \times 32 \times 32 \times 64$	
	BatchNorm + LeakyReLU		
	Conv2DTranspose		
	Spatial Dropout		
	Concatenate ( $z_m^1$ )		
	Conv2D + Sigmoid		$t \times 64 \times 64 \times 2 (z_{d_t}^1)$

250 porosity, permeability, facies, and injection well(s) location, represented as a matrix  $m$  of dimensions  $64 \times$   
 251  $64 \times 4$ . The outputs are the results from the numerical reservoir simulation, namely pressure and saturation  
 252 distributions over time, represented as a matrix  $d$  of dimensions  $64 \times 64 \times 60 \times 2$ . This yields an ill-posed  
 253 and under-determined estimation problem, which is extremely difficult to resolve [citation]. To improve the  
 254 training efficiency and performance, we subsample in time from 60 timesteps to 11. In other words, instead  
 255 of monthly monitoring, we predict the dynamic outputs at the initial step and every 6 months afterward;  
 256 therefore the output matrix  $d$  a final dimension of  $64 \times 64 \times 11 \times 2$ . We also perform min-max normalization  
 257 so that the input and output features are in the range of  $[0, 1]$ , which greatly improves the performance of  
 258 the nonlinear activation functions. Furthermore, we perform data augmentation by  $90^\circ$  rotation, making  
 259 the network agnostic to orientation and effectively learning the flow physics in the system rather than  
 260 memorizing spatial distribution patterns. The total amount of training data is therefore 2,000 realizations  
 261 (after augmentation), which is split into 1,500 realizations for training and 500 realizations for testing. To  
 262 improve model generalizability, at each epoch, each minibatch is split into 80/20 for training and validation  
 263 sets, respectively.

264 A custom three-part loss function is used to accurately predict pixel-wise and perceptual information in  
 265 the predictions. The mean squared error (MSE) is used to reconstruct the pixel-wise intensity values, while  
 266 the mean absolute error (MAE) is used to optimize for the pressure and saturation plume edges. The third  
 267 part is the structural similarity index metric (SSIM), which provides a perceptual image-to-image comparison



**Figure 12:** Architecture of the Stochastic pix2vid model. The input data,  $X \equiv m$ , is encoded through a series of convolutional layers to capture the spatial dependencies in the geologic models. The latent representation,  $z_m$ , is then recursively passed through a spatiotemporal decoder with convolutional-recurrent layers, and concatenated with the residuals of the encoder to reconstruct iteratively the frames of the output (video) data,  $y \equiv d$ .

268 of luminance, contrast, and structure [2]. For optimal training, the aim is to minimize the MSE and MAE  
 269 while maximizing the SSIM for the true versus predicted outputs,  $d$  and  $\hat{d}$ , such that the total loss is given  
 270 by:

$$\mathcal{L} = \alpha(1 - SSIM) + (1 - \alpha)[\beta MSE + (1 - \beta) MAE] \quad (7)$$

271 where  $\alpha$  and  $\beta$  are weighting coefficients obtained empirically as 0.33 and 0.66, respectively.

272 The model is trained using the AdamW optimizer [citation]. This variant of the well-known adaptive  
 273 momentum (Adam) optimizer [citation] includes an added method to decay weights for the adaptive esti-  
 274 mation of first-order and second-order moments. We implement a learning rate of  $1 \times 10^{-3}$  with a weight  
 275 decay term of  $1 \times 10^{-5}$ .

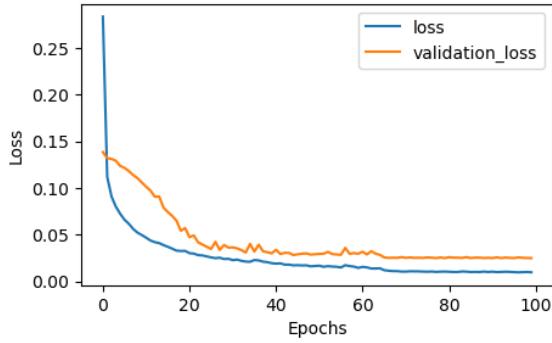
## 276 3 Results

277 This section describes the Stochastic pix2vid model training performance and discusses the results for various  
 278 training and testing realizations.

### 279 3.1 Training Performance

280 Using an NVIDIA Quadro M6000 GPU, we train for 100 epochs with a batch size of 50. The model has  
 281 a total of 97,523,370 parameters, and the training time required is approximately 88 minutes for all 1,500

282 training realizations. The training and validation performance per epoch is shown in Figure 13. We observe  
 283 minimal overfit in the validation set, corresponding to good model generalizability and prediction accuracy.  
 284 Using physics-based numerical simulation, each realization requires approximately 30 seconds to obtain the  
 285 dynamic pressure and saturation predictions from the static geologic models. Our Stochastic pix2vid model  
 286 obtains the same results in approximately 4.59 milliseconds, corresponding to a  $6,500 \times$  speedup. The average  
 287 MSE for the ensemble is  $9.21 \times 10^{-4}$  and  $9.70 \times 10^{-4}$  for training and testing, respectively. Similarly, the  
 288 average SSIM for the ensemble is 98.97 and 97.91 for training and testing, respectively.



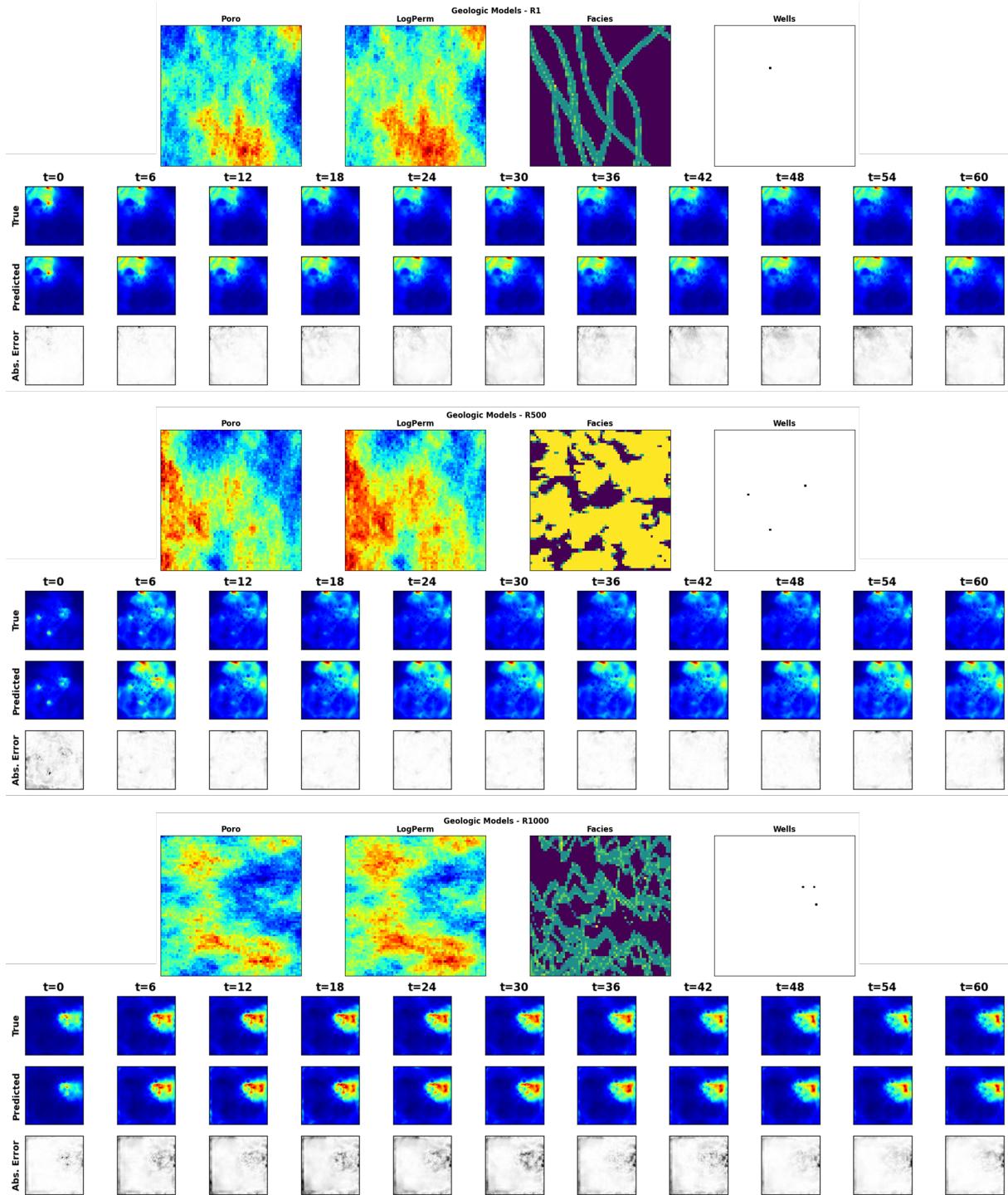
**Figure 13:** The total training and validation losses,  $\mathcal{L}$ , as a function of epoch number.

### 289 3.2 Prediction Results

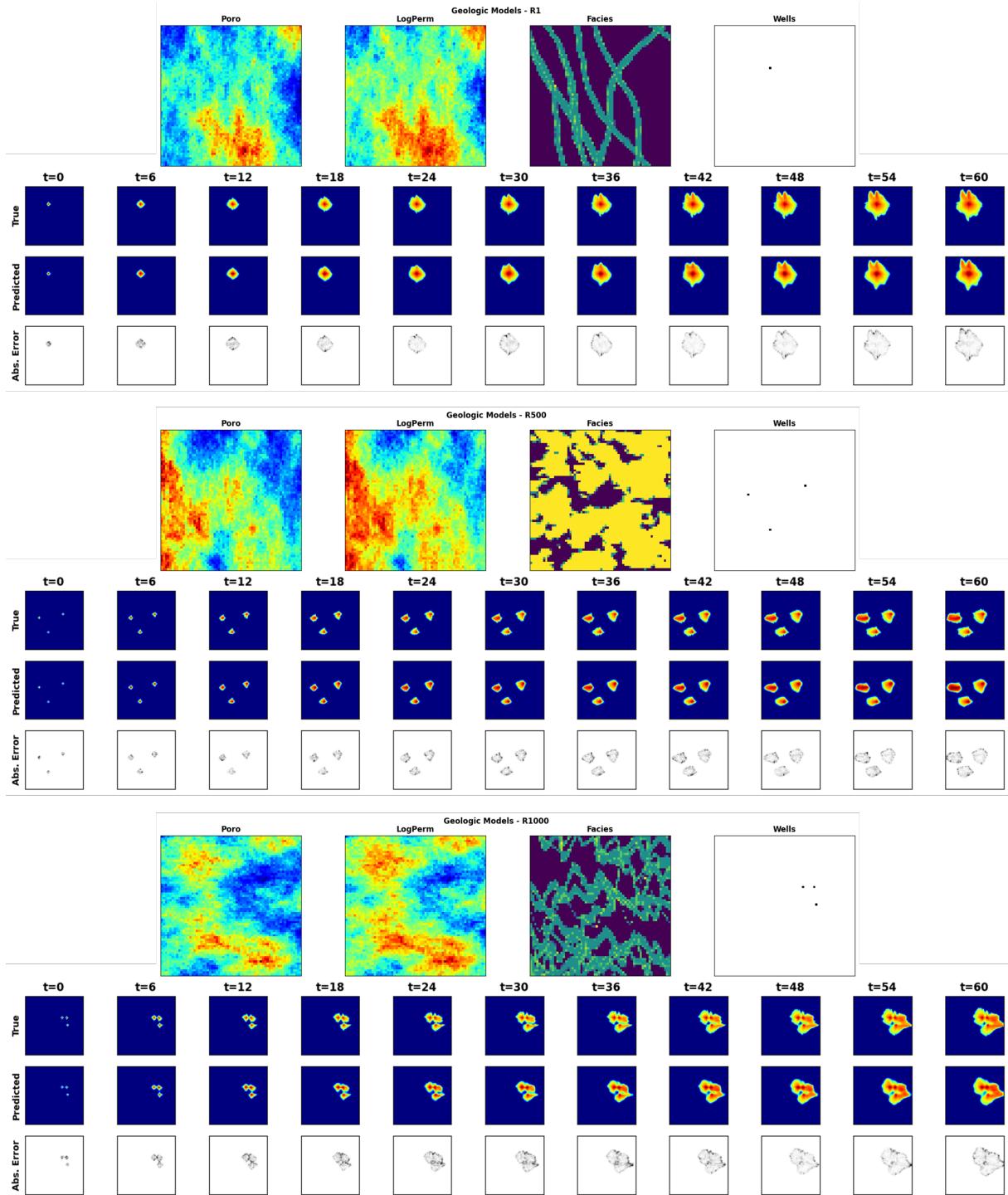
290 The Stochastic pix2vid model is capable of predicting dynamic reservoir response from a static geologic  
 291 model as an image-to-video data mapping operator. Figures 14 and 15 show the predicted pressure and  
 292 saturation distributions, respectively, along with the absolute difference to HFS for 3 training realizations.  
 293 We observe reasonable agreement between the true and predicted CO<sub>2</sub> pressure and saturation plumes over  
 294 time, with an average MSE of  $3.25 \times 10^{-4}$  and SSIM of 98.59% for pressure predictions and MSE of  $1.50 \times 10^{-4}$   
 295 and SSIM of 97.31% for saturation predictions.

296 Similarly, Figures 16 and 17 show the pressure and saturation distributions predictions along with the  
 297 absolute difference to HFS for 3 testing realizations. We observe a similar performance, with an average MSE  
 298 of  $3.71 \times 10^{-4}$  and SSIM of 97.55% for pressure predictions and MSE of  $1.61 \times 10^{-3}$  and SSIM of 96.19% for  
 299 saturation predictions. This indicates that the Stochastic pix2vid model has excellent generalization ability  
 300 and achieves on par performance with HFS at a fraction of the computational cost.

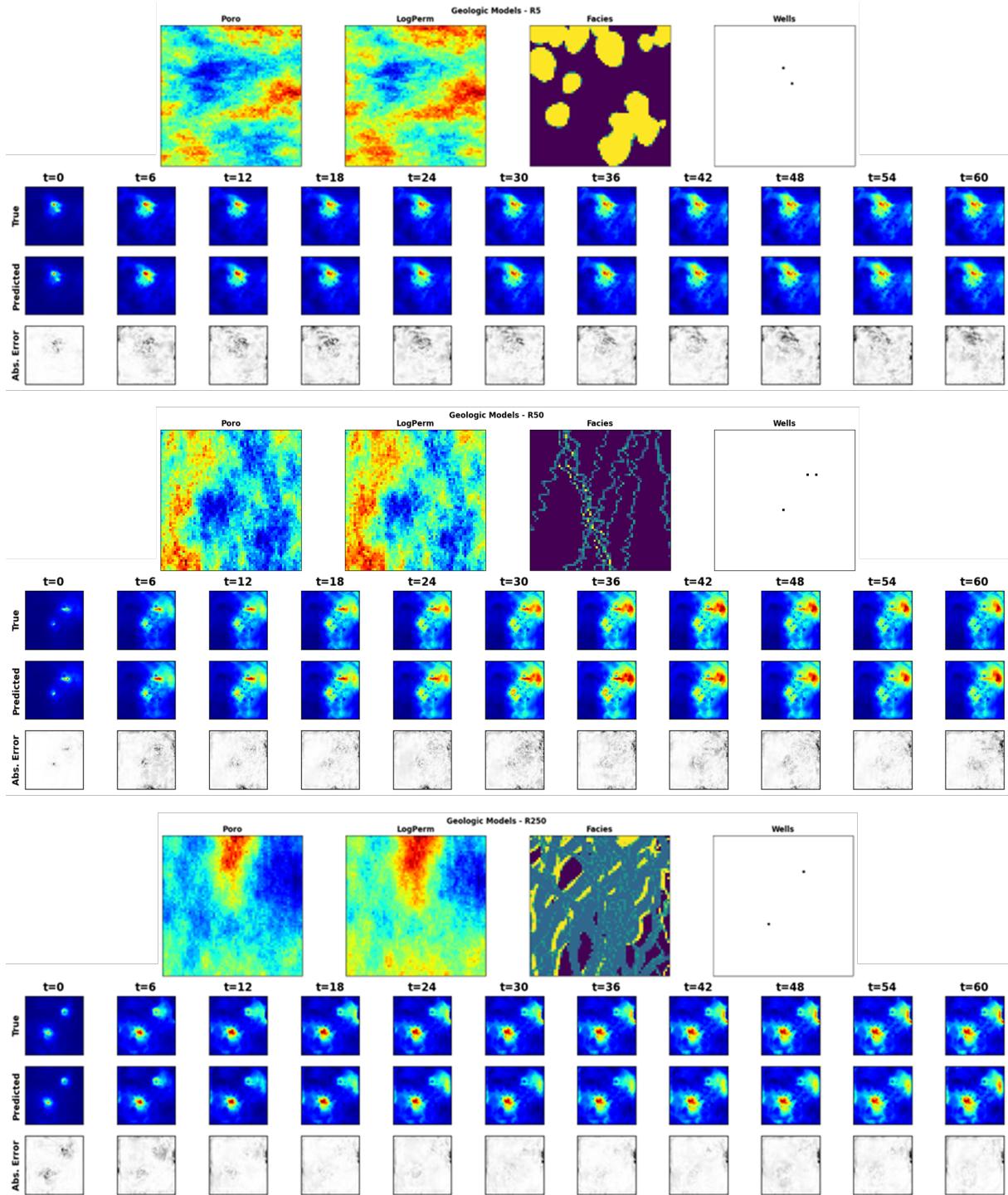
301 These results imply that our Stochastic pix2vid is capable of learning the spatiotemporal relationship  
 302 between the static geologic models and the dynamic reservoir response. Thus, our image-to-video architecture  
 303 can outperform current image-to-image and encoder-recurrent-decoder architectures for improved reservoir  
 304 behavior prediction. To quantify the uncertainty in predictions, a comparison of true ( $d$ ) versus predicted ( $\hat{d}$ )  
 305 response distributions for pressure and saturation for the testing data is shown in Figure 18. The average  $R^2$



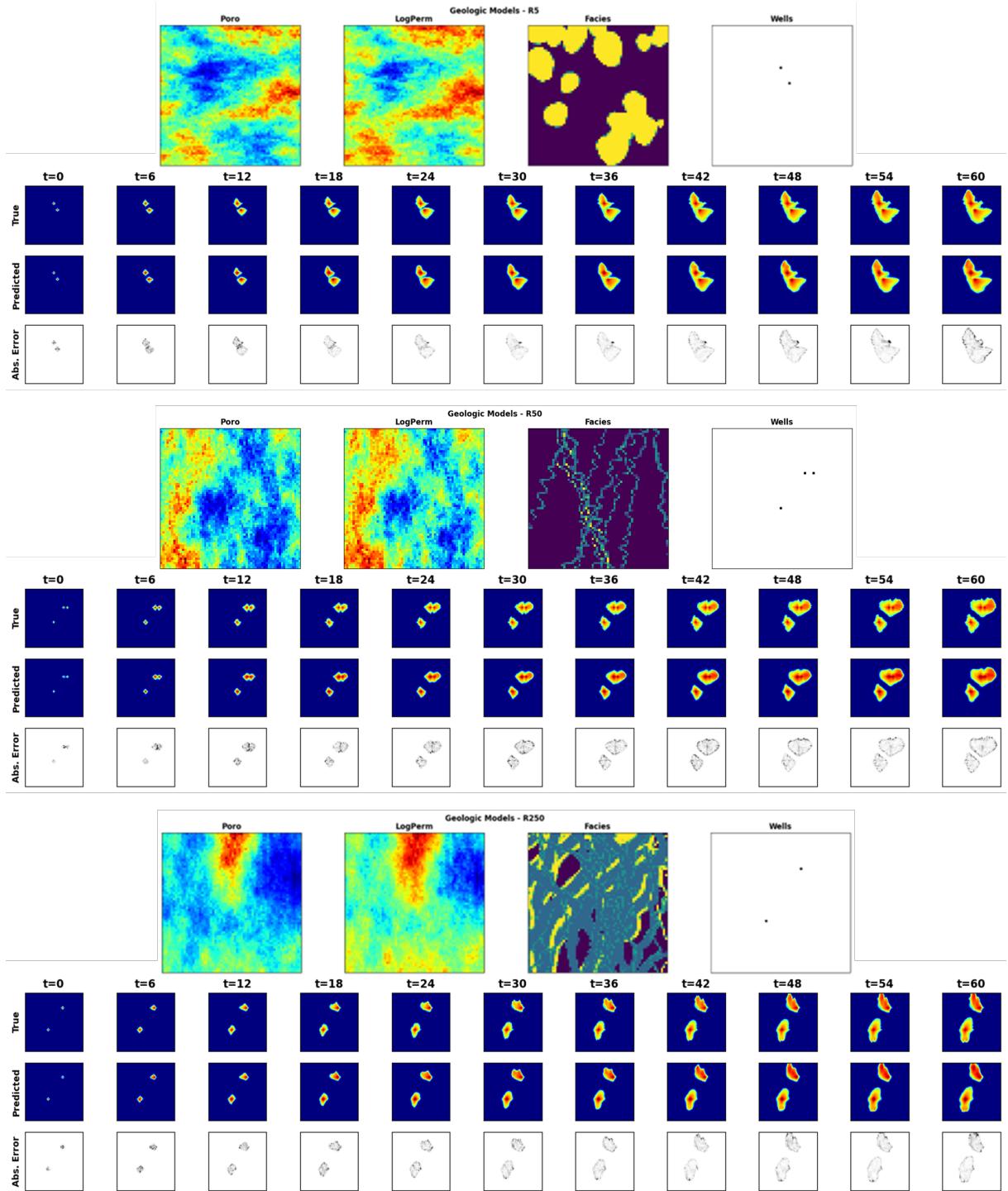
**Figure 14:** (Normalized) pressure distribution over time for 3 random training realization. For each panel, the top row is the ground truth from the HFS, the middle row is the Stochastic pix2vid prediction, and the bottom row is the absolute difference.



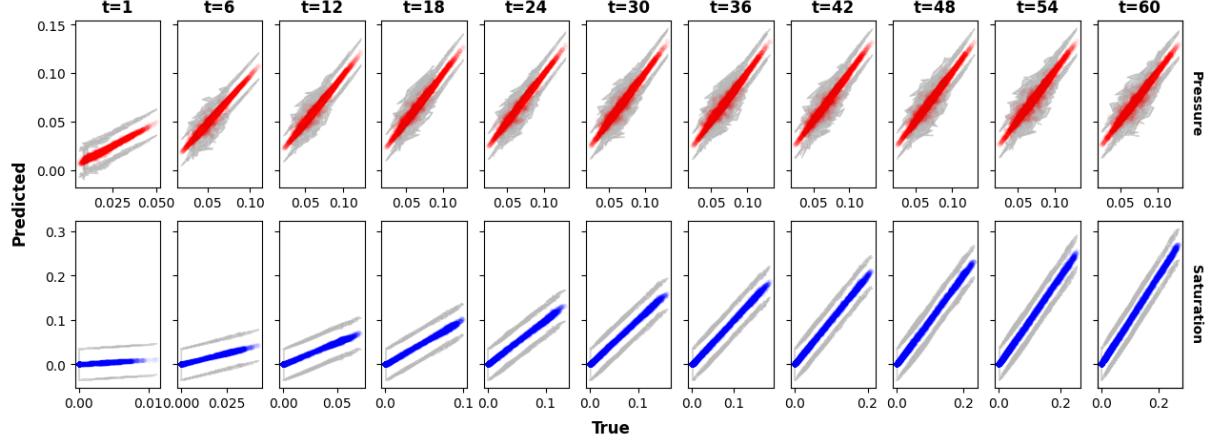
**Figure 15:** Saturation distribution over time for 3 random training realization. For each panel, the top row is the ground truth from the HFS, the middle row is the Stochastic pix2vid prediction, and the bottom row is the absolute difference.



**Figure 16:** (Normalized) pressure distribution over time for 3 random testing realization. For each panel, the top row is the ground truth from the HFS, the middle row is the Stochastic pix2vid prediction, and the bottom row is the absolute difference.



**Figure 17:** Saturation distribution over time for 3 random testing realization. For each panel, the top row is the ground truth from the HFS, the middle row is the Stochastic pix2vid prediction, and the bottom row is the absolute difference.



**Figure 18:** True versus predicted average (normalized) pressure (top) and saturation (bottom) over time for the testing data. The gray portion represents the 95% confidence bands, which become narrower over time.

over time is approximately 99% with narrow 95% prediction bands that recursively narrow over time. From Figure 18 we observe the Stochastic pix2vid model’s performance at recursively refining the predictions over time due to the residual connections in the spatiotemporal decoder network, reducing the spatiotemporal uncertainty in the predictions.

This shows that despite some minor inaccuracies in the plume front prediction, the overall shape and intensity of the CO<sub>2</sub> pressure and saturation plumes are accurately recovered, and the model can be used as a reliable replacement for expensive numerical reservoir simulations, especially in cases where large number of runs are required to obtain dynamic estimates (e.g., well placement and control optimization, history matching, uncertainty quantification).

### 3.3 Discussion

The results shown above suggest that our Stochastic pix2vid is an efficient and accurate predictor of dynamic reservoir response from static geologic models, serving as a reasonable replacement for physics-based numerical reservoir simulation.

CO<sub>2</sub> saturation and pressure buildup fronts are important quantities for geologic CO<sub>2</sub> storage projects and are often used for regulatory oversight [citation], monitoring metrics or history matching purposes [citation]. The distance between the injection well(s) and the saturation fronts represents the maximum extent of the CO<sub>2</sub> plume. However, these are often very difficult to capture accurately with data-driven proxy models. Our Stochastic pix2vid model shows greater absolute error on and around the plume fronts compared to within the plumes. However, the overall shape and intensity of the pressure and saturation plumes is very well captured for all realizations despite being highly heterogeneous.

The encoder block is composed of separable convolutions, squeeze and excite layers, and instance normal-

ization. These three special implementations allows for precise parameterization of the geologic realization into a latent representation, without mixing the effects of Gaussian distributed properties against binary or binomial distributed properties. Using recursive ConvLSTM layers, the recurrent decoder block recursively predicts each dynamic state, or frame, from the concatenation of the previous latent representation and the intermediate encoding parameterizations. Thus, this architecture presents the proxy as an image-to-video prediction formulation for dynamic reservoir states from a static geologic realization.

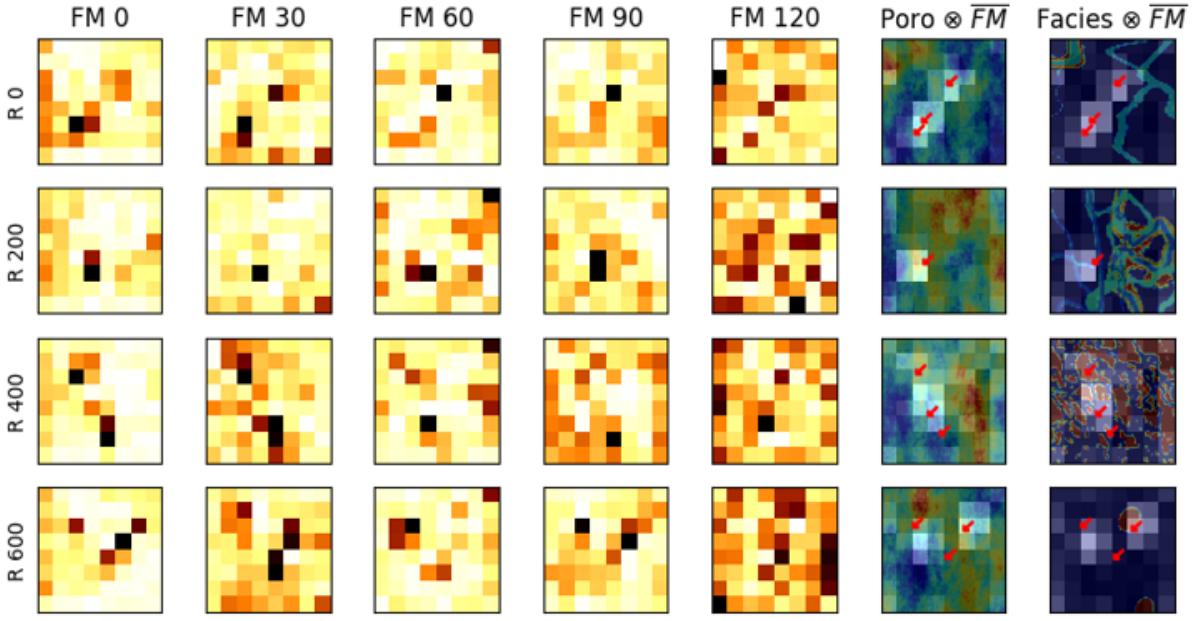
By using GPU-enabled computations, we can significantly accelerate the training and prediction time of the pix2vid model. Each HFS run was performed on an Intel ®i9-10900KF processor with 10 cores. The 1,000 realizations are parallelized equally among all cores and the total simulation time accounting for parallelization for all realizations is about 8.33 hours. Dynamic prediction using the pix2vid model on an NVIDIA Quadro M6000 GPU require a total of 4.6 seconds, or 0.001275 hours, with an accuracy of 99% and 98% for training and testing, respectively. This provides a sustainable argument for the usage of our Stochastic pix2vid model as a replacement for HFS when computational time is a constraint.

As described in Section 2.3, the Stochastic pix2vid model takes the static geologic realizations,  $m$ , and compresses them into a latent space representation,  $z_m$ . Here we provide a visualization for a random selection of latent feature maps, along with their superposition on the porosity and facies distribution, as shown in Figure 19. This can be interpreted as an analog to the attention head mechanisms recently developed in transformer-based architectures [citation]. We observe that the latent feature maps are essentially learning the injection location(s) and direction of flow based on the geologic distributions. Thus, proving that the Stochastic pix2vid model is learning multiphase flow physics and dynamic reservoir behavior appropriately.

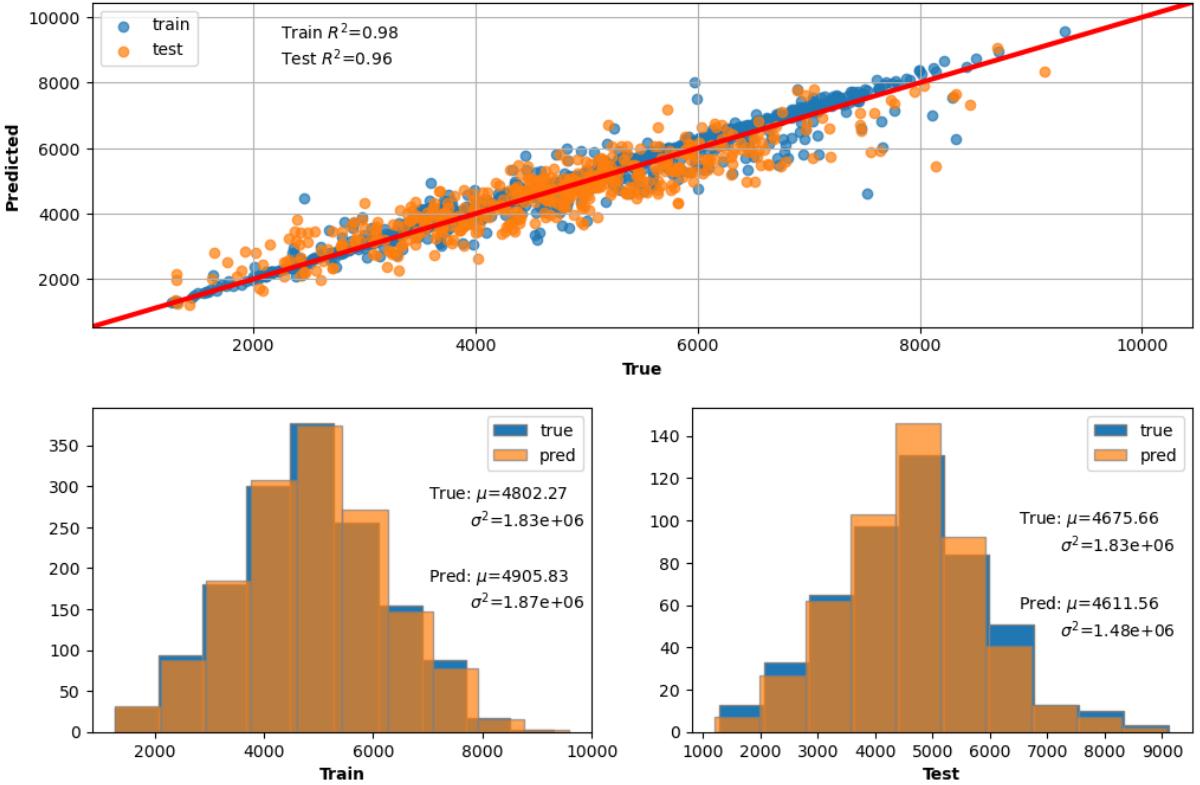
To further demonstrate the effectiveness of our Stochastic pix2vid model for geologic CO<sub>2</sub> storage operations, we plot the cumulative pixel-wise CO<sub>2</sub> saturation as a surrogate for the cumulative CO<sub>2</sub> volume injected. For all training and testing realizations, Figure 20 shows the sum of pixel-wise CO<sub>2</sub> saturation and the probability density function (PDF) of the true versus predicted saturations. We observe an  $R^2$  of 98% for training and 96% for testing in the cumulative CO<sub>2</sub> saturation of true versus predicted results, and a conformable PDFs for both training and testing.

## 4 Conclusions

In this study, we developed a deep learning-based spatiotemporal proxy model to provide flow predictions for a large-scale GCS operation. The key extension introduced is the use of a spatiotemporal convolutional-recurrent architecture for dynamic predictions of CO<sub>2</sub> pressure and saturation distributions over time from an uncertain static geologic realization. The framework was developed as an image-to-video prediction, which is



**Figure 19:** Five random feature maps (FM) of  $z_m^3$  for 4 random realizations. Their average is overlaid on top of the porosity and facies distributions to show the attention mechanism of the encoder. Bright colors represent higher attention and dark colors represent lower attention.



**Figure 20:** (Top) True vs. predicted cumulative CO<sub>2</sub> volume injected via pixel-wise saturation. (Bottom) True vs. predicted distributions of cumulative CO<sub>2</sub> saturation for training (left) and testing (right).

358 a noteworthy under-determined estimation problem. Specifically, the implementation extends the architec-  
359 tures of current encoder-recurrent-decoder models and provides a fast and accurate proxy as a replacement  
360 for high-fidelity numerical reservoir simulation.

361 The spatiotemporal proxy was applied to a synthetic 2D GCS project with multiple uncertain geologic  
362 scenarios and random number and location of injection well(s). A total of 1,000 geologic models were obtained  
363 from a variety of possible geologic scenarios including fluvial, turbidite, and deepwater lobe systems. The  
364 spatial distribution of porosity, permeability and facies, and the spatial location of the injector well(s) were  
365 used as the input data. The proxy then predicts the dynamic reservoir response over time, namely the video  
366 frames, corresponding to the dynamic CO<sub>2</sub> pressure and saturation distributions, which are obtained offline  
367 for training using HFS. The total training time is 88 minutes on a single NVIDIA Quadro M6000 GPU,  
368 and predictions are obtained with 98-99% accuracy within approximately 4.6 milliseconds, compared to the  
369 approximate 30 seconds required for HFS – a 6,500× speedup.

370 There are several possible directions that could be considered for future work. Firstly, an extension to  
371 3D geologic models and their corresponding dynamic predictions is key to extending this method to real-  
372 world applications. Similarly, although the spatiotemporal convolutional-recurrent proxy was only trained  
373 for CO<sub>2</sub> sequestration, it should be applicable for a range of processes such as compositional, geothermal,  
374 or conventional oil and gas systems. The proxy could also be applied to several subsurface energy resource  
375 workflows such as optimization and history matching. Moreover, it would be interesting to extend the  
376 proxy from a data-driven mapping operator to a PINN by including the discretized form of the governing  
377 PDE in the loss function and minimizing the residuals. Another future direction would be to test the  
378 performance of the Stochastic pix2vid model on unseen timesteps, either interpolating the training timesteps  
379 or extrapolating beyond the training timesteps. Furthermore, the proxy is robust to uncertain geology and  
380 variable number and placement of injector wells but could be extended to variable well controls and applied  
381 to robust optimization and closed-loop reservoir management workflows.

## 382 Reproducibility

383 The code will be made publicly available on the author’s repository ([github.com/misaelmmorales](https://github.com/misaelmmorales) and  
384 [github.com/GeostatsGuy](https://github.com/GeostatsGuy)).

## 385 Funding

386 This research did not receive any specific grant from funding agencies in the public, or not-for-profit sectors.

387 **Declarations**

388 The authors declare no conflict of interests.

389 **Acknowledgements**

390 The authors thank the Digital Reservoir Characterization Technology (DIRECT) and Formation Evaluation  
391 (FE) Industry Affiliate Program at the University of Texas at Austin for supporting this work.

392 **References**

- 393 [1] Knut-Andreas Lie. *An introduction to reservoir simulation using MATLAB/GNU Octave: User guide*  
394 *for the MATLAB Reservoir Simulation Toolbox (MRST)*. Cambridge University Press, 2019.
- 395 [2] Zhou Wang, Alan C. Bovik, Hamid R. Sheikh, and Eero P. Simoncelli. Image quality assessment: from  
396 error visibility to structural similarity. *IEEE Transactions on Image Processing*, 13:600–612, 4 2004.  
397 ISSN 1941-0042. doi: doi.org/10.1109/TIP.2003.819861.